

OptiTrust: Producing Trustworthy High-Performance Code via Source-to-Source Transformations

GUILLAUME BERTHOLON, ARTHUR CHARGUÉRAUD, THOMAS KÉHLER, BEGATIM BYTYQI, and DAMIEN ROUHLING, Inria & ICube lab, CNRS, Université de Strasbourg, France

Developments in hardware have delivered formidable computing power. Yet, the increased hardware complexity has made it a real challenge to develop software that exploits hardware to its full potential. Numerous approaches have been explored to help programmers turn naive code into high-performance code, finely tuned for the targeted hardware. However, these approaches have inherent limitations, and it remains common practice for programmers seeking maximal performance to follow the tedious and error-prone route of writing optimized code by hand.

This paper presents OptiTrust, an interactive source-to-source optimization framework. The programmer develops a script describing a series of code transformations. The framework provides continuous feedback in the form of human-readable *diffs* over conventional C syntax. OptiTrust supports advanced code transformations, including transformations exploited by the state-of-the-art DSL tools Halide and TVM, and transformations beyond the reach of existing tools. OptiTrust also supports user-defined transformations, as well as defining complex transformations by composition of simpler transformations.

Crucially, to check the validity of code transformations, OptiTrust leverages a *static resource analysis* in a simplified form of Separation Logic. Our analysis exploits user-provided annotations on functions and loops, and deduces precise resource usage throughout the code. Through three representative case studies, we demonstrate how OptiTrust can be employed to produce state-of-the-art, high-performance programs.

1 INTRODUCTION

1.1 Motivation

Performance matters in numerous fields of computer science, and in particular in applications from machine learning, computer graphics, and numerical simulation. Massive speedups can be achieved by fine-tuning the code to best exploit the available hardware [Kelefouras and Keramidas 2022]. Between a naive implementation and an optimized implementation, it is common to see a speedup of the order of 50×, on a single core. For many applications, the code can then be accelerated further by one or two orders of magnitude by exploiting multicore parallelism or GPUs.

Yet, producing high performance code is hard. Over the past decades, nontrivial mechanisms with subtle interactions were integrated into hardware architectures. Reasoning about performance requires reasoning about the effects of multiple levels of caches, the limitations of memory bandwidth, the intricate rules of atomic operations, and the diversity of vector instructions (SIMD). These aspects and their interactions make it challenging to build cost models. For example, the cost of a memory access can range from one CPU cycle to hundreds of CPU cycles, depending on whether the corresponding data is already in cache. In the general case, accurately modeling cache behavior requires a deep understanding of the algorithm and hardware at play.

Accurately predicting runtime behavior is challenging for expert programmers, and appears beyond the capabilities of automated tools. Therefore, compilers struggle to navigate the exponentially large search space of all possible code candidates [Vachharajani et al. 2003], resorting to best-effort heuristics, and often failing to produce competitive code [Barham and Isard 2019].

Today, it remains common practice in industry for programmers to write optimized code *by hand* [Amaral et al. 2020; Evans et al. 2022]. However, manual code optimization is unsatisfactory for at least three reasons. First, manually implementing optimized code is time-consuming. Second, the optimized code is hard to maintain through hardware and software evolutions. Third, the

Authors' address: Guillaume Bertholon; Arthur Charguéraud, arthur.chargueraud@inria.fr; Thomas Koehler; Begatim Bytyqi; Damien Rouhling, Inria & ICube lab, CNRS, Université de Strasbourg, France.

rewriting process is error-prone: not only every manual code edition might introduce a bug, but the code complexity also increases, especially when introducing parallelism. These three factors are exacerbated by the fact that optimizations typically make code size grow by an order of magnitude (Section 2 contains examples).

In summary, neither fully automatic nor fully manual approaches are satisfying for generating high performance code. *Semi-automatic code optimization* aims at combining the benefits of machine automation with the strength of human insight. Before reviewing tools for semi-automatic code optimization, let us introduce a number of qualitative properties on which to evaluate these tools.

- **Generality:** How large is the domain of applicability of the tool? In particular, is it restricted to a domain-specific language (DSL)?
- **Expressiveness:** How advanced are the code transformations supported by the tool? Is it possible to express state-of-the-art code optimizations?
- **Control:** How much control over the final code is given to the user by the tool? In particular, is there a monolithic code generation stage?
- **Feedback:** Does the tool provide easily readable intermediate code after each transformation?
- **Composability:** Is it possible to define transformations as the composition of existing transformations? Can transformations be higher-order, i.e., parameterized by other transformations?
- **Extensibility** of transformations: Does the tool facilitate defining custom transformations that are not expressible as the composition of built-in ones?
- **Modularity** of analyses: for transformations whose correctness depends on a code analysis, can the tool deal with specifications that summarize the effects of each function, or are all functions inlined during the analyses?
- **Trustworthiness:** Does the tool ensure that user-requested transformations preserve the semantics of the code? Can it moreover provide mechanized proofs?

There exists other properties for optimization tools, such as the ease of integration in an existing code base, the maintainability of optimized code, or the steepness of the learning curve for new users. These are certainly important aspects, yet they are even harder to evaluate objectively. Hence, we omit them from the discussion, and focus on the aforementioned technical properties.

1.2 Closely Related Work

Table 1 summarizes the properties of existing approaches, highlighting their diversity. For the tools considered, generality appears negatively correlated with expressiveness, i.e., with how advanced the supported transformations are. For each property considered, at least two tools show strengths on that property. However, even if we leave out the ambition of achieving mechanized proofs, each tool considered shows weaknesses on several properties. Hence, it appears that there remains a lot of room for improvement. Before presenting the contribution of the OptiTrust framework, we first describe the tools listed in the table.

Halide [Ragan-Kelley et al. 2013] is an industrial-strength domain-specific compiler for image processing, used e.g. to optimize code running in Photoshop and YouTube. Halide popularized the idea of separating an *algorithm* describing what to compute from a *schedule* describing how to optimize the computation. This separation makes it easy to try different schedules. TVM [Chen et al. 2018] is a tool directly inspired by Halide, but tuned for machine learning applications; it is used by most of the major CPU/GPU manufacturers. Other tools inspired by Halide include Fireiron [Hagedorn et al. 2020a], used at Nvidia, as well as PartIR [Alabed et al. 2024], used at Google. All these tools are inherently limited to the domains (DSLs) that they target. They do

	Halide/TVM	Elevate+Rise	Exo	Clay/LoopOpt	ATL	Alpinist	Clava+LARA
Generality	◐	◑	◑	◑	◑	◑	◑
Expressiveness	●	●	●	◑	◑	◑	◑
Control	◑	◑	◑	◑	◑	●	●
Feedback	◑	◑	●	●	◑	●	◑
Composability	○	●	●	◑	●	○	●
Extensibility	○	●	●	○	●	●	●
Modularity	○	(not applicable)	○	○	●	●	○
Trustworthiness	◑	◑	◑	◑	●	●	○

Table 1. Overview of user-guided tools for high-performance code generation. Darker is better.

not support higher-order composition of transformations, and are not extensible [Barham and Isard 2019; Ragan-Kelley 2023]. Moreover, understanding their output is difficult as the applied transformations are not detailed to the user, even though interactive scheduling systems have been proposed to mitigate this difficulty [Ikarashi et al. 2021].

Elevate [Hagedorn et al. 2020b] is a functional language for describing *optimization strategies* as composition of simple *rewrite rules*. Advanced optimizations from TVM and Halide can be reproduced using Elevate. One key benefit is extensibility: adding rewrite rules is much easier than changing complex and monolithic compilation passes [Ragan-Kelley 2023]. Elevate strategies are applied on programs expressed in a functional array language named Rise, followed by compilation to imperative code. The use of a functional array language greatly simplifies rewriting, however it restricts applicability and makes controlling imperative aspects difficult (e.g. memory reuse).

Exo [Ikarashi et al. 2022] is an imperative DSL embedded in Python, geared towards the development of high-performance libraries for specialized hardware. The strength of Exo lies in externalizing target-specific code generation to user-level code instead of compilation backends. Exo programs can be optimized by applying a series of source-to-source transformations. These transformations are described in a Python script, with a cursor mechanism for targeting code points. The user can add custom transformations, possibly defined by (higher-order) composition. A major limitation of Exo is that it is restricted to static control programs with linear integer arithmetic. Another important limitation of Exo is that the transformations are performed on code in which all functions are inlined. This approach, which lacks modularity, may harm scalability to larger or more complex programs.

Clay [Bagnères et al. 2016a] is a framework to assist in the optimization of loop nests that can be described in the *polyhedral model* [Feautrier 1992]. The polyhedral model only covers a specific class of loop transformations, with restriction over the code contained in the loop bodies, however it has proved extremely powerful for optimizing code falling in that fragment. Clay provides a decomposition of polyhedral optimizations as a sequence of basic transformations with integer arguments. The corresponding transformation script can then be customized by the programmer. Clint [Zinenko et al. 2018b] adds visual manipulation of polyhedral schedules through interactive 2D diagrams. LoopOpt [Chelini et al. 2021] provides an interactive interface that helps users design optimization sequences (featuring unrolling, tiling, interchange, and reverse of iteration order) that can be bound in a declarative fashion to loop nests satisfying specific patterns.

ATL [Liu et al. 2022] is a purely functional array language for expressing Halide-style programs. Its particularity is to be embedded into the Coq proof assistant. ATL programs can be transformed through the application of rewrite rules expressed as Coq theorems. With this approach, transformations are inherently accompanied by machine-checked proofs of correctness. The set of rules includes expressive transformations, some beyond the scope of Halide, and can be extended by the user. Once optimized, ATL programs are then compiled into imperative C code. Like Rise, generality and control are restricted by the functional array language nature of ATL.

148 Alpinist [Sakar et al. 2022] is a *pragma*-based tool for optimizing GPU-level, array-based code. It
149 is able to apply basic transformations such as loop tiling, loop unrolling, data prefetching, matrix
150 linearization, and kernel fusion. The key characteristic of Alpinist is that it operates over code
151 formally verified using the VerCors framework [Blom et al. 2017]. Concretely, Alpinist transforms
152 not only the code but also its formal annotations. If Alpinist were to leverage transformation scripts
153 instead of pragmas, it might be possible to chain and compose transformations; yet, this possibility
154 remains to be demonstrated.

155 Clava [Bispo and Cardoso 2020] is a general-purpose C++ source-to-source analysis and trans-
156 formation framework implemented in Java. The framework has been instantiated mainly for code
157 instrumentation purpose and auto-tuning of parameters. Clava can also be used in conjunction
158 with a DSL called LARA [Silvano et al. 2019] for optimizing specific programs. LARA allows ex-
159 pressing user-guided transformations by combining declarative queries over the abstract syntax
160 tree and imperative invocations of transformations, with the option to embed JavaScript code. The
161 application paper on the Pegasus tool [Pinto et al. 2020] illustrates this approach on loop tiling and
162 interchange operations.

163

164 1.3 Contribution

165 This paper introduces OptiTrust, the first interactive optimization framework that operates, from
166 the perspective of the user, at the level of C syntax, and that supports and validates state-of-the-art
167 optimizations. OptiTrust is open-source, and available from: <https://github.com/charguer/optitrust>.

168

169 *Overview.* In OptiTrust, the user starts from an unoptimized code in C syntax, and develops a
170 *transformation script* describing a series of optimization steps. Each step consists of an invocation
171 of a specific transformation at specified *targets*. OptiTrust provides an expressive target mechanism
172 for describing, in a concise and robust manner, one or several code locations. On any step of the
173 transformation script, the user can press a key shortcut to view the *diff* associated with that step,
174 in the form of a comparison between two human-readable programs in C syntax. Concretely,
175 a transformation script consists of an OCaml program linked against the OptiTrust library of
176 transformations.

177 A central aspect of OptiTrust is that it guarantees that the code transformations requested by
178 the programmer preserve the semantics of the program. To that end, OptiTrust leverages our *static*
179 *resource analysis*, which concretely takes the form of a type checking algorithm in a type system
180 featuring linear resources. Technically, OptiTrust’s type system consists of a scaled down version
181 of Separation Logic [Reynolds 2002].

182 For type-checking resources, functions and loops need to be equipped with *contracts* describing
183 their resource usage. These contracts may be inserted either directly as annotations (in the form of
184 no-op instructions) in the C source code, or they may be inserted by dedicated commands as part
185 of the transformation script. OptiTrust is able to automatically infer simple loop contracts, thus not
186 all loops need to be annotated manually. Crucially, every OptiTrust transformation takes care of
187 updating contracts in order to reflect changes in the code. In other words, a well-typed program
188 must remain well-typed after a successful transformation. This property is essential to ensure that
189 subsequent transformations in the optimization chain can be validated by exploiting information
190 from our resource analysis.

191 The implementation of OptiTrust distinguishes between *basic* transformations and *combined*
192 transformations. On the one hand, a basic transformation applies minimalistic changes to the
193 abstract syntax tree (AST). The validity of a basic transformation is checked by leveraging the
194 resource analysis. On the other hand, a combined transformation is implemented as a composition
195 of basic transformations. Combined transformations aim to implement high-level strategies, that
196

196

197 may trigger the execution of dozens of basic transformation. These more complex combined trans-
198 formations need not be accompanied with code for checking validity: their validity is guaranteed
199 by the validity checks performed by the basic transformations. This two-layer approach enables us
200 to minimize the size of the trusted code base (TCB) of OptiTrust.

201 OptiTrust operates on a subset of the C language, with a slightly simplified semantics, and
202 augmented with typing annotations—function and loop *contracts*, as well as *ghost code*. We call this
203 user-level language *OptiC*. OptiTrust does not directly manipulate an abstract syntax tree (AST)
204 for *OptiC*. Instead, it operates on an intermediate representation that essentially consists of an
205 imperative λ -calculus. We call this internal language *Opti λ* . Concretely, the OptiC code, expressed
206 in C syntax, is first parsed using Clang. Then, the OptiC code is translated into Opti λ . In particular,
207 our translation eliminates mutable variables and operations involving *l*-values. Importantly, our
208 translation is bidirectional, allowing to print back C syntax after transformations are applied on
209 the internal AST. Considering a syntax and semantics simpler than that of C considerably helps to
210 tame the complexity of the design and implementation of typing rules, code transformations, and
211 correctness criteria associated with transformations.

212 *Limitations.* In the long term, our aim is for OptiTrust to perform full-score on all the aforemen-
213 tioned evaluation criteria. On the way towards this highly ambitious goal, we have considered four
214 simplifications that apply to the work described in the present paper.

- 215 (1) We restrict ourselves to OptiC, which includes a subset of the C language. As our case studies
216 show, this subset nevertheless suffices to express numerous practical, high-performance
217 programs, in an idiomatic programming style both for the unoptimized and for the optimized
218 code. For simplicity, we currently ignore complications related to arithmetic overflows, and
219 we treat floating point numbers as reals. Besides, due to the complexity of the semantics of
220 the C language, we have not yet formalized the relationship between OptiC and C.
- 221 (2) We have already implemented dozens of transformations, among the most standard ones.
222 We believe that these transformations suffice to assess the interest of the OptiTrust approach
223 to code optimization. However, for production usage, dozens of additional transformations
224 remain to implement.
- 225 (3) We have so far restricted ourselves to a subset of Separation Logic. Our resource-based type
226 system is able to describe the ownership of arrays, matrices, or individual cells, however it
227 does not allow specifying properties about the values stored in data structures. Nevertheless,
228 as our case studies show, shape-based resources suffice to justify the correctness of many
229 practical code optimization patterns.
- 230 (4) We present formal definitions and theorems for our typing judgment, and describe a proof
231 strategy for justifying semantic-preservation for individual transformations. However, we
232 do not present correctness proofs for the transformations that we have implemented. Such
233 correctness proofs are extremely tedious and error-prone, thus it would only make sense
234 to carry them out using a proof assistant. Yet, completing such mechanized proofs will
235 presumably require a couple years of additional work. Note that state-of-the-art compilers
236 such as Halide have been described in publications that did not include correctness proofs.

237
238 In the long term, our resource-based system aims to be similar in spirit to RefinedC [Sammler
239 et al. 2021], a Separation Logic-based type system for C code. The effectiveness of Separation Logic
240 has been successfully demonstrated across a broad range of applications, both for low-level and
241 high-level code [Charguéraud 2020a; O’Hearn 2019]. By building OptiTrust on Separation Logic,
242 we are confident that our framework has the potential to be generally applicable.

243 In summary, we present a framework that can readily be exploited to optimize certain classes of
244 programs, and acknowledge that future work remains necessary to achieve full generality. Note
245

that we have taken great care in our design and implementation to anticipate for the extensions to a richer programming language and to a richer Separation Logic.

1.4 Contents of the Paper

We first present the features of OptiTrust by means of example, in Section 2. Then, we present the construction of OptiTrust. In Section 3, we present Opti λ , the language used internally by OptiTrust, and describe at a high level the bidirectional translation between OptiC and Opti λ . In Section 4, we explain the core of our resource-based typechecker. This part presents relatively standard Separation Logic concepts, but following an algorithmic rather than a declarative presentation of the reasoning rules. In Section 5, we explain a key addition to the typechecker, namely the computation of *usage information* for every resource and for every subterm. In Section 6, we present a set of representative code transformations, illustrating in particular how usage information is exploited to guide transformations and to justify their correctness. Finally, we discuss additional related work in Section 7, then conclude in Section 8.

2 OPTITRUST IN PRACTICE

Let us present the features of OptiTrust through three case studies. In Section 2.1, we reproduce a manually written code from OpenCV—a very popular, optimized computer vision library. In Section 2.2, we consider a physics simulation program featuring a kernel typical of particle simulations; we demonstrate how to apply, using OptiTrust, several optimizations that are ubiquitous in this kind of applications. In Section 2.3, we reproduce an optimized implementation of matrix-multiply, similar to the one produced by TVM, the state-of-the-art specialized compiler for machine learning applications. Then, in Section 2.4, we evaluate OptiTrust against the desirable properties for semi-automatic code optimization frameworks.

2.1 The OpenCV Row-Based Blur Case Study

In image processing, a *blur* is typically used to remove noise and smoothen images. A two-dimensional blur can be decomposed as a combination of *column-based blur*, *row-based blur*, and (optionally) the application of a normalization pass. Our case study focuses on a *row-based blur* function, as implemented in the state-of-the-art OpenCV library [Bradski et al. 2000].

Unoptimized Code. If performance was not a concern at all, the row-based blur function would be implemented as shown in Fig. 1. The output is a single-row image, stored in an array named D , made of n pixels. The input is a single-row image, stored in an array named S , made of $n+w-1$ pixels, where the parameter w corresponds to the width of the blur. The input pixels in S are encoded on c_n integers of type τ , whereas the output pixels in D are encoded on c_n integers of type $\sigma\tau$. Typically, the type $\sigma\tau$ is represented on more bits than the type τ . The output pixel $D[i]$ is computed as the sum of the values of the input pixels in the range from $S[i]$ to $S[i+w-1]$. This sum is computed independently for each of the c_n color channels. The code accommodates any value of c_n , but practical values include $c_n=1$ for grayscale, $c_n=3$ for RGB, $c_n=4$ for RGBA.

Optimized Code. The handwritten OpenCV library includes an implementation of row-sum blur structured like the code shown in Fig. 2. The original OpenCV code may be viewed online.¹ The code from Fig. 2 corresponds to the code that we produce using OptiTrust.

¹https://github.com/opencv/opencv/blob/4.10.0/modules/imgproc/src/box_filter.simd.hpp#L75: The OpenCV code is implemented as a class with the types S and $\sigma\tau$ as template arguments, whereas for the moment our code refers to fixed yet unspecified integer types; we look forward to add support for template polymorphism in the future. The OpenCV code also traverses certain arrays by incrementing pointers, whereas we use explicit array indexing everywhere. In general, this choice is not performance critical and we leave OptiTrust support for pointer shifting to future work.

```

295 void rowSum(const int n, const int cn, const int w, const T S[n+w-1][cn], ST D[n][cn]) {
296     for (int i = 0; i < n; i++) { // for each target pixel in the row described by D
297         for (int c = 0; c < cn; c++) { // for each channel (e.g., red, green, and blue)
298             ST s = 0;
299             for (int k = i; k < i+w; k++) // for each source pixel nearby to the right
300                 s += (ST) S[k][c];
301             D[i][c] = s;
302         } } }

```

Fig. 1. Unoptimized C code for the OpenCV case study, using multidimensional arrays.

```

303 void rowSum(const int n, const int cn,
304             const int w, const T* S, ST* D) {
305     if (w == 3) {
306         for (int ic = 0; ic < cn * n; ic++) {
307             D[ic] = (ST) S[ic]
308                 + (ST) S[ic + cn]
309                 + (ST) S[ic + 2 * cn];
310         }
311     } else if (w == 5) {
312         for (int ic = 0; ic < cn * n; ic++) {
313             D[ic] = (ST) S[ic]
314                 + (ST) S[ic + cn]
315                 + (ST) S[ic + 2 * cn]
316                 + (ST) S[ic + 3 * cn]
317                 + (ST) S[ic + 4 * cn];
318         }
319     } else if (cn == 1) {
320         ST s = (ST) 0;
321         for (int i = 0; i < w; i++) {
322             s += (ST) S[i];
323         }
324         D[0] = s;
325     } else if (cn == 3) {
326         ST s0 = (ST) 0;
327         ST s1 = (ST) 0;
328         ST s2 = (ST) 0;
329         for (int i = 0; i < n - 1; i++) {
330             s += (ST) S[i + w] - (ST) S[i];
331             D[i + 1] = s;
332         }
333     } else if (cn == 4) {
334         for (int i = 0; i < 3 * n - 3; i += 3) {
335             s0 += (ST) S[3 * w + i] - (ST) S[i];
336             s1 += (ST) S[3 * w + i + 1] - (ST) S[i + 1];
337             s2 += (ST) S[3 * w + i + 2] - (ST) S[i + 2];
338             D[i + 3] = s0;
339             D[i + 4] = s1;
340             D[i + 5] = s2;
341         }
342     } else if (cn == 1) {
343         ST s = (ST) 0;
344         for (int i = 0; i < w; i++) {
345             s += (ST) S[i];
346         }
347         D[0] = s;
348     } else if (cn == 3) {
349         for (int i = 0; i < n - 1; i++) {
350             s += (ST) S[i + w] - (ST) S[i];
351             D[i + 1] = s;
352         }
353     } else if (cn == 4) {
354         // [...] similar to cn == 3, with one more variable
355         for (int c = 0; c < cn; c++) {
356             ST s = (ST) 0;
357             for (int i = 0; i < cn * w; i += cn) {
358                 s += (ST) S[c + i];
359             }
360             D[c] = s;
361             for (int i = c; i < cn * n - cn + c; i += cn) {
362                 s += (ST) S[cn * w + i] - (ST) S[i];
363                 D[cn + i] = s;
364             }
365         }
366     }
367 }

```

Fig. 2. Our optimized C code for the OpenCV case study, showing the body of the rowSum function. This code exploits essentially the same optimizations as the original OpenCV code.

This optimized implementation is a *multi-versioned* code, with dedicated execution paths for handling specific values of the parameters. The branches $w = 3$ and $w = 5$ correspond to values of the width that are commonly used by library users. For these small constant values of w , the inner loop on k from Fig. 2 is unfolded. Otherwise, the loop on k is not unfolded and a standard algorithmic optimization called *sliding window* is applied. Note that Halide, the state-of-the-art specialized compiler for image processing, does not support the introduction of sliding windows—and the developers of Halide do not plan to lift this limitation.²

The branch of the code that uses the sliding window optimization is then further specialized with branches for commonly used parameters: $cn = 1$ and $cn = 3$ and $cn = 4$. For these small constant values of cn , the outer loop on c is unfolded, then the multiple occurrences of the loop on i that result from this unfolding are fused into a single loop. The final else-branch in the code from Fig. 2

²Halide does not support sliding windows for reasons explained on: <https://github.com/halide/Halide/issues/180>. Hence, the programmer either needs to manually refine the code to introduce the sliding window before scheduling; or needs to exploit other transformation tools specialized in sliding window optimizations [Chaurasia et al. 2015; Kanetaka et al. 2024].

```

344 void rowSum(const int n, const int cn, const int w, const T* S, ST* D) {
345     __requires("w >= 0, n >= 1, cn >= 0");
346     __reads("S ~> Matrix2(n+w-1, cn)");
347     __modifies("D ~> Matrix2(n, cn)");
348     for (int i = 0; i < n; i++) { // for each pixel
349         __xmodifies("for c in 0..cn -> &D[MINDEX2(n, cn, i, c)] ~> Cell");
350         for (int c = 0; c < cn; c++) { // for each channel
351             __xmodifies("&D[MINDEX2(n, cn, i, c)] ~> Cell");
352             __ghost(assume, "is_subrange(i..i + w, 0..n + w - 1)");
353             ST s = 0;
354             for (int k = i; k < i+w; k++) {
355                 __ghost(in_range_extend, "k, i..i+kn, 0..n+kn-1");
356                 __ghost_begin(focus, matrix2_ro_focus, "S, k, c");
357                 s += (ST) S[MINDEX2(n+w-1, cn, k, c)];
358                 __ghost_end(focus);
359             }
360             D[MINDEX2(n, cn, i, c)] = s;
361         }
362     }
363 }

```

Fig. 3. Unoptimized C code for the OpenCV case study, using flat arrays and resource annotations.

corresponds to the generic implementation. Moreover, in the last three branches, the loops are reindexed to augment the counter i by steps of cn , thereby saving multiplication operations.

Multidimensional vs Flat Arrays. The code from Fig. 1 is presented using C syntax for multidimensional arrays, for the sake of improved readability. However, the optimized code from Fig. 2 and our contract-annotated code from Fig. 3 instead use a flat array representation. The flat representation is frequently used in high-performance code: it allows performing simplifications in array accesses, moreover it allows for compatibility with C++ parsers. For technical reasons, and to anticipate for extensions of OptiTrust, OptiTrust relies on a C++ parser. We leave to future work the parsing of multidimensional arrays and their elimination via a source-to-source transformation.

Annotated Unoptimized Code. Before we can start optimizing the code from Fig. 1 using OptiTrust, we need to annotate the code with *function contracts*, *loop contracts*, as well as *ghost instructions*. A contract consists of a description of the assumptions and guarantees associated with a function or a loop, as well as a description of the side-effects that may be performed. A ghost instruction behaves, semantically, as a no-op. Its purpose is to guide the typechecker of OptiTrust, typically by altering the way the memory state is described in the Separation Logic invariants. These invariants may be exploited for guiding code transformations, and for checking their correctness.

Ghost instructions may also be used to keep track of nontrivial arithmetic reasoning involved in the typechecking process. Typically, we need to derive arithmetic inequalities, to justify that a certain range falls within the bounds of an array. The mathematical implications are recorded in the source code, e.g., $\forall i k n w. (0 \leq i < n) \wedge (i \leq k < i + w) \rightarrow (0 \leq k < n + w)$. They can be validated at any point during the optimization process using, e.g., an off-the-shelf decision procedure or SMT solver.

Besides, to ease the manipulation and typechecking of multidimensional arrays, all accesses are assumed to be written using a family of functions called `MINDEX`. For example, `D[MINDEX2(n, cn, i, c)]` denotes an access in the flat array D , of dimensions $n \times cn$, at the coordinates (i, c) . For the purpose of readability of generated programs, OptiTrust offers the option to print the same access in the form `D[i;c]`. This syntax is purposely out of the syntax of C. The form `D[i;c]` remains non-ambiguous

393 because the size information appears in the description of the Separation Logic resources at hand.
 394 We leave it to future work to support input programs written without explicit dimensions on array
 395 accesses.

396 Fig. 3 shows the same C code as in Fig. 1 augmented with contracts, relevant ghost instructions,
 397 and MINDEX accesses. The clause `__requires` contains assumptions about the input parameters. The
 398 clause `__reads` asserts that the input array `s` can be accessed in read-only mode. The clause `__modifies`
 399 asserts that the output array `D` can be modified in place. The clause `__xmodifies` describes a *loop*
 400 *contract*: it indicates not only that the i -th iteration can modify certain cells, but also that the other
 401 iterations do not access these cells. In other words, the i -th iteration has exclusive access to that
 402 cell. The “x” prefix in `__xmodifies` stands for “exclusive”.

403 In particular, the outer loop on i is annotated with a clause involving an iteration construct:
 404 `__xmodifies("for c in 0..cn -> &D[MINDEX2(n, cn, i, c)] ~ Cell")`. This clause indicates that the i -th
 405 iteration of that outer loop requires exclusive access to all the cells in the i -th row of the destination
 406 array `D`. Further in the paper, this same resource may also be written using the corresponding math
 407 notation, as: $\star_{c \in 0..cn} (\&D[i;c] \sim \text{Cell})$, where the star symbol is called *iterated separating conjunction*
 408 in Separation Logic. The iteration construct is also used to define the `Matrix2` predicate, which
 409 describes a 2D range of individual cells. Concretely, the resource $D \sim \text{Matrix2}(n, cn)$ is equivalent to
 410 $\star_{i \in 0..n} \star_{c \in 0..cn} (\&D[i;c] \sim \text{Cell})$, which covers all the $n \times cn$ cells of the matrix `D`.

411 The lines introduced by `__ghost_begin`, `__ghost_end`, or sometimes just `__ghost` correspond to *ghost*
 412 *instructions*: no-ops whose purpose is to change the view on the resources. The need for ghost
 413 instructions is standard in Separation Logic frameworks. The specialized keywords `__ghost_begin`
 414 and `__ghost_end` materialize a pair of ghost instructions that are the reciprocal of one another.
 415 For example, the ghost *focus* operation allows recovering a single memory cell from the array `s`,
 416 isolating $\&S[k;c] \sim \text{Cell}$ from $\star_{j \in 0..n+w-1} \star_{c \in 0..cn} (\&S[j;c] \sim \text{Cell})$. Technically, the focus involves
 417 read-only fractions and a “magic wand” describing the remaining cells. The matching `__ghost_end`
 418 pseudo-instruction applies the symmetrical operation, recovering the original resource. In the
 419 future, we could try to rely on heuristics for automatically inferring certain ghost operations, and
 420 reduce the number of such ghost operations that need to be explicitly provided by the programmer.
 421

422 *Optimization Script Syntax.* Fig. 4 shows our script for generating the optimized code of Fig. 2
 423 starting from the annotated unoptimized code of Fig. 3. In OptiTrust, optimizations are dictated
 424 by means of a script written in the OCaml programming language. For the reader not familiar
 425 with OCaml, $f \ x \ y$ denotes the call of f on the arguments x and y ; the symbol `-` is used to provide
 426 optional (or named) arguments; $[x; y; z]$ denotes a list; (x, y, z) denotes a tuple; $x \ ^ \ y$ denotes
 427 a string concatenation; and `let f x = t in` introduces a local function.

428 A transformation script consists of a series of calls to functions from the OptiTrust library.
 429 Each call may depend on a number of arguments controlling the transformations. By convention,
 430 the last argument of a transformation always denotes a *target*. Before explaining the working
 431 of targets, we first present the transformations involved in our script from Fig. 2. `Reduce.intro`
 432 introduces a map-reduce operation for computing the sum over a segment. `Reduce.elim` eliminates
 433 a map-reduce into an explicit summation. `Reduce.slide` performs a sliding window optimization
 434 on a map-reduce computation. `Specialize.variable_multi` introduces a cascade of if-statements
 435 for testing specific variable values. `Loop.collapse` takes two nested loops and replaces them with
 436 a single loop that iterates over the product space. `Loop.swap` takes two nested loops and swaps them.
 437 `Variable.elim_reuse` takes two variables with equal values and eliminates the second variable.
 438 `Loop.shift_range` and `Loop.scale_range` allow altering the iteration range of a loop. `Loop.unroll`
 439 unrolls a loop with a statically known number of iterations. `Loop.fusion_targets` fuses targeted
 440 loops into a single one. `Instr.gather_targets` reorders instructions in a sequence to make the
 441

```

442 Reduce.intro [cVarDef "s"];
443 Specialize.variable_multi ~mark_then:fst ~mark_else:"anyw"
444   ["w", int 3; "w", int 5] [cFunBody "rowSum"; cFor "i"];
445 Reduce.elim ~inline:true [nbMulti; cMark "w"; cCall "reduce_spe1"];
446 Loop.collapse [nbMulti; cMark "w"; cFor "i"];
447 Loop.swap [nbMulti; cMark "anyw"; cFor "i"];
448 Reduce.slide ~mark_alloc:"acc" [nbMulti; cMark "anyw"; cArrayWrite "D"];
449 Reduce.elim [nbMulti; cMark "acc"; cCall "reduce_spe1"];
450 Variable.elim_reuse [nbMulti; cMark "acc"];
451 Reduce.elim ~inline:true [nbMulti; cMark "anyw"; cFor "i"; cCall "reduce_spe1"];
452 Loop.shift_range (StartAtZero) [nbMulti; cMark "anyw"; cFor "i"];
453 Loop.scale_range ~factor:(trm_find_var "cn" []) [nbMulti; cMark "anyw"; cFor "i"];
454 Specialize.variable_multi ~mark_then:fst ~mark_else:"anycn" ~simpl:custom_specialize_simpl
455   ["cn", int 1; "cn", int 3; "cn", int 4] [cMark "anyw"; cFor "c"];
456 Loop.unroll [nbMulti; cMark "cn"; cFor "c"];
457 Target.foreach [cMark "cn"] (fun c ->
458   Loop.fusion_targets ~into:FuseIntoLast [nbMulti; c; cFor "i" ~body:[cArrayWrite "D"]];
459   Instr.gather_targets [c; cStrict; cArrayWrite "D"];
460   Loop.fusion_targets ~into:FuseIntoLast [nbMulti; c; cFor ~stop:[cVar "w"] "i"];
461   Instr.gather_targets [c; cFor "i"; cArrayWrite "D"]; );
462 Loop.shift_range (ShiftBy (trm_find_var "c" [cMark "anycn"]))
463   [cMark "anycn"; cFor ~body:[cArrayWrite "D"] "i"];
464 Cleanup.std ();

```

Fig. 4. Optimization script for the OpenCV case study.

targeted instructions consecutive. `Cleanup.std` eliminates all dependencies on the OptiTrust header file and performs arithmetic simplifications in order to produce conventional C syntax as final output.

Targets. A target provides a way to concisely and robustly refer to one or several code locations, at which to apply a transformation. The construct `Target.foreach`, visible in Fig. 2, can also be used to explicitly iterate over several code locations. A target consists of a list of constraints (prefixed by “c”) that is satisfied by code paths that go through nodes satisfying each constraint, in the given order. For example, the constraint `cFunBody "rowSum"` requires visiting a function definition with the name “rowSum”. The constraint `cFor "c"` requires visiting a for loop over an index with the name “c”. The constraint `cMark "cn"` requires visiting an AST node that carries the mark “cn”. Such marks are introduced by transformations, on demand of the programmer.

Constraints may also take targets as arguments: `cFor "i"~body:[cArrayWrite "D"]` requires visiting a for loop over an index with the name “i”, whose body also contains a write on the array `D`. Besides, targets may include special modifiers. The modifier `nbMulti` indicates that the programmer expects to find not one but multiple AST nodes that match this target. The modifier `tBefore`, which appears in the other two case studies, allows targeting the interstice before an instruction.

Interactive Visualization. Each step of an evaluation script may be executed interactively: with the cursor on a line, the OptiTrust user can press a shortcut key in their code editor to visualize the *diff* associated with the transformation on that line. Fig. 5 shows the *diff* associated with the `Loop.scale_range` transformation that appears near the middle of the script from Fig. 4. This transformation reindexes a loop. In the present example, it modifies the indexing from `for (int i = 0; i < w; i++)` to `for (int i = 0; i < cn*w; i+=cn)`, and replaces every occurrence of the index `i` with the expression `exact_div(i,cn)`. In particular, the array access `s[i,c]`, which is

```

491 | 29     for (int c = 0; c < cn; c++) {
492 | 30         ST s = (ST)0;
493 | 31     for (int i = 0; i < w; i++) {
494 | 32         s = s + (ST)S[i; c];
495 | 33     }
496 | 34     D[0; c] = s;
497 | 35     for (int i = 0; i < n - 1; i++) {
498 | 36         s = s + (ST)S[i + w; c] - (ST)S[i; c];
499 | 37     }
500 | 38     D[i + 1; c] = s;
501 | 39 }
502 | 39 }
503 |
504 | 29     for (int c = 0; c < cn; c++) {
505 | 30         ST s = (ST)0;
506 | 31     for (int i = 0; i < cn * w; i += cn) {
507 | 32         s = s + (ST)S[exact_div(i, cn); c];
508 | 33     }
509 | 34     D[0; c] = s;
510 | 35     for (int i = 0; i < cn * (n - 1); i += cn) {
511 | 36         s = s + (ST)S[exact_div(i, cn) + w; c] -
512 | 37             (ST)S[exact_div(i, cn); c];
513 | 38         D[exact_div(i, cn) + 1; c] = s;
514 | 39     }
515 | 40 }

```

Fig. 5. Diff for the `Loop.scale_range` transformation that appears near the middle of the script from Fig. 4. The printer uses the abbreviation for `MINDEX`, e.g., `S[i; c]` corresponds to `S[MINDEX2(n+w-1, cn, i, c)]`. OptiTrust can also produce a more verbose diff that includes contracts, ghost instructions, and all `MINDEX` arguments.

an abbreviation for `S[MINDEX2(n + w - 1, cn, i, c)]`, becomes `S[exact_div(i, cn); c]`, which is an abbreviation for `S[MINDEX2(n + w - 1, cn, exact_div(i, cn), c)]`. The final cleanup step of our script unfolds the definition of `MINDEX2` to obtain `S[cn * exact_div(i, cn) + c]`, then applies an arithmetic simplification to obtain the index `S[c + i]`. The latter expression appears in the final code presented in Fig. 2. Additionally, OptiTrust can produce a complete *execution trace* in the form of an interactive tree. This tree reports the diff not only for every transformation visible in the script, but also for all the internal transformations that are leveraged in the process.³

Validity Checks. The transformation script from Fig. 4 consists of *combined* transformations, whose evaluation triggers the application of a chain of *basic* transformations. As said earlier, basic transformations are those that directly modify the AST, whereas combined transformations are defined as the composition of basic transformations. For every basic transformation being applied, OptiTrust checks that this transformation preserves the semantics of the program, by leveraging resource typing information. Because the checks performed by OptiTrust depend on resource typing, every intermediate program must typecheck. In particular, if a transformation modifies the code, it may need to also modify the annotations, such as the loop contracts and the ghost instructions. Correctness criteria and preservation of typing are discussed in details in Section 6.

2.2 The Particle Simulation Case Study

Particle-In-Cell (PIC) is a technique commonly used to simulate plasma, where charged particles are in motion, by approximating the charge distribution using a grid. Our case study is inspired by the work from Barsamian et al. [2018], who present a PIC implementation featuring state-of-the-art optimizations. In the present case study, we consider a simplified PIC simulation, focusing on the computations associated with one particular cell of the grid. Our goal is to illustrate how OptiTrust can be used to derive a certain number of transformations ubiquitous in particle simulation as well as other physics simulation code.

Unoptimized Code. Fig. 6 shows the unoptimized simulation kernel that we consider. A number of particles, all with the same mass and charge, move inside a cubic cell. For simplicity, we assume in this case study that the particles do not leave the cube. The initial position and speed of every particle is given. Positions are described with values in the range $[0, 1]$, for each axis. We assume that the particles do not affect each other, and that an external electric field affects the acceleration of the particles.⁴ The electric field is described by 8 vectors, one per corner of the cell. The electric

³The traces showing the diff for every major step of the script can be browsed online at:

<https://www.chargueraud.org/softs/optitrust/traces/index.html>. Due to their large size, the traces that include all the substeps are only available by constructing them using a local installation of OptiTrust.

⁴Note that this is a simplification compared to Barsamian et al. [2018], as they also optimize code for the “charge deposit”.

```

540 void simulate(const vect* fieldAtCorners,
541             const int nbSteps, const double deltaT,
542             const double pCharge, const double pMass,
543             const int nbPart, particle* part) {
544     __reads("fieldAtCorners ~> Matrix1(nbCorners)");
545     __modifies("part ~> Matrix1(nbPart)");
546     for (int idStep = 0; idStep < nbSteps; idStep++) {
547         for (int idPart = 0; idPart < nbPart; idPart++) {
548             // Each particle is updated at each time step
549             __xmodifies("&part[MINDEX1(nbPart, idPart)] ~> Cell");
550             __ghost_begin(part, particle_open, "&part[MINDEX1(nbPart, idPart)]");
551             // Interpolate the field based on the position relative to the corners of the cell
552             double* const coeffs = MALLOC1(double, nbCorners);
553             compute_corner_interpolation_coeffs(part[MINDEX1(nbPart, idPart)].pos, coeffs);
554             const vect fieldAtPos = matrix_vect_mul(coeffs, fieldAtCorners);
555             free(coeffs);
556             // Compute the acceleration: F = m*a and F = q*E gives a = q/m*E
557             const vect accel = vect_mul(pCharge / pMass, fieldAtPos);
558             // Compute the new speed and position for the particle
559             const vect speed2 = vect_add(part[MINDEX1(nbPart, idPart)].speed, vect_mul(deltaT, accel));
560             const vect pos2 = vect_add(part[MINDEX1(nbPart, idPart)].pos, vect_mul(deltaT, speed2));
561             // Update the particle
562             part[MINDEX1(nbPart, idPart)].speed = speed2;
563             part[MINDEX1(nbPart, idPart)].pos = pos2;
564             __ghost_end(part);
565         } } }

```

Fig. 6. Unoptimized code for the particle simulation case study, with resource annotations.

field that applies at a given position inside the cubic cell is obtained by linearly interpolating the vectors associated with the corners—a standard technique in particle-in-cell (PIC) simulations.

The simulation proceeds as follows. At each time step, all the particles are updated. For a given particle, its speed is first updated, based on the value of the acceleration at the position of this particle. Then, the position of the particle is updated, based on its speed. Observe how, in Fig. 6, these updates are described at a high-level of abstraction, using vector operations, as well as a matrix-vector product for computing the interpolation. The auxiliary function `compute_corner_interpolation_coeffs` computes the interpolation coefficients associated with the position of the particle.

Optimized Code. Fig. 7 shows our optimized code for the function `simulate`. Two preliminary transformations are applied. First, auxiliary functions are inlined. In particular, the first 14 lines of the loop on `idPart` visible in the optimized code (involving the variables `rX`, `rY`, `rZ`, as well as `cX`, `cY`, `cZ`) correspond to the code inlined from `compute_corner_interpolation_coeffs`, whose implementation was not shown in Fig. 6. Second, the allocation of the array `coeffs`, used to store the interpolation coefficients, is moved outside the loop.⁵ Then, two key optimizations are applied.

First, the vector and matrix operations are replaced with operations over individual fields (named `pos.x`, `pos.y`, `pos.z`, `speed.x`, `speed.y`, and `speed.z`). Moreover, local vector variables are replaced with families of variables (e.g., `fieldAtPos_x`, `fieldAtPos_y`, and `fieldAtPos_z`).

Second, a *scaling* transformation is applied on the data in order to simplify the arithmetic computations that need to be performed at every time step. To understand how this scaling optimization works, consider a particle. For simplicity, let us focus on its behavior on the x -coordinate. At a given time step, its speed, written v , and its position, written x , are updated according to the formulae: $a = qE/m$ and $v += a\Delta_t$ and $x += v\Delta_t$. Here, E denotes the electric field interpolated at the location of this particle. The constants q , m , and Δ_t corresponds to the program variables `pCharge`, `pMass`, and `deltaT`, respectively. The idea is to store not the values of E and v , but

⁵In the full-featured Particle-in-Cell code Barsamian et al. [2018], the array `coeffs` is entirely eliminated by further optimizations, which generate large-size arithmetic expressions that may then be processed by vector instructions.

```

589 void simulate(const vect* fieldAtCorners,
590             const int nbSteps, const double deltaT,
591             const double pCharge, const double pMass,
592             const int nbPart, particle* part) {
593     const double fieldFactor = deltaT * deltaT * pCharge / pMass;
594     vect* const lFieldAtCorners = (vect*) malloc(nbCorners * sizeof(vect));
595     for (int i = 0; i < nbCorners; i++) {
596         lFieldAtCorners[i].x = fieldAtCorners[i].x * fieldFactor;
597         lFieldAtCorners[i].y = fieldAtCorners[i].y * fieldFactor;
598         lFieldAtCorners[i].z = fieldAtCorners[i].z * fieldFactor;
599     }
600     for (int i = 0; i < nbPart; i++) {
601         part[i].speed.x *= deltaT;
602         part[i].speed.y *= deltaT;
603         part[i].speed.z *= deltaT;
604     }
605     double* const coeffs = (double*) malloc(nbCorners * sizeof(double));
606     for (int idStep = 0; idStep < nbSteps; idStep++) {
607         for (int idPart = 0; idPart < nbPart; idPart++) {
608             const double rX = part[idPart].pos.x;
609             const double rY = part[idPart].pos.y;
610             const double rZ = part[idPart].pos.z;
611             const double cX = 1. - rX;
612             const double cY = 1. - rY;
613             const double cZ = 1. - rZ;
614             coeffs[0] = cX * cY * cZ;
615             coeffs[1] = cX * cY * rZ;
616             coeffs[2] = cX * rY * cZ;
617             coeffs[3] = cX * rY * rZ;
618             coeffs[4] = rX * cY * cZ;
619             coeffs[5] = rX * cY * rZ;
620             coeffs[6] = rX * rY * cZ;
621             coeffs[7] = rX * rY * rZ;
622             double fieldAtPos_x = 0.;
623             double fieldAtPos_y = 0.;
624             double fieldAtPos_z = 0.;
625             for (int k = 0; k < nbCorners; k++) {
626                 fieldAtPos_x += coeffs[k] * lFieldAtCorners[k].x;
627                 fieldAtPos_y += coeffs[k] * lFieldAtCorners[k].y;
628                 fieldAtPos_z += coeffs[k] * lFieldAtCorners[k].z;
629             }
630             const double speed2_x = part[idPart].speed.x + fieldAtPos_x;
631             const double speed2_y = part[idPart].speed.y + fieldAtPos_y;
632             const double speed2_z = part[idPart].speed.z + fieldAtPos_z;
633             part[idPart].pos.x += speed2_x;
634             part[idPart].pos.y += speed2_y;
635             part[idPart].pos.z += speed2_z;
636             part[idPart].speed.x = speed2_x;
637             part[idPart].speed.y = speed2_y;
638             part[idPart].speed.z = speed2_z;
639         }
640     }
641     free(coeffs);
642     for (int i = 0; i < nbPart; i++) {
643         part[i].speed.x /= deltaT;
644         part[i].speed.y /= deltaT;
645         part[i].speed.z /= deltaT;
646     }
647     free(lFieldAtCorners);
648 }

```

Fig. 7. Optimized code for the particle simulation case study.

```

638 let ctx = cFunBody "simulate_single_cell" in
639 let find_var n = trm_find_var n [ctx] in
640 let vect = typ_find_var "vect" [ctx] in
641 let particle = typ_find_var "particle" [ctx] in
642 let dims = ["x"; "y"; "z"] in
643 Matrix.local_name_tile ~var:"fieldAtCorners"
644 ~elem_ty:vect ~uninit_post:true ~mark_load:"loadField"
645 ~local_var:"lFieldAtCorners" [ctx; cFor "idStep"];
646 Function.inline_multi [ctx; cCalls ["cornerInterpolationCoeff"; "matrix_vect_mul"; "vect_add"; "vect_mul"]];
647 Variable.inline_and_rename [ctx; cVarDef "fieldAtPos"];
648 Record.split_fields ~tys:[particle; vect] [tSpanSeq [ctx]];
649 Record.to_variables [ctx; cVarDefs ["fieldAtPos"; "pos2"; "speed2"; "accel"]];
650 let deltaT = find_var "deltaT" in
651 Variable.insert ~name:"fieldFactor" ~value:(trm_mul (trm_mul deltaT deltaT) (trm_exact_div (find_var "pCharge")
652 (find_var "pMass"))) [ctx; tBefore; cVarDef "lFieldAtCorners"];
653 let scaleFieldAtPos d =
654   Accesses.scale_var ~factor:(find_var "fieldFactor") [nbMulti; ctx; cVarDef ("fieldAtPos_" ^ d)] in
655 List.iter scaleFieldAtPos dims;
656 let scaleSpeed2 d = Accesses.scale_immutable ~factor:deltaT [nbMulti; ctx; cVarDef ("speed2_" ^ d)] in
657 List.iter scaleSpeed2 dims;
658 let scaleFieldAtCorners d =
659   let address_pattern = Trm.(struct_access (array_access (find_var "lFieldAtCorners") (pattern_var "i"))) d in
660   Accesses.scale ~factor:(find_var "fieldFactor") ~address_pattern ~uninit_post:true
661     [ctx; tSpan [tBefore; cMark "loadField"] [tAfter; cFor "idStep"]] in
662 List.iter scaleFieldAtCorners dims;
663 let scaleParticles d =
664   let address_pattern =
665     Trm.(struct_access (struct_access (array_access (find_var "part") (pattern_var "i")) "speed") d) in
666   Accesses.scale ~factor:deltaT ~address_pattern ~mark_preprocess:"partsPrep" ~mark_postprocess:"partsPostp"
667     [ctx; tSpanAround [cFor "idStep"]]; in
668 List.iter scaleParticles dims;
669 List.iter Loop.fusion_targets [[cMark "partsPrep"]; [cMark "partsPostp"]];
670 Variable.unfold ~at:[cFor "idStep"] [cVarDef "fieldFactor"];
671 Variable.inline [ctx; cVarDefs (Tools.cart_prod (^) ["accel_"; "pos2_"] dims)];
672 Arith.(simpls_rec [expand; gather_rec]) [ctx];
673 Loop.hoist_alloc ~indep:["idStep"; "idPart"] ~dest:[tBefore; cFor "idStep"] [cVarDef "coeffs"];
674 Cleanup.std ();

```

Fig. 8. Optimization script for the particle simulation case study.

667 instead the values E' and v' defined as: $E' = qE\Delta_t^2/m$ and $v' = \Delta_t v$. The interest is that the speed and
668 position updates at a given time step are now described using much simpler formulae that avoid the
669 need for computing multiplications: $v' += E'$ and $x += v'$. To implement this scaling transformation,
670 the components of the field speed of the array `part` are multiplied, in-place, by a factor Δ_t before
671 starting the simulation; symmetrically, at the end of the simulation, the values are divided by Δ_t . For
672 the electric field array, which is read-only, the scaling factor is applied on an auxiliary array named
673 `lFieldAtCorners`, obtained by multiplying the values of `fieldAtCorners` by $q\Delta_t^2/m$. (An in-place update
674 would be disallowed because the array `fieldAtCorners` is described using a read-only permission.)
675 By linearity of the interpolation computations, this scaling propagates to the values computed for
676 the electric field at the particle location (`fieldAtPos_x`, `fieldAtPos_y`, and `fieldAtPos_z`). Note that we
677 currently treat floating-point numbers as real numbers during such transformations—reasoning
678 about precision in the optimized code is an orthogonal challenge, which we leave to future work.
679

680 *Optimization Script.* Fig. 8 shows our optimization script. Let us describe the keys steps. The
681 transformation `Function.inline_multi` inlines auxiliary functions, in particular vector operations.
682 `Record.split_fields` turns record assignments operations into per-field assignment operations.
683 `Variable.insert` inserts a definition for the multiplicative factor $q\Delta_t^2/m$, which is applied to the
684 electric field. `Accesses.scale` (as well as `scale_var` and `scale_mut`) apply the relevant multiplicative
685 factors on the values stored in the various data structures at hand. Crucially, the correctness
686

of the scaling transformation relies on the knowledge that the same arrays are not accessed by means of other (aliased) pointers. The verification of this property relies on the Separation Logic information computed during typechecking. `Loop.fusion_targets` fuses the several loops that applied per-field scaling. `Variable.unfold` reveals the definition of a variable at certain of its occurrences. `Variable.inline` eliminates a variable definition, replacing all occurrences with the definition. `Loop.hoist_alloc` pulls the allocation of the `coeffs` array outside the loop. `Cleanup.std` applies final simplifications, as previously explained.

All these transformations refer to targets, whose purpose is to match AST subtrees. In future work, we look forward to improve the conciseness of certain targets.⁶

Benefits of using OptiTrust. Applied mathematicians commonly write optimized code such as that of Fig. 7 by hand. Revealing the x , y and z coordinates triples the size of the code, and applying a scaling transformation by hand is a highly error-prone task. The aim of OptiTrust is to provide them with an alternative route, more productive and more trustworthy. As we have already explained, for each of the transformations being applied, OptiTrust exploits the Separation Logic invariants to check criteria that guarantee that the transformations preserve the semantics of the code.

2.3 The Matrix-Multiply Case Study

TVM [Chen et al. 2018] is the state-of-the-art, industrial-strength, semi-automatic compiler for machine learning. The TVM tutorial presents an optimization script⁷ (a.k.a. *schedule*) for optimizing a matrix multiplication function, specialized for square matrices of size 1024. This script has been carefully tuned to produce code optimized for specific Intel CPUs. On a 4-core Intel i7-8665U CPU with AVX2 support, the TVM experts thereby achieve a speedup of 150× over a totally naive, sequential implementation of matrix multiply.⁸ The aim of this third case study is to demonstrate the ability of OptiTrust to produce code that matches the performance delivered by TVM. More precisely, we show that we are able to generate code that features the exact same optimization patterns as in the TVM case study, with a reasonably short transformation script.

Unoptimized Code. Fig. 9 shows the unoptimized and annotated matrix-multiply code that we take as input. Note that some annotations could be inferred automatically with additional tooling.

Optimized Code. TVM output code is expressed not as C code, but directly in the intermediate representation of LLVM. We manually inspected the TVM schedule, intermediate representation, and LLVM IR output to infer what C code we should generate. The code we produce using OptiTrust is shown in Fig. 10. Compared with the naive code from Fig. 9, the optimized code from Fig. 10 integrates 7 key optimizations:

- (1) The body of the generic matrix multiply function `mm` is specialized to the size 1024.
- (2) An auxiliary matrix named `pB` is allocated to store the transpose of the matrix `B`. The introduction of this auxiliary matrix induces a cost for the initial copy, but then greatly improves the memory access patterns.
- (3) The matrices are processed by blocks of size 32: each loop over a range of size 1024 is replaced with 2 loops each of range 32. Blocking improves locality in matrix-multiply.

⁶For example, our script evaluates, for each dimension d , the target `struct_access (array_access (find_var "lFieldAtCorners"))(pattern_var "i"))d`. With a concrete syntax for expressing patterns based on their string representation, we could presumably shorten the target to `"lFieldAtCorners[?i].{x|y|z}"`.

⁷https://github.com/apache/tvm/blob/v0.19.0/gallery/how_to/optimize_operators/opt_gemm.py

⁸The 150× speed up achieved using TVM does not quite match the 204× speedup achieved by the proprietary Intel’s MKL, a library manually optimized by Intel’s experts. Yet, keep in mind that the MKL provides optimized implementation for a fixed set of functions, whereas the TVM compiler can be used to optimize entire classes of functions. We leave it to future work to investigate the extent to which OptiTrust could be used to derive code that matches the performance of MKL.

```

736 void mm(float* C, float* A, float* B, int m, int n, int p) { // naive matrix-multiply
737   __reads("A ~> Matrix2(m, p), B ~> Matrix2(p, n)");
738   __modifies("C ~> Matrix2(m, n)");
739   for (int i = 0; i < m; i++) {
740     __xmodifies("for j in 0..n -> &C[MINDEX2(m, n, i, j)] ~> Cell");
741     for (int j = 0; j < n; j++) {
742       __xmodifies("&C[MINDEX2(m, n, i, j)] ~> Cell");
743       float sum = 0.0f;
744       for (int k = 0; k < p; k++) {
745         __ghost_begin(focusA, matrix2_ro_focus, "A, i, k");
746         __ghost_begin(focusB, matrix2_ro_focus, "B, k, j");
747         sum += A[MINDEX2(m, p, i, k)] * B[MINDEX2(p, n, k, j)];
748         __ghost_end(focusA);
749         __ghost_end(focusB);
750       }
751       C[MINDEX2(m, n, i, j)] = sum;
752     }
753   }
754 }
755 void mm1024(float* C, float* A, float* B) { // specialization to 1024x1024 matrices
756   __reads("A ~> Matrix2(1024, 1024), B ~> Matrix2(1024, 1024)");
757   __modifies("C ~> Matrix2(1024, 1024)");
758   mm(C, A, B, 1024, 1024, 1024);
759 }

```

Fig. 9. Unoptimized C code for the matrix-multiply case study, using flat arrays and resource annotations.

- (4) Results are not accumulated into a scalar accumulator, but instead into a stack-allocated array named `sum` of size 32×32 that contains all scalar accumulators for a block.
- (5) Around the inner vectorized loops, the locally relevant row of `sum` is promoted to a smaller array `s` that can be mapped onto a few 256-bit vector registers. On every `i` iteration, two `memcpy` operations are used for synchronizing `s` with `sum`.
- (6) The various loops are reordered in a particular manner, both to improve cache locality and to enable parallelization. The outermost loops are executed in parallel by several cores. The instructions of the inner loop are parallelized by means of SIMD operations.
- (7) The 4 loops tagged as `#pragma omp simd` in Fig. 10 are very similar. However, if we attempt to factorize them into a loop with 4 iterations, then Intel’s compiler (ICX) produces slower code. Unfolding the loops as shown makes relying on unrolling heuristics unnecessary.

Again, this particular set of optimizations directly comes from the TVM case study. We demonstrate how to reproduce the same optimizations using OptiTrust.

Optimization Script. Fig. 11 shows our optimization script, which consists of only 10 lines. Internally, though, the high-level transformations mentioned in the script trigger the application of 55 basic transformations. An illustrative example is the call to `Loop.reorder_at` on Line 4 of Figure 11. This combined transformation takes as argument a specific instruction (referred to as “an instruction of the form +=”) as well as a description of the desired order for the loops that surround this instruction (the list `["bi"; "bj"; "bk"; "i"; "k"; "j"]`). The `reorder` transformation iteratively “brings down” the loops that need to be swapped closer to the instruction, starting from the innermost loops, and processing the loops until the outermost one. The call to `reorder_at` in our script triggers a total of 4 *loop swaps*, 6 *loop fissions*, and 2 *loop hoist* operations. In particular, the effect of these 2 hoist operations is to turn local variable named `sum` in Fig. 9 into the 2D-array named `sum` in Fig. 10.


```

785 void mm1024(const float* A, const float* B, float* C) {
786   float* pB = (float*)malloc(1048576 * sizeof(float));
787   #pragma omp parallel for
788   for (int bj = 0; bj < 32; bj++) {
789     for (int bk = 0; bk < 256; bk++) {
790       for (int k = 0; k < 4; k++) {
791         for (int j = 0; j < 32; j++) {
792           pB[32768 * bj + 128 * bk + 32 * k + j] = B[32 * bj + 4096 * bk + 1024 * k + j]; }}}}
793   #pragma omp parallel for
794   for (int bi = 0; bi < 32; bi++) {
795     for (int bj = 0; bj < 32; bj++) {
796       float* sum = (float*)malloc(1024 * sizeof(float));
797       for (int i = 0; i < 32; i++) {
798         for (int j = 0; j < 32; j++) {
799           sum[32 * i + j] = 0.f; }}
800       for (int bk = 0; bk < 256; bk++) {
801         for (int i = 0; i < 32; i++) {
802           float s[32];
803           memcpy(&s[0], &sum[32 * i], 32 * sizeof(float));
804           #pragma omp simd
805           for (int j = 0; j < 32; j++) {
806             s[j] += A[32768 * bi + 4 * bk + 1024 * i] * pB[32768 * bj + 128 * bk + j]; }
807           #pragma omp simd
808           for (int j = 0; j < 32; j++) {
809             s[j] += A[32768 * bi + 4 * bk + 1024 * i + 1] * pB[32768 * bj + 128 * bk + j + 32]; }
810           #pragma omp simd
811           for (int j = 0; j < 32; j++) {
812             s[j] += A[32768 * bi + 4 * bk + 1024 * i + 2] * pB[32768 * bj + 128 * bk + j + 64]; }
813           #pragma omp simd
814           for (int j = 0; j < 32; j++) {
815             s[j] += A[32768 * bi + 4 * bk + 1024 * i + 3] * pB[32768 * bj + 128 * bk + j + 96]; }
816           memcpy(&sum[32 * i], &s[0], 32 * sizeof(float)); }}
817       for (int i = 0; i < 32; i++) {
818         for (int j = 0; j < 32; j++) {
819           C[32768 * bi + 32 * bj + 1024 * i + j] = sum[32 * i + j]; }}
820       free(sum);
821     } }
822   free(pB);
823 }

```

Fig. 10. Our optimized C code for the matrix-multiply case study. This code features the same optimization patterns as the reference output of TVM.

```

816 Function.inline_def [cFunDef "mm"];
817 let tile (id, tile_size) =
818   Loop.tile (int tile_size) ~index:("b" ^ id) ~bound:TileDivides [cFor id] in
819   List.iter tile [("i", 32); ("j", 32); ("k", 4)];
820   Loop.reorder_at ~order:["bi"; "bj"; "bk"; "i"; "k"; "j"] [cPlusEq ()];
821   Loop.hoist_expr ~dest:[tBefore; cFor "bi"] "pB" ~indep:["bi"; "i"] [cArrayRead "B"];
822   Matrix.stack_copy ~var:"sum" ~copy_var:"s" ~copy_dims:1 [cFor ~body:[cPlusEq ()] "k"];
823   Loop.simd [cFor ~body:[cPlusEq ()] "j"];
824   Loop.parallel [cFunBody "mm1024"; cStrict; cFor ""];
825   Loop.unroll [cFor ~body:[cPlusEq ()] "k"];
826   Cleanup.std ();
827

```

Fig. 11. Optimization script for the matrix-multiply case study.

Comparison Against TVM. The TVM matrix-multiply case study appears in Fig. 12. We only comment on specific aspects and refer to TVM’s tutorial for further details. In TVM, input programs are written in a domain-specific language embedded in Python. Ideally, the matrix-multiply program shown on the left-hand side of Fig. 12 would replace the definitions of `pB` and `C` with a simpler definition of `C`:

833

```

834     k = tvm.reduce_axis((0, P))
835     A = tvm.placeholder((M, P))
836     B = tvm.placeholder((P, N))
837
838     pB = tvm.compute((N / 32, P, 32),
839                     lambda bj, k, j:
840                         B[k, bj * 32 + j])
841
842     C = tvm.compute((M, N),
843                   lambda i, j:
844                       sum(A[i, k] * pB[j // 32, k, j % 32],
845                           axis=k))
846
847     CC = s.cache_write(C, "global")
848     bi, bj, i, j = s[C].tile(
849         C.op.axis[0], C.op.axis[1], 32, 32)
850     s[CC].compute_at(s[C], bj)
851     i2, j2 = s[CC].op.axis
852     (kaxis,) = s[CC].op.reduce_axis
853     bk, k = s[CC].split(kaxis, factor=4)
854     s[CC].reorder(bk, i2, k, j2)
855     s[CC].vectorize(j2)
856     s[CC].unroll(k)
857     s[C].parallel(bi)
858     bj3, _, j3 = s[pB].op.axis
859     s[pB].vectorize(j3)
860     s[pB].parallel(bj3)

```

Fig. 12. TVM case study for matrix-multiply. On the left, input code in TVM’s domain specific language. On the right, TVM optimization script (a.k.a. *schedule*). Both use Python syntax.

```

847     C = tvm.compute((M, N), lambda i, j: sum(A[i, k] * B[k, j], axis=k))

```

Yet, TVM is unable to express the introduction of the transposed matrix of B , named pB , as a code transformation. The programmer therefore needs to introduce this auxiliary structure manually in the input code. Likewise, the blocking strategy needs to be hardwired in the source code on the left-hand side of Fig. 12. In contrast, our input program for matrix multiply shown in Fig. 9 builds upon standard C syntax and, most importantly, includes no optimization. Starting from a totally unoptimized reference code improves readability, trustworthiness, and maintainability. Besides, although our input code for matrix-multiply is currently expressed using explicit loops, in the future we could alternatively express it using higher-order combinators as well.

The right-hand side of Fig. 12 shows TVM’s optimization script. Our optimization script shown in Fig. 11 is not much more verbose than that of TVM. We have carefully checked that the C code produced using OptiTrust features the same optimizations as the LLVM IR code produced using TVM. To the best of our knowledge, OptiTrust is the first general-purpose optimization framework to demonstrate the ability to reproduce a case study from a state-of-the-art, specialized compiler such as TVM.

Finally, let us comment on interactivity. Guided by all the contents from Fig. 12, TVM applies a monolithic compilation pass to produce optimized code. TVM does not provide interactive, easily-readable feedback for the transformations performed. In contrast, OptiTrust applies a series of local, source-to-source transformations, manipulating programs expressed in conventional C syntax. Moreover, it provides human-readable diffs for every step and every substep involved in the optimization process. Although after the final cleanup step the optimized code contains somewhat-obfuscated flat array indices (recall Fig. 10), all the previous steps from the optimization script result in diffs where array accesses are presented as multidimensional accesses.

2.4 Evaluation of OptiTrust

Now that we have given a tour of the features of the OptiTrust framework, let us try to evaluate it against the set of desirable properties for semi-automatic code optimization listed in Section 1.1.

Generality. As pointed out in Section 1.3, this first release of OptiTrust has a number of limitations: it includes a subset of the C language, it applies to a simplified version of Separation Logic, and there remains many useful transformations to implement. Thus, OptiTrust in its current certainly form does not yet demonstrate full generality. Yet, every aspect of OptiTrust has been designed towards that goal.

Expressiveness and Control. OptiTrust supports a number of basic transformations that, taken individually, might appear relatively straightforward. However, by chaining such transformations

883 in the desired manner, the OptiTrust user is able to achieve state-of-the-art high-performance code,
884 similar to what an expert might have written by hand. Moreover, the many basic transformations
885 involved need not be explicitly invoked by the user: the use of high-level combined transformations
886 allows us to achieve expressiveness via concise scripts—recall, e.g., the call to `Loop.reorder` in the
887 matrix-multiply case study. A key feature of OptiTrust is that, at any stage in the optimization
888 process, the user remains fully in control.

889 Expressiveness also depends on the generality of the correctness criteria associated with every
890 transformation. In practice, there could be situations where the user may want to legitimately apply
891 a basic transformation, yet OptiTrust’s implementation is unable to recognize this application as
892 correct. In the short term, one option is for the user to treat this particular step as “user-trusted”, and
893 to rely on human review of the diff associated with that step. In the long term, users might be able
894 to replace a piece of code with any other piece of code that provably satisfies the same specification,
895 by leveraging a full-blown Separation Logic, possibly combined with the use of interactive proofs.

896 *Feedback.* For each step in the transformation script, OptiTrust delivers feedback in the form of
897 human-readable C syntax. The user usually only needs to read the diff against the previous code.
898 Interestingly, OptiTrust also records a trace that allows investigating all the substeps triggered by
899 a *combined* transformation. This information is critically useful when the result of a high-level
900 transformation does not match the user’s intention. Full traces can also be very useful for third-party
901 reviewing of an optimization process. Besides, a key feature of OptiTrust is its fast feedback loop.
902 The production of fast, human-readable feedback in a system with significant control is reminiscent
903 of interactive proof assistants, and of the aforementioned ATL tool [Liu et al. 2022].
904

905 *Composability.* OptiTrust transformation scripts are expressed as OCaml programs, and each
906 transformation from our library consists of an OCaml function. Because OCaml is a full-featured
907 programming language, OptiTrust users may define additional transformations at will by combining
908 existing transformations. User-defined transformations may query the abstract syntax tree (AST),
909 allowing to perform analyses before deciding what transformations to apply. Furthermore, because
910 OCaml is a higher-order programming language, transformation can take other transformations as
911 argument. We use this programming pattern for example to customize the arithmetic simplifications
912 to be performed after certain transformations.
913

914 *Extensibility.* If in need of a transformation that is not expressible as a combination of trans-
915 formations from the OptiTrust library, the user may devise a custom transformation. Because
916 OptiTrust does not rely on heuristics, adding a new transformation to OptiTrust does not impact
917 in any way the behavior of existing scripts. To define relatively simple custom transformations,
918 OptiTrust provides a term-rewriting facility based on pattern matching. For more complicated
919 transformations, one can follow the patterns employed in the OptiTrust’s library. For all custom
920 transformations, it is the programmer’s responsibility to work out the criteria under which applying
921 the transformation preserves the semantics of the code, and to adapt contracts if necessary in order
922 to produce well-typed code.

923 *Modularity.* The Separation Logic contract provided by the programmer for a function f con-
924 stitutes a complete summary of the side effects that this function may perform. Hence, when a
925 transformation operates on a piece of code that contains a call to f , the analysis involved in checking
926 the correctness of that transformation needs not traverse the implementation of f . In that sense,
927 all our analyses, including the typechecking process, are modular. This modularity has numerous
928 benefits. First, it implies that one may change the implementation of f without invalidating the
929 optimization script associated with another function g , provided that the optimization of g was not
930 relying on an inlining of the function f . Second, it means that analyses can much more easily scale
931

up to larger and more complex programs, without computation costs blowing up. Third, it makes it easier to devise clearer, more concise error messages. Indeed, in a modular system, errors depend solely on local information.

In compiler design in general, there exists a tension between modularity and optimizations, because certain key optimizations need to be applied across abstraction barriers. OptiTrust handles this tension by leaving it up to the user to decide where functions should be inlined—thereby deciding on a per-need basis where modularity should be given up to the benefits of performance.

Trustworthiness. Compilers are well-known to be incredibly hard to get 100% correct [Yang et al. 2011]. Like compilers, interactive optimization tools are highly subject to bugs. OptiTrust mitigates the risks of producing incorrect code in two ways. First, the diff of every step can be thoroughly scrutinized. Secondly, as explained in Section 1.3, we have organized the OptiTrust code base in such a way as to isolate the implementation of the *basic* transformations, which consists of transformations that directly modify the AST. Only basic transformations need to be trusted. We have been careful to systematically minimize the complexity of the interface and of the implementation of our basic transformations. All other transformations—the *combined* transformations—are *not* part of the trusted computing base (TCB).

This completes our high-level presentation of the OptiTrust framework. The remaining sections present the implementation of OptiTrust: its internal AST, its typechecking algorithm, and its transformations.

3 OPTITRUST’S INTERNAL AST

In OptiTrust, input programs written in OptiC (the targeted subset of C, augmented with annotations) are encoded into Opti λ (OptiTrust’s internal imperative λ -calculus). All code transformations are performed on that internal language Opti λ . Then, programs are decoded back into OptiC. As explained in the introduction, this approach enables OptiTrust to report the diff associated with every transformation in terms of a concise syntax familiar to the programmer.

The purpose of this section is to present Opti λ , whose constructs appear throughout the rest of the paper, from the statement of the typing rules to the description of the transformations. In Section 3.1, we present the grammar of Opti λ . In Section 3.2, we describe, at a high-level, OptiTrust’s translation between OptiC and Opti λ . Such a translation is relatively standard: C compilers generally include a phase that eliminates mutable variables and *l*-values. The specificity of our translation is that it attaches annotations on certain subterms to allow computing the reciprocal translation.

3.1 OptiTrust’s Internal AST

Fig. 13 gives the grammar of Opti λ . In this language, variables are bound by let-bindings and function definitions, and they are always immutable. Immutable variables allow for a straightforward implementation of substitution: variables may be substituted with values without concern on whether occurrences appear as right- or left-values. We next describe the grammar, starting with the less common features. The standard, call-by-value semantics, may be found in Appendix A.

Sequences. A sequence is a term that consists of a list of subterms with side effects or let-bindings, to be executed in order, and of a return value. A sequence is written $\{t_1; \dots; t_n; r\}$, where each t_i could be of the form **let** $x = t$, and where r denotes a return value for the sequence. This return value may be just the unit value (**void** in C), written \emptyset . (This presentation of sequences is similar to that found in, e.g., the Rust language.) We enforce that the expression r does not perform side-effects. In our current implementation, the result value r is syntactically restricted to be either unit or a

981	$R :=$	range ($t_{\text{start}}, t_{\text{stop}}, t_{\text{step}}$)	range for simple for-loops
982	$\pi :=$	par \cdot	optional parallel flag on simple for-loops
983	$r :=$	\emptyset x	result of a sequence
984	$t :=$	x	variables
985		b n	boolean values, and number values
986		$\{t_1; \dots; t_n; r\}$	sequence with result r
987		let $x = t$	variable definition
988		fun (a_1, \dots, a_n) $\mapsto t$	function definition
989		$t_0(t_1, \dots, t_n)$	function call
990		for $^\pi(i \in R) t$	possibly parallel, simple for-loop
991		if t_0 then t_1 else t_2	conditional
992		$\{f_1 = t_1; \dots; f_n = t_n\}$ $[t_1; \dots; t_n]$	structure and array as values
993		$t_1[t_2]$ $t.f$	projection from array/struct values
994		$t_1 \boxplus t_2$ $t \boxminus f$	address computation
995		$t_1 \boxplus t_2$ $t \boxminus f$	address computation
996		$t_1 \boxplus t_2$ $t \boxminus f$	address computation

Fig. 13. Grammar of Opti λ , the internal λ -calculus of OptiTrust. The actual abstract syntax tree moreover features placeholders for carrying type information, as well as annotations used to guide the reverse translation.

variable. We translate a statement of the form **return** t that appears in terminal position of a C function into “**let** $x = t; x$ ” where x is a fresh variable name.

A sequence $\{t_1; \dots; t_n; r\}$ introduces a lexical scope. If t_i is of the form **let** $x = t$, then the variable x may occur in any t_j for $j > i$. The variable x does not scope beyond the closing brace. We impose in Opti λ the invariant that every function body consists of a sequence block, even if the sequence contains a single instruction.

Moreover, in Opti λ , we enforce that all the instructions in a sequence have type unit. To do so, we insert calls to the built-in function “ignore” around instructions that are not of type **void** in the C code. Eliminating the feature of *implicitly ignored returned value* coming from the C language helps to simplify typechecking and transformations.

Sequences in Opti λ may also include *ghost instructions*. A ghost instruction behaves, semantically, as a no-op. It guides, however, the typechecker of OptiTrust, typically by altering the way the memory state is described in the Separation Logic invariants. These invariants may be exploited for guiding code transformations, and for checking their correctness. A key interest of our design is that it allows placing instructions *after* the point at which the return value is computed. Doing so is specifically useful for ghost instructions that depend on the result value. From the perspective of our bidirectional translation, ghost instructions are treated exactly like regular function calls.

Manipulation of Heap and Stack Cells. To account for heap-allocated data, OptiTrust provides the following standard primitive functions: $\text{heapAlloc}_{\text{Cell}_{\hat{\tau}}}$ for allocating an uninitialized cell of type $\hat{\tau}$ on the heap, get for reading a cell, set for writing a cell, and free for freeing allocated cells. As usual, a read in an uninitialized memory cell is undefined behavior. More generally, heapAlloc can be used for matrix allocation. For example $\text{heapAlloc}_{\text{Matrix}_{2,\text{int}}(5,8)}$ allocates an uninitialized matrix of 5×8 integers. Additionally, to account for stack-allocated variables, OptiTrust includes special functions. The operation $\text{stackAlloc}_{\text{Cell}_{\hat{\tau}}}()$ allocates a memory cell of type $\hat{\tau}$ on the stack without initializing its contents. The corresponding space is automatically reclaimed at the end of the surrounding sequence. Like for heapAlloc , stackAlloc can also be used to allocate matrices on the stack. The operation $\text{stackRef}(t)$ also allocates a memory cell on the stack but initializes it with t . These two special operations are meant to occur as part of a let-binding, for example **let** $x = \text{stackRef}(3)$, occurring directly within a sequence. Note that a binding **let** $x = \text{stackRef}(t)$

is equivalent to $\mathbf{let} \ x = \mathbf{stackAlloc}_{\text{Cell}_\tau}(\); \ \mathbf{set}(x, t)$ where $\hat{\tau}$ is the type of t . The two stack-allocation operators, apart from their implicit-free behavior, are treated like other primitive functions.

Possibly Parallel, Simple For Loops. The construct $\mathbf{for}^\pi(i \in \mathbf{range}(t_{\text{start}}, t_{\text{stop}}, t_{\text{step}})) \ t_{\text{body}}$ describes a *simple-for-loop*. In such a loop, the immutable variable i denotes the loop index. The loop range consists of the loop bounds and the per-iteration step, that are evaluated only once before starting the loop. Following the convention used by Python and other languages, the index goes from the *start* value inclusive to the *stop* value exclusive. If the *step* value is negative, the loop index iterates downwards. Optionally, the loop may be tagged with a *parallel* flag (i.e., setting π to **par**), thereby asserting that the loop should be treated as a parallel loop by the compiler and the runtime. This flag corresponds to the directive: `#pragma openmp parallel`. The restrictions imposed by OpenMP on the ranges of parallel for-loops essentially constraint them to fit the format $\mathbf{range}(t_{\text{start}}, t_{\text{stop}}, t_{\text{step}})$, which is the format that we use for our simple-for-loops.

Structured Data. The constructs $\{f_1 = t_1; \dots; f_n = t_n\}$ and $[t_1; \dots; t_n]$ build records and arrays as constant values. Mutable record and arrays are allocated by means of a call to the `stackAlloc` or `heapAlloc` functions. OptiTrust features 4 operations to manipulate structured data. If a corresponds to a constant array value, then the operation $a[i]$ reads the i -th cell of the array a . If, however, a corresponds to the address of a heap-allocated or a mutable stack-allocated array, then the memory address of i -th cell of the array a can be computed by the operation $t \boxplus_{\hat{\tau}} i$, where $\hat{\tau}$ denotes the type of the elements of t . This operation corresponds to the C pointer arithmetic operation $t+i$. The contents of that cell may be retrieved by evaluating $\mathbf{get}(t \boxplus_{\hat{\tau}} i)$. Likewise, reading the field f of a constant record r is described by the operation $r.f$, whereas the memory address of the field f of a record r allocated in memory is described by the operation $r \boxminus_{\hat{\tau}} f$, where $\hat{\tau}$ denotes the type of r . This operation corresponds to shifting the pointer r by the offset associated with the field f from the type $\hat{\tau}$. All these projection and address-shifting operations are here presented as constructs of the grammar. From the perspective of typechecking, however, we treat these operations like function applications.

Other Language Constructs. The other language constructs of Opti λ are standard. They include function abstraction, function calls, and conditionals. Our implementation accounts for a diversity of literal types. For simplicity, we consider in this paper only two kinds of literals: the metavariable b denotes a boolean literal (either true or false), and the metavariable n denotes an integer literal.

Other Primitive Operations. Besides the aforementioned primitive operations for manipulating heap and stack cells, Opti λ provides primitive functions that correspond to the arithmetic and boolean operators of the C language. One notable exception are the short-circuiting operators `&&` and `||` from C. We encoded them in Opti λ using conditionals, carrying annotations for guiding the reverse translation as detailed further on. Indeed, we wish to keep the simplest possible semantics for Opti λ .

Annotations. In addition to the ghost instructions presented earlier, each subterm of an Opti λ program can carry a number of extra information that do not affect the semantics in the form of annotations. Currently, our internal AST carries the following information:

- the location of the subterm in the initial source code;
- user-placed marks allow referring to subterms by name in transformation script's targets;
- Separation Logic contracts for functions and loops;
- type information for all bindings, operators, and for every subterm;
- style annotations to guide the reverse translation from Opti λ to OptiC, as described in more details in the next subsection.

1079 *Unsupported Language Features.* As mentioned earlier, the present paper aims at demonstrating
 1080 the interest of OptiTrust’s approach to code optimization. It does not aim at covering all the features
 1081 of the C language. Let us nevertheless comment on three features that we look forward to support
 1082 in the near future.

1083 For while-loops and general forms of for-loops, we plan to use an encoding into a single form of
 1084 repeat-loop. Observe that, despite the absence of general loops, the language that we consider in
 1085 this paper is Turing-complete thanks to our support for general recursive functions.

1086 To handle abrupt termination, as triggered by **break**, **continue**, and non-final **return** statements,
 1087 we need a generalization of our type system. The treatment of abrupt termination in Separation
 1088 Logic is well-understood—they are handled, for example, in the VST program verification framework
 1089 for C programs [Cao et al. 2018]. Yet, its support introduces a fair amount of additional complexity,
 1090 explaining why we have not included them in the present paper.

1091 The C language allows mutation of function arguments, whereas OptiC features only immutable
 1092 arguments. Even though mutating function arguments in C is sometimes considered bad practice,
 1093 we could support this pattern in OptiTrust by introducing an auxiliary fresh local mutable variable,
 1094 and turning the mutated argument into a constant argument.

1095 *Implementation of the AST.* The Opti λ abstract syntax tree (AST) is represented as an immutable
 1096 tree data structure. A program transformation takes as input such an immutable AST, and produces
 1097 as output another AST, which may share subtrees with the input AST. There are two major
 1098 benefits to following a purely functional programming style using immutable trees. First, this
 1099 approach avoids numerous bugs typically associated with inadvertent sharing of subtrees when
 1100 modifying data structures in-place. Second, this approach, by enabling sharing, can lead to a more
 1101 compact construction of complete execution traces, which are used for reporting to the user all the
 1102 intermediate ASTs constructed during the evaluation of the user’s transformation script.

1104 3.2 Bidirectional Translation between OptiC and Opti λ

1105 OptiTrust *encodes* OptiC input programs into the internal Opti λ language. Then, after one or several
 1106 transformations on the internal AST, OptiTrust *decodes* the program back into OptiC. In the rest of
 1107 this section, we first explain how C syntax is parsed to produce the OptiC AST. We then present the
 1108 key ideas of the encoding from OptiC to Opti λ by means of example, and explain how annotations
 1109 in Opti λ are used to ensure that a *round-trip* property holds. Additional details on our translation
 1110 may be found in a separate workshop article [Bertholon and Charguéraud 2025].

1112 *Initial Parsing.* The implementation of OptiTrust currently relies on Clang for parsing C syntax.
 1113 The Clang AST is then translated into the OptiC AST. During this initial translation, some amount
 1114 of semantically-irrelevant information may be lost. In particular, we currently do not attempt
 1115 to keep in our ASTs information about comments and spacing. Beyond this initial translation,
 1116 no more information is lost. Indeed, a *round-trip* property holds: encoding an OptiC program is
 1117 the reciprocal to decoding an Opti λ program. Crucially, a lot of style information is preserved in
 1118 normalized OptiC programs through annotations, as described next.

1119 *Annotations.* To deal with the fact that several OptiC expressions might admit the same encoding
 1120 in the Opti λ , our translation attaches annotation on certain Opti λ terms. For example, $(\star r) . f$ and $r \rightarrow f$
 1121 are both encoded as $\text{get}(r \sqcap f)$, therefore we instrument the encoding to attach a “dont-use-arrow”
 1122 annotation on the get term when translating $(\star r) . f$. As another example, the two terms $e1 \ \&\& \ e2$
 1123 and $e1 \ ? \ e2 \ : \ \text{false}$ have the same encoding, therefore we attach a “use- $\&\&$ ” on the conditional when
 1124 translating $e1 \ \&\& \ e2$. If a transformation step modifies the else-branch from false into something
 1125 else, then the annotation “use- $\&\&$ ” is ignored by the decoding operation.

1127

1128 *Encoding Scheme and Pure Variables.* The core of OptiTrust’s encoding consists of eliminating
 1129 *l*-values. For a heap allocated piece of data, a read operation $*p$ is encoded as the function call $\text{get}(p)$,
 1130 and an assignment $*p = v$ is encoded as $\text{set}(p, v)$. For stack-allocated C variables, the encoding
 1131 distinguishes two cases, *pure* and *non-pure*, depending on whether the address of the variable
 1132 needs to be manipulated. A variable x can only be *pure* if there is no assignment operation on x
 1133 and if the address of the variable x is never computed via the address-of operator.⁹ Equivalently, a
 1134 variable x can be *pure* if and only if x could have been declared with the modifiers `const register`,
 1135 in the terminology of the C standard. For such a pure variable, its definition, say `const int x = 3`,
 1136 is encoded simply as `let x = 3`. For a non-pure variable, its encoding involves a stack-allocation.
 1137 For example, the definition `int x = 3` is encoded as `let x = stackRef(3)`, the assignment `x = 4` is
 1138 encoded as `set(x, 4)`, and an occurrence of $\&x$ is encoded as x .

1139 The programmer may want to translate variables that can be pure into stack-allocated cells, to
 1140 enable further code transformations. Hence, we need to rely on a keyword (or attribute) to indicate
 1141 which variables should be translated without stack allocation. We could rely on `const register`, yet
 1142 for brevity we decided that in OptiC the keyword `const` alone would indicate the intention of the
 1143 programmer to introduce a pure variable.

1144 Fig. 14 provides an example translation.

1146	<code>const int x = 3;</code>	\longleftrightarrow	<code>let_{int} x = 3;</code>
1147	<code>f(x);</code>	\longleftrightarrow	<code>f(x);</code>
1148			
1149	<code>int z;</code>	\longleftrightarrow	<code>let_{ptr_{int}} z = stackAlloc_{int}();</code>
1150	<code>z = 6;</code>	\longleftrightarrow	<code>set(z, 6);</code>
1151	<code>const int v = z;</code>	\longleftrightarrow	<code>let_{int} v = get(z);</code>
1152			
1153	<code>int* const a = malloc(sizeof(int));</code>	\longleftrightarrow	<code>let_{ptr_{int}} a = heapAlloc_{int}();</code>
1154	<code>*a = *a + 2;</code>	\longleftrightarrow	<code>set(a, get(a) + 2);</code>
1155	<code>free(a);</code>	\longleftrightarrow	<code>free(a);</code>
1156			
1157	<code>int y = 5;</code>	\longleftrightarrow	<code>let_{ptr_{int}} y = stackRef_{int}(5);</code>
1158	<code>f(y);</code>	\longleftrightarrow	<code>f(get(y));</code>
1159	<code>y = y + 2;</code>	\longleftrightarrow	<code>set(y, get(y) + 2);</code>
1160	<code>y += 4;</code>	\longleftrightarrow	<code>inplaceAdd(y, 4);</code>
1161	<code>y++;</code>	\longleftrightarrow	<code>ignore(getThenIncr(y));</code>
1162			
1163	<code>int* const p = &y;</code>	\longleftrightarrow	<code>let_{ptr_{int}} p = y;</code>
1164	<code>*p = *p + 2</code>	\longleftrightarrow	<code>set(p, get(p) + 2);</code>
1165			
1166	<code>int* q = &y;</code>	\longleftrightarrow	<code>let_{ptr_{ptr_{int}}} q = stackRef_{ptr_{int}}(y);</code>
1167	<code>q = &z;</code>	\longleftrightarrow	<code>set(q, z);</code>
1168	<code>*q = *q + 2;</code>	\longleftrightarrow	<code>set(get(q), get(get(q)) + 2);</code>

1170 Fig. 14. Example translation from OptiC into the Opti λ . The functions `ignore`, `inplaceAdd`, and `getThenIncr`
 1171 are provided by OptiTrust’s library. The example assumes a function `void f(int)` to be defined.

1173 ⁹For example, a variable x cannot be *pure* if the code includes an occurrence of $\&x$, or an expression of the form $\&x.f$ or
 1174 $\&x[i]$. That said, a variable x could be *pure* despite occurring below an address-of operator. For example, x could be a pure
 1175 variable and appear as part of the expression $\&(x \rightarrow f)$, which is encoded as $x \square f$.

4 COMPUTING PROGRAM RESOURCES: CONTEXTS

Traditional typecheckers have a typing judgment of the form $\Gamma \vdash t : \tau$. Yet, the OptiTrust type-checker needs to account also for linear resources. Following the presentation of Separation Logic, OptiTrust’s typing judgment is written as a *triple* of the form $\{\{\Gamma\}\} t \{\{\Gamma'\}\}$. The input context Γ decomposes as $\langle E \mid F \rangle$, where E consists of *pure resources* and F consists of *linear resources*. Symmetrically, the output context Γ' contains both pure and linear resources. The pure resources from Γ' typically correspond to ghost return values and to pure postconditions. We qualify as *ghost*, any entity that is useful during program typechecking but is erased in the final executable code. Triples will be later extended in Section 5 to the form $\{\{\Gamma\}\} t^\Delta \{\{\Gamma'\}\}$, where Δ denotes a *usage map*, providing a summary explaining which resources are used by every subterm, and how they are used. This section presents the typing entities and the algorithmic typing rules, ignoring usage maps.

The section is organized as follows. Section 4.1 presents the grammar of *pure resources* and *linear resources*. Section 4.2 presents the grammar of *contexts*. Section 4.3 presents the grammar of *function contracts* and *loop contracts*. Section 4.4 presents the *entailment relation*. Section 4.5 presents the *subtraction procedure*, which corresponds to an algorithmic implementation of the entailment relation. Section 4.6 presents the typing judgment for *logical expressions*. Section 4.7 presents our algorithmic typing rules, which define the judgment $\{\{\Gamma\}\} t \{\{\Gamma'\}\}$. Finally, Section 4.8 presents soundness results.

Throughout the section, we assume a substitution operator for every entity. Concretely, given a map σ associating variable names to values, we write $\sigma(X)$ the substitution of the bindings from σ throughout X .

4.1 Grammar of Resources

As mentioned earlier, a context Γ decomposes as $\langle E \mid F \rangle$, where E contains pure resources and F contains linear resources. A pure resource describes a fact that remains true until the end of the program, or describes a variable permanently bound to a given value. Pure resources may be freely duplicated during typechecking. Linear resources describe the ownership of a given subset of the memory. Each linear resource describes a fragment of memory. Two *full* linear resources that appear in a same context must describe disjoint parts of the memory. A given full linear resource may be split into *fractional* resources, in which case several fractional linear resources may cover the same parts of memory. Subsequently, resources that were split may be joined back together. In any case, a linear resource cannot be duplicated and cannot be silently dropped. We next describe the grammar of pure and of linear resources.

Pure Resources. The pure part of a typing context contains resources that are bindings of the form “ $x : \tau$ ”, where τ corresponds either to a C type or to a *logical type*. A C type is denoted by the meta-variable $\hat{\tau}$. A logical type corresponds to a type from higher-order logic. Thus, intuitively, the pure part of a typing context Γ can be thought of as an interleaving of a traditional program typing context, which binds immutable program variables to C types, and a Coq context, which binds ghost variables to Coq types. Let us give examples of bindings that may appear in a pure context—that is, in the pure part of a context.

- “ $\tau : \text{Type}$ ” quantifies a type variable, useful for expressing polymorphism in Opti λ .
- “ $x : \tau$ ” quantifies a variable of type τ ; and “ $x : \hat{\tau}$ ” quantifies a variable with the C type $\hat{\tau}$.
- “ $f : \tau_1 \xrightarrow{\text{logic}} \tau_2$ ” quantifies a logical function, which corresponds to functions that are pure and terminating.

Syntax in C	Syntax in the theory	Description
$p \rightsquigarrow \text{Cell}$	$p \rightsquigarrow \text{Cell}_{\hat{\tau}}$	permission to access the cell at address p of type $\hat{\tau}$
$p \rightsquigarrow \text{Matrix1}(n)$	$p \rightsquigarrow \text{Matrix1}_{\hat{\tau}}(n)$	permission on an array of length n
$p \rightsquigarrow \text{Matrix2}(m, n)$	$p \rightsquigarrow \text{Matrix2}_{\hat{\tau}}(m, n)$	permission on a $m \times n$ matrix
for i in $r \rightarrow H(i)$	$\star_{i \in r} H(i)$	union of resources $H(i)$, for i in the range r
$_RO(\alpha, H)$	αH	read-only permission on H with fraction α
$_Uninit(H)$	$\text{Uninit}(H)$	permission on H that disallows reads before a write
$_Freeable(p, H)$	$\text{Freeable}(p, H)$	permission to free p by giving away the resource H

Fig. 15. Grammar of heap predicates. User-defined representation predicates are left to future work.

- “ $P : \text{Prop}$ ” quantifies an abstract proposition; and “ $Q : \tau \xrightarrow{\text{logic}} \text{Prop}$ ” quantifies an abstract logical predicate over values of type τ .
- “ $p : P$ ” quantifies a proof witness of a proposition P ; for example “ $p : i > 0$ ” captures the assumption that i is positive.
- “ $p : \text{Spec}(f, [a_1, \dots, a_n], \gamma)$ ” describes a *function specification*¹⁰ asserting that the function f expects arguments named a_i and admits the *function contract* γ .
- “ $H : \text{Hprop}$ ” quantifies an abstract heap predicate¹¹, and “ $I : \text{int} \xrightarrow{\text{logic}} \text{Hprop}$ ” quantifies an abstract invariant parameterized by a loop index.

Thereafter, to avoid confusion between the separating conjunction operation \star from Separation Logic and the star-symbol that denotes a C pointer type, we use the alternative syntax ptr_A to denote the C type A^* .

Linear Resources. The linear part of a typing context contains *resources*. A resource is described by a binding of the form “ $y : H$ ”, where H is a *heap predicate*, and where y is a name. For example, “ $y : p \rightsquigarrow \text{Cell}$ ” is a resource. This name y is used in particular to refer to resources in usage maps. A heap predicate H describes “ownership” of part of the memory. When a linear context contains several resources, each resource must describe a disjoint part of the memory. Interestingly, heap predicates guarantee the absence of hidden aliasing.

Fig. 15 summarizes the most common heap predicates, which have already been discussed in Section 2, but for which we here introduce math notation, moreover making the type annotations explicit. The resource $p \rightsquigarrow \text{Cell}_{\hat{\tau}}$ corresponds to the ownership of a single cell of type $\hat{\tau}$, located at address p . The resource $p \rightsquigarrow \text{Matrix1}_{\hat{\tau}}(n)$ is syntactic sugar for $\star_{i \in 0..n} p[i] \rightsquigarrow \text{Cell}_{\hat{\tau}}$. This resource corresponds to the ownership of the set of all the cells in the array. The big-star symbol corresponds to the *iterated separating conjunction* of Separation Logic. Likewise, $p \rightsquigarrow \text{Matrix2}_{\hat{\tau}}(n, m)$ denotes $\star_{i \in 0..n} \star_{j \in 0..m} p[i][j] \rightsquigarrow \text{Cell}_{\hat{\tau}}$. We leave it to future work to provide mechanisms allowing the user to define *representation predicates* [Reynolds 2002] for custom data types. The three heap predicates listed at the bottom of Fig. 15 are explained in the following paragraphs.

Read-Only Fractions. Following standard Separation Logic, we represent read-only resources using *fractional resources* [Boyland 2003; Jung et al. 2018a]. Intuitively, possessing a non-zero fraction of a linear resource gives read-only access to this resource. Possessing the full fraction (i.e., 1) of a resource gives read-write exclusive access to this resource. Possessing both αH and βH is equivalent to possessing the single resource $(\alpha + \beta)H$. Said differently, if we have αH at hand

¹⁰Function contracts may appear in typing contexts, while typing contexts are involved in the statement function contracts. This form of *impredicativity* is standard in higher-order Separation Logic [Charguéraud 2020b].

¹¹In formalizations of Separation Logic, Hprop is typically defined as $\text{state} \xrightarrow{\text{logic}} \text{Prop}$, where state denotes the type of a memory state, however this definition needs not be revealed to the OptiTrust user.

in the context, we can *carve out* a subfraction βH , leaving as remainder $(\alpha - \beta)H$. This splitting operation can be performed for any fraction β such that $0 < \beta < \alpha$.

Every time our typechecker requires a read-only permission on H in a context containing αH , it carves out a subfraction βH out of αH . This strategy ensures that we always keep around a fraction of the read-only resources initially available. These fractions may be useful for typing subsequent terms. When a read-only permission is returned after being used, our typing algorithm eagerly merges back βH and $(\alpha - \beta)H$ into the original form αH . Interestingly, carve-out operations may be performed in cascade, and merge-back operations can be performed in any order. To support this general pattern, we introduce the operation `CloseFracs`, which appears in our typing rules. The operation `CloseFracs` repeatedly applies the following rewrite rule:

$$(\alpha - \beta_1 - \dots - \beta_n)H \star (\beta_i - \gamma_1 - \dots - \gamma_m)H \longrightarrow (\alpha - \beta_1 - \dots - \beta_{i-1} - \gamma_1 - \dots - \gamma_m - \beta_{i+1} - \dots - \beta_n)H.$$

In general, if we start with a full permission H , that is $1H$, then whatever the order in which we carve out and merge back all the fractions of H , we ultimately recover $1H$.

Resources for Uninitialized Cells. Separation Logic can guarantee that a program never reads from an uninitialized memory cell. The traditional way to formalize this approach is as follows.

- (A1) Allocation of a memory cell at address p is specified as producing the heap predicate $p \rightsquigarrow \perp$, where \perp is a special token denoting uninitialized content.
- (A2) The specification of the read operation requires not only a fraction of a permission of the form $p \rightsquigarrow v$, but also requires the property $v \neq \perp$.

OptiTrust operates not on predicates of the form $p \rightsquigarrow v$, but on less precise predicates of the form $p \rightsquigarrow \text{Cell}$. Hence, we follow a slightly different approach for handling uninitialized cells.

- (B1) Our heap predicate $p \rightsquigarrow \text{Cell}$ denotes not only the ownership of the cell at location p but also the information that its contents is previously initialized.
- (B2) Our heap predicate $\text{Uninit}(p \rightsquigarrow \text{Cell})$ denotes the ownership of the cell p , yet without the permission to read its contents before it is initialized.
- (B3) We specify a write operation on p as consuming $\text{Uninit}(p \rightsquigarrow \text{Cell})$ and producing $p \rightsquigarrow \text{Cell}$.
- (B4) We allow a permission $p \rightsquigarrow \text{Cell}$ to be downgraded into $\text{Uninit}(p \rightsquigarrow \text{Cell})$ at any time.

The combination of (B3) and (B4) means that a write operation can also be typechecked as an operation that consumes and returns the permission $p \rightsquigarrow \text{Cell}$. More generally, as detailed further on (in Section 4.5), when our typechecker encounters a term that requires $\text{Uninit}(H)$ in a context where the plain resource H is available, it weakens H into $\text{Uninit}(H)$ on-the-fly.

We generalize the predicate to the form $\text{Uninit}(H)$ to describe uninitialized arrays and matrices. Concretely, for a matrix, $\text{Uninit}(p \rightsquigarrow \text{Matrix2}(m, n))$ corresponds to $\star_{i \in 0..n} \star_{j \in 0..m} \text{Uninit}(p[i][j] \rightsquigarrow \text{Cell})$. We do not attempt to provide a definition of $\text{Uninit}(H)$ for arbitrary H : like for read-only resources, we use uninitialized resources only for cells and groups of cells.

Permission to free. A permission of the form $\text{Freeable}(p, H)$ is obtained when p is the address returned by an allocation operation, these cells being described by the heap predicate H . The predicate H must be given back in order to invoke the free function on p .

4.2 Construction and Operations on Typing Contexts

Construction of Contexts. A context Γ takes the form $\langle E \mid F \rangle$, where E consists of a list of *pure resources* and F consists of a set of *linear resources*. In its expanded form, a context is written $\langle x_0 : \tau_0, \dots, x_n : \tau_n \mid y_0 : H_0, \dots, y_n : H_n \rangle$, where x_i denotes a pure resource of type τ_i , and y_i denotes a linear resource with heap predicate H_i . The names x_i and y_i must all be distinct. The pure part E is a

1324 *telescope*: the variable x_i may occur in any τ_j where $i < j$. Moreover, all the pure variables x_i scope
 1325 over the linear formulas H_j . The order of the linear resources is irrelevant.

1326 The pure part E of a context Γ may contain bindings of a special form, called *alias bindings*.
 1327 Such a binding takes the form “ $x_i : \tau_i := v_i$ ”. The intention is that, in presence of such an alias, our
 1328 typechecker eagerly replaces x_i with v_i during internal unification operations. An alias binding
 1329 corresponds exactly to a *local definition* in Coq. An alias binding “ $x_i : \tau_i := v_i$ ” may also be interpreted
 1330 as a conventional binding that associates x_i to a singleton type whose sole inhabitant is v_i .

1331 Following standard practice in proof assistants, variable names that are nowhere mentioned may
 1332 be hidden. For example the context $\langle p : \text{ptr}_{\text{int}}, n : \text{int}, n > 0 \mid p \rightsquigarrow \text{Cell}_{\text{int}} \rangle$ contains two anonymous
 1333 resources: $n > 0$ and $p \rightsquigarrow \text{Cell}_{\text{int}}$. Internally, though, all context items are identified by a variable
 1334 name.

1335 *Bindings of the Special Result Variable*. In contexts, we use a special variable **res** as a canonical
 1336 name to denote the result value of an expression. Therefore, if t has a non-void type τ then, in
 1337 the triple $\{\{\Gamma\}\} t \{\{\Gamma'\}\}$, this variable **res** may be bound in Γ' as an alias. The variable **res** also
 1338 appears in function contracts, to specify properties about the return value of the function. The
 1339 use of a dedicated name such as **res** is common practice in program verification tools, such as
 1340 ESC/Java [Flanagan et al. 2002] or Why3 [Filliâtre 2003].

1342 *Projection of Context Components*. We define two projection functions. For a context $\Gamma = \langle E \mid F \rangle$,
 1343 the projection “ $\Gamma.\text{pure}$ ” returns E , and the projection “ $\Gamma.\text{linear}$ ” returns F .

1345 *Syntax for Contexts with One Component*. As syntactic sugar, we define $[x_0 : \tau_0, \dots, x_n : \tau_n]$ as
 1346 $\langle x_0 : \tau_0, \dots, x_n : \tau_n \mid \emptyset \rangle$, for contexts that are entirely pure. Furthermore, we allow ourselves to write
 1347 F to mean $\langle \emptyset \mid F \rangle$, where F denotes a set of linear resources.

1348 *Separated Conjunction of Two Contexts*. We write $F_1 \star F_2$ the disjoint union of two sets of linear
 1349 resources. Furthermore, for two contexts Γ_1 and Γ_2 , we define $\Gamma_1 \otimes \Gamma_2$ as $\langle \Gamma_1.\text{pure}, \Gamma_2.\text{pure} \mid \Gamma_1.\text{linear} \star$
 1350 $\Gamma_2.\text{linear} \rangle$, assuming the variables in this result are well-scoped (that is, Γ_1 and Γ_2 have disjoint
 1351 domains and the formulas in Γ_2 are well-scoped in $\Gamma_1.\text{pure}$). Observe that $[E] \otimes F = \langle E \mid F \rangle$.

1352 *Pointwise Operators Over Linear Contexts*. Consider a linear context F of the form $(y_0 : H_0,$
 1353 $\dots, y_n : H_n)$. We define $\star_{i \in r} F$ as $(y_0 : \star_{i \in r} H_0, \dots, y_n : \star_{i \in r} H_n)$, that is, the iterated separating
 1354 conjunction distributes pointwise over the set of linear resources. Similarly, we define αF as
 1355 $(y_0 : \alpha H_0, \dots, y_n : \alpha H_n)$.

1357 *Filtering on Contexts*. We define a filtering operation, written $G \vdash X$, where G is a set of resources
 1358 (linear or pure) and X is a set of variable names. This operation computes a set of resources where
 1359 only the entries from G whose name belongs to the set X are kept. Filtering also applies to contexts:
 1360 $\langle E \mid F \rangle \vdash X$ is defined as $\langle E \vdash X \mid F \vdash X \rangle$.

1362 *Specialization of Contexts*. The specialization operation is used for example to specialize the
 1363 contract of a function for a specific call to that functions. The contract is then specialized on
 1364 the arguments, as well as on the ghost arguments, on which the function is applied. In case of a
 1365 polymorphic function, type arguments are specialized as well. The specialization operation takes
 1366 the form $\text{Specialize}_{\Gamma_0} \{\sigma\}(\Gamma)$. The definition of this operation is fairly technical, yet it is a direct
 1367 generalization of the process of typechecking function applications in higher-order logics. Rather
 1368 than presenting technical definitions, let us illustrate the specialization operation on an example.

1369 Consider a function f whose input is described by a context $\Gamma \equiv \langle A : \text{Type}, C : \text{Type}, n : \text{int},$
 1370 $p : \text{ptr}_A, b : A, c : C \mid p \rightsquigarrow \text{Matrix}1_A(n) \rangle$, where A and C are type arguments, where p and n
 1371 denote physical arguments, and where b and c are ghost arguments. Consider a function call of
 1372

1373 $\{ \{ \text{heapAlloc}_{C_{\hat{\tau}}}() \} \{ [\mathbf{res} : \text{ptr}_{\hat{\tau}}] \otimes \text{Uninit}(\mathbf{res} \rightsquigarrow C_{\hat{\tau}}) \otimes \text{Freeable}(\mathbf{res}, \mathbf{res} \rightsquigarrow C_{\hat{\tau}}) \}$
 1374 $\{ [\hat{\tau} : \text{Type}, a : \text{ptr}_{\hat{\tau}}, \alpha : \text{frac}] \otimes \alpha(a \rightsquigarrow \text{Cell}_{\hat{\tau}}) \} \{ \{ \text{get}(a) \} \{ [\mathbf{res} : \hat{\tau}] \otimes \alpha(a \rightsquigarrow \text{Cell}_{\hat{\tau}}) \}$
 1375 $\{ [\hat{\tau} : \text{Type}, a : \text{ptr}_{\hat{\tau}}, b : \hat{\tau}] \otimes \text{Uninit}(a \rightsquigarrow \text{Cell}_{\hat{\tau}}) \} \{ \{ \text{set}(a, b) \} \{ a \rightsquigarrow \text{Cell}_{\hat{\tau}} \}$
 1376 $\{ [\hat{\tau} : \text{Type}, a : \text{ptr}_{\hat{\tau}}, H : \text{Hprop}] \otimes \text{Freeable}(a, H) \otimes \text{Uninit}(H) \} \{ \{ \text{free}(a) \} \}$
 1377
 1378

Fig. 16. Contracts assigned to key primitive functions; $\hat{\tau}$ denotes a C type; a and b denote program variables. $C_{\hat{\tau}}$ is either $\text{Cell}_{\hat{\tau}}$, $\text{Matrix1}_{\hat{\tau}}(n)$ or $\text{Matrix2}_{\hat{\tau}}(m, n)$, for size expressions m and n .

1381 the form $f(7, q)$, where q is a program variable of type ptr_{int} in scope at the call site. This call
 1382 specializes n to 7 and p to q , hence it is described by a substitution $\sigma \equiv (n := 7, p := q)$. Let Γ_0 be the
 1383 context describing the pure variables bound at the call site. In particular, we have $(q : \text{ptr}_{\text{int}}) \in \Gamma_0$.
 1384 For the example considered, the specialization operation yields the context: $\langle C : \text{Type}, b : \text{int},$
 1385 $c : C \mid q \rightsquigarrow \text{Matrix1}_A(7) \rangle$. Observe how the types and arguments being specialized (namely A, n
 1386 and p) are eliminated from the pure part of the context, and the corresponding values (namely $\text{int},$
 1387 7 and q) are substituted in the entities that remain.
 1388

1389 *Renaming on Contexts.* A renaming operation is involved when the programmer explicitly spec-
 1390 ifies the names to assign to the ghost variables obtained as part of the result of a function call.
 1391 The operation $\text{Rename}\{\rho\}(\Gamma)$ renames certain keys from Γ . Here, ρ denotes a map that associates
 1392 resource names to other resource names. The keys from ρ may or may not be bound in Γ . The
 1393 values from ρ must be fresh from Γ . For example, $\text{Rename}\{x := x', y := y'\}(\langle E_1, x : \tau, E_2 \mid F \rangle)$,
 1394 where y has no occurrence in E_1, E_2 or F , evaluates to $\langle E_1, x' : \tau, x := x'(E_2) \mid x := x'(F) \rangle$. As
 1395 another example, $\text{Rename}\{y := y'\}(\langle E \mid F_1, y : H, F_2 \rangle)$ evaluates to $\langle E \mid F_1, y' : H, F_2 \rangle$.
 1396

1397 4.3 Grammar of Contracts

1398 Every function and every loop carries a contract to guide the typechecker. We next detail the
 1399 grammar of contracts.
 1400

1401 *Function Contracts.* A function definition annotated with a *function contract* γ takes the form
 1402 $\text{fun}(a_1, \dots, a_n)_{\gamma} \mapsto t$. The contract γ consists of two contexts, one for the *precondition*, written $\gamma.\text{pre}$,
 1403 and one for the *postcondition*, written $\gamma.\text{post}$. Intuitively, a function f with arguments named a_i and
 1404 with contract γ satisfies the Separation Logic triple $\{ \gamma.\text{pre} \} f(a_1, \dots, a_n) \{ \gamma.\text{post} \}$. This property is
 1405 formally captured by the proposition $\text{Spec}(f, [a_1, \dots, a_n], \gamma)$, which may appear in contexts.
 1406

1407 Technically, a function contract γ takes the form $\{ \text{pre} = \Gamma_{\text{pre}} ; \text{post} = \Gamma_{\text{post}} \}$. The precondition
 1408 Γ_{pre} must contain all the formal parameters a_i , and may refer to any of the free variables in scope.
 1409 The postcondition Γ_{post} may also refer to all these variables, as well as to the pure variables bound
 1410 in the precondition Γ_{pre} .
 1411

1412 *Contracts for Primitive Functions.* Fig. 16 gives the contracts that we axiomatize for the opera-
 1413 tions on heap cells—technically, we present not their contracts but the triples derived from their
 1414 contracts, to improve readability. These contracts illustrate key mechanisms of the formalism. A
 1415 heap allocation produces an uninitialized permission and a permission to free the allocated cells. A
 1416 write requires an uninitialized permission and returns a full permission. A read requires a read-only
 1417 permission and returns it. A free operation requires a permission to free, the associated uninitialized
 1418 permission and returns nothing. Recall that a full permission can be split into read-only resources,
 1419 and that it may be downgraded at any time into an uninitialized permission. Additionally, we can
 1420 see that bindings on \mathbf{res} appear in output contexts.
 1421

Contracts for arithmetic operations are described later on, in Section 4.6.

1422 *Contracts for Ghost Functions.* In addition to contracts for primitive heap-manipulating func-
 1423 tions, OptiTrust provides contracts for primitive ghost functions. For example, the ghost function
 1424 `swap_groups` allows swapping two iterators (iterated separating conjunctions). It is involved for
 1425 example in the loop-swap operation, which is used in our case studies (Section 2), and which is
 1426 presented further on in Section 6.5. The transformation is specified as shown below, where H is a
 1427 heap predicate that depends on the two indices i and j . The type range corresponds to a triple of
 1428 integers.

$$1430 \{[R_i : \text{range}, R_j : \text{range}, H : (\text{int}, \text{int}) \xrightarrow{\text{logic}} \text{Hprop}] \otimes_{i \in R_i, j \in R_j} (\star \star H(i, j))\} \text{swap_groups} \{ \star \star_{j \in R_j, i \in R_i} H(i, j) \}$$

1432 The OptiTrust user can define custom ghost functions to factorize repetitive resource-manipulation
 1433 patterns. Ghost functions are written and typechecked like regular C functions whose body com-
 1434 poses other ghost functions, typically through sequences and for-loops. Importantly, the body of a
 1435 ghost function does not need to be executed, and simply serves as a proof witness.

1436 *Loop Contracts.* A for-loop annotated with a *loop contract* χ takes the form **for** $(i \in R)_{\chi} \{t\}$. The
 1437 loop contract χ consists of a record structured as follows.

$$1439 \left\{ \begin{array}{l} \text{vars} = E \quad \text{Pure variables, scoping over the other contract components} \\ \text{excl} = \left\{ \begin{array}{l} \text{pre} = F_{pre} \quad \text{Resources consumed exclusively by one iteration} \\ \text{post} = F_{post} \quad \text{Resources produced exclusively by one iteration} \end{array} \right. \\ \text{shrd} = \left\{ \begin{array}{l} \text{reads} = F_{reads} \quad \text{Read only resources shared between iterations} \\ \text{inv} = F_{inv} \quad \text{Sequential invariant, threaded through iterations} \end{array} \right. \end{array} \right.$$

1440 We call E the *loop ghost variables*. The variables from E scope over F_{pre} , F_{post} , F_{reads} and F_{inv} .
 1441 We call F_{pre} the *consumed per-iteration resources* and F_{post} the *produced per-iteration resources*.
 1442 Resources in F_{pre} and in F_{post} may (and typically do) refer to the loop index. We call F_{reads} the
 1443 *shared reads*, because in practice this context consists of read-only resources. Resources in F_{reads}
 1444 cannot refer to the loop index. We call F_{inv} the *sequential invariant*. It corresponds to a standard
 1445 loop invariant in sequential Separation Logic. In this paper, we consider for simplicity that F_{inv}
 1446 does not depend on the loop index.

1447 *Parallel Loop Contracts.* A loop is parallelizable if it can be typechecked with an empty sequential
 1448 invariant F_{inv} . Hence, we say that a loop contract χ is *parallelizable*, and write $\text{parallelizable}(\chi)$,
 1449 when $\chi.\text{shrd}.\text{inv} = \emptyset$.

1457 4.4 Entailment

1458 We next introduce the *entailment* judgment, written $\Gamma \Rightarrow \Gamma'$. The entailment judgment can be used
 1459 to assert that a context Γ obtained at a given program point corresponds to a context Γ' expected
 1460 at that same point. For example, the context at the end of a function body must entail the context
 1461 described by the postcondition of this function. The entailment judgment $\Gamma \Rightarrow \Gamma'$ is a declarative
 1462 judgment, for which we will present our algorithmic implementation in the next section.

1463 The literature on Separation Logic includes two types of entailment: *linear* and *affine* entailment
 1464 relations. OptiTrust is based on a linear entailment relation, disallowing resources to be silently
 1465 “dropped”. The benefits of using linear entailment is that it allows checking the absence of memory
 1466 leaks—every piece of heap allocated data must eventually be freed.

1467 The OptiTrust entailment between two contexts:

$$1469 \langle x_0 : \tau_0, \dots, x_n : \tau_n \mid y_0 : H_0, \dots, y_m : H_m \rangle \Rightarrow \langle x'_0 : \tau'_0, \dots, x_{m'} : \tau_{m'} \mid y'_0 : H'_0, \dots, y_{m'} : H'_{m'} \rangle$$

is defined as

$$\forall x_0 : \tau_0, \dots, \forall x_n : \tau_n, \quad \text{SL}(H_0 \star \dots \star H_m) \stackrel{\text{SL}}{\Rightarrow} (\exists x'_0 : \tau'_0, \dots, \exists x'_m : \tau'_m, \text{SL}(H'_0 \star \dots \star H'_{m'}))$$

where $\stackrel{\text{SL}}{\Rightarrow}$ denotes the standard Separation Logic entailment, and where the SL function converts resources into standard Separation Logic heap predicates.

For example, the following entailments hold:

- $\langle x : \text{loc} \mid x \rightsquigarrow \text{Cell} \rangle \Rightarrow \langle \alpha : \text{frac} \mid \alpha(x \rightsquigarrow \text{Cell}), (1 - \alpha)(x \rightsquigarrow \text{Cell}) \rangle$
- $\langle y : \text{loc} \mid y \rightsquigarrow \text{Cell} \rangle \Rightarrow \langle \mid \text{Uninit}(y \rightsquigarrow \text{Cell}) \rangle$
- $\langle A : \text{loc}, n : \text{int}, m : \text{int} \mid \star_{i \in 0..n} \star_{j \in 0..m} A[i][j] \rightsquigarrow \text{Cell} \rangle \Rightarrow \langle \mid \star_{j \in 0..m} \star_{i \in 0..n} A[i][j] \rightsquigarrow \text{Cell} \rangle$
- $\langle n : \text{int}, n \text{ even} \mid \rangle \Rightarrow \langle m : \text{int}, n = 2m \mid \rangle$.

However, the entailment $\langle x : \text{loc} \mid x \rightsquigarrow \text{Cell} \rangle \Rightarrow \langle \mid \rangle$ does not hold because linear resources cannot be dropped, and the entailment $\langle x : \text{loc} \mid x \rightsquigarrow \text{Cell} \rangle \Rightarrow \langle x : \text{loc} \mid x \rightsquigarrow \text{Cell}, x \rightsquigarrow \text{Cell} \rangle$ does not hold because linear resources cannot be duplicated.

As a shorthand, we write $\Gamma \Leftrightarrow \Gamma'$ to assert that entailment holds both ways, that is, to assert that the conjunction $(\Gamma \Rightarrow \Gamma') \wedge (\Gamma' \Rightarrow \Gamma)$ holds.

4.5 Subtraction

The *subtraction* operation provides a sound (yet incomplete) algorithmic implementation of the entailment judgment. The subtraction operation not only allows checking the validity of an entailment, it also enables a certain amount of inference. At a high level, given Γ and Γ' , the subtraction operation computes the *frame*, written F , which denotes the set of linear resources such that $\Gamma \Rightarrow \Gamma' \star \langle \emptyset \mid F \rangle$. The subtraction operation also infers the instantiation map σ providing the witnesses for the instantiations of the variables that are bound (and therefore existentially quantified) in Γ' . Such a subtraction operator is found in most—if not all—practical verification frameworks based on Separation Logic.

The typing rules of OptiTrust actually make use of two variants of the subtraction operation. The *core subtraction operation*, written $\Gamma \boxminus \Gamma'$, is able to convert uninitialized resources into full resources on-the-fly, however it does not support splitting read-only resources on-the-fly. The *carving subtraction operation*, written $\Gamma \ominus \Gamma'$, extends the former with the feature of carving out a fraction of a read-only permission from Γ every time a corresponding read-only permission is requested in Γ' . (Carving was described in Section 4.1.)

The core subtraction operation $\Gamma \boxminus \Gamma'$ is formally specified as a partial operation. It may fail (that is, return \perp) if a resource in Γ' cannot be matched against a corresponding resource in Γ . Otherwise, the operation returns a result of the form (σ, F) . When $\Gamma \boxminus \Gamma' = (\sigma, F)$, then the entailment $\Gamma \Rightarrow \text{Specialize}_\Gamma\{\sigma\}(\Gamma') \star \langle \emptyset \mid F \rangle$ holds. In particular, the subtraction operation can be used to prove an entailment $\Gamma \Rightarrow \Gamma'$, by checking that $\Gamma \boxminus \Gamma'$ evaluates to (σ, \emptyset) for some σ .

The subtraction operation is implemented following a standard scheme.

- (1) The substitution map σ is initialized with bindings that associates each of the pure variables of Γ' to a fresh unification variable.
- (2) Each of the linear resources from Γ' are syntactically matched against a corresponding resource from Γ . This process may trigger unifications, resulting in partial or total resolution of certain unification variables.
- (3) If Γ' requests a linear resource of the form $\text{Uninit}(H)$, and if Γ contains the resource H , then our algorithm applies an on-the-fly weakening from H to $\text{Uninit}(H)$.
- (4) The items from Γ that remains at the end are assigned to the frame F .

The carving subtraction operation $\Gamma \ominus \Gamma'$ behaves almost like $\Gamma \boxminus \Gamma'$ but outputs a triple $(E_{\text{frac}}, \sigma, F)$ where σ and F are the same as in core subtraction and E_{frac} is a pure context for generated

	VAR	INT	BOOL	INTTYPE	BOOCTYPE
1520	$(x : \tau) \in E$				
1521	$\frac{}{E \vdash x : \tau}$	$\frac{}{E \vdash n : \text{int}}$	$\frac{}{E \vdash b : \text{bool}}$	$\frac{}{E \vdash \text{int} : \text{Type}}$	$\frac{}{E \vdash \text{bool} : \text{Type}}$
1522					
1523					
1524	PROP		HPROP		PTRTYPE
1525	P is a logical proposition		H is a heap predicate		$E \vdash A : \text{Type}$
1526	with free variables in E		with free variables in E		$\frac{}{E \vdash \text{ptr}_A : \text{Type}}$
1527	$\frac{}{E \vdash P : \text{Prop}}$		$\frac{}{E \vdash H : \text{Hprop}}$		
1528					
1529	LOGICFUN			LOGICAPP	
1530	$(E, x_1 : \tau_1, \dots, x_n : \tau_n) \vdash t : \tau$			$E \vdash t_0 : (\tau_1, \dots, \tau_n) \xrightarrow{\text{logic}} \tau$	
1531	$\frac{}{E \vdash \text{fun}(x_1 : \tau_1, \dots, x_n : \tau_n) \mapsto t : (\tau_1, \dots, \tau_n) \xrightarrow{\text{logic}} \tau}$			$\frac{\forall i \in [1..n], E \vdash t_i : \tau_i}{E \vdash t_0(t_1, \dots, t_n) : \tau}$	
1532					

Fig. 17. Selected rules defining the typing judgment for logical expressions, written $E \vdash t : \tau$. Arithmetic operations such as $t_1 + t_2$ are viewed as functions calls and are therefore handled by the rule PUREAPP.

fractions containing only bindings of the form $\alpha : \text{frac}$. At step (1), E_{frac} is initialized as an empty environment. Compared to the core subtraction, the carving subtraction refines step (2) as follows. If Γ' requests a fractional resource αH , if α is an unconstrained unification variable that denotes a fraction, and if Γ contains a fractional resource $\beta H'$ for some fraction β and where H unifies with H' , then our algorithm applies an on-the-fly splitting operation to convert $\beta H'$ into the conjunction of $\alpha' H'$ and $(\beta - \alpha') H'$ for a fresh α' added to E_{frac} . Our algorithm then adds the binding $\alpha := \alpha'$ into σ . The interest of extracting a carved fraction from βH rather than consuming the whole read-only permission βH is that the left-over fraction remains available in Γ , allowing to match other resources of the form $\alpha'' H$ that might appear in the other elements from Γ' .

4.6 Typechecking of Logical Expressions

A *logical expression* is an expression that may appear in specifications and invariants; technically, a logical expression is an expression whose evaluation terminates and does not depend on the memory state. Logical expressions include program variables (which are always immutable in the Opti λ), constant literals, logical propositions, heap predicates, C types, logical types, pure functions definitions, and pure function calls. Figure 17 shows the main typing rules for logical expressions (we omitted technical details for the treatment of dependent types). The judgment is written $E \vdash t : \tau$, where E is a pure context.

An arithmetic expression $t_1 + t_2$ can be considered as a logical expression if its two arguments are pure. The contract for addition is: $\{[a : \text{int}, b : \text{int}]\} (a + b) \{[\text{res} := a \hat{+} b : \text{int}]\}$, where $+$ denotes the addition operator from the programming language, and where $\hat{+}$ denotes the corresponding addition operator from the logic. Partial functions may also be treated as logical expressions, simply with an additional precondition. The contract for division is: $\{[a : \text{int}, b : \text{int}, b \neq 0]\} (a/b) \{[\text{res} := a \hat{/} b : \text{int}]\}$ where $\hat{/}$ denotes the logical integer division operator. Following standard practice in proof assistants, the operator $\hat{/}$ is defined in the logic as a total function that returns unspecified results when the divisor is equal to zero.

4.7 Typechecking of Terms

Our typing judgment takes the form $\{\{\Gamma\}\} t^\Delta \{\{\Gamma'\}\}$, capturing the fact that, in context Γ , the term t is well typed and produces a context Γ' with a *usage map* Δ . We are interested in describing the *algorithmic* typing rules exploited by OptiTrust. Our typing algorithm takes Γ and t as input, and

1569 produces Γ' and Δ as output. The refinement with usage maps will be discussed further in Section 5.3.
 1570 For now, we focus on describing typing rules for the judgment $\{\{\Gamma\}\} t \{\{\Gamma'\}\}$.

1571 In general, in a valid triple $\{\{\Gamma\}\} t \{\{\Gamma'\}\}$, variables from the postcondition Γ' may refer to variables
 1572 from the precondition Γ . For the purpose of the algorithmic typechecking, however, we design the
 1573 typing rules in such a way that Γ' is always *closed*, meaning that variable occurrences in Γ' refer to
 1574 variables that are all previously bound in Γ' . The purpose of this design decision is to maximize the
 1575 amount of information that is propagated forward during the typechecking.

1576 In particular, in the algorithmic typechecking, all the logical bindings (ghost variables and pure
 1577 facts) from Γ are reproduced in Γ' . The pure bindings that appear in Γ' but not in Γ correspond
 1578 either (1) to the binding for **res**, which denotes the result value produced by t , as explained in
 1579 Section 4.3; or (2) to logical bindings (ghost variables and pure facts) that correspond to existentially
 1580 quantified variables and pure postconditions.

1581 The linear bindings of Γ' , compared with those in Γ , reflect the side effects performed by t . Linear
 1582 resources that are bound with the same name in Γ' as in Γ necessarily correspond to resources that
 1583 have not been modified by t .

1584 Figure 18 presents our typing rules. The typing rule for applications handles the particular case
 1585 where the subterms are program variables (i.e., functions calls in A-normal form)—the processing of
 1586 effectful subterms depends on resource usage, and is explained further in Section 5.5. The soundness
 1587 of these rules stems from the fact that they correspond to an algorithmic reformulation of the
 1588 standard reasoning rules from Separation Logic. We next describe the rules individually.

1589 *Literals and Variables.* Consider a term t that corresponds either to a program variable or to a
 1590 literal. In its triple, of the form $\{\{\Gamma\}\} t \{\{\Gamma'\}\}$, the output context Γ' is obtained by extending Γ with
 1591 an alias binding from **res** to t itself. Alias bindings were defined in Section 4.2. This is possible since
 1592 for literals and variables, t is a logical expression and therefore can directly appear in contexts.
 1593 The type of t is computed by means of the typing judgment for logical expressions, defined in
 1594 Section 4.6.

1596 *Let-Bindings.* Consider an instruction of the form **let** $x = t$. Recall from Section 3.1 that such
 1597 instructions only appear in sequences. The subexpression t produces a value, hence the output
 1598 context Γ_1 associated with t binds the special variable **res**. The expression **let** $x = t$ itself does not
 1599 produce a value, hence its output context Γ_2 does not bind **res**. However, the output context Γ_2 is
 1600 extended with a binding on x . Concretely, Γ_2 is obtained by replacing in Γ_1 the bound name **res**
 1601 with the bound name x .

1602 *Sequence of Instructions.* We decompose the treatment of sequences in two rules: a first rule
 1603 named SEQ for handling the sequence of instructions per se, and a second rule named BLOCK for
 1604 handling the disposal of stack-allocated variables. The rest of this paragraph describes the SEQ rule.
 1605 Consider a sequence $(t_1; \dots; t_n; r)$. Starting from an input context Γ_0 , each subterm t_i makes the
 1606 context evolves from Γ_{i-1} to Γ_i . Recall from Section 3.1 that each subterm t_i must have unit type
 1607 (a.k.a. void type), else it would have been wrapped into a call to the “ignore” function. The sequence
 1608 itself may return a value identified by the optional result variable r . If such a result variable is set,
 1609 the final context is patched to include a **res** binding instead of the original r binding.

1611 *Scope Blocks.* The typing rule BLOCK is responsible for collecting the resources that corresponds
 1612 to stack-allocated variables, when reaching the end of a sequence, that is, the end of their scope.
 1613 Recall from Section 3.1 that stack allocation takes the form **let** $x = \text{stackRef}(T)$ or **let** $x =$
 1614 $\text{stackAlloc}_C()$, with such instructions occurring directly within a sequence. The auxiliary function
 1615 $\text{StackAllocCells}(t_1, \dots, t_n)$ synthesizes, based on the syntax of the terms t_i that appear in the se-
 1616 quence at hand, a conjunction of resources, each of the form $\text{Unit}(p \rightsquigarrow \text{Cell}_\tau)$. These resources
 1617

$$\begin{array}{c}
1618 \\
1619 \\
1620 \\
1621 \\
1622 \\
1623 \\
1624 \\
1625 \\
1626 \\
1627 \\
1628 \\
1629 \\
1630 \\
1631 \\
1632 \\
1633 \\
1634 \\
1635 \\
1636 \\
1637 \\
1638 \\
1639 \\
1640 \\
1641 \\
1642 \\
1643 \\
1644 \\
1645 \\
1646 \\
1647 \\
1648 \\
1649 \\
1650 \\
1651 \\
1652 \\
1653 \\
1654 \\
1655 \\
1656 \\
1657 \\
1658 \\
1659 \\
1660 \\
1661 \\
1662 \\
1663 \\
1664 \\
1665 \\
1666
\end{array}$$

$$\begin{array}{c}
\frac{\Gamma.\text{pure} \vdash t : \tau \quad t \text{ is a literal or a variable}}{\{\{\Gamma\}\} t \{\{\Gamma \otimes [\text{res} : \tau := t]\}\}} \text{LITORVAR} \\
\\
\frac{\{\{\Gamma_0\}\} t \{\{\Gamma_1\}\} \quad \Gamma_2 = \text{Rename}\{\text{res} := x\}(\Gamma_1)}{\{\{\Gamma_0\}\} \mathbf{let} x = t \{\{\Gamma_2\}\}} \text{LET} \\
\\
\frac{\forall i \in [1, n]. \quad x_i \text{ fresh} \quad \wedge \quad \{\{\Gamma_{i-1}\}\} t_i \{\{\Gamma_i\}\} \quad \Gamma_r = \begin{cases} \text{Rename}\{r := \text{res}\}(\Gamma_n) & \text{if } r \neq \emptyset \\ \Gamma_n & \text{if } r = \emptyset \end{cases}}{\{\{\Gamma_0\}\} (t_1; \dots; t_n; r) \{\{\Gamma_r\}\}} \text{SEQ} \\
\\
\frac{\{\{\Gamma_0\}\} (t_1; \dots; t_n; r) \{\{\Gamma_r\}\} \quad (\emptyset, F) = \Gamma_r \boxminus \text{StackAllocCells}(t_1, \dots, t_n)}{\{\{\Gamma_0\}\} \{t_1; \dots; t_n; r\} \{\{\Gamma_r.\text{pure} \mid F\}\}} \text{BLOCK} \\
\\
\frac{\begin{array}{c} \{\{\Gamma_0.\text{pure}\} \otimes \gamma.\text{pre}\} t \{\{\Gamma_1\}\} \quad (_, \emptyset) = \Gamma_1 \boxminus \gamma.\text{post} \\ (\text{res} : \hat{t}_r) \in \gamma.\text{post} \quad \hat{t}_f = (\hat{t}_1, \dots, \hat{t}_n) \rightarrow \hat{t}_r \end{array}}{\{\{\Gamma_0\}\} (\mathbf{fun}(a_1 : \hat{t}_1, \dots, a_n : \hat{t}_n)_\gamma \mapsto t) \{\{\Gamma_0 \otimes [\text{res} : \hat{t}_f, \text{Spec}(\text{res}, [a_1, \dots, a_n], \gamma)]\}\}} \text{FUN} \\
\\
\frac{\begin{array}{c} \text{Spec}(x_0, [a_1, \dots, a_n], \gamma) \in \Gamma_0 \\ (E_{\text{frac}}, \sigma', F) = \Gamma_0 \ominus \text{Specialize}_{\Gamma_0} \{\overline{a_i} := x_i^{i \in [1, n]}, \sigma\}(\gamma.\text{pre}) \\ \text{dom}(\rho) = \text{dom}(\gamma.\text{post}) \quad \text{im}(\rho) \cap \text{dom}(\Gamma_0) = \emptyset \\ \Gamma_q = \text{Rename}\{\rho\}(\overline{a_i} := x_i^{i \in [1, n]}, \sigma, \sigma'(\gamma.\text{post})) \\ \Gamma_r = \text{CloseFrac}([\Gamma_0.\text{pure}, E_{\text{frac}}] \otimes F \otimes \Gamma_q) \end{array}}{\{\{\Gamma_0\}\} x_0(x_1, \dots, x_n)_{\sigma, \rho} \{\{\Gamma_r\}\}} \text{APP} \\
\\
\frac{\begin{array}{c} \Gamma_p = [\chi.\text{vars}] \otimes (\star_{i \in R} \chi.\text{excl}.\text{pre}) \otimes \chi.\text{shrd}.\text{reads} \otimes \chi.\text{shrd}.\text{inv} \\ (E_{\text{frac}}, \sigma', F) = \Gamma_0 \ominus \Gamma_p \\ \Gamma'_p = [i : \text{int}, i \in R] \otimes [\chi.\text{vars}] \otimes \chi.\text{excl}.\text{pre} \otimes \frac{1}{R.\text{len}} \chi.\text{shrd}.\text{reads} \otimes \chi.\text{shrd}.\text{inv} \\ \{\{\Gamma_0.\text{pure}\} \otimes \Gamma'_p\} t \{\{\Gamma'_q\}\} \\ (_, \emptyset) = \Gamma'_q \boxminus \chi.\text{excl}.\text{post} \otimes \frac{1}{R.\text{len}} \chi.\text{shrd}.\text{reads} \otimes \chi.\text{shrd}.\text{inv} \\ \Gamma_q = \sigma'((\star_{i \in R} \chi.\text{excl}.\text{post}) \otimes \chi.\text{shrd}.\text{reads} \otimes \chi.\text{shrd}.\text{inv}) \\ \Gamma_r = \text{CloseFrac}([\Gamma_0.\text{pure}, E_{\text{frac}}] \otimes F \otimes \Gamma_q) \end{array}}{\{\{\Gamma_0\}\} \mathbf{for} (i \in R)_\chi t \{\{\Gamma_r\}\}} \text{FOR} \\
\\
\frac{\text{parallelizable}(\chi) \quad \{\{\Gamma_0\}\} \mathbf{for} (i \in R)_\chi t \{\{\Gamma_r\}\}}{\{\{\Gamma_0\}\} \mathbf{for}^{\text{PAR}} (i \in R)_\chi t \{\{\Gamma_r\}\}} \text{FORPAR} \\
\\
\frac{\{\{\Gamma_0\}\} t_1 \{\{\Gamma_1\}\} \quad \{\{\text{Learn}\{\text{res} = \text{true}\}(\Gamma_1)\}\} t_2 \{\{\Gamma_2\}\} \quad \{\{\text{Learn}\{\text{res} = \text{false}\}(\Gamma_1)\}\} t_3 \{\{\Gamma_3\}\} \\ (_, \emptyset) = \Gamma_2 \boxminus \Gamma_r \quad (_, \emptyset) = \Gamma_3 \boxminus \Gamma_r}{\{\{\Gamma_0\}\} \mathbf{if}_{\Gamma_r} t_1 \mathbf{then} t_2 \mathbf{else} t_3 \{\{\Gamma_r\}\}} \text{IF}
\end{array}$$

Fig. 18. Algorithmic typing rules for establishing triples of the form $\{\{\Gamma\}\} t \{\{\Gamma'\}\}$. These rules are generalized in Section 5.3 to derive triples the form $\{\{\Gamma\}\} t^\Delta \{\{\Gamma'\}\}$, where Δ describes the resource usage.

are subtracted from the context available at the end of the sequence. Crucially, the subtraction operation checks that the resources indeed appear in the current resource set. Doing so ensures, in particular, that the address of a stack-allocated piece of data was not subject to a prior call to free.

1667 *Function Definition.* Consider a function definition $\mathbf{fun}(a_1 : \hat{\tau}_1, \dots, a_n : \hat{\tau}_n)_\gamma \mapsto t$, with arguments
 1668 a_i of type $\hat{\tau}_i$, with body t , and with contract γ . Recall from Section 4.3 that the function contract
 1669 consists of a precondition $\gamma.\text{pre}$ and a postcondition $\gamma.\text{post}$, both described as contexts. The function
 1670 is a closure that may capture free variables from the current context. In the rule, the pure variables
 1671 from the current context are described as $\Gamma_0.\text{pure}$. Note, however, that the function is not allowed
 1672 to capture linear resources. Hence, the body of the function is typechecked in an environment that
 1673 consists of the conjunction of $\Gamma_0.\text{pure}$ and $\gamma.\text{pre}$. Ultimately, the body of the function must produce
 1674 a context Γ_r that entails the postcondition $\gamma.\text{post}$. The postcondition of the function definition itself
 1675 binds \mathbf{res} with the correct function type ($\mathbf{res} : \hat{\tau}_f$) and gives its specification hypothesis ($\text{Spec}(\mathbf{res},$
 1676 $[a_1, \dots, a_n], \gamma)$). As explained earlier in Section 4.3, this hypothesis captures $\{\gamma.\text{pre}\} \mathbf{res}(a_1, \dots,$
 1677 $a_n) \{\gamma.\text{post}\}$, which is indeed the triple intended for the function named \mathbf{res} .

1678 *Function Applications.* Consider a function application of the form $x_0(x_1, \dots, x_n)$, where the x_i are
 1679 program variables. (The general form will be discussed in section 5.5.) To typecheck it, the input
 1680 context Γ_0 must contain an entry of the form $\text{Spec}(x_0, [a_1, \dots, a_n], \gamma)$ for the function x_0 . This same
 1681 context Γ_0 must entail the precondition $\gamma.\text{pre}$, specialized for the arguments x_i by means of the
 1682 Specialize operations defined in Section 4.2. This entailment is checked by means of the carving
 1683 subtraction operation defined in Section 4.5. The subtraction produces a frame F that contains the
 1684 resources from Γ_0 that are not used by the function call, and produces a substitution named σ'
 1685 that describes the instantiation of the ghost arguments and resources. The final postcondition Γ_q is
 1686 obtained by considering the postcondition $\gamma.\text{post}$, adding the frame F and $\Gamma_0.\text{pure}$, then invoking the
 1687 CloseFrac operation described in Section 4.1 for eagerly recombining carved-out fractions.

1688 Two additional technicalities are involved in the statement of the APP rule. They correspond to the
 1689 handling of optional user-provided annotations, named σ and ρ , that may guide the typechecking
 1690 of an application. Such annotations are commonly found both in proof assistants and in program
 1691 verification frameworks. The map σ allows instantiating a subset of the ghost arguments. Indeed,
 1692 there could be situations where the subtraction operation would fail to infer a unique possible
 1693 instantiation, by the only means of the unification process. Hence, user annotations are required to
 1694 resolve the instantiation. In all our case studies, σ is only used on ghost calls. The map ρ corresponds
 1695 to a renaming map. Its purpose is to rename all the resources that are produced by the postcondition
 1696 to avoid name conflicts. Typically, the map ρ is initialized with fresh variable names during the first
 1697 typechecking of each function application. In the future, we might let the user explicitly provide
 1698 some entries of ρ to manipulate the produced resources by name.

1700 *Simple for-loops.* Consider a possibly parallel, simple for-loop of the form $\mathbf{for}^\pi(i \in R)_\chi t$. The
 1701 typechecking of such a loop is driven by the loop contract annotation χ . The loop body t is
 1702 typechecked in a context that binds an index i of type int , a hypothesis of type $i \in R$, the variables
 1703 from E , the resources F_{pre} , (subfractions of) the resources in F_{reads} , and the resources in F_{inv} .
 1704 The loop body needs to produce the resources F_{post} , and it needs to give back the resources that
 1705 it had received from F_{reads} and from F_{inv} . There are three complications. First, the shared-read
 1706 resources, described by $\chi.\text{shrd.reads}$, are split into $\frac{1}{R.\text{len}}$ subfractions, where $R.\text{len}$ denotes the
 1707 number of iterations associated with the range R . Note that, when typechecking the body of the loop
 1708 for a particular iteration $i \in R$, the denominator $R.\text{len}$ can be assumed to be nonzero—indeed, $i \in R$
 1709 is equivalent to $0 \leq i < R.\text{len}$. Second, like for function calls, the instantiation of the contract using
 1710 the resources from the input environment Γ_0 is computed using a subtraction, involving a frame F
 1711 as well as an instantiation map σ' . Also, like for function calls, the output context is obtained by
 1712 invoking the CloseFrac operation. Third, loops, like functions calls, feature optional annotations
 1713 σ and ρ , which we have omitted from the statement of the rule, for simplicity. The map σ guides
 1714 how the contract is instantiated in the input environment Γ_0 . The map ρ can be used to explicit the
 1715

names associated with the resources produced by the loop. The two maps are handled in a similar way as in the `APP` rule.

Conditionals. Consider a conditional of the form `if t_1 then t_2 else t_3` . The condition t_1 is evaluated in the input context Γ_0 and produces a context Γ_1 . Then, both branches t_2 and t_3 need to typecheck in the context Γ_1 . This context needs to be patched to reflect the knowledge that t_1 evaluated to either true or false, depending on the branch. The patch is implemented by means of the operation $\text{Learn}\{\mathbf{res} = b\}(\Gamma)$. This operation applies the following three steps.

- (1) If an aliasing binding of the form $\mathbf{res} := v : \text{bool}$ appears in Γ , then the operation replaces this binding with a conventional binding $\mathbf{res} : \text{bool}$, and extends Γ with an equality $[\mathbf{res} = v]$.
- (2) It specializes the variable \mathbf{res} with b , that is, it removes the binding $\mathbf{res} : \text{bool}$, and replaces all occurrences of \mathbf{res} with the boolean value b .
- (3) It applies basic simplifications on the expressions in which \mathbf{res} has been substituted with b .

For example, assume t_1 is a test of the form $x == y$, and consider the evaluation of $\text{Learn}\{\mathbf{res} = \text{true}\}(\Gamma_1)$. The output context of t_1 contains the alias binding $\mathbf{res} := (x==y) : \text{bool}$. At step (1), this binding is replaced with an equality $\mathbf{res} = (x==y)$. At step (2), \mathbf{res} is replaced with true, hence the equality becomes $\text{true} = (x==y)$. At step (3), this hypothesis is rewritten as the logical equality $x = y$.

The then-branch t_2 produces an output context Γ_2 , and likewise the else-branch t_3 produce an output context Γ_3 . What should be the output context of the entire conditional `if t_1 then t_2 else t_3` ? It must be a context, call it Γ_r , that both Γ_2 and Γ_3 entail. This context Γ_r is usually called the *join* context in program logics. In general, there is no way to automatically infer join contexts—it is almost as hard as inferring contracts for loops. Therefore, typechecking and verification tools must resort to a combination of user-provided annotations and heuristics. For now, we assume join contexts to be provided by the user. In our box-blur case study (Section 2.1), the conditionals appear in terminal position in the body of a function, hence our typechecker can simply instantiate the join context using the (user-provided) postcondition of that function. We leave it to future work to devise heuristics well-suited for our typesystem, in order to reduce the number of situations where `OptiTrust` users need to provide annotations.

4.8 Type Soundness

The purpose of this section is to present formal statements that reflect the design principles of our type system. This section may be safely skipped for a first read. A number of auxiliary definitions, such as the evaluation rules or the satisfaction of a linear resource by a heap fragment, may be found in the appendix.

We follow the standard approach of justifying soundness of a separation logic by providing a semantic interpretation of triples. The general pattern asserts that: “a triple holds if and only if, in any input state satisfying the precondition (i.e., the input context), the evaluation of the term terminates and produces an output state satisfying the postcondition (i.e., the output context)”. This statement relies on two central ingredients. First, a definition of the semantics of a term. Second, a definition of what it means for a program state to satisfy a context.

We formalize the semantics using an *omni-big-step* evaluation judgment [Charguéraud et al. 2023]. This judgment has been shown to simplify proofs of the frame rule of separation logic, and proofs of compiler correctness results. Concretely, the judgment $t/(s, m) \Downarrow Q$ asserts that the term t , in an input program state (s, m) , evaluates to output program states that belong to the set Q . A program state, written (s, m) , consists of an immutable stack s and a store m . If t produces an output value, then this value is bound in the output program state to the dedicated name \mathbf{res} . For

1765 simplicity, we focus on total correctness: $t/(s, m) \Downarrow Q$ asserts that all possible evaluations of the
 1766 term t do terminate, without error.¹² The evaluation rules may be found in Section A.

1767 Let us now focus on context satisfaction. As usual in separation logic that involves fractional
 1768 permissions (or more general forms of ghost state), one asserts that a program state satisfies a
 1769 context if and only if there exists a *logic state*, which consists of this program state augmented with
 1770 additional (“ghost”) information, such that this logical state satisfies the context. A *logical state* is
 1771 one that may *satisfy* a context Γ . We define further on an elision function that extracts a program
 1772 state from a logical state. We first describe the representation of a logical state.

1773 A logical state consists of a logical stack, written σ , and a logical store, written μ . A logical stack
 1774 is similar to a program stack except that it includes additional bindings for ghost variables. A logical
 1775 store is similar to a program store except that every memory location is tagged with a fraction,
 1776 written α , in the range $(0, 1]$. As standard in realizations of separation logic, a fraction less than
 1777 one corresponds to a read-only permission.

1778 As said, a context Γ corresponds to a specification of a logical state. We say that a logical state
 1779 (σ, μ) satisfies a context Γ of the form $\langle E \mid F \rangle$, and write $(\sigma, \mu) \in \Gamma$, if the bindings in σ have types
 1780 that correspond to the bindings in E , and if the memory cells described by μ correspond to the
 1781 linear resources described in F . The technical details of the definition of $(\sigma, \mu) \in \Gamma$ are given in
 1782 Section B.

1783 To state the semantic interpretation of triples, we need a projection function for extracting a
 1784 program state out of a logical state. We write $\sigma|_{\text{prog}}$ the operation that converts a logical stack σ
 1785 into a program stack s by restricting the entries to program variables, or, equivalently said, by
 1786 removing entries associated with ghost variables. We write $\mu|_{\text{prog}}$ the operation that turns a logical
 1787 store μ into a program store m by removing all fractions. By leveraging the two operations, we
 1788 define $(\sigma, \mu)|_{\text{prog}}$ as $(\sigma|_{\text{prog}}, \mu|_{\text{prog}})$, to convert a logical state into a program state.

1789 Before defining triples, we introduce $\text{AcceptableStates}(\sigma, \mu, \Gamma')$ to denote the set of program
 1790 output states satisfying the postcondition Γ' and satisfying certain constraints with respect to the
 1791 input state (σ, μ) . The set $\text{AcceptableStates}(\sigma, \mu, \Gamma')$ corresponds to the set of states that are the
 1792 projection of a logical state (σ', μ') such that: (1) the logical state (σ', μ') satisfies the specification
 1793 Γ' , and (2) the read-only restriction of μ' is identical to the read-only restriction of μ , and (3) the
 1794 stacks in σ' and σ agree on the intersection of their domain. To formalize the second constraint,
 1795 we let $\text{OnlyRO}(\mu)$ denote the restriction of the logical store μ to the cells that are tagged with a
 1796 fraction strictly less than 1, that is, as $\{l \mapsto (\alpha, v) \mid (l \mapsto (\alpha, v)) \in \mu \wedge \alpha < 1\}$. We then define:

$$1797 \text{AcceptableStates}(\sigma, \mu, \Gamma') := \left\{ (\sigma', \mu')|_{\text{prog}} \left| \begin{array}{l} (\sigma', \mu') \in \Gamma' \\ \wedge \text{OnlyRO}(\mu) = \text{OnlyRO}(\mu') \\ \wedge \forall x \in \text{dom}(\sigma) \cap \text{dom}(\sigma'), \sigma(x) = \sigma'(x) \end{array} \right. \right\}.$$

1801 We are now ready to define logical triples, written $\{\Gamma\} t \{\Gamma'\}$. Such a triple asserts that for any
 1802 logical state satisfying Γ , starting in the program state that corresponds to this logical state, all
 1803 executions of t terminate and produce output states that belong to the set $\text{AcceptableStates}(\sigma, \mu,$
 1804 $\Gamma')$. The latter means that an output state must satisfy Γ' , must preserve read-only entries, and
 1805 must feature an output stack that agrees with the input stack.

1807 ¹²As explained in the omniseantics paper [Charguéraud et al. 2023], the omni-big-step evaluation judgment is related to
 1808 the standard big-step judgment via the following equivalence.

$$1809 t/(s, m) \Downarrow Q \iff \begin{array}{l} \text{all possible executions of } t \text{ terminate without error} \\ \wedge (\forall (s', m'), t/(s, m) \Downarrow (s', m') \implies (s', m') \in Q) \end{array}$$

1814 *Definition 4.1 (Logical triples).*

$$1815 \quad \{\Gamma\} t \{\Gamma'\} := \forall (\sigma, \mu) \in \Gamma, \quad t/(\sigma, \mu)|_{\text{prog}} \Downarrow \text{AcceptableStates}(\sigma, \mu, \Gamma')$$

1817 The fundamental property of separation logic is the frame rule, which we prove correct for our
 1818 logical triples in Appendix C. The contexts involved here are dependently-typed, hence we need
 1819 additional assumptions to ensure that the composed contexts are *well-typed*, in the sense that every
 1820 variable that appears in a type or a resource is properly bound earlier in the context, and that all
 1821 the types that appear in the context are themselves well-typed. (Well-typed contexts are formalized
 1822 by Definition C.3 in Section C.) The statement of the frame rule is thus as follows.

1823 **THEOREM 4.2 (FRAME RULE FOR LOGICAL TRIPLES).**

$$1824 \quad \{\Gamma\} t \{\Gamma'\} \wedge \Gamma * \Gamma'' \text{ is well-typed} \wedge \Gamma' * \Gamma'' \text{ is well-typed} \implies \{\Gamma * \Gamma''\} t \{\Gamma' * \Gamma''\}$$

1826 Our typing rules presented earlier on in this section are designed as algorithmic variants of the
 1827 standard typing rules of separation logic. The soundness of our algorithmic typing rules stems
 1828 from the soundness of the standard typing rules of separation logic. Soundness is formally stated
 1829 as follows.

1831 **PROPOSITION 4.3 (SOUNDNESS OF THE ALGORITHMIC TYPING RULES).** $\{\{\Gamma\}\} t \{\{\Gamma'\}\} \implies \{\Gamma\} t \{\Gamma'\}$

1832 We leave to future work the completion of a mechanized proof of this statement.

1834 5 COMPUTING PROGRAM RESOURCES: USAGE MAPS

1835 The first goal of this section is to formalize the usage maps, written Δ , and to generalize triples
 1836 from the form $\{\{\Gamma\}\} t \{\{\Gamma'\}\}$ to the form $\{\{\Gamma\}\} t^\Delta \{\{\Gamma'\}\}$. Section 5.1 presents the grammar of usage
 1837 maps. Section 5.2 presents operations on usage maps. Section 5.3 explains how usage maps are
 1838 computed by our typing algorithm.

1839 The second goal of this section is to formalize the *triple minimization* operations, which plays
 1840 a central role in the typechecking of function calls involving effectful subexpressions. Triple
 1841 minimization will also be useful later on to minimize the loop contracts produced by transformations.
 1842 Section 5.4 presents the triple minimization procedure. Section 5.5 presents the typing rule for
 1843 subexpressions—this typing rule applies as a preprocessing before the `APP` rule presented earlier.
 1844 Section 5.6 presents formal statements about the contents of usage maps.

1846 5.1 Grammar of Usage Maps

1847 A *usage map*, written Δ , is an association map that binds resource names to *usage kinds*. For a
 1848 *pure* resource name, there are 2 possible usage kinds: required and ensured. For a *linear* resource
 1849 name, there are 5 possible usage kinds: full, uninit, splittedFrac, joinedFrac and produced. In a
 1850 triple $\{\{\Gamma\}\} t^\Delta \{\{\Gamma'\}\}$, the usage map Δ binds names of resources that can be bound in Γ or Γ' , or
 1851 possibly in both. The usage map Δ only binds names of resources that are effectively manipulated
 1852 by t . (In other words, the *framed* resources are omitted from usage maps.) Let us now explain the
 1853 meaning of each possible binding in a usage map Δ associated with the triple $\{\{\Gamma\}\} t^\Delta \{\{\Gamma'\}\}$.

- 1854 • “ x : required” means that x is a pure resource in Γ that was used during the typing of t .
- 1855 • “ x : ensured” means that x is a pure resource added to the context Γ' during the typing of t .
 1856 In such a situation, x is not bound in Γ .
- 1857 • “ y : full” can arise when Γ contains a linear resource “ $y : H$ ”, for some predicate H . The
 1858 usage “ y : full” means that this resource is consumed during the typing of t . As a result y is
 1859 not bound in Γ' . Even if t produces a linear resource with the same predicate H , this new
 1860 occurrence of H is assigned a fresh name, distinct from y .

- “ y : uninit” is similar to “ y : full” but moreover captures the information that t needs not read the original contents of the memory cells associated with the resource named y . In particular, if t performs a write operation in a cell y before any read operation on y , then the usage of y is uninit.
- “ y : splittedFrac” can arise when Γ contains a splittable linear resource “ y : H ”, for some predicate H . The usage “ y : splittedFrac” means that t uses an unspecified subfraction of this resource. In such a situation, the name y is bound both in Γ and in Γ' . It may be the case, however, that the resource named y carries different fractions in Γ and Γ' .
- “ y : joinedFrac” can arise when Γ contains a linear resource of the form “ y : $(\alpha - \beta_1 - \dots - \beta_n)H$ ”. The usage “ y : joinedFrac” means that: (1) the linear resource named y is not used by t , and (2) t produced a resource of the form $(\beta_i - \gamma_1 - \dots - \gamma_m)H$, and (3) these two resources are merged and the result appears in Γ' under the name y . If a single merge operation is applied, then the resulting resource is y : $(\alpha - \beta_1 - \dots - \beta_{i-1} - \gamma_1 - \dots - \gamma_m - \beta_{i+1} - \dots - \beta_n)H$. (Recall Section 4.1.)
- “ y : produced” means that the linear resource y has been produced by t . In this case, y is the name of a linear resource in Γ' , and does not occur in Γ .
- If a resource name is bound in Γ but not in Δ , then its absence indicates that the corresponding resource is not touched by t . Such a resource is bound under the same name in Γ and Γ' .

5.2 Operations on Usage Maps

Projections of Usage Maps. We define $\Delta.\text{full}$ as the set of names y such that “ y : full” appears in Δ . Likewise, we define $\Delta.\text{required}$, $\Delta.\text{ensured}$, $\Delta.\text{splittedFrac}$, $\Delta.\text{joinedFrac}$ and $\Delta.\text{produced}$. In addition, we define the following operations.

$$\begin{aligned}
 \Delta.\text{consumed} &= \Delta.\text{full} \cup \Delta.\text{uninit} \\
 \Delta.\text{read} &= \Delta.\text{splittedFrac} \cup \Delta.\text{joinedFrac} \\
 \Delta.\text{alter} &= \Delta.\text{consumed} \cup \Delta.\text{produced} \cup \Delta.\text{ensured}
 \end{aligned}$$

Intersection and Filtering. We define:

$$\begin{aligned}
 \Delta_1 \cap \Delta_2 &= \text{dom}(\Delta_1) \cap \text{dom}(\Delta_2) \\
 \Gamma \upharpoonright \Delta &= \Gamma \upharpoonright \text{dom}(\Delta)
 \end{aligned}$$

Sequential Composition of Usage Maps. This section defines the *usage composition operator*, written $\Delta_1; \Delta_2$. This operator plays a central role in computing the usage of a sequence of terms. Let us begin with an example.

Consider the sequence “ $(\text{let } r = \text{heapAlloc}_C())^{\Delta_1}; (\text{let } k = \text{get}(r))^{\Delta_2}; \text{free}(r)^{\Delta_3}$ ”. In Δ_1 , we have a binding “ y : produced” because the first instruction produces the resource “ y : $r \rightsquigarrow \text{Cell}$ ”. In Δ_2 , we have a binding “ y : splittedFrac” because the instruction only reads with y (thus it accepts any subfraction). In Δ_3 , we have a binding “ y : uninit” because the third instruction destroys the resource y without caring about the value of the Cell.

Let us give three examples of compositions. First, the usage map $\Delta_1; \Delta_2$ contains a binding “ y : produced” because, taken together, the sequential composition of the those two instructions still creates the resource y . Second, the usage map $\Delta_2; \Delta_3$ contains a binding y : full because, taken together, the second and the third instruction consume the Cell, and they read the value that was contained inside. Third, the usage map $\Delta_1; \Delta_2; \Delta_3$ contains no binding for y because the Cell cannot be seen from outside the sequence of instruction.

$\Delta_1; \Delta_2$	\emptyset	required	ensured
\emptyset	\emptyset	required	ensured
required	required	required	\perp
ensured	ensured	ensured	\perp

$\Delta_1; \Delta_2$	\emptyset	full	uninit	splittedFrac	joinedFrac	produced
\emptyset	\emptyset	full	uninit	splittedFrac	joinedFrac	produced
full	full	\perp	\perp	\perp	\perp	\perp
uninit	uninit	\perp	\perp	\perp	\perp	\perp
splittedFrac	splittedFrac	full	full	splittedFrac	splittedFrac	\perp
joinedFrac	joinedFrac	full	uninit	splittedFrac	joinedFrac	\perp
produced	produced	\emptyset	\emptyset	produced	produced	\perp

Fig. 19. Tables for sequential composition of two usage maps, for pure and for linear resources. For example, in the second table, the cell on the row “splittedFrac” and on the column “full” expresses that if “ $x : splittedFrac$ ” is a binding from Δ_1 and “ $x : full$ ” is a binding from Δ_2 , then “ $x : full$ ” is a binding in $\Delta_1; \Delta_2$. The input or output \emptyset corresponds to cases where the usage map contains no binding for the resource name considered. The output \perp corresponds to cases that cannot arise according to our typechecking rules.

Formally, the usage composition operation $\Delta_1; \Delta_2$ is defined by merging the two usage maps pointwise by resource name, using the table shown in Fig. 19 to compute the “combined usage” in case a same resource name is bound both in Δ_1 and Δ_2 .

The input or output \emptyset corresponds to cases where there is no binding for a resource name in the usage map. Note that a resource produced in Δ_1 and then fully used in Δ_2 will be absent from $\Delta_1; \Delta_2$. As illustrated in the earlier example, a usage map abstracts away intermediate resources not present in the final triple.

The output \perp corresponds to cases that cannot arise. For example, it is not possible to have a linear resource used as full and used again afterwards, since usage full corresponds to a removal from the context. Similarly, the same resource name cannot be produced or ensured twice.

Finally, let us comment on the naming policy. If a linear resource is entirely consumed, its name disappears. If a resource $y : \beta H$ is split as αH and $(\beta - \alpha)H$, the $(\beta - \alpha)H$ part keeps the initial resource name y (and αH takes a fresh resource name). If `CloseFrac` merges the fractions $y : (\beta - \alpha)H$ and $y' : \alpha H$, it produces a resource βH with the name y (and the name y' disappears).

Let us illustrate how these rules play out on a concrete example. Assume a term t_1 uses a full resource named y to only perform a read operation, and subsequently a term t_2 uses the same resource to perform a write operation. Then, thanks to the fact that the name y was preserved during the carve-out and subsequent `CloseFrac` operation, the usage map of the sequence $t_1; t_2$ contains, as one would naturally expect, the binding $y : full$.

5.3 Computing Usage Maps

Usage of a context subtraction. Each time a typing rule performs a subtraction, we add entries to the usage map of the term invoking this rule. This paragraph explains the usage map associated with a subtraction. The usage map of a subtraction $(\sigma, F) = \Gamma_1 \boxminus \Gamma_2$ contains:

- One entry required for each pure variable of Γ_1 mentioned in σ .
- One entry uninit or full for each linear resource of Γ_1 that was unified with a resource of Γ_2 . The entry is uninit if the resource in Γ_2 is of the form `Uninit(H)`. Otherwise, it is a full.

For a subtraction performing read-only carving $\Gamma_1 \ominus \Gamma_2$, the usage map is defined in the same way as $\Gamma_1 \boxminus \Gamma_2$ except that if a linear resource from Γ_2 is found by carving a resource of Γ_1 , the entry

for that resource from Γ_1 has kind `splittedFrac`, and we also add an ensured entry for each newly generated fraction.

Usage of a CloseFrac. When closing fractions, we need to add entries to the usage map to account for the modifications on the context. We try to do so in a way that preserves as much information as possible. When `CloseFrac` finds a possible reduction on two resources $y_1 : (\alpha - \beta_1 - \dots - \beta_n)H$ and $y_2 : (\beta_i - \gamma_1 - \dots - \gamma_m)H$ it keeps the name y_1 from the carved part for the generated closed resource $(\alpha - \beta_1 - \dots - \beta_{i-1} - \gamma_1 - \dots - \gamma_m - \beta_{i+1} - \dots - \beta_n)H$. On the one hand, the resource y_2 disappears from the context. Therefore, we have to put the usage $y_2 : \text{full}$ in the usage map. On the other hand, the resource y_1 remains in the context. Since the absence of y_1 would not have blocked the typechecking, it gets the usage $t_1 : \text{joinedFrac}$. Note this is currently the only way `joinedFrac` usage are generated. Note also that the order of reduction does not matter for the final usage map (all the fractions that disappear will have a usage `full`, and all the fractions that got bigger will have a usage `joinedFrac`).

Computing Usage During Term Typing. In order to produce triples of the form $\{\{\Gamma\}\} t \{\{\Gamma'\}\}$, we need to patch our typing rules to record the usage information.

Here is the full version of the rules `LITORVAR` and `LET` described earlier:

$$\frac{\Gamma.\text{pure} \vdash t : \tau \quad t \text{ is a literal or a variable} \quad \Delta = \{\mathbf{res} : \text{ensured}\}}{\{\{\Gamma\}\} t^\Delta \{\{\Gamma \otimes [\mathbf{res} : \tau := t]\}\}} \text{LITORVAR}$$

$$\frac{\{\{\Gamma_0\}\} t^\Delta \{\{\Gamma_1\}\} \quad \Gamma_2 = \text{Rename}\{\mathbf{res} := x\}(\Gamma_1) \quad \Delta' = \text{Rename}\{\mathbf{res} := x\}(\Delta)}{\{\{\Gamma_0\}\} (\mathbf{let} \ x = t)^\Delta \{\{\Gamma_2\}\}} \text{LET}$$

For the rule `VAL`, the usage map contains a single binding $\mathbf{res} : \text{ensured}$ to account for the alias added in the context. For the rule `LET`, the typechecker uses the operator $\text{Rename}\{x := x'\}(\Delta)$, that renames the key x into x' inside the map Δ . This renaming is applied on the usage map of the body to follow the renaming in the context.

For the interested reader, we now explain how usage maps are computed in practice. Instead of rewriting each typing rules with explicit usage maps, which would be quite verbose, we simply explain how the rules are extended. We reuse the variables names of the rules described in figure 18.

- For the rule `SEQ`, if each instruction t_i has a usage map Δ_i , the usage map of the sequence Δ is given by:

$$\Delta = \begin{cases} \text{Rename}\{r := \mathbf{res}\}(\Delta_1; \dots; \Delta_n) & \text{if } r \neq \emptyset \\ (\Delta_1; \dots; \Delta_n) & \text{if } r = \emptyset \end{cases}$$

- For the rule `BLOCK`, if we name Δ_r the usage map of the sequence, and Δ_c the usage map of the subtraction of `StackAllocCells`, the usage map of the whole block is $(\Delta_r; \Delta_c)$.
- For the rule `FUN`, if we name Δ_1 the usage map of the function body, Δ_2 the usage map of the subtraction, and S the generated specification hypothesis, then the usage map of the function definition is $((\Delta_1; \Delta_2) \vdash \Gamma_0) \cup \{\mathbf{res} : \text{ensured}, S : \text{ensured}\}$. Indeed, viewed from outside the only dependencies of the function definition are the pure resources captured from the surrounding context.
- For the rule `APP`, if Δ_σ is a usage map containing an entry required for each x_i and each pure resource from Γ_0 mentioned in σ , Δ_p is the usage map of the subtraction on Γ_0 , Δ_q is a usage map containing one produced (resp. ensured) for each linear (resp. pure) resource in Γ_q , and Δ_f the usage map of the `CloseFrac` operation, the usage map of the application is $(\Delta_\sigma; \Delta_p; \Delta_q; \Delta_f)$.

- For the rule **FOR**, only the outer contract instantiation and the required pure variables needed to typecheck the loop body are considered for computing the usage map. If Δ_p is the usage map of the subtraction $\Gamma_0 \ominus \Gamma_p$, Δ_b is the usage of the body of the loop, Δ'_q is the usage of the subtraction on Γ'_q , Γ_q is a usage map containing one produced (resp. ensured) for each linear (resp. pure) resource in Γ_q , and Γ_r is the usage of the `CloseFrac`s operation, then $(\Delta_p; ((\Delta_b; \Delta'_q) \vdash \Gamma_0); \Delta_q; \Delta_r)$ is the usage map of the for-loop. Note that the $((\Delta_b; \Delta'_q) \vdash \Gamma_0)$ part of this usage map correspond to the usage of pure resources from outside the loop in the body of the loop (they all have a required usage kind).
- For the rule **IF**, applied to a conditional **if** _{Γ} t_1 **then** t_2 **else** t_3 , it is always sound (though possibly imprecise) to combine the usage map Δ_0 of the condition expression t_1 to another usage map Δ_1 that gives a full usage to each linear resource in Γ_1 (the output context of t_1) and a usage map Δ_r that contains a produced usage for each linear resource of Γ_r . For the usage of pure resources, we name Δ_2 (resp. Δ_3) the required usage from t_2 (resp. t_3). Then, we take all the pure facts from Γ_r that are not in Γ_1 as ensured in a usage map Δ'_r . In summary, we compute the usage map of the whole conditional as $(\Delta_0; \Delta_1; \Delta_2; \Delta_3; \Delta_r; \Delta'_r)$.

5.4 Minimization of Triples

The *triple minimization operation* is used for typing function calls with effectful arguments and for minimizing loop contracts produced by transformations. The operation $\text{Minimize}(\Gamma, \Gamma', \Delta)$ is defined when its input corresponds to a valid triple $\{\{\Gamma\}\} t^\Delta \{\{\Gamma'\}\}$. The output of the operation is a quadruplet $(E^{\text{fracs}}, \hat{F}, \hat{F}', F^{\text{framed}})$.

- \hat{F} is the *minimized linear precondition*: a linear context containing resources from Γ .linear that are needed to typecheck t .
- \hat{F}' is the *minimized linear postcondition*: a linear context produced after typechecking t if we give only \hat{F} as the linear precondition.
- F^{framed} is the *maximal frame*: a linear context of resources from Γ .linear that were superfluous in the typechecking of t . It means resources in F^{framed} can be framed during the typechecking of t . Since these resources are not touched by t , they must also occur in Γ' .linear.
- E^{fracs} is the *generated fraction set*: a set of pure fractions that are created by the `Minimize` algorithm to give only an arbitrary subfraction of the resource in Γ .linear in \hat{F} when such a fractional resource suffices to typecheck t .

Concretely, the result of `Minimize` is guided by the entries in the usage map Δ , which is computed when typechecking t .

- If t can typecheck without a linear resource H , then H should be put in the frame F^{framed} .
- If t can typecheck with only the uninitialized version of H (because, for instance, it starts by overwriting the data accessible through H), then $\text{Uninit}(H)$ should be placed in \hat{F} .
- If t can typecheck with only an arbitrary subfraction of H (because, for instance, t only reads using H), then a fresh fraction α should be created and placed in E^{fracs} , the resource αH should be placed in \hat{F} , and $(1 - \alpha)H$ should remain in F^{framed} .

Detailed examples and an algorithmic description of `Minimize` can be found in Appendix D. From the perspective of establishing soundness results, the following three properties about the quadruplet $(E^{\text{fracs}}, \hat{F}, \hat{F}', F^{\text{framed}})$ are useful.

- $\{\{\Gamma.\text{pure}, E^{\text{fracs}} \mid \hat{F}\}\} t \{\{\Gamma'.\text{pure}, E^{\text{fracs}} \mid \hat{F}'\}\}$, which corresponds to the minimized triple.
- $\Gamma \Rightarrow \langle \Gamma.\text{pure}, E^{\text{fracs}} \mid \hat{F} \star F^{\text{framed}} \rangle$, which describes the decomposition of Γ .
- $\langle \Gamma'.\text{pure}, E^{\text{fracs}} \mid \hat{F}' \star F^{\text{framed}} \rangle \Rightarrow \Gamma'$, which describes the decomposition Γ' .

2059 5.5 Typechecking of Order-Irrelevant Subexpressions

2060 We next explain how to leverage the minimization procedure for typechecking functions calls
 2061 that are not in A-normal form, but possibly include effectful subexpressions. In C, and in our
 2062 subset OptiC, the arguments of a function call may be evaluated in an arbitrary order. The fact
 2063 that the order is not fixed is interesting because it leaves additional flexibility for optimizations.
 2064 Our typesystem checks that, for well-typed OptiTrust programs, the order of evaluation is indeed
 2065 irrelevant. To that end, we consider a sufficient condition: that the arguments can be evaluated in
 2066 parallel, in the sense that the side-effects performed by the arguments either should be disjoint.

2067 Remark: there exists valid C programs that fail to typecheck in OptiTrust because our condition
 2068 is slightly more restrictive. However, such programs may be easily rewritten by binding variables
 2069 to arguments before the function call.

2070 The rule SUBEXPR reduces the typechecking of a term with possibly effectful subexpressions to
 2071 the typechecking of a term whose subexpressions are program variables. In particular, the rule may
 2072 be used to compute the output context associated with a call of the form $f(t_1, \dots, t_n)$ in an input
 2073 context Γ_0 , by reducing the problem to the typechecking of a call of the form $f(x_1, \dots, x_n)$, in an
 2074 input context Γ_p that binds the fresh variables x_i .

2075 The rule SUBEXPR, shown below, applies to a term of the form $\hat{\mathcal{E}}[t_0, \dots, t_n]$, where $\hat{\mathcal{E}}$ denotes a
 2076 multi-evaluation-context and where the t_i variables denote the subterms in evaluation position.
 2077 A multi-evaluation-context is a term with ordered holes that are all in evaluation position. We
 2078 write $\hat{\mathcal{E}}[t_0, \dots, t_n]$ the operation that fills the holes with terms t_0 to t_n . For example, if $\hat{\mathcal{E}}$ denotes the
 2079 multi-evaluation-context $\square(\square, \dots, \square)$, then the application $\hat{\mathcal{E}}[f, t_1, \dots, t_n]$ produces the function call
 2080 $f(t_1, \dots, t_n)$.

2081 The goal of the rule SUBEXPR is to distribute the linear resources from the input context Γ_0 between
 2082 the subterms t_i . If several subterms read the same resource, then this resource needs to be split. If
 2083 one subterm reads a resource and another subterm modifies that same resource, the rule must fail
 2084 to apply. The key idea is to typecheck the subterms one after the other, taking advantage of the
 2085 Minimize operation to remove the minimal amount of resources from the input context, thereby
 2086 leaving as many resources as possible for the remaining subterms.

2088 SUBEXPR

$$\begin{array}{l}
 2089 \forall i \in [1, n]. \quad \{\{\Gamma_{i-1}\}\} t_i^{\Delta_i} \{\{\Gamma'_i\}\} \wedge (E_i^{\text{fracs}}, \hat{F}_i, \hat{F}'_i, F_i^{\text{framed}}) = \text{Minimize}(\Gamma_{i-1}, \Gamma'_i, \Delta_i) \wedge x_i \text{ fresh} \\
 2090 \quad \forall i \in [1, n]. \quad \Gamma_i = \langle \Gamma_i.\text{pure}, E_i^{\text{fracs}} \mid F_i^{\text{framed}} \rangle \wedge \hat{\Gamma}'_i = \langle \Gamma'_i.\text{pure} \uparrow \Delta_i.\text{ensured} \mid \hat{F}'_i \rangle \\
 2091 \quad \Gamma_p = \text{CloseFrac}^{\text{Ap}}(\Gamma_n \otimes \star_{i \in [0, n]} \text{Rename}\{\mathbf{res} := x_i\}(\hat{\Gamma}'_i)) \\
 2092 \quad \{\{\Gamma_p\}\} \hat{\mathcal{E}}[x_1, \dots, x_n]^{\Delta_q} \{\{\Gamma_q\}\} \\
 2093 \quad \Delta = \text{Rename}\{\mathbf{res} := x_1\}(\Delta_1); \dots; \text{Rename}\{\mathbf{res} := x_n\}(\Delta_n); \Delta_p; \Delta_q \\
 2094 \hline
 2095 \quad \{\{\Gamma_0\}\} \hat{\mathcal{E}}[t_1, \dots, t_n]^{\Delta} \{\{\Gamma_q\}\}
 \end{array}$$

2096 Appendix E presents an example application of this rule.

2098 5.6 Formal Properties of Usage Maps

2099 To conclude this section, we present three propositions that specify the contents of usage maps
 2100 computed by our typing algorithm. These propositions have guided all our definitions. We claim
 2101 that these propositions hold by design; we leave to future work a thorough mechanized proof of
 2102 the claims.
 2103

2104 Consider an algorithmic triple $\{\{\Gamma\}\} t^{\Delta} \{\{\Gamma'\}\}$, where Γ decomposes as $\langle E \mid F \rangle$ and Γ' decomposes
 2105 as $\langle E' \mid F' \rangle$. The first proposition explains how F and F' are partitioned by the usage map.

2108 PROPOSITION 5.1 (DECOMPOSITION BY USAGE).

$$2109 \quad F = F \vdash \Delta. \text{full} * F \vdash \Delta. \text{uninit} * F \vdash \Delta. \text{splittedFrac} * F \vdash \Delta. \text{joinedFrac} * F \setminus \Delta$$

$$2110 \quad F' = F' \vdash \Delta. \text{produced} * F' \vdash \Delta. \text{splittedFrac} * F' \vdash \Delta. \text{joinedFrac} * F' \setminus \Delta$$

2112 The second proposition explains how E' extends E , and how the frame resources from F are
 2113 preserved in F' . Besides, the proposition captures the fact that a resource with usage `splittedFrac`
 2114 or `joinedFrac` in F must also appear in F' , albeit with a possibly different fraction.

2115 PROPOSITION 5.2 (PRESERVED PARTS OF TYPING CONTEXTS).

$$2116 \quad \begin{aligned} & \{\{E \mid F\}\} t^\Delta \{\{E' \mid F'\}\} \\ \implies & E' = E, (E' \vdash \Delta. \text{ensured}) \\ & \wedge F' \setminus \Delta = F \setminus \Delta \\ & \wedge (\forall y, \forall H, (\exists \alpha, (y : \alpha H) \in F \vdash \Delta. \text{splittedFrac}) \iff (\exists \beta, (y : \beta H) \in F' \vdash \Delta. \text{splittedFrac})) \\ & \wedge (\forall y, \forall H, (\exists \alpha, (y : \alpha H) \in F \vdash \Delta. \text{joinedFrac}) \iff (\exists \beta, (y : \beta H) \in F' \vdash \Delta. \text{joinedFrac})) \end{aligned}$$

2122 The third proposition explains that the entries of the usage map Δ imply that the term t may
 2123 be typed in a context with smaller footprint. If a resource H appears in F but not used, then it is
 2124 omitted. If a resource H appears in F but used only as `uninit` (i.e., the corresponding cells are written
 2125 before being read), then the resource is replaced with `Uninit(H)`. If a resource H is only read, then it
 2126 is replaced with a fractional resource αH , where α is a constant that can be chosen arbitrarily small.
 2127 These operations are formally captured in the following statement, which also covers additional
 2128 complications related to the case where a set of input resources are splitted or merged together for
 2129 producing certain output resources. Below, $\{\hat{\Gamma}\} t \{\hat{\Gamma}'\}$ corresponds to a semantic triple, a notion
 2130 introduced in Section 4.8; and the \hat{F} variables are explained afterwards.

2132 PROPOSITION 5.3 (MINIMIZATION WITH USAGE MAPS).

$$2133 \quad \begin{aligned} & \{\{E \mid F\}\} t^\Delta \{\{E' \mid F'\}\} \\ \implies & \forall \alpha, \exists \hat{F}^{SP}, \exists \hat{F}^{ST}, \exists \hat{F}^{JS}, \exists \hat{F}^{JF}, \\ & \text{let } \hat{\Gamma} := \langle E \vdash \Delta. \text{required} \mid F \vdash \Delta. \text{full} * \text{IntoUninit}(F \vdash \Delta. \text{uninit}) * \alpha(F \vdash \Delta. \text{splittedFrac}) \rangle \text{ in} \\ & \text{let } \hat{\Gamma}' := \langle E' \vdash \Delta. \text{required}, E' \vdash \Delta. \text{ensured} \mid F' \vdash \Delta. \text{produced} * \hat{F}^{SP} * \hat{F}^{JS} * \hat{F}^{JF} \rangle \text{ in} \\ & \quad \{\hat{\Gamma}\} t \{\hat{\Gamma}'\} \\ & \wedge \alpha(F \vdash \Delta. \text{splittedFrac}) \iff \hat{F}^{SP} * \hat{F}^{ST} \\ & \wedge F' \vdash \Delta. \text{splittedFrac} \iff (1 - \alpha)(F \vdash \Delta. \text{splittedFrac}) * \hat{F}^{SP} * \hat{F}^{JS} \\ & \wedge F' \vdash \Delta. \text{joinedFrac} \iff (F \vdash \Delta. \text{joinedFrac}) * \hat{F}^{JF} \end{aligned}$$

2142 We explain the role of the \hat{F} variables at a high level, by means of example.

- 2144 • Assume a resource $y : (\beta - \gamma)H$ from F with usage `joinedFrac` in Δ meaning that t does not
 2145 read this resource. It must be the case that t produces (directly or indirectly) a resource
 2146 γH that is immediately merged into y . This produced resource appears in \hat{F}^{JF} , short for
 2147 *joined-framed*.
- 2148 • Assume a resource $y : (\beta - \gamma)H$ from F with usage `splittedFrac` in Δ meaning that t reads
 2149 this resource. Assume moreover t produces a resource γH that is immediately merged into
 2150 y . This produced resource appears in \hat{F}^{JS} , short for *joined-split*.
- 2151 • Assume a resource $y : \beta H$ from F with usage `splittedFrac` in Δ . In the minimized triple
 2152 for t , which takes an arbitrarily-small fraction α of the `splittedFrac` resources, (up to) two
 2153 subfractions of $\alpha\beta H$ may be involved. A first subfraction corresponds to a subresource of
 2154 y that t does not alter; this subfraction appears in \hat{F}^{SP} , short for *split-preserved*. A second
 2155 subfraction corresponds to a subresource of y that t alters; this subfraction appears in \hat{F}^{ST} ,

short for *split-transformed*. The line $\alpha(F \cdot \Delta.\text{splittedFrac}) \Leftrightarrow \hat{F}^{\text{SP}} \star \hat{F}^{\text{ST}}$ captures that the `splittedFrac` resources from F divide between \hat{F}^{SP} and \hat{F}^{ST} .

Again, we leave it to future work to carry out a mechanized proof of these propositions.

6 JUSTIFYING TRANSFORMATION CORRECTNESS

In this section, we explain how OptiTrust leverages resource typing information to check the correctness of the transformations requested by the programmer. The aim of this section is not to cover all the transformations implemented in OptiTrust, but to present a representative subset thereof. We focus in particular on transformations that leverage the resource information in an interesting way. All the transformations presented in this section are invoked multiple times in our case studies from Section 2.

Recall that we only need to check the correctness of *basic* transformations, because *combined* transformations are defined as composition of basic transformations. For each basic transformation considered, we present a generally applicable *sufficient condition* for the transformation to be correct. For certain transformations, this sufficient condition includes the property that the produced program successfully typechecks. For other transformations, typechecking is not required to ensure correctness. Nevertheless, OptiTrust re-typechecks the program after every transformation, for the purpose of allowing the application of subsequent transformations.

A number of *basic* transformation might seem “simple” to the reader. This simplicity is precisely a strength of OptiTrust. As explained in the introduction, we aim to minimize the trusted code base, by considering the simplest possible *basic* transformations and by implementing as many transformations as possible as composition of *basic* transformations. Other transformations are more involved. In fact, for certain loop transformations, we have considered only simplified sufficient conditions, which we could further generalize in future work.

Before presenting the key aspects of specific transformations, we introduce notation for describing transformations. Transformations apply to instructions or group of instructions; they depend on typing context and usage information; and they produce code with possibly updated loop contracts, and possibly including ghost instructions. Hence, we need a convenient way to visualize all these entities.

Notation for Well-Typed Programs. Transformations leverage typing information, not only for checking correctness, but also for guiding the generation of the output code. Recall from the previous section that our typechecking algorithm computes, for every subterm t , its input context Γ_1 , its output context Γ_2 , and its usage map Δ , establishing triples of the form $\{\{\Gamma_1\}\} t^\Delta \{\{\Gamma_2\}\}$. In this section, we use an alternative syntax, better-suited for describing the input of transformations. If t denotes an instruction, we write $\boxed{\Gamma_1 \ t; \Delta \ \Gamma_2}$ as straight-line syntax for $\{\{\Gamma_1\}\} t^\Delta \{\{\Gamma_2\}\}$.

Groups of Instructions. Some transformations operate on groups of consecutive instructions. We let the meta-variable T range over a (possibly empty) group of instructions. We generalize our alternative syntax by writing $\boxed{\Gamma_1 \ T; \Delta \ \Gamma_2}$, where Γ_1 and Γ_2 denotes the initial and final contexts, and Δ denotes the *composition* of the usages from the group of instructions, as defined in Section 5.3:

$$\boxed{\Gamma_0 \ T; \Delta \ \Gamma_n} \equiv \boxed{\Gamma_0 \ t_1; \Delta_1 \Gamma_1 \ t_2; \Delta_2 \Gamma_2 \ \dots t_n; \Delta_n \ \Gamma_n} \quad \begin{array}{l} \text{where } T \equiv t_1; t_2; \dots; t_n \\ \text{and } \Delta \equiv \Delta_1; \Delta_2; \dots; \Delta_n \end{array}$$

Program Contexts. Transformations generally apply to a program subterm, that is, apply under a *program context*. Unlike evaluation contexts, program contexts can reach subterms that are not in evaluation position. We let the meta-variable \mathcal{E} range over program contexts. For example, evaluating a subexpression $1 + 1$ that appears in a program context \mathcal{E} is described as the transition

from $\mathcal{E}[1 + 1]$ to $\mathcal{E}[2]$. We also allow program contexts to denote a hole in the middle of a sequence. For example, swapping two instructions that appear inside a sequence is described as the transition from $\mathcal{E}[t_1; t_2]$ to $\mathcal{E}[t_2; t_1]$, to be interpreted as a transition from $\mathcal{E}'[\{T_0; t_1; t_2; T_3\}]$ to $\mathcal{E}'[\{T_0; t_2; t_1; T_3\}]$, where \mathcal{E}' denotes the program context associated with the outer sequence that contains $t_1; t_2$. We will only explicitly mention the surrounding program context \mathcal{E} for the first few transformations.

Evaluation Contexts. Some transformations operate on subexpressions that appear inside an instruction. For those, we may need to restrict the form of the program contexts in which the subexpression may appear, to avoid nontrivial control-flow arising from, e.g., a conditional. Recall from Section 5.5 that an *evaluation context*, written $\hat{\mathcal{E}}$, denotes a program context whose holes (possibly just one) are in evaluation position. For example, $f(g(\square, 2), g(3, a + 4))$ is an evaluation context with a single hole written \square . One key property is that the rewrite $\hat{\mathcal{E}}[t] \mapsto \mathbf{let} \ x = t; \hat{\mathcal{E}}[x]$ is correct for any evaluation context $\hat{\mathcal{E}}$. The reciprocal rewrite holds for programs well-typed in our type system.

The validity of this rewrite rule, and more generally the interest of evaluation contexts for transformations, crucially relies on the hypothesis that the input code typechecks against our typing rules. Indeed, the `SUBEXPR` rule ensures that, if a function has multiple arguments, then the available resources are distributed across the arguments—only read-only resources can be distributed onto several arguments. For example, $f(g_1(), g_2(), g_3())$ is equivalent to $\mathbf{let} \ x = g_2(); f(g_1(), x, g_3())$ because, if the former term is well-typed, then the effects of $g_2()$ do commute with the effects of $g_1()$ and $g_3()$.

Notation for Introducing Ghost Calls. Recall that a call to a ghost function is an instruction that semantically behaves as a no-op, yet updates the context available. In the output of transformations, we write $\mathbf{ghost}(\Gamma \rightarrow \Gamma')$ to mean the insertion of an appropriate ghost call $g()$, such that g admits Γ as precondition and Γ' as postcondition. Concretely, the effect of $\mathbf{ghost}(\Gamma \rightarrow \Gamma')$ is to consume the resources Γ then to produce the resources Γ' .

We are now ready to present transformations. We begin with transformations on instructions and variable bindings, then move on to transformations on storage, and transformations on loops.

6.1 Transformations on Sequences of Instructions

Moving Instructions. The basic transformation `Instr.move` allows to move a group of instructions to a given destination within the same sequence. Doing so amounts to swapping a group of instructions T_1 with an adjacent group of instructions T_2 . The *move* transformation turns a program of the form $\mathcal{E}[T_1; T_2]$ into $\mathcal{E}[T_2; T_1]$, where \mathcal{E} denotes a program context. The transformation is formalized as shown below. The variables Δ_1 and Δ_2 denote the usage associated with T_1 and T_2 . The correctness criterion, stated on the right-hand-side, is explained next.

$$\boxed{\mathcal{E}[T_1; \Delta_1; T_2; \Delta_2]} \mapsto \boxed{\mathcal{E}[T_2; T_1]} \quad \text{correct if: } \begin{cases} \Delta_1.\text{alter} \cap \Delta_2 = \emptyset \\ \Delta_2.\text{alter} \cap \Delta_1 = \emptyset \end{cases}$$

The expression $\Delta_1.\text{alter}$ denotes the resources that T_1 adds or removes (consumes, produces, or ensures). It excludes resources that remained unaltered (carving or merging a fraction does not count). The property $\Delta_1.\text{alter} \cap \Delta_2 = \emptyset$ captures the idea that if a resource is altered by T_1 , then T_2 must not use it (this includes “Write After Read” dependencies), otherwise swapping T_1 and T_2 might not be correct. (The resource intersection operator \cap was defined in Section 5.2.) The second property, namely $\Delta_2.\text{alter} \cap \Delta_1 = \emptyset$, captures the symmetrical property: if a resource is altered by T_2 , then T_1 must not use it (this includes “Read After Write” dependencies). When both conditions

are met, the only resources that both T_1 and T_2 depend on are accessed in read-only mode, and T_1 and T_2 may be safely swapped without impacting their evaluation result.

Deleting Instructions. The basic transformation `Instr.delete` allows deleting a group of instructions T from a sequence. It therefore maps a program $\mathcal{E}'[\{T_0; T; T_2\}]$ to a program $\mathcal{E}'[\{T_0; T_2\}]$, for a program context \mathcal{E}' . Following our convention that program contexts may describe subsequences, we may also describe the transformation as mapping $\mathcal{E}[T]$ to $\mathcal{E}[\emptyset]$, for a program context \mathcal{E} .

Intuitively, the deletion operation preserves program semantics if the resources altered by T are not observed by the rest of the program. More precisely, if T has been typechecked as $\Gamma T; \Delta$, then we start with the resources Γ corresponding to not executing T , then forget the contents of the resources that might be different when not executing T . The resources to “uninitialize” Γ_m are computed by the filtering operation $\Gamma \vdash \Delta.\text{alter}$. (Filtering was defined in Section 5.2.) Finally, we typecheck the auxiliary program $\mathcal{E}[G]$, in which the T is replaced with a ghost instruction G casting the Γ_m resources into their corresponding “uninitialized form”, as performed by the `IntoUninit` operator. If a resource H is consumed by T , then G consumes H and produces `Uninit(H)`.

The transformation can therefore be formalized as follows.

$$\boxed{\mathcal{E}[\Gamma T; \Delta]} \mapsto \boxed{\mathcal{E}[\emptyset]} \quad \text{correct if } \mathcal{E}[\mathbf{ghost}(\Gamma_m \longrightarrow \text{IntoUninit}(\Gamma_m))] \text{ typechecks,} \\ \text{where } \Gamma_m \equiv \Gamma \vdash \Delta.\text{alter.}$$

If the auxiliary program $\mathcal{E}[G]$ typechecks, then we can discard this program, and safely replace the original program $\mathcal{E}[T]$ with $\mathcal{E}[\emptyset]$. Note that this pattern of introducing an auxiliary program for the purpose of evaluating a correctness criterion will appear again for other transformations.

Inserting Instructions. The transformation `Instr.insert` refines a program from $\mathcal{E}[\emptyset]$ to $\mathcal{E}[T]$, where T denotes the group of inserted instructions. The correctness criterion, described below, is essentially the same as that for instruction deletion. Indeed, for $\mathcal{E}[T]$ to admit the same semantics as $\mathcal{E}[\emptyset]$, it suffices that $\mathcal{E}[\emptyset]$ admits the same semantics as $\mathcal{E}[T]$.

$$\boxed{\mathcal{E}[\emptyset]} \mapsto \boxed{\mathcal{E}[T]} \quad \text{correct if:} \\ \begin{array}{l} (1) \text{ the program } \mathcal{E}[T] \text{ typechecks as } \mathcal{E}[\Gamma T; \Delta] \text{ for some } \Gamma \text{ and } \Delta; \\ (2) \text{ the program } \mathcal{E}[\mathbf{ghost}(\Gamma_m \longrightarrow \text{IntoUninit}(\Gamma_m))] \text{ type-} \\ \text{checks, where } \Gamma_m \equiv \Gamma \vdash \Delta.\text{alter, for the above values of } \Gamma \text{ and } \Delta. \end{array}$$

Idempotent Terms. A number of transformations depend on the notion of idempotence. In the C23 standard, an expression is said to be “idempotent” if, intuitively, evaluating this expression multiple times in immediate sequence produces the same results. In OptiTrust, we leverage our resource analysis to capture a practical over-approximation of idempotence.¹³ A term can be considered idempotent if the resources that this term produces correspond: (1) either to uninitialized resources that were consumed by this term; or (2) to read-only resources that the term consumes and returns with the exact same fraction. These criteria may be formalized as follows.

$$\text{A term } T \text{ that appears in a program } \mathcal{E}[\Gamma_1 T; \Delta \Gamma_2] \text{ is idempotent iff:} \quad \left\{ \begin{array}{l} \Delta.\text{full} = \emptyset \\ (\Gamma_2 \vdash \Delta.\text{produced}) \boxplus (\Gamma_1 \vdash \Delta.\text{uninit}) = (\sigma, \emptyset) \\ \Gamma_1 \vdash \Delta.\text{reads} = \Gamma_2 \vdash \Delta.\text{reads} \end{array} \right. \quad \text{for some } \sigma$$

¹³The C23 standard defines a number of related notions. In particular, an expression is said to be “effectless” iff “if any store operation that is sequenced during the execution is the modification of an object that synchronizes with the call”. An expression is said to be “reproducible” iff it is both effectless and idempotent. Reproducibility is essentially equivalent to the notion of pure expression in GCC’s terminology [ale 2022]. Due to our resource typing discipline, all OptiTrust terms can be considered “effectless”. Hence, in the context of OptiTrust, “idempotent” and “reproducible” are equivalent.

In particular, these criteria rule out terms that consume full resources, or produce resources they did not consume. For example, $x = y+1$, which reads y and assigns x is idempotent; however $x++$, which modifies x , is *not* idempotent. One key property that holds for an idempotent term t is that the following program equivalence holds: $\boxed{\text{let } x = t; \text{let } y = t; \mathcal{E}[x, \dots, y, \dots]} \leftrightarrow \boxed{\text{let } x = t; \mathcal{E}[x, \dots, x, \dots]}$.

Duplicating and Deduplicating Instructions. If an instruction T (or possibly a group of instructions) is idempotent, then after a first instruction T , a second instruction T may be inserted or removed without affecting the semantics. The transformation `Instr.dup` and its reciprocal `Instr.dedup` are formalized, for the general case of groups of instructions, as follows.

$$\boxed{\mathcal{E}[T]} \leftrightarrow \boxed{\mathcal{E}[T; T]} \quad \text{where } T \text{ is idempotent.}$$

Similarly, if a term t is idempotent, then after the instruction `let x = t`, an instruction `let y = t` may be inserted or removed, for a fresh variable y . The corresponding transformations are named `Instr.dup_let` and `Instr.dedup_let`. Thereafter, for brevity, we omit the program context surrounding the code snippet, previously written \mathcal{E} .

$$\boxed{\text{let } x = t;} \leftrightarrow \boxed{\text{let } x = t; \text{let } y = t;} \quad \text{where } t \text{ is idempotent and } y \text{ fresh.}$$

Deduplicating expressions is a building block for common subexpression elimination, which is detailed further on. Duplicating expressions can also improve performance in certain situations: recomputing a simple expression may be cheaper than storing its value in memory and subsequently retrieving this value, especially if the redundant computations are scattered in distinct loops.

6.2 Exploiting Equalities

Read after Write. The transformation `Eq.read_after_write` captures the fact that reading immediately after a write yields the value that was written. On its own, this transformation may seem of little interest; however, it is useful when combined with moves of the read or the write instruction.

$$\boxed{\text{set}(p, v); \hat{\mathcal{E}}[\text{get}(p)];} \mapsto \boxed{\text{set}(p, v); \hat{\mathcal{E}}[v];} \quad \text{correct if } \hat{\mathcal{E}} \text{ is an evaluation context and } v \text{ is a logical expression.}$$

Results of Idempotent Expressions. The transformation `Eq.idempotent` captures the fact that evaluating an idempotent expression twice yields equal results.

$$\boxed{\text{let } x = t; \text{let } y = t; \mathcal{E}[x]} \mapsto \boxed{\text{let } x = t; \text{let } y = t; \mathcal{E}[y]} \quad \text{correct if } \mathcal{E} \text{ is a program context and } t \text{ is idempotent.}$$

6.3 Transformations on Bindings

Inlining/Binding for Logical Expressions. The basic transformation `Variable.inline_pure` eliminates a binding of the form `let x = v`, where v is a logical expression, by replacing all occurrences of x with v . This transformation is always correct and requires no check. The reciprocal transformation, `Variable.bind_pure`, introduces a binding for one or several occurrences of a logical expression v . Likewise, it is always correct.

Inlining a Binding with a Single Occurrence, in the Next Instruction. The basic transformation `Variable.inline_one` eliminates a binding `let x = t` in programs where x has exactly one occurrence,

and this occurrence is contained in the immediately succeeding instruction, under an evaluation context $\hat{\mathcal{E}}$. As mentioned earlier, the correctness of this inlining transformation critically relies on the fact that our typing rules ensure that the order of evaluation of subexpressions is irrelevant.

$$\boxed{\text{let } x = t; \hat{\mathcal{E}}[x];} \longrightarrow \boxed{\hat{\mathcal{E}}[t];}$$
 correct if $\hat{\mathcal{E}}$ contains no other occurrence of x than the one in its hole, and the output program typechecks.

Inlining a Binding with Multiple Occurrences, in the Next Instruction. The transformation `Variable.inline_dup` expands a binding at one of its occurrences, without removing the binding. Here again, we consider an occurrence appearing in an immediately succeeding evaluation context. This transformation is implemented as a *combined* transformation, decomposed as shown below. Recall that we do not need to devise correctness criteria for combined transformations.

$$\boxed{\text{let } x = t; \hat{\mathcal{E}}[x];} \xrightarrow{\text{Instr.dup_let}} \boxed{\text{let } x = t; \text{let } y = t; \hat{\mathcal{E}}[x];} \xrightarrow{\text{Eq.idempotent}} \boxed{\text{let } x = t; \text{let } y = t; \hat{\mathcal{E}}[y];} \xrightarrow{\text{Variable.inline_one}} \boxed{\text{let } x = t; \hat{\mathcal{E}}[t];}$$

Inlining a Binding in the Scope of a Sequence. The combined transformation `Variable.inline` eliminates a binding `let x = t` in the general case. If t is a logical expression, then `Variable.inline_pure` is invoked. Otherwise, we implement the inlining as a combination of several of the aforementioned transformations. Indirectly, our combined transformation enforces the minimal checks required for eliminating a binding `let x = t` without affecting the semantics.

- If x has no occurrences, the effects of t need to be irrelevant to the rest of the program.
- If x has exactly one occurrence, then the effects of t needs to commute with all the instructions located between the binding on x and the occurrence of x .
- If x has several occurrences, then, in addition to the requirement from the previous case, t moreover needs to be idempotent.

Concretely, our transformation proceeds as follows. If there are no occurrences of x , it invokes the transformation `Instr.delete`. If there is exactly one occurrence of x , it attempts to move, using `Instr.swap`, the binding on x just in front of this binding, then invoke `Variable.inline_one`. If there are several occurrences of x in the sequence, then it moves the binding to the front of the first instruction that contains occurrences of x ; then it applies the transformation `Variable.inline_dup`; then it repeats the process until reaching the last occurrence of x . We show below an example decomposition of `Variable.inline`, where t is assumed to be idempotent.

$$\begin{aligned} & \text{let } x = t; g(); \text{set}(a, x); \text{set}(b, x); \\ \longmapsto & g(); \text{let } x = t; \text{set}(a, x); \text{set}(b, x); && (\text{Instr.swap}) \\ \longmapsto & g(); \text{let } x = t; \text{set}(a, t); \text{set}(b, x); && (\text{Variable.inline_dup}) \\ \longmapsto & g(); \text{set}(a, t); \text{let } x = t; \text{set}(b, x); && (\text{Instr.swap}) \\ \longmapsto & g(); \text{set}(a, t); \text{set}(b, t); && (\text{Variable.inline_one}) \end{aligned}$$

We leave to future work the support, in a combined transformation, of more complex patterns where occurrences of a non-pure binding appear in depth under control flow constructs.

Binding Introduction. The basic transformation `Variable.bind_one` is essentially the reciprocal of `Variable.inline_one`.

$$\boxed{\hat{\mathcal{E}}[t];} \mapsto \boxed{\mathbf{let } x = t; \hat{\mathcal{E}}[x];}$$

Folding for Additional Occurrences. The combined transformation `Variable.bind_dup` is essentially the reciprocal of `Variable.inline_dup`. (We implement it as a combination of `Variable.bind_one`, `Instr.swap`, `Eq.idempotent`, and `Instr.delete`.)

$$\boxed{\mathbf{let } x = t; T; \hat{\mathcal{E}}[t];} \mapsto \boxed{\mathbf{let } x = t; T; \hat{\mathcal{E}}[x];}$$

Common Subexpression Elimination. The combined transformation `Variable.bind` is essentially the reciprocal of `Variable.inline`. Internally, it exploits the transformations `Variable.bind_one` and `Variable.bind_dup` to introduce a binding that factorizes the evaluation of common subexpressions. For example, if e is idempotent and commutes with $g()$, the program “ $g(); \text{set}(a, t); \text{set}(b, y)$ ” can be transformed into “ $\mathbf{let } x = t; g(); \text{set}(a, x); \text{set}(b, x)$ ”.

6.4 Transformations on Storage

The purpose of this section is to present transformations for introducing, eliminating, and converting between various forms of storage. We present transformations operating on single cells, and omit from the discussion the generalizations to arrays and N -dimensional matrices.

Recall from Section 3 that a pure program variable `const int x = 3` is represented in the OptiTrust AST as `let x = 3`, that a non-pure stack-allocated variable `int x = 3` is represented as `let x = stackRef(3)`, and that an uninitialized variable `int x` is represented as `let x = stackAllocCell()`. For stack-allocated data, the resources produced by `stackAlloc` are automatically reclaimed at the end of the scope. For heap-allocated data, the resources produced by `heapAlloc` are consumed by the matching call to `free`.

Separating Declaration from Initialization. For a stack-allocated variable, the basic transformation `Variable.init_detach` separates its declaration from its initialization. This transformation is useful as a preliminary step for the combined transformation that hoists a variable declaration appearing inside a loop into an array allocated outside that loop. The basic transformation `Variable.init_attach` applies the reciprocal operation.

$$\boxed{\mathbf{let } x = \text{stackRef}(t);} \leftrightarrow \boxed{\mathbf{let } x = \text{stackAllocCell}(); \text{set}(x, t);}$$

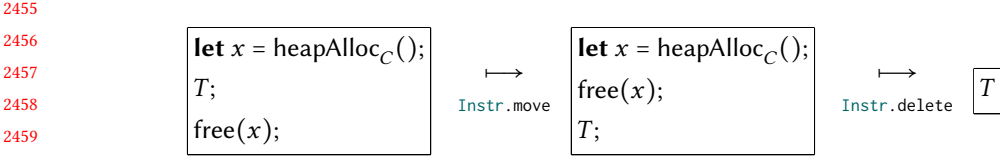
Converting between Stack and Heap Allocation. The basic transformation `Variable.to_heap` transforms an uninitialized stack-allocated storage into a corresponding heap-allocated storage. The transformation takes as optional argument the target at which the free instruction should be inserted; by default, it is placed at the end of the scope. The reciprocal transformation is named `Variable.to_stack`.

$$\boxed{\{T_1; \mathbf{let } x = \text{stackAlloc}_C(); T_2;\}} \leftrightarrow \boxed{\{T_1; \mathbf{let } x = \text{heapAlloc}_C(); T_2; \text{free}(x);\}}$$

Removal of Unused Storage. If a stack-allocated storage is never used, it may be removed by means of the operation `Instr.delete`. Concretely, the instruction `let x = stackAllocCell()` may be deleted if x has no occurrences, and the instruction `let x = stackRef(e)` may be deleted if moreover the effects performed by e are not observed by the rest of the program.

If a heap-allocated space is never used, then it may also be removed. To that end, one needs to delete both the `heapAlloc` and the corresponding `free` instructions. Neither of them can be removed

2451 independently, because both depend on each other. However, if we move using `Instr.move` the
 2452 `heapAlloc` instruction next to the `free` instruction, or vice-versa, then the group made of the two
 2453 instructions may be removed at once by means of `Instr.delete`. The combined transformation
 2454 `Variable.delete`, described below, performs this task.



2461 *Temporary Alternative Storage.* The transformation `Variable.local_name` is the most complex that
 2462 we have implemented in terms of operations on plain sequences of instructions. The transformation
 2463 `Variable.local_name` operates over a specified group of instructions, say T , for a specified storage,
 2464 say x . Over this scope, a fresh storage, call it y , is allocated. Just before executing T , the contents
 2465 of x are copied into y . All instructions from T are updated to use y instead of x . Just after these
 2466 instructions, the possibly-updated contents of y is copied into x . Depending on the situation, the
 2467 initial copy from x to y , or the final copy from y into x might be unnecessary—and even ill-typed.
 2468 Such unnecessary copy operations are omitted.

2469 The variable x may be allocated either on the stack or on the heap. The user may choose to
 2470 allocate y on the stack or on the heap. Moreover, our implementation supports the general case
 2471 where x is not just a variable but an N -dimensional matrix. In case where x is a matrix, y may
 2472 correspond to only a subset (i.e., a tile) of the matrix. The interest of the `local_name` transformation
 2473 is to enable the program to operate on a local piece of data. Crucially, the memory layout of this
 2474 data may be refined by subsequent transformations, for example to store the transposed of a matrix
 2475 in a cache friendly way (as in Section 2.3), or to enable vectorization.

2476 The transformation is described in Figure 20. There, the group of instructions T is represented as
 2477 $\mathcal{E}[x, \dots, x]$, i.e., as a program context with multiple occurrences of x . The typing context Γ_1 describes
 2478 the resources available before T , and Γ_2 the resources available after T . This typing information
 2479 is used not only for checking the correctness criterion, but also for guiding the generation of the
 2480 output code.

2481 The correctness criterion appears at the bottom of Figure 20. An essential aspect of this criterion
 2482 is to check that, during the execution of T , the resource H_x corresponding to the full permission on
 2483 x is “frozen” (i.e., made unavailable) in order to ensure that no operation may be performed on x
 2484 via potential aliases of this pointer. The first ghost call uses a standard technique for enforcing such
 2485 a “freeze” in Separation Logic: introducing a magic wand operator (\multimap), guarded by a token named
 2486 H in the postcondition of the ghost operation, and bound as H' in the rest of the sequence. The
 2487 heap predicate H' admits the type `Hprop`, which is the type of all heap predicates in Separation
 2488 Logic. This heap predicate H' serves the role of a *key* for unfreezing H_x at the desired point—here,
 2489 the end of the scope on which y is used in place of x , where the second ghost call is placed. As far
 2490 as the present paper is concerned, the magic wand operator can be viewed as a binary operator on
 2491 heap predicates whose definition needs not be revealed to the user.

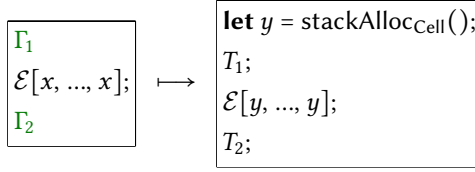
2492

2493 6.5 Transformations on Loops

2494 Loop transformations depend on the contracts associated with the loops from the input code. For
 2495 every loop being modified or introduced, the transformations also need to produce appropriate
 2496 contracts. In what follows, we present details for loop tiling, loop interchange, loop fission, and
 2497 loop hoisting. We then list other loop transformations that we have implemented.

2498

2499



where E is a multi-hole program context with one hole per occurrence of x , and where:

$$\begin{cases} T_1 \equiv \text{set}(y, \text{get}(x)); & \text{if } x \rightsquigarrow \text{Cell} \text{ or } \alpha(x \rightsquigarrow \text{Cell}) \text{ appears in } \Gamma_1 \\ T_1 \equiv \emptyset & \text{if } \text{Uninit}(x \rightsquigarrow \text{Cell}) \text{ appears in } \Gamma_1 \\ T_2 \equiv \text{set}(x, \text{get}(y)); & \text{if } x \rightsquigarrow \text{Cell} \text{ appears in } \Gamma_2 \\ T_2 \equiv \emptyset & \text{if } \text{Uninit}(x \rightsquigarrow \text{Cell}) \text{ or } \alpha(x \rightsquigarrow \text{Cell}) \text{ appears in } \Gamma_2 \end{cases}$$

correct if the program to the right typechecks successfully, where H_x is:

- $(x \rightsquigarrow \text{Cell})$
- $\alpha(x \rightsquigarrow \text{Cell})$
- or $\text{Uninit}(x \rightsquigarrow \text{Cell})$

depending on what appears in Γ_1 .

```

let y = stackAllocCell();
T₁;
ghost(⟨∅ | Hₓ⟩ → ⟨H : Hprop | H, (H * Hₓ)⟩); binding H as H'
E[y, ..., y];
ghost(⟨∅ | H', (H' * Hₓ)⟩ → ⟨∅ | Hₓ⟩);
T₂;

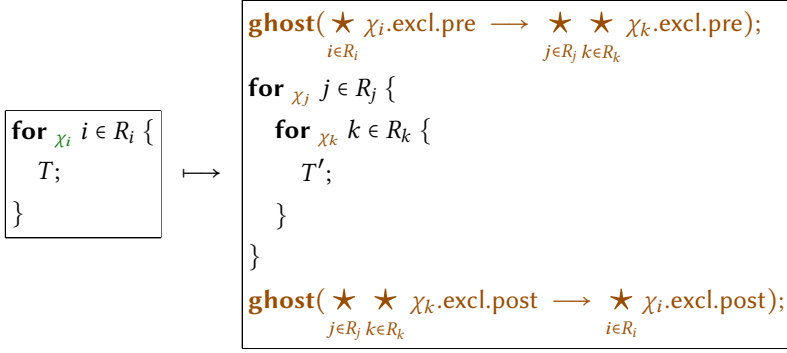
```

Fig. 20. Description of the basic transformation `Variable.local_name`.

Loop Tiling. The basic transformation `Loop.tile` allows tiling (a.k.a. strip-mining) a loop. Concretely, it transforms a loop, say with index i , into two nested loops, with indices j and k . Intuitively, the outer loop on j iterates over the *blocks*, whereas the inner loop on k iterates inside every block. Depending on the form of the input range, and on whether the block size divides the width of the loop range, the transformation is able to generate different ranges for the output loops. For each kind of output, the expression `RecoverIndex(j, k)` indicates how to compute the original index i in terms of the two new indices j and k .

The 4 variants supported by `Loop.tile` are described in Figure 21. The ranges of the three loops are written R_i , R_j and R_k , respectively. A range is of the form **range**(*start, stop, step*). The notation *start..stop* is a shorthand for **range**(*start, stop, 1*). In particular, $0..n$ describes the range of values from 0 inclusive to n exclusive. The contracts for the three loops involved are written χ_i , χ_j and χ_k , respectively. To typecheck the output code, *ghost tiling* operations need to be inserted, materialized before and after the produced loops in the figure. Indeed, the loop on i consumes, in particular, the resource: $\star_{i \in R_i} \chi_i.\text{excl.pre}$ whereas the loop on j consumes instead: $\star_{j \in R_j} \star_{k \in R_k} \chi_k.\text{excl.pre}$.

Loop Interchange. The basic transformation `Loop.swap` allows interchanging (i.e. swapping) two loops. It is described at the top of Figure 22. There exists a general criterion capturing when two loops may be swapped, however this criterion requires reasoning about the resources required by specific iterations, e.g. when executing T for iterations i, j and i', j' with $i' > i$ and $j > j'$. Instead, we focus on two conditions that are simpler yet sufficient for many practical situations: if at least one of the outer loop or the inner loop is parallelizable, then swapping the two loops is correct. Figure 22 describes the case where the outer loop is parallelizable. The case where the inner loop is parallelizable, not shown, is treated with just a few changes.



where:

$$\begin{aligned}
 T' &\equiv [i \mapsto \text{RecoverIndex}(j, k)](T) \\
 \chi_k &\equiv [i \mapsto \text{RecoverIndex}(j, k)](\chi_i) \\
 \chi_j &\equiv \begin{cases} \text{vars} \equiv \chi_i.\text{vars} \\ \text{shrd} \equiv \chi_i.\text{shrd} \\ \text{excl} \equiv \{\text{pre} \equiv \star_{k \in R_k} \chi_k.\text{excl.pre}; \text{post} \equiv \star_{k \in R_k} \chi_k.\text{excl.post}\} \end{cases}
 \end{aligned}$$

with the following possible instantiations for the ranges:

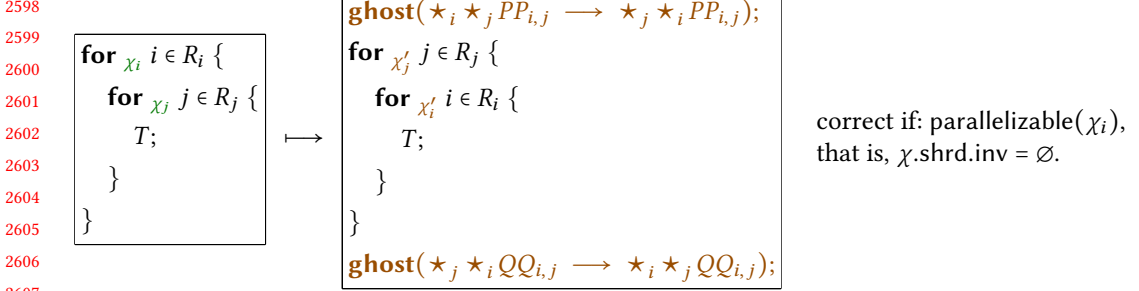
Variant	Range R_i	Range R_j	Range R_k	Formula for recovering i : $\text{RecoverIndex}(j, k)$
A	$0..(m \times b)$	$0..m$	$0..b$	$j * b + k$
B	$0..n$ where b divides n	$0..(n/b)$	$0..b$	$j * m + k$
C	$0..n$ where b divides n	range (0, n , b)	$j..j + b$	k
D	$0..n$	range (0, n , b)	$j..\min(j + b, n)$	k

Fig. 21. Description of the 4 variants of the basic transformation `Loop.tile`.

The first step is to partition the resources from the inner loop contract depending on where they come from relative to the resources from the outer loop. We name partitions by using the first letter to denote its inner loop origin, and the second letter to denote its outer loop origin. We use I for invariant, R for shared reads, P for exclusive precondition and Q for exclusive postcondition. For example, the inner shared reads are partitioned into RP_i that comes from the outer precondition, and RR that comes from the outer shared reads.

Then, we appropriately place the resources obtained from the partitioning in the contracts χ'_i and χ'_j associated with the swapped loops. Compared with χ_j , the new contract χ'_j essentially adds a \star_i operator to certain components. Compared with χ_i , the new contract χ'_i removes occurrences of the \star_j operators. Note that the loop on i remains parallelizable. Around the new loop nest, a pair of ghost operations is inserted for swapping groups of resources—a necessary step to match the resources required by the new loop nest.

Loop Fission. The transformation `Loop.fission`, in its *basic* version, breaks a loop with body $T_1; T_2$ into two loops over the same range, a first loop with body T_1 , and a second loop with body T_2 . The transformation is described in Figure 23. As for loop swapping, there exists a general correctness criterion expressed using inequalities on indices, but for now we focus on a simpler yet practical criterion.



2608 The contracts from the input code are decomposed as follows:

2609 $\chi_i.\text{shrd} = \{\text{inv} = \emptyset, \text{reads} = (\star_j PR_j \star IR \star RR)\}$

2610 $\chi_i.\text{excl} = \{\text{pre} = (\star_j PP_{i,j} \star IP_i \star RP_i), \text{post} = (\star_j QQ_{i,j} \star IP_i \star RP_i)\}$

2611 $\chi_j.\text{shrd} = \{\text{inv} = (IP_i \star IR), \text{reads} = (RP_i \star RR)\}$

2612 $\chi_j.\text{excl} = \{\text{pre} = (PP_{i,j} \star PR_j), \text{post} = (QQ_{i,j} \star PR_j)\}$

2614 The contracts for the output code are built as follows:

2615 $\chi'_j \equiv \begin{cases} \text{vars} \equiv \chi_j.\text{vars} \\ \text{shrd} \equiv \{\text{inv} \equiv (\star_i IP_i \star IR), \text{reads} \equiv (\star_i RP_i \star RR)\} \\ \text{excl} \equiv \{\text{pre} \equiv (\star_i PP_{i,j} \star PR_j), \text{post} \equiv (\star_i QQ_{i,j} \star PR_j)\} \end{cases}$

2616

2617

2618

2619

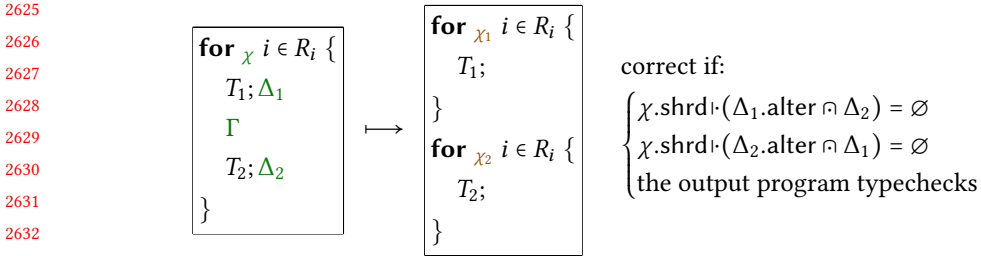
2620 $\chi'_i \equiv \begin{cases} \text{vars} \equiv \chi_j.\text{vars} \\ \text{shrd} \equiv \{\text{inv} \equiv \emptyset, \text{reads} \equiv (PR_j \star IR \star RR)\} \\ \text{excl} \equiv \{\text{pre} \equiv (PP_{i,j} \star IP_i \star RP_i), \text{post} \equiv (QQ_{i,j} \star IP_i \star RP_i)\} \end{cases}$

2621

2622

2623

2624 Fig. 22. The basic transformation `Loop.swap`, in the particular case where the outer loop is parallelizable.



2633 with:

2634 $\rightarrow F \equiv \Gamma \boxplus \chi.\text{shrd.inv} \quad \rightarrow F_{\text{cut}} \equiv F \boxplus \text{StackAllocCells}(T_1)$

2635

2636 $\chi_1 \equiv \begin{cases} \text{vars} \equiv \chi.\text{vars} \\ \text{shrd} \equiv \chi.\text{shrd} \cdot \Delta_1 \\ \text{excl} \equiv \{\text{pre} \equiv \chi.\text{excl.pre}, \text{post} \equiv F_{\text{cut}}\} \end{cases}$

2637 $\chi_2 \equiv \begin{cases} \text{vars} \equiv \chi.\text{vars} \\ \text{shrd} \equiv \chi.\text{shrd} \cdot \Delta_2 \\ \text{excl} \equiv \{\text{pre} \equiv F_{\text{cut}}, \text{post} \equiv \chi.\text{excl.post}\} \end{cases}$

2638

2639

2640 where $\chi.\text{shrd} \cdot X$ is a shorthand for $\{\text{inv} = (\chi.\text{shrd.inv} \cdot X), \text{reads} = (\chi.\text{shrd.reads} \cdot X)\}$.

2641 Fig. 23. The basic transformation `Loop.fission`.

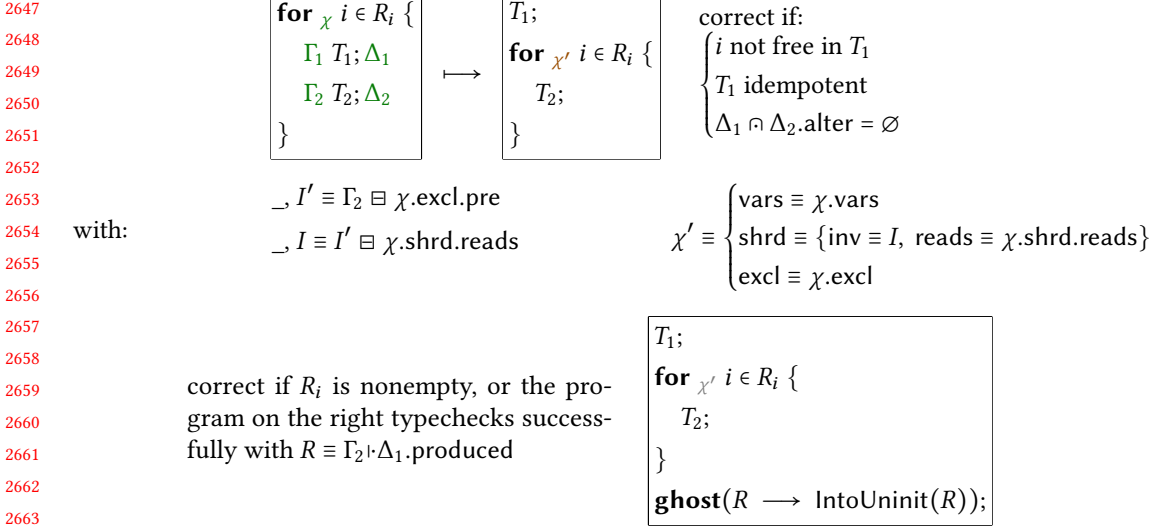
2642

2643

2644

2645

2646

Fig. 24. The basic transformation `Loop.move_out`.

Our criterion asserts that loop fission is correct if the resources altered by T_1 at any iteration i do not interfere with the resources altered by T_2 at any other iteration $i' \neq i$. To implement this check, we inspect the usage maps Δ_1 and Δ_2 associated with T_1 and T_2 , respectively. If T_1 alters one resource from $\chi.\text{shrd}$, then T_2 must not use this same resources; symmetrically, if T_2 alters a resource, then T_1 must not use it. Note, however, that T_1 and T_2 are allowed to both read the same resource; moreover, the resources exclusively consumed or produced by T_1 at the i -th iteration of the first loop may be consumed by T_2 at the i -th iteration of the second loop.

There remains to explain how to synthesize the contracts χ_1 and χ_2 , associated with the two generated loops, from the original contract χ . For `shrd` resources, we project the subsets of $\chi.\text{shrd}$ resources used by T_1 and T_2 . For `excl` resources, we need to synthesize the resources at the cut point, written F_{cut} . The first loop takes the exclusive resources from $\chi.\text{excl.pre}$ to F_{cut} , whereas the second loop takes the exclusive resources from F_{cut} to $\chi.\text{excl.post}$. At a high level, F_{cut} is computed by subtracting the shared resources as well as the local allocations from T_1 , described by $\chi.\text{shrd}$ and $\text{StackAllocCells}(T_1)$, from the typing context Γ computed by our typechecker at the location just between T_1 and T_2 .

Observe that the loop contracts χ_1 and χ_2 generated by the loop fission transformation may contain a larger typing context than strictly necessary. We describe further on, in Section 6.6, a procedure for minimizing loop contracts.

Loop Invariant Code Motion. The basic transformation `Loop.move_out` applies to a loop with body $T_1; T_2$, where T_1 performs instructions that are redundant at every iteration. It produces as output a code that first executes T_1 , exactly once, then executes a loop with body T_2 . The transformation is formalized in Figure 24. We assume for simplicity the loop range to be provably nonempty, or T_1 to be provably deletable. Alternatively, T_1 could be wrapped into a conditional.

The key properties to check are that T_1 is the same for all iterations (it does not depend on i), can be safely deduplicated (it is idempotent as required by `Instr.dedup`), and does not interfere with the remaining instructions of the loop, described by T_2 (that is, the condition $\Delta_1 \cap \Delta_2.\text{alter} = \emptyset$). Note that, contrarily to the `Instr.move` criterion, it is safe for T_2 to read resources modified by T_1 .

Other Loop Transformations. There are other important loop transformations that we support.

- `Loop.fusion` (reciprocal of `Loop.fission`): fuse two consecutive loops into a single one.
- `Loop.collapse` (reciprocal of `Loop.tile`): collapse two nested loops into a single one.
- `Loop.hoist_alloc`: hoist a variable allocated inside a loop into an array allocated outside the loop; more generally, it hoists a matrix of dimension N allocated inside a loop into a matrix of dimension $N + 1$ allocated outside the loop.
- `Loop.shift_range`: reindex a loop by applying a positive or negative offset to its values.
- `Loop.scale_range`: reindex a loop using an index that takes either smaller or larger steps.
- `Loop.extend_range`: extend the range of a loop by wrapping its body in a conditional.
- `Loop.unroll`: unroll a loop whose range is statically known.
- `Loop.parallel`: set (or unset) a parallel flag on a loop using our parallelizable criterion.

6.6 Transformations on Annotations

Modification to Contracts and Ghost Code. The semantics of a program is fully determined by its *proper* OptiC code: it does not depend in any way on the *ghost* code nor on the function and loop contracts. Therefore, contracts may be freely modified, and ghost instructions may be freely inserted, deleted, or modified. The requirement is to reach, after one or several updates, a set of annotations for which the typechecking of the same code succeeds. These modifications can be applied either directly by the programmer, or during the evaluation of transformations.

Minimization of Loop Contracts. The aforementioned loop transformations produce correct resource annotations, yet these annotations might be suboptimal for later transformations. Typically, the generated loop contracts would include clauses covering a set of resources possibly larger than strictly necessary. For example, after the basic loop fission transformation, the contract of the first loop would typically mention resources that are in fact only used by the instructions from the second loop. Mentioning unnecessary resources in a contract may impede the applicability of further transformations. OptiTrust therefore includes a procedure, implemented as a basic transformation, to *minimize* loop contracts. OptiTrust’s combined transformations for loops systematically include a call to this procedure.

The loop contract minimization procedure takes as input a loop with contract χ , and updates this contract to χ' , without modifying the code. The procedure depends on the usage map Δ computed for the instructions T that constitute the loop body.

$$\boxed{\text{for } \overset{\pi}{\chi} i \in R_i \{T; \Delta\}} \mapsto \boxed{\text{for } \overset{\pi}{\chi'} i \in R_i \{T\}}$$

Intuitively, the contract χ' is obtained by filtering out and by weakening resources from χ , depending on their usage in Δ . First, if a resource is unused by T and thus is absent from Δ or has usage `joinedFrac`, then it is excluded from χ' . As a result, certain variables that were quantified in χ might no longer have occurrence in χ' , hence they can be removed as well. Second, if a resource appears with fraction 1 in χ , yet this resource is marked as `splittedFrac` in Δ , then this resource is replaced with a read-only version of it. Technically, an additional fraction variable must be quantified in χ' , and this fraction variable is used for describing the resource as read-only. Internally, the implementation of contract minimization reuses our *minimization of triple* procedure (Section 5.4 and appendix D). Details may be found in appendix F.

Moving and Cancelling Ghost Instructions. OptiTrust includes a transformation that attempts to remove pairs of ghost transformations that cancel each other. Indeed, the sequence `ghost`($H \rightarrow H'$); `ghost`($H' \rightarrow H$) is equivalent to a no-op. More generally, the user as well as combined transformations may request a ghost instruction to be moved so as to be (logically) executed as early as possible in the program; or, symmetrically, to be executed as late as possible. Moving ghost

instructions in such a way may lead to the apparition of cancellable pairs of ghost instructions; and, even when ghost instructions do not disappear, moving them away from, e.g., a loop kernel, may unlock certain transformations.

7 RELATED WORK

The most closely related frameworks were discussed in the introduction. In this section, we comment on the remaining related work, focusing in turn on each of the ingredients that constitute OptiTrust.

General-Purpose Compilers. General purpose compilers such as GCC or ICC are able to apply a large class of program optimizations, from the classic ones such as inlining, dead code elimination, move of instructions to more advanced ones such as loop fission, loop fusion, or loop reordering. The same transformations are available in OptiTrust, yet with three major differences. First, general-purpose compilers apply these transformations on intermediate representations that are not suitable for producing feedback to the user. In contrast, OptiTrust operates on an intermediate representation that has been designed not only to simplify transformations, but also to support its translation back into a conventional program syntax. Second, general-purpose compilers relies on fully-automated procedures, often guided by heuristics, to determine what transformations to apply. In contrast, OptiTrust transformations are fully controlled by the programmer, either directly via basic transformations, or indirectly via combined transformations. Third, general-purpose compilers rely on static analysis applied to plain C code to determine whether certain transformations are applicable, and as a result may lack information to trigger a transformation. In contrast, OptiTrust leverages expressive resource typing information deduced from annotations to justify the correctness of transformations, significantly enlarging the set of applicable transformations.

Guidance in General-Purpose Compilers. To introduce human guidance in general-purpose compilers, a common approach is to insert *pragmas* into the code. For example, Scout [Krzikalla et al. 2011] is a pragma-based tool for guiding source-to-source transformations that introduce vector instructions. As another example, Radtke and Weinzierl [2024] makes use of C++ attributes for switching between array-of-structures and structures-of-arrays over the scope of specific computation kernels; the compiler automatically inserts instructions for copying the data before and after the loop. A similar approach could be expressed as an OptiTrust transformation, by composing the `local_name` transformation for arrays (discussed in Section 6.4) with the `aos_to_soa` transformation (not discussed in this paper). The main limitation of *pragma*-based approaches is that they are ill-suited for describing sequences of optimizations. Indeed, there is no easy way to attach a pragma to a line of code that is generated by a first optimization. Kruse and Finkel [2018] suggest the possibility to stack up pragmas, by providing labels as additional pragma arguments: a second pragma may refer to the labels introduced by the transformation described in a first pragma. Yet, this approach does not scale up well beyond a handful of successive transformations. OptiTrust, in contrast, supports chains of dozens of transformations.

Domain-Specific Compilers. Another possible approach to overcome the limitations of general-purpose compilers is to leverage *domain specific languages* (DSL), such as Halide [Ragan-Kelley et al. 2013], TVM [Chen et al. 2018], Fireiron [Hagedorn et al. 2020a] (used at Nvidia), PartIR [Alabed et al. 2024] (used at DeepMind). Specialized compilers can benefit from carefully tuned heuristics. Yet, even for programs expressed in a specific DSL, the optimization search space remains vast, hence programmer guidance is key to achieve good performance. Halide and its descendants makes use of a script, called a *schedule*, for guiding the compilation strategy.

For DSLs, the language restriction is also their Achilles' heel: as soon as the user's application requires a single feature that falls outside of what the DSL can express, the programmer loses

2794 most if not all of the benefits of the DSL. In practice, DSLs typically support the possibility to
2795 include foreign functions (or, inlined general-purpose code), however these foreign functions must
2796 be treated as black box by the DSL compiler, preventing the applications of any domain-specific
2797 optimization across this black box.

2798 In contrast to DSLs, OptiTrust sticks to a standard, general-purpose language. At the same time,
2799 OptiTrust retains the ability to manipulate domain-specific operations and to exploit transforma-
2800 tions that are specific to these operations, as illustrated with the *reduce* function in our OpenCV
2801 case study. At any point in the transformation script, an occurrence of a domain-specific operation
2802 may be lowered into standard C code, thereby enabling further lower-level optimizations.

2803
2804 *Code Transformations via Rewrite Rules.* A rewrite rule maps a code pattern to another code
2805 pattern. A number of tools exploit rewrite rules to perform source-to-source transformations. For
2806 example, TXL [Cordy 2006] is a multi-language rewrite system, whose patterns are expressed at
2807 the level of syntax, using grammars. Coccinelle [Lawall and Muller 2018] allows the programmer
2808 to describe *semantic patches* in C code. CodeBoost [Bagge et al. 2003] applies the Stratego program
2809 transformation language [Bravenboer et al. 2008] to C++ code. CodeBoost can be used to turn
2810 high-level operations on matrices and vectors into typical high-performance source code.

2811 OptiTrust relies on OCaml to provide a very expressive language for describing transformations,
2812 going beyond rewrite rules. Although many transformations *can* be encoded as rewrite rules,
2813 the encoding involved can be cumbersome or inefficient. For example, reconstructing a for-loop
2814 for a series of similar blocks of instructions can be encoded via rewrite rules, yet the blocks
2815 must be merged into the for-loop one by one. Other transformations, especially those involving
2816 contracts, would be challenging to express as rewrite rules. For example, *loop contract minimization*
2817 (Section 6.6) would require the rewrite rule to depend on side-conditions and meta-operations that
2818 involve resources and usage maps.

2819
2820 *Intermediate Languages.* The use of an intermediate language with simpler semantics is com-
2821 monplace, both in the domain of compilation and in the domain of program verification. Let us cite
2822 a few examples. The Common Intermediate Language (CIL) serves as intermediate compilation
2823 language for the whole .NET ecosystem [Gough and Gough 2001]. Why3 [Filliâtre and Paskevich
2824 2013] serves as intermediate verification language for C, Java, and Ada programs. Viper [Müller
2825 et al. 2017] serves as intermediate verification language for Java, Rust, Go, OpenCL, etc. Although
2826 intermediate languages are commonplace, we are not aware of any framework that leverages a
2827 translation into an intermediate language *and* provides a reciprocal translation back to the source
2828 language, with a round-trip property such as that provided by OptiTrust.

2829
2830 *Source Code Manipulation Frameworks.* Frameworks that offer more expressiveness than rewrite
2831 rules generally give access to the abstract syntax tree (AST) of the source code. Traditional compilers
2832 employ an AST, but they are not designed for synthesizing pieces of AST at the source level.
2833 Moreover, traditional compilers operate on intermediate representations, and lose the structure
2834 of the original code. These two limitations of general-purpose compilers have motivated the
2835 development of frameworks that are specifically designed to support code transformations (and
2836 code analyses) at the level of C code. ROSE [Quinlan 2000; Quinlan and Liao 2011] and Cetus [Bae
2837 et al. 2013; Dave et al. 2009] are two such frameworks that provide facilities for manipulating C ASTs.
2838 Source-to-source transformation frameworks have also been employed to produce code targeting
2839 GPUs [Amini 2012; Konstantinidis 2013; Lebras 2019]. These frameworks implement generic
2840 optimization strategies, in a similar fashion as general-purpose compilers. In contrast, OptiTrust
2841 leverages transformation scripts to guide the optimization of a specific program. Moreover, the
2842

2843 OptiTrust infrastructure supports resource typing, which provides much more precise information
2844 than the classic static code analyses implemented in the frameworks such as ROSE and Cetus.

2845 *Transformation Scripts.* Expressing a series of source-level transformations for a specific program
2846 can be done by means of a transformation script. Such scripts have appeared in particular in the
2847 context of polyhedral transformations [Bagnères et al. 2016b; Bondhugula et al. 2008], for example
2848 in Loopy [Namjoshi and Singhanian 2016] and in work by Zinenko et al. [2018a]. CHiLL [Chen et al.
2849 2008; Rudy et al. 2011] includes transformations that go beyond the polyhedral model. It has been
2850 applied to generate finely tuned CUDA code from high-level linear algebra kernels. POET [Yi and
2851 Qasem 2008; Yi et al. 2014] is a scripting language for performing program transformations, for
2852 C/C++ as well as other languages. POET has been employed to generate optimized code for linear
2853 algebra kernels, including semi-automated exploration of a search space of possible optimizations.

2854 Several pieces of work already discussed in the introduction exploit transformation scripts.
2855 Halide [Ragan-Kelley et al. 2013], TVM [Chen et al. 2018] feature schedules that can be viewed as
2856 transformation scripts. Elevate [Hagedorn et al. 2020b] expresses the transformation script in the
2857 form of a composition of functions. ATL [Liu et al. 2022] leverages “tactic”-based proof scripts as
2858 support for expressing transformations scripts. LARA consists of a transformation script featuring
2859 declarative queries as well as arbitrary JavaScript instructions.

2860 MLIR (Multi-Level Intermediate Representation) [Lattner et al. 2021] is a framework for building
2861 reusable and extensible compiler infrastructure. MLIR aims in particular at improving compilation
2862 for heterogeneous hardware, and at improving support for DSL constructs. To that end, MLIR
2863 provides *dialects*, which enables expressing extended language constructs. For example, the *tensor*
2864 *dialect* helps representing multidimensional arrays and operations on them. Recently, a *transform*
2865 *dialect* [Lücke et al. 2024] was added to MLIR to express transformation scripts. This extension
2866 confirms the interest for finer-grained control, going beyond the simple ordering of global opti-
2867 mizations passes. A major limitation of MLIR is that its dialects and passes do not share a common
2868 specification language that could be used to exploit loop invariants and summaries of function
2869 effects across different analyses. We believe that Separation Logic, as implemented in OptiTrust,
2870 could offer such a common language for expressing invariants. We leave it to future work to explore
2871 how the OptiTrust AST and transformations could be extended to support user-defined language
2872 constructs.

2873 All this related work demonstrates a strong interest in leveraging transformation scripts for
2874 putting control of optimizations in the hand of the programmer. Systems differ in what language
2875 they target, and what transformations they support. None of the aforementioned systems support
2876 in their transformation scripts a system for targeting program points with the expressiveness and
2877 conciseness offered by OptiTrust targets. Moreover, as far as we know, LARA [Silvano et al. 2019]
2878 and OptiTrust are the only two frameworks making use of transformation scripts for applying
2879 general-purpose transformations at the level of C syntax. OptiTrust is the first to demonstrate the use
2880 of transformation scripts to produce high-performance code for state-of-the-art benchmarks. Most
2881 importantly, unlike LARA, OptiTrust checks that the transformations requested by the programmer
2882 preserve the semantics of the code.

2883 *Proof-Transforming Compilation.* The notion of *Proof Carrying Code* [Necula 1998] refers to the
2884 idea that compilers could be instrumented to carry invariants from high-level source code down to
2885 low-level code. The original line of work on Proof Carrying Code did not aim at full functional
2886 correctness properties, but rather focused on simpler invariants capturing safety properties, such
2887 as the absence of out-of-bound accesses.

2888 Subsequent work introduced the notion of *Proof-Transforming Compilation* to refer to a compiler
2889 that takes as input a formally-verified program and produces as output compiled code accompanied
2890

2891

2892 by a formal proof (i.e., a proof tree in a program logic) that the compiled code satisfies the same
2893 functional correctness specification as the input program. In particular, the PhD work of César
2894 Kunz [Barthe et al. 2009; Kunz 2009] shows how to realize proof-transforming compilation for
2895 standard compiler optimizations, applied at the level of the RTL intermediate language. More
2896 recently, the work on Alpinist [Sakar et al. 2022] demonstrates the feasibility, for a small number
2897 of GPU-oriented optimizations, of transforming GPU code while preserving logical invariants.

2898 Our results obtained so far with OptiTrust demonstrates the feasibility, for a fair number of
2899 general-purpose code optimizations, of transforming C code while preserving resource-based
2900 invariants. In future work, we look forward to extending OptiTrust in order to handle richer logical
2901 invariants and to produce optimized programs accompanied by formal proofs of correctness.

2902
2903 *Separation Logic.* OptiTrust leverages a standard Separation Logic. The most closely related
2904 program logics are VST [Cao et al. 2018], a program verification tool for C, and RefinedC [Sammler
2905 et al. 2021], a very expressive type system for C. Both these systems are grounded on the Iris
2906 framework [Jung et al. 2018a,b], at this day the most advanced formalization of Concurrent Sep-
2907 aration Logic. Other tools, such as Alpinist [Sakar et al. 2022] leverage Viper’s *dynamic frames*
2908 technique [Müller et al. 2017], a cousin of Separation Logic. Fractional resources [Boyland 2003]
2909 are a standard ingredient of Separation Logic [Jung et al. 2018a]. Following common practice,
2910 OptiTrust leverages fractional resources to describe read-only resources. The technique of making
2911 fractions essentially transparent to the end-user is directly inspired by the work by Heule et al.
2912 [2013] implemented in the Chalice verification tool.

2913 OptiTrust is, as far as we know, the first transformation framework based on separation logic
2914 to compute and leverage *usage information*. This information describes how the Separation Logic
2915 resources available for typechecking of a subterm are *actually* exploited for typechecking this
2916 subterm.

2917
2918 *Contract Inference.* OptiTrust currently requires the programmer to annotate the input program
2919 with ghost operations as well as function and loop contracts. One may wonder the extent to which
2920 such contracts could be automatically inferred, at least for reasonably simple programs.

2921 The experience from other practical Separation Logic frameworks (e.g., [Müller et al. 2017]) is
2922 that heuristics can be devised to significantly reduce the number of ghost operations that need to
2923 be explicitly provided by the programmer. For example, if we have at hand no other permissions
2924 on an array than a permission covering a range of its cells, then when facing a read operations on
2925 a particular cell from this array, isolating this cell from the range at hand is the only way in which
2926 typechecking could succeed.

2927 Inference is not limited to ghost operations: certain contracts may also be automatically inferred.
2928 For example, Journault and Miné [2018] show that, by leveraging abstract interpretation, for func-
2929 tions such as matrix-multiplication or similar linear algebra operations, full functional correctness
2930 specifications can be automatically computed. Besides, *bi-abduction* [Calcagno et al. 2019; Spies
2931 et al. 2024] is a technique for inferring function contracts, at the heart of the *infer* automated
2932 program analysis tool [Calcagno et al. 2019].

2933 We leave it to future work to integrate techniques for inferring ghost operations and contracts,
2934 for decreasing the amount of user annotations required.

2935

2936 8 CONCLUSION

2937 In this paper, we have presented OptiTrust, the first modular tool for programmer-guided opti-
2938 mization that demonstrates both a high degree of control and a high degree of generality. We have
2939 demonstrated the benefits of OptiTrust on 3 realistic case studies, comparing against manually
2940

2941 optimized code from image processing and from scientific computing applications, and comparing
2942 against the state-of-the-art specialized compiler TVM.

2943 OptiTrust leverages in a crucial way Separation Logic. As we have shown, for numerous opti-
2944 mizations, shape-based assertions are sufficient. To support more ambitious transformations, we
2945 plan to extend OptiTrust’s specification language to a full-blown Separation Logic. We would also
2946 like to develop mechanized proofs of the metatheory of OptiTrust, by first formalizing the type
2947 system, then formalizing the correctness criteria of the basic transformations.

2948 Besides mechanized proofs, there are numerous directions for future work on OptiTrust. Let us
2949 mention a few. First, we will work on improving the user experience, in particular by making the
2950 typechecker incremental for improved performance, and by augmenting the amount of inference
2951 for contracts and ghost operations. Second, we plan to make the language extensible with DSL
2952 constructs (like MLIR), to complete our library of transformations, and to provide support for
2953 reasoning about numerical accuracy. Third, we would like OptiTrust to support a diversity of
2954 hardware targets, including exotic accelerators. Finally, we would like to integrate in OptiTrust
2955 tools for reporting performance feedback from benchmarks, tools for autotuning parameters, and
2956 tools for providing suggestions for the possible next step in a transformation script.

2957
2958
2959
2960
2961
2962
2963
2964
2965
2966
2967
2968
2969
2970
2971
2972
2973
2974
2975
2976
2977
2978
2979
2980
2981
2982
2983
2984
2985
2986
2987
2988
2989

2990 A SEMANTICS

2991 As said in Section 4.8, we formalize the semantics of $\text{Opti}\lambda$ using an omni-big-step evaluation
 2992 judgment in call by value style. The judgment $t/(s, m) \Downarrow Q$ asserts that the term t , in a program
 2993 stack s and in a program store m , evaluates to result states that belong to the set Q . The result
 2994 states in Q are of the form (s', m') where s' is a program stack and m' a program store. Program
 2995 stacks maps program variables to values, and program stores maps locations to values. The values,
 2996 denoted v , can be logical expressions, locations, function closures of the form $\mathbf{fun}^s(x_1, \dots, x_n) \mapsto t$,
 2997 and the special uninitialized value \perp .

2998 Fig. 25 gives the semantic rules of $\text{Opti}\lambda$. The evaluation contexts consist of function arguments
 2999 and ranges of for-loops.

3000 This semantics rules are standard except maybe for the rule SEQ to handle sequences with optional
 3001 result value. The rule SEQ encode the fact that a sequence creates a lexical scope by restoring the
 3002 program stack after its execution. The result value (if there is one) is bound in the output stacks.

3003 By design, like all omni-big-step judgments, the judgment $t/(s, m) \Downarrow Q$ is preserved when
 3004 enlarging Q . This property named consequence will be used in the proof of the frame rule.

3005 **THEOREM A.1 (CONSEQUENCE PROPERTY FOR OMNISEMANTICS).**

$$3006 \quad t/(s, m) \Downarrow Q \wedge Q \subseteq Q' \implies t/(s, m) \Downarrow Q'$$

3007 We refer to the omnisemantics paper [Charguéraud et al. 2022] for the inductive proof pattern.

3010 B ASSERTION AND CONTEXT SATISFACTION

3011 In Section 4.8, we introduced the judgment $(\sigma, \mu) \in \Gamma$ to assert that a logical state (σ, μ) satisfies a
 3012 context Γ of the form $\langle E \mid F \rangle$. This section formally defines this judgment. Doing so involves two
 3013 auxiliary judgments $\sigma : E$ and $F \models \mu$ that we define below.

3014 First, $\sigma : E$ is a characterization of the fact that bindings in σ have types that correspond to the
 3015 bindings in E . Recall that the operator *Specialize* described in Section 4.2 checks that each binding
 3016 $x : v$ in σ corresponds to a binding $x : T$ in E such that v is of type T . The operator *Specialize* returns
 3017 the subset of E that is not instantiated by σ . Here we enforce that this subset is empty.

3018 *Definition B.1.*

$$3019 \quad \sigma : E \quad := \quad \text{Specialize}_{\emptyset}\{\sigma\}([E]) = \emptyset$$

3020 $F \models \mu$ is a characterization of the fact that memory cells described by μ correspond to the linear
 3021 resources described in F . Before giving its formal definitions, we need to introduce additional
 3022 operators on logical stores. These definitions are essentially standard in Separation Logic.

3023 We denote by $\mu_1 \uplus \mu_2$ the *compatible* union between two logical store. We denote by $\mu_1 \# \mu_2$ the
 3024 fact that two logical stores are compatible. Two logical stores are compatible if and only if, on their
 3025 intersection, all the bindings have the same value, and the sum of the fractions does not exceed one.

3026 *Definition B.2.*

$$3027 \quad \mu_1 \# \mu_2 \quad := \quad \forall l \in \text{dom}(\mu_1) \cap \text{dom}(\mu_2), \exists \alpha_1, \exists \alpha_2, \exists v, \\ 3028 \quad \mu_1(l) = (\alpha_1, v) \wedge \mu_2(l) = (\alpha_2, v) \wedge \alpha_1 + \alpha_2 \leq 1$$

3029 When two logical stores are compatible, their compatible union is defined as follows:

3030 *Definition B.3.* Assume $\mu_1 \# \mu_2$. Then:

$$3031 \quad \mu_1 \uplus \mu_2 \quad := \quad \left\{ l \mapsto (\alpha, v) \left| \begin{array}{l} \mu_1(l) = (\alpha, v) \quad \wedge \quad l \notin \text{dom}(\mu_2) \\ \vee \quad \mu_2(l) = (\alpha, v) \quad \wedge \quad l \notin \text{dom}(\mu_1) \\ \vee \quad \mu_1(l) = (\alpha_1, v) \in \mu_1 \quad \wedge \quad \mu_2(l) = (\alpha_2, v) \quad \wedge \quad \alpha = \alpha_1 + \alpha_2 \end{array} \right. \right\}$$

$$\begin{array}{c}
3039 \\
3040 \\
3041 \\
3042 \\
3043 \\
3044 \\
3045 \\
3046 \\
3047 \\
3048 \\
3049 \\
3050 \\
3051 \\
3052 \\
3053 \\
3054 \\
3055 \\
3056 \\
3057 \\
3058 \\
3059 \\
3060 \\
3061 \\
3062 \\
3063 \\
3064 \\
3065 \\
3066 \\
3067 \\
3068 \\
3069 \\
3070 \\
3071 \\
3072 \\
3073 \\
3074 \\
3075 \\
3076 \\
3077 \\
3078 \\
3079 \\
3080 \\
3081 \\
3082 \\
3083 \\
3084 \\
3085 \\
3086 \\
3087
\end{array}$$

$$\begin{array}{c}
\frac{}{v/(s, m) \Downarrow \{(s[\mathbf{res} \mapsto v], m)\}} \text{VAL} \qquad \frac{}{x/(s, m) \Downarrow \{(s[\mathbf{res} \mapsto s(x)], m)\}} \text{VAR} \\
\frac{}{(\mathbf{fun}(a_1, \dots, a_n) \mapsto t_f)/(s, m) \Downarrow (s[\mathbf{res} \mapsto (\mathbf{fun}^s(a_1, \dots, a_n) \mapsto t_f)], m)} \text{FUN} \\
\frac{}{(\mathbf{let } x = \mathbf{stackAlloc}())/(s, m) \Downarrow \{(s[x \mapsto l], m[l \mapsto \perp]) \mid l \notin \text{dom}(m)\}} \text{STACKALLOC} \\
\frac{t/(s, m) \Downarrow Q}{(\mathbf{let } x = t)/(s, m) \Downarrow \{(s[x \mapsto s'(\mathbf{res})], m') \mid (s', m') \in Q\}} \text{LET} \\
\frac{t/(s, m) \Downarrow Q' \quad \forall (s', m') \in Q', \mathcal{E}[s'(\mathbf{res})]/(s, m') \Downarrow Q \quad \mathcal{E} \text{ is an evaluation context}}{\mathcal{E}[t]/(s, m) \Downarrow Q} \text{BIND} \\
\frac{s_c = s_f[\overline{a_i} \mapsto \overline{v_i}] \quad t_f/(s_c, m) \Downarrow Q}{(\mathbf{fun}^{s_f}(a_1, \dots, a_n) \mapsto t_f)(v_1, \dots, v_n)/(s, m) \Downarrow Q} \text{CALL} \\
\frac{\forall (s_c, m_c) \in Q_c, (s_c(\mathbf{res}) = \mathbf{true} \implies t_t/(s, m_c) \Downarrow Q) \wedge (s_c(\mathbf{res}) = \mathbf{false} \implies t_f/(s, m_c) \Downarrow Q)}{(\mathbf{if } t_c \mathbf{ then } t_t \mathbf{ else } t_f)/(s, m) \Downarrow Q} \text{IF} \\
\frac{Q_0 = \{(s_0, m_0)\} \quad \forall i \in [1, n], \forall (s, m) \in Q_{i-1}, t_i/(s, m) \Downarrow Q_i}{Q_A = \{(s, m \setminus A(s)) \mid (s, m) \in Q_n\} \text{ where } A(s) = \{s(x_i) \mid t_i \text{ is of the form } \mathbf{let } x_i = \mathbf{stackAlloc}()\}} \\
Q = \begin{cases} \{(s_0, m) \mid (s, m) \in Q_A\} & \text{if } r = \emptyset \\ \{(s_0[\mathbf{res} \mapsto s(x)], m) \mid (s, m) \in Q_A\} & \text{if } r = x \end{cases} \\
\frac{}{\{t_1; \dots; t_n; r\}/(s_0, m_0) \Downarrow Q} \text{SEQ} \\
\frac{m(l) \neq \perp}{\mathbf{get}(l)/(s, m) \Downarrow \{(s[\mathbf{res} \mapsto m(l)], m)\}} \text{GET} \qquad \frac{l \in \text{dom}(m)}{\mathbf{set}(l, v)/(s, m) \Downarrow \{(s, m[l \mapsto v])\}} \text{SET} \\
\frac{}{\mathbf{ignore}(v)/(s, m) \Downarrow \{(s \setminus \mathbf{res}, m)\}} \text{IGNORE} \qquad \frac{}{\mathbf{add}(v_1, v_2)/(s, m) \Downarrow \{(s[\mathbf{res} \mapsto v_1 + v_2], m)\}} \text{ADD} \\
\frac{m(l_1) \neq \perp}{\mathbf{inplaceAdd}(l_1, v_2)/(s, m) \Downarrow \{(s, m[l_1 \mapsto m(l_1) + v_2])\}} \text{INPLACEADD} \\
\frac{}{\mathbf{heapAlloc}()/s, m) \Downarrow \{(s[\mathbf{res} \mapsto l], m[l \mapsto \perp]) \mid l \notin \text{dom}(m)\}} \text{HEAPALLOC} \\
\frac{l \in \text{dom}(m)}{\mathbf{free}(l)/(s, m) \Downarrow \{(s, m \setminus l)\}} \text{FREE} \\
\frac{n_{\text{start}} \leq n_{\text{stop}} \quad t/(s[i \mapsto n_{\text{start}}], m) \Downarrow Q_1}{\forall (s_1, m_1) \in Q_1, (\mathbf{for } (i \in \mathbf{range}(n_{\text{start}} + n_{\text{step}}, n_{\text{stop}}, n_{\text{step}})) t)/(s_1, m_1) \Downarrow Q} \text{FORITER} \\
\frac{}{(\mathbf{for } (i \in \mathbf{range}(n_{\text{start}}, n_{\text{stop}}, n_{\text{step}})) t)/(s, m) \Downarrow Q} \\
\frac{n_{\text{start}} > n_{\text{stop}}}{(\mathbf{for } (i \in \mathbf{range}(n_{\text{start}}, n_{\text{stop}}, n_{\text{step}})) t)/(s, m) \Downarrow \{(s, m)\}} \text{FOREND}
\end{array}$$

Fig. 25. Semantics of the Opti λ internal language in omni-big-step style as explained in Section 4.8. Other arithmetic built-in functions follow the pattern of ADD or INPLACEADD.

In program and logical stores, we allow a special value \perp for variables that are not initialized. We define the fact that a linear resource H models a logical store μ recursively as follows:

Definition B.4. We define $H \models \mu$ with the following rules:

$$\begin{array}{llll}
l \rightsquigarrow v & \models & \{l \mapsto (1, v)\} & \\
l \rightsquigarrow \text{Cell} & \models & \{l \mapsto (1, v)\} & \text{when } v \neq \perp \\
\star_{i \in r} H_i & \models & \uplus_{i \in r} \mu_i & \text{when } \forall i \in r, H_i \models \mu_i \\
\alpha H & \models & \uplus_{l \mapsto (\beta, v) \in \mu} \{l \mapsto (\alpha \cdot \beta, v)\} & \text{when } H \models \mu \\
\text{Uninit}(H) & \models & \uplus_{l \mapsto (\alpha, v) \in \mu} \{l \mapsto (\alpha, v')\} & \text{when } H \models \mu \\
H_1 \multimap H_2 & \models & \mu & \text{when } \forall \mu_1, \mu_1 \# \mu \wedge H_1 \models \mu_1 \\
& & & \implies H_2 \models \mu_1 \uplus \mu
\end{array}$$

Above, all the occurrences of the operator \uplus must be well-defined.

We say that a linear context F models a logical store μ and write $F \models \mu$ if and only if the disjoint union of all resources in F models μ . Formally:

Definition B.5. Consider F of the form H_1, \dots, H_n . The predicate $F \models \mu$ holds if and only if $(\star_{i \in [1, n]} H_i) \models \mu$.

With the two relations $\sigma : E$ and $F \models \mu$ defined above, we define $(\sigma, \mu) \in \Gamma$ in the following way:

Definition B.6 (Context satisfaction).

$$(\sigma, \mu) \in \langle E \mid F \rangle \quad := \quad \sigma : E \quad \wedge \quad \sigma(F) \models \mu$$

C PROOF OF THE FRAME RULE

This section gives a proof of the frame rule for logical triples. This proof is divided in two steps. First, we need to prove correct the frame property with respect to the semantics of our language for omni-big-step evaluation judgments. Then, we can use this property along with the other omnisemantics properties described in the previous section to show that the frame rule for logical triples holds.

Before formally stating the frame property for omni-big-step evaluation judgment, we need one technical definition to take the compatible union of two sets of program states. Two program stacks (respectively program stores) are compatible, and we write $s \# s'$ (resp. $m \# m'$), if their domain is disjoint. In that case, we write $s \uplus s'$ (resp. $m \uplus m'$) their disjoint union. We can define the compatible union of two set of program states as follows:

Definition C.1.

$$Q \star Q' \quad := \quad \{(s \uplus s', m \uplus m') \mid (s, m) \in Q \wedge (s', m') \in Q' \wedge s \# s' \wedge m \# m'\}$$

Then, the frame property for omni-big-step evaluation judgments reads as follows:

THEOREM C.2 (FRAME PROPERTY FOR OMNISEMANTICS).

$$t/(s, m) \Downarrow Q \implies \forall s' \# s, \forall m' \# m, t/(s \uplus s', m \uplus m') \Downarrow (Q \star \{(s', m')\})$$

The proof sketch of this property is given in the omnisemantics paper [Charguéraud et al. 2023, §5.4].

Before expressing the frame property for logical triples, we need one last technical definition to characterize typing contexts that are well-typed. Recall that since E is a telescope, bindings defined in E can be used in the following bindings of the typing context.

Definition C.3 (Well-typed contexts). A typing context $\Gamma = \langle E \mid F \rangle$ is *well-typed* iff there is no name conflict in E or in F and for any $x : \tau$ in E , τ is of type `Type` and for any $y : H$ in F , H is of type `Hprop`.

We can now state and prove the frame property for logical triples:

THEOREM C.4 (FRAME PROPERTY FOR LOGICAL TRIPLES).

$$\{\Gamma\} t \{\Gamma'\} \wedge \Gamma \star \Gamma'' \text{ is well-typed} \wedge \Gamma' \star \Gamma'' \text{ is well-typed} \implies \{\Gamma \star \Gamma''\} t \{\Gamma' \star \Gamma''\}$$

PROOF.

- Suppose we have $\{\Gamma\} t \{\Gamma'\}$. Let $(\sigma_0, \mu_0) \in \Gamma \star \Gamma''$. By definition of triples, we have to prove $t/(\sigma_0, \mu_0)|_{\text{prog}} \Downarrow \text{AcceptableStates}(\sigma_0, \mu_0, \Gamma' \star \Gamma'')$.
- There is a decomposition $\sigma_0 = \sigma \uplus \sigma''$ and $\mu_0 = \mu \uplus \mu''$ such that $(\sigma, \mu) \in \Gamma$ and $(\sigma'', \mu'') \in \sigma(\Gamma'')$.
- By definition of $\{\Gamma\} t \{\Gamma'\}$ applied to $(\sigma, \mu) \in \Gamma$, we have $t/(\sigma, \mu)|_{\text{prog}} \Downarrow \text{AcceptableStates}(\sigma, \mu, \Gamma')$.
- We have $(\sigma \uplus \sigma'')|_{\text{prog}} = \sigma|_{\text{prog}} \cup \sigma''|_{\text{prog}}$ and $\sigma''|_{\text{prog}} \perp \sigma|_{\text{prog}}$.
- We pose $m = \mu|_{\text{prog}}$, and $m'' = \mu''|_{\text{prog}} \setminus \text{dom}(m)$. Therefore $m'' \perp m$. By well-definedness of $\mu \uplus \mu''$, we also have $\mu_0|_{\text{prog}} = (\mu \uplus \mu'')|_{\text{prog}} = m \cup m''$.
- By the frame property for Omni-big-step applied on $t/(\sigma, \mu)|_{\text{prog}} \Downarrow \text{AcceptableStates}(\sigma, \mu, \Gamma')$, $\sigma''|_{\text{prog}} \perp \sigma|_{\text{prog}}$ and $m'' \perp m$, we obtain $t/(\sigma|_{\text{prog}} \cup \sigma''|_{\text{prog}}, m \cup m'') \Downarrow (\text{AcceptableStates}(\sigma, \mu, \Gamma') \star \{(\sigma''|_{\text{prog}}, m'')\})$.
- Since $(\sigma|_{\text{prog}} \cup \sigma''|_{\text{prog}}, m \cup m'') = (\sigma_0, \mu_0)|_{\text{prog}}$, by the consequence property of Omni-big-step, it suffices to show $\text{AcceptableStates}(\sigma, \mu, \Gamma') \star \{(\sigma''|_{\text{prog}}, m'')\} \subseteq \text{AcceptableStates}(\sigma_0, \mu_0, \Gamma' \star \Gamma'')$.
- Take $(s_r, m_r) \in \text{AcceptableStates}(\sigma, \mu, \Gamma') \star \{(\sigma''|_{\text{prog}}, m'')\}$. We need to show that $(s_r, m_r) \in \text{AcceptableStates}(\sigma_0, \mu_0, \Gamma' \star \Gamma'')$.
- There is a decomposition $s_r = s'_r \uplus s''_r$ and $m_r = m'_r \uplus m''_r$ such that $(s'_r, m'_r) \in \text{AcceptableStates}(\sigma, \mu, \Gamma')$ and $(s''_r, m''_r) \in \{(\sigma''|_{\text{prog}}, m'')\}$. Since $\{(\sigma''|_{\text{prog}}, m'')\}$ is a singleton, we have $s''_r = \sigma''|_{\text{prog}}$ and $m''_r = m''$.
- By definition of AcceptableStates , there is $(\sigma', \mu') \in \Gamma'$ such that $(s'_r, m'_r) = (\sigma', \mu')|_{\text{prog}}$ and $\forall x \in \text{dom}(\sigma) \cap \text{dom}(\sigma')$, $\sigma(x) = \sigma'(x)$ and $\text{OnlyRO}(\mu) = \text{OnlyRO}(\mu')$.
- $(s'_r \uplus \sigma''|_{\text{prog}}) = (\sigma'|_{\text{prog}} \uplus \sigma''|_{\text{prog}}) = (\sigma' \uplus \sigma'')|_{\text{prog}}$
- We know that $\Gamma' \star \Gamma''$ is well-scoped and that $\sigma' : \Gamma'$.pure. Therefore, for any x free in Γ'' , $x \in \text{dom}(\Gamma') = \text{dom}(\sigma')$. Similarly, we know that for any x free in Γ'' , $x \in \text{dom}(\Gamma) = \text{dom}(\sigma)$. Therefore, since $\forall x \in \text{dom}(\sigma) \cap \text{dom}(\sigma')$, $\sigma(x) = \sigma'(x)$, for any x free in Γ'' , $\sigma(x) = \sigma'(x)$. This implies $\sigma(\Gamma'') = \sigma'(\Gamma'')$.
- We have $(\sigma'', \mu'') \in \sigma(\Gamma'')$ and $\sigma(\Gamma'') = \sigma'(\Gamma'')$. Therefore, $(\sigma'', \mu'') \in \sigma'(\Gamma'')$ and thus $(\sigma' \uplus \sigma'', \mu' \uplus \mu'') \in \Gamma' \star \Gamma''$.
- $\forall x \in \text{dom}(\sigma \uplus \sigma'') \cap \text{dom}(\sigma' \uplus \sigma'')$, $(\sigma \uplus \sigma'')(x) = (\sigma' \uplus \sigma'')(x)$ directly follows from $\forall x \in \text{dom}(\sigma) \cap \text{dom}(\sigma')$, $\sigma(x) = \sigma'(x)$.
- Let us show that $\mu' \uplus \mu''$ is well-defined. Take $l \in \text{dom}(\mu') \cap \text{dom}(\mu'')$. We need to find α, α'', v such that $\mu'(l) = (\alpha, v)$ and $\mu''(l) = (\alpha'', v)$ and $\alpha + \alpha'' \leq 1$. Let us consider two cases: either $l \in \text{dom}(\mu)$ or $l \notin \text{dom}(\mu)$.
 - If $l \notin \text{dom}(\mu)$, $l \in \text{dom}(m'')$. Since $m'_r \uplus m''$ is well-defined, $l \notin \text{dom}(m'_r)$. Since $\text{dom}(\mu') = \text{dom}(\mu'|_{\text{prog}}) = \text{dom}(m'_r)$, we conclude $l \notin \text{dom}(\mu')$. This is absurd therefore this case is not possible.

- 3186 – If $l \in \text{dom}(\mu)$, we have $l \in \text{dom}(\mu) \cap \text{dom}(\mu'')$, by well-definedness of $\mu \uplus \mu''$ we have
3187 v, α and α'' such that $\mu(l) = (\alpha, v)$ and $\mu''(l) = (\alpha'', v)$ and $\alpha + \alpha'' \leq 1$. In particular,
3188 we have $\alpha < 1$ and therefore $l \in \text{dom}(\text{OnlyRO}(\mu))$. Since $\text{OnlyRO}(\mu) = \text{OnlyRO}(\mu')$,
3189 $\mu'(l) = (\alpha, v)$. This is the last fact needed to conclude.
- 3190 • Now we want to show that $(\mu' \uplus \mu'')|_{\text{prog}} \subseteq (m'_r \uplus m'')$. Take $l \in \text{dom}((\mu' \uplus \mu'')|_{\text{prog}}) =$
3191 $\text{dom}(\mu' \uplus \mu'')$. We need to show that $((\mu' \uplus \mu'')|_{\text{prog}})(l) = (m'_r \uplus m'')(l)$. Consider two
3192 cases: either $l \in \text{dom}(\mu')$ or $l \notin \text{dom}(\mu')$.
- 3193 – If $l \in \text{dom}(\mu')$, then there is α' and v such that $\mu'(l) = (\alpha', v)$. Therefore, $\mu'|_{\text{prog}}(l) = v$
3194 and thus since $m'_r = \mu'|_{\text{prog}}$, $m'_r(l) = v$. By definition of $\mu' \uplus \mu''$, there is α such that
3195 $(\mu' \uplus \mu'')(l) = (\alpha, v)$. Therefore, $((\mu' \uplus \mu'')|_{\text{prog}})(l) = v$. By definition of $m'_r \uplus m''$,
3196 $(m'_r \uplus m'')(l) = v$.
- 3197 – If $l \notin \text{dom}(\mu')$, then since $l \in \text{dom}(\mu' \uplus \mu'')$ we have $l \in \text{dom}(\mu'')$. Consider two cases:
3198 either $l \in \text{dom}(\mu)$ or $l \notin \text{dom}(\mu)$.
- 3199 * If $l \in \text{dom}(\mu)$, we have $l \in \text{dom}(\mu) \cap \text{dom}(\mu'')$, by well-definedness of $\mu \uplus \mu''$ we
3200 have v, α and α'' such that $\mu(l) = (\alpha, v)$ and $\mu''(l) = (\alpha'', v)$ and $\alpha + \alpha'' \leq 1$.
3201 In particular, we have $\alpha < 1$ and therefore $l \in \text{dom}(\text{OnlyRO}(\mu))$. Thus, since
3202 $\text{OnlyRO}(\mu) = \text{OnlyRO}(\mu')$, we get $l \in \text{dom}(\mu')$ which is a contradiction.
- 3203 * If $l \notin \text{dom}(\mu)$, then $l \in \text{dom}(m'')$. There is α'' and v such that $\mu''(l) = (\alpha'',$
3204 $v)$. Since $m'' = \mu''|_{\text{prog}} \setminus \text{dom}(m)$, we have $m''(l) = v$. By definition of $\mu' \uplus \mu''$,
3205 there is α such that $(\mu' \uplus \mu'')(l) = (\alpha, v)$. Therefore, $((\mu' \uplus \mu'')|_{\text{prog}})(l) = v$. By
3206 definition of $m'_r \uplus m''$, $(m'_r \uplus m'')(l) = v$.
- 3207
- 3208 • We want to show that $(\mu' \uplus \mu'')|_{\text{prog}} \supseteq (m'_r \uplus m'')$. Take $l \in \text{dom}(m'_r \uplus m'')$. Since we
3209 already have the equality of values from the previous point, we only need to show that
3210 $l \in \text{dom}((\mu' \uplus \mu'')|_{\text{prog}}) = \text{dom}(\mu' \uplus \mu'')$. There are two cases: either $l \in \text{dom}(m'_r)$ or
3211 $l \in \text{dom}(m'')$.
- 3212 – If $l \in \text{dom}(m'_r)$, then since $m'_r = \mu'|_{\text{prog}}$, we have $l \in \text{dom}(\mu'|_{\text{prog}}) = \text{dom}(\mu')$. Therefore, $l \in$
3213 $\text{dom}(\mu' \uplus \mu'')$.
- 3214 – If $l \in \text{dom}(m'')$, then since $m'' = \mu''|_{\text{prog}} \setminus \text{dom}(m)$, we have $l \in \text{dom}(\mu''|_{\text{prog}}) =$
3215 $\text{dom}(\mu'')$. Therefore, $l \in \text{dom}(\mu' \uplus \mu'')$.
- 3216 • We have $(\mu' \uplus \mu'')|_{\text{prog}} \subseteq (m'_r \uplus m'')$ and $(\mu' \uplus \mu'')|_{\text{prog}} \supseteq (m'_r \uplus m'')$ Therefore, $(\mu' \uplus \mu'')|_{\text{prog}} =$
3217 $(m'_r \uplus m'')$.
- 3218 • $\text{OnlyRO}(\mu \uplus \mu'') = \text{OnlyRO}(\mu' \uplus \mu'')$ directly follows from $\text{OnlyRO}(\mu') = \text{OnlyRO}(\mu)$.
- 3219 • We conclude by choosing $\sigma' \uplus \sigma''$ and $\mu' \uplus \mu''$ and instantiate the definition of AcceptableStates
3220 with $(s'_r \cup \sigma''|_{\text{prog}}) = (\sigma' \uplus \sigma'')|_{\text{prog}}$ and $(m'_r \cup m'') = (\mu' \uplus \mu'')|_{\text{prog}}$ and $(\sigma' \uplus \sigma'', \mu' \uplus$
3221 $\mu'') \in \Gamma' \star \Gamma''$ and $\forall x \in \text{dom}(\sigma \uplus \sigma'') \cap \text{dom}(\sigma' \uplus \sigma'')$, $(\sigma \uplus \sigma'')(x) = (\sigma' \uplus \sigma'')(x)$ and
3222 $\text{OnlyRO}(\mu \uplus \mu'') = \text{OnlyRO}(\mu' \uplus \mu'')$

□

3225 D DETAILS OF TRIPLE MINIMIZATION

3226 In the description of the triple minimization operator in Section 5.4, we did not explain how it is
3227 computed. This section gives the missing implementation details.

3228 Minimize is computed by looking at the usage of each resource:

- 3229
- 3230 • For a resource H that appear as `uninit` in the usage map, if H is not already of the form
3231 `Uninit(H')`, we can wrap it as `Uninit(H)` in \hat{F} .
- 3232 • For resources that appear as `splittedFrac` in the usage map, we can give an arbitrarily small
3233 fraction to t , and keep the rest in the frame.
- 3234

- For resources that appear as `joinedFrac` in the usage map, we can completely place them in the frame.

These two last points, some care is needed for the minimized postcondition \hat{F}' , because new subfractions might have been created by t and were immediately merged into a resource that is now not given anymore. You can find below some examples of minimization made by our version of `Minimize`:

$\Gamma(y)$	$\Gamma'(y)$	$\Delta(y)$	E^{fracs}	\hat{F}	\hat{F}'	F^{framed}
H	H	$y \notin \Delta$	\emptyset	\emptyset	\emptyset	$y:H$
H	\emptyset	full	\emptyset	$y:H$	\emptyset	\emptyset
$\text{Uninit}(H)$	\emptyset	uninit	\emptyset	$y:\text{Uninit}(H)$	\emptyset	\emptyset
H	\emptyset	uninit	\emptyset	$y:\text{Uninit}(H)$	\emptyset	\emptyset
H	H	splittedFrac	$\alpha:\text{frac}$	$y':\alpha H$	$y':\alpha H$	$y:(1-\alpha)H$
αH	αH	splittedFrac	$\beta:\text{frac}$	$y':\beta H$	$y':\beta H$	$y:(\alpha-\beta)H$
$(\alpha-\beta)H$	αH	splittedFrac	$\gamma:\text{frac}$	$y':\gamma H$	$y':\gamma H, y_2:\beta H$	$y:(\alpha-\beta-\gamma)H$
αH	$(\alpha-\beta)H$	splittedFrac	$\gamma:\text{frac}$	$y':\gamma H$	$y':(\gamma-\beta)H$	$y:(\alpha-\gamma)H$
$(\alpha-\beta_1-\beta_2)H$	$(\alpha-\beta_1-\beta_3)H$	splittedFrac	$\gamma:\text{frac}$	$y':\gamma H$	$y':(\gamma-\beta_3)H, y_2:\beta_2 H$	$y:(\alpha-\gamma)H$
$(\alpha-\beta)H$	αH	joinedFrac	\emptyset	\emptyset	$y':\beta H$	$y:(\alpha-\beta)H$
$(\alpha-\beta_1-\beta_2-\beta_3)H$	$(\alpha-\beta_2)H$	joinedFrac	\emptyset	\emptyset	$y_1:\beta_1 H, y_3:\beta_3 H$	$y:(\alpha-\beta_1-\beta_2-\beta_3)H$

Algorithmically, `Minimize` can be defined by iterating over its first argument.

Start with $E^{\text{fracs}} = \emptyset$, $\hat{F} = \emptyset$, $\hat{F}' = \Gamma'$.linear and $F^{\text{framed}} = \emptyset$.

For each binding $y : H$ in Γ , lookup y in Δ :

- If y is not in Δ , add $y : H$ in F^{framed} and remove it from \hat{F}' (it must exist there by the invariants of triples).
- If $y : \text{full}$ is in Δ , add $y : H$ in \hat{F} .
- If $y : \text{uninit}$ is in Δ , add $y : \text{Uninit}(H)$ in \hat{F} or $y : H$ if H already is of the form $\text{Uninit}(H')$.
- If $y : \text{splittedFrac}$ is in Δ , decompose H as $(\alpha - \beta_1 - \dots - \beta_n)H'$. It is always possible since $\alpha = 1$ is allowed and the list of β_i can be empty. Create a fresh fraction $\gamma \leq \alpha - \beta_1 - \dots - \beta_n$ and place it in E^{fracs} . Add $y' : \gamma H'$ in \hat{F} and $y : (\alpha - \beta_1 - \dots - \beta_n - \gamma)H'$ in F^{framed} . Replace the binding $y : (\alpha - \delta_1 - \dots - \delta_m)H'$ in \hat{F}' by the following: try to pair δ_i with a matching β_j . For each unmatched β_i , add a binding $y'_i : \beta_i H'$ to \hat{F}' . These correspond to the subfractions that were created by t and merged into y . Let the unmatched δ_i form the list of $\tilde{\delta}_i$. These correspond to the subfractions consumed by t and not given back. Add the binding $y' : (\gamma - \tilde{\delta}_1 - \dots - \tilde{\delta}_m)H$ to \hat{F}' .
- If $y : \text{joinedFrac}$ is in Δ , decompose H as $(\alpha - \beta_1 - \dots - \beta_n)H'$. Add $y : H$ in F^{framed} . Remove the binding $y : (\alpha - \delta_1 - \dots - \delta_m)H'$ in \hat{F}' . Given that `joinedFrac` usage are only created by `CloseFrac`s and given how the `CloseFrac`s algorithm works, each δ_i will necessarily match one of the β_j , however there will be some β_i that are not matched. For each unmatched β_i , add the binding $y'_i : \beta_i H'$ in \hat{F}' .

The next two sections give details for two applications of this `Minimize` operator: typechecking subexpressions and loop contract minimization.

E EXAMPLE TYPECHECKING OF SUBEXPRESSIONS

This section presents an example application of the `SUBEXPR` from Section 5.5 and repeated below.

i	0	1	2	3
$\Gamma_i.\text{pure}$	$p, q, c : \text{ptr}_{\text{int}}$	$p, q, c : \text{ptr}_{\text{int}}$	$p, q, c : \text{ptr}_{\text{int}}$	$p, q, c : \text{ptr}_{\text{int}}$ $\alpha : \text{frac}$
$\Gamma_i.\text{linear}$	$Hp : p \rightsquigarrow \text{Cell}$ $Hq : q \rightsquigarrow \text{Cell}$ $Hc : c \rightsquigarrow \text{Cell}$	$Hp : p \rightsquigarrow \text{Cell}$ $Hq : q \rightsquigarrow \text{Cell}$ $Hc : c \rightsquigarrow \text{Cell}$	$Hp : p \rightsquigarrow \text{Cell}$ $Hq : q \rightsquigarrow \text{Cell}$	$Hp : (1 - \alpha)(p \rightsquigarrow \text{Cell})$ $Hq : q \rightsquigarrow \text{Cell}$
t_i	q	$\text{get_incr}(c)$	$\text{get}(p)$	$\text{get}(p)$
Δ_i	$q : \text{required}$ res : ensured	$c : \text{required}$ $Hc : \text{full}$ $Hc' : \text{produced}$ res : ensured	$p : \text{required}$ $Hp : \text{splittedFrac}$ res : ensured	$p : \text{required}$ $Hp' : \text{splittedFrac}$ res : ensured
E_i^{fracs}	\emptyset	\emptyset	$\alpha : \text{frac}$	$\beta : \text{frac}$
\hat{F}_i	\emptyset	$Hc : c \rightsquigarrow \text{Cell}$	$\alpha(p \rightsquigarrow \text{Cell})$	$\beta(p \rightsquigarrow \text{Cell})$
\hat{F}'_i	\emptyset	$Hc' : c \rightsquigarrow \text{Cell}$	$\alpha(p \rightsquigarrow \text{Cell})$	$\beta(p \rightsquigarrow \text{Cell})$
F_i^{framed}	$Hp : p \rightsquigarrow \text{Cell}$ $Hq : q \rightsquigarrow \text{Cell}$ $Hc : c \rightsquigarrow \text{Cell}$	$Hp : p \rightsquigarrow \text{Cell}$ $Hq : q \rightsquigarrow \text{Cell}$	$Hp : (1 - \alpha)(p \rightsquigarrow \text{Cell})$ $Hq : q \rightsquigarrow \text{Cell}$	$Hp :$ $(1 - \alpha - \beta)(p \rightsquigarrow \text{Cell})$ $Hq : q \rightsquigarrow \text{Cell}$
$\hat{\Gamma}'_i.\text{pure}$	res := $q : \text{ptr}_{\text{int}}$	res : int	res : int	res : int
$\Gamma_p.\text{pure}$	$p, q, c : \text{ptr}_{\text{int}}, x_0 := q : \text{ptr}_{\text{int}}, x_1, x_2, x_3 : \text{int}$			
$\Gamma_p.\text{linear}$	$Hp : p \rightsquigarrow \text{Cell}, Hq : q \rightsquigarrow \text{Cell}, Hc' : c \rightsquigarrow \text{Cell}$			

Fig. 26. Example of an application of the SUBEXPR rule on an expression $\hat{\mathcal{E}}[q, \text{get_incr}(c), \text{get}(p), \text{get}(p)]$, in a context $\langle p, q, c : \text{ptr}_{\text{int}} \mid Hp : p \rightsquigarrow \text{Cell}, Hq : q \rightsquigarrow \text{Cell}, Hc : c \rightsquigarrow \text{Cell} \rangle$.

SUBEXPR

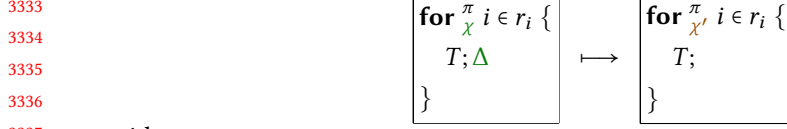
$$\begin{array}{l}
\forall i \in [1, n]. \quad \{\Gamma_{i-1}\} t_i^{\Delta_i} \{\Gamma'_i\} \wedge (E_i^{\text{fracs}}, \hat{F}_i, \hat{F}'_i, F_i^{\text{framed}}) = \text{Minimize}(\Gamma_{i-1}, \Gamma'_i, \Delta_i) \wedge x_i \text{ fresh} \\
\forall i \in [1, n]. \quad \Gamma_i = \langle \Gamma_i.\text{pure}, E_i^{\text{fracs}} \mid F_i^{\text{framed}} \rangle \wedge \hat{\Gamma}'_i = \langle \Gamma'_i.\text{pure} \wedge \Delta_i.\text{ensured} \mid \hat{F}'_i \rangle \\
\Gamma_p = \text{CloseFracs}^{\Delta_p} (\Gamma_n \otimes \star_{i \in [0, n]} \text{Rename}\{\mathbf{res} := x_i\}(\hat{\Gamma}'_i)) \\
\{\Gamma_p\} \hat{\mathcal{E}}[x_1, \dots, x_n]^{\Delta_q} \{\Gamma_q\} \\
\Delta = \text{Rename}\{\mathbf{res} := x_1\}(\Delta_1); \dots; \text{Rename}\{\mathbf{res} := x_n\}(\Delta_n); \Delta_p; \Delta_q \\
\hline
\{\Gamma_0\} \hat{\mathcal{E}}[t_1, \dots, t_n]^{\Delta} \{\Gamma_q\}
\end{array}$$

The example consists of a multi-evaluation-context $\hat{\mathcal{E}}$, which could be a function call, featuring 4 subexpression holes: $\hat{\mathcal{E}}[q, \text{get_incr}(c), \text{get}(p), \text{get}(p)]$. This expression is typechecked in a typing context: $\langle p, q, c : \text{ptr}_{\text{int}} \mid Hp : p \rightsquigarrow \text{Cell}, Hq : q \rightsquigarrow \text{Cell}, Hc : c \rightsquigarrow \text{Cell} \rangle$.

Figure 26 shows the typechecking steps. The figure includes 4 columns, describing the steps associated with each of the 4 subterms. The rows explain how the metavariables from the rule SUBEXPR are instantiated. In particular, observe how the two subexpressions $\text{get}(p)$ both have read-only access to the same resource H . As the details in the Figure show, the first $\text{get}(p)$, according to its minimized precondition, only needs a fraction of H . This fraction is carved out, obtaining a subfraction αH and leaving $(1 - \alpha)H$ for the second $\text{get}(p)$.

F DETAILS OF LOOP MINIMIZATION

Figure 27 describes the loop minimization transformation. Essentially, it uses the Minimize operator to minimize the exclusive part of the loop contract, and it tries to reduce the footprint of the shared part of the loop contract by taking arbitrary subfractions and using shared reads whenever possible.



3337 with:

3338 $(E^{\text{frac}}, \hat{F}, \hat{F}', _) = \text{Minimize}(\chi.\text{excl.pre}, \Sigma^{-1}(\chi.\text{excl.post}), \Delta)$

3339 $\langle E_{RO} \mid F_{RO} \rangle = \text{IntoRO}((\chi.\text{shrd.inv} \vdash \Delta.\text{read}) * (\chi.\text{shrd.reads} \vdash \Delta.\text{read}))$

3340 $\left\{ \begin{array}{l} \text{shrd.reads} \equiv F_{RO} * (\chi.\text{shrd.reads} \vdash \Delta.\text{alter}) \\ \text{shrd.inv} \equiv \chi.\text{shrd.inv} \vdash \Delta.\text{alter} \end{array} \right.$

3341 $\chi' \equiv \left\{ \begin{array}{l} \text{excl.pre} \equiv \langle (\chi.\text{excl.pre.pure} \vdash \Delta.\text{required}) \mid \hat{F} \rangle \\ \text{excl.post} \equiv \Sigma(\langle \chi.\text{excl.post.pure} \mid \hat{F}' \rangle) \end{array} \right.$

3342 $\left\{ \begin{array}{l} \text{vars} \equiv (\chi.\text{vars} \vdash (\text{usedVars}(\chi'.\text{shrd}) \cup \text{usedVars}(\chi'.\text{excl}) \cup \Delta.\text{required})), E_{RO}, E^{\text{frac}} \end{array} \right.$

3343

3344

3345

3346

3347

3348 Fig. 27. The basic transformation `Loop.minimize`. The `Minimize` operation is that defined for triples in

3349 Section 5.4. `usedVars(X)` is the set of all variables used in X at least once.

3350

3351 For that we use a new operator `IntoRO` that operates on linear contexts and is defined recursively

3352 as follows:

3353
$$\text{IntoRO}(F) := \begin{cases} \langle \mid \rangle & \text{if } F = \emptyset \\ \langle \alpha : \text{frac} \mid y : \alpha H \rangle \otimes \text{IntoRO}(F') & \text{if } F = (y : H) :: F' \end{cases}$$

3354

3355

3356 For pure variables, it simply removes those that are not used and adds the new arbitrary fractions

3357 generated during the previous steps.

3358 One technical difficulty: postcondition of a loop contract uses names for linear resources, and

3359 these names must be matched to corresponding resources at the end of the body of the loop. In

3360 fact our typechecker had to prove an entailment there. We can remember the map Σ from linear

3361 resource names at the end of the loop body to linear resources names in the postcondition, and use

3362 it in the loop minimization transformation.

3363

3364 G PARTIAL SUBTRACTION AND PARTITIONING OF RESOURCES

3365 A concrete way to compute the partitions involved in loop swap is by using the `PartialSub`(Γ_1, Γ_2)

3366 operator similar to the context subtraction operator from Section 4.5. Instead of failing like $\Gamma_1 \boxminus \Gamma_2$

3367 when a resource in Γ_2 cannot be found in Γ_1 , `PartialSub`(Γ_1, Γ_2) returns both the resources that were

3368 found, and the ones that could not be found, including pure resources. Intuitively, F represents

3369 resources that are common between Γ_1 and Γ_2 , Γ'_1 represents resources that remain from Γ_1 , and

3370 Γ'_2 represents resources that remain from Γ_2 . More formally, if $F, \Gamma'_1, \Gamma'_2 = \text{PartialSub}(\Gamma_1, \Gamma_2)$, then:

3371 $\exists \sigma. \Gamma_1 \Rightarrow \Gamma'_1 \otimes F \wedge \text{Specialize}_{\Gamma_1} \{\sigma\}(\Gamma_2) \Rightarrow \Gamma'_2 \otimes F$, with Γ'_2 containing as few resources as possible.

3372 Using this operator, our example partition can be computed as follows:

3373

3374
$$RR, RP_i, _ \equiv \text{PartialSub}(\chi_j.\text{shrd.reads}, \chi_i.\text{shrd.reads})$$

3375

3376 In practice, we can sometimes avoid this computation altogether in our implementation by

3377 instead leveraging information left by our type checker regarding realized contract instantiations.

3378 REFERENCES

3379 2022. Unsequenced functions. <https://www.open-std.org/jtc1/sc22/wg14/www/docs/n2956.htm> ISO/IEC JTC1/SC22/WG14

3380 document N2956.

3381

- 3382 Sami Alabed, Daniel Belov, Bart Chrzaszcz, Juliana Franco, Dominik Grewe, Dougal Maclaurin, James Molloy, Tom Natan,
3383 Tamara Norman, Xiaoyue Pan, Adam Paszke, Norman A. Rink, Michael Schaarschmidt, Timur Sitdikov, Agnieszka
3384 Swietlik, Dimitrios Vytiniotis, and Joel Wee. 2024. PartIR: Composing SPMD Partitioning Strategies for Machine Learning.
arXiv:2401.11202 [cs.LG] <https://arxiv.org/abs/2401.11202>
- 3385 Vasco Amaral, Beatriz Norberto, Miguel Goulão, Marco Aldinucci, Siegfried Benkner, Andrea Bracciali, Paulo Carreira, Edgars
3386 Celms, Luís Correia, Clemens Grelck, Helen Karatza, Christoph Kessler, Peter Kilpatrick, Hugo Martiniano, Ilias Mavridis,
3387 Sabri Pllana, Ana Respício, José Simão, Luís Veiga, and Ari Visa. 2020. Programming languages for data-Intensive HPC
3388 applications: A systematic mapping study. *Parallel Comput.* 91 (2020), 102584. <https://doi.org/10.1016/j.parco.2019.102584>
- 3389 Mehdi Amini. 2012. *Source-to-source automatic program transformations for GPU-like hardware accelerators*. Ph.D. Disserta-
3390 tion. Ecole Nationale Supérieure des Mines de Paris.
- 3391 Hansang Bae, Dheya Mustafa, Jae-Woo Lee, Aurangzeb, Hao Lin, Chirag Dave, Rudolf Eigenmann, and Samuel P. Midkiff.
3392 2013. The Cetus Source-to-Source Compiler Infrastructure: Overview and Evaluation. *Int. J. Parallel Program.* 41, 6
3393 (2013), 753–767. <https://doi.org/10.1007/S10766-012-0211-Z>
- 3394 O.S. Bagge, K.T. Kalleberg, M. Haverlaen, and E. Visser. 2003. Design of the CodeBoost transformation system for domain-
3395 specific optimisation of C++ programs. In *Proceedings Third IEEE International Workshop on Source Code Analysis and*
3396 *Manipulation*. 65–74. <https://doi.org/10.1109/SCAM.2003.1238032>
- 3397 Léo Bagnères, Oleksandr Zinenko, Stéphane Huot, and Cédric Bastoul. 2016a. Opening Polyhedral Compiler’s Black Box.
3398 In *IEEE/ACM International Symposium on Code Generation and Optimization (CGO)*.
- 3399 Léo Bagnères, Oleksandr Zinenko, Stéphane Huot, and Cédric Bastoul. 2016b. Opening Polyhedral Compiler’s Black Box.
3400 In *IEEE/ACM International Symp. on Code Generation and Optimization*.
- 3401 Paul Barham and Michael Isard. 2019. Machine Learning Systems are Stuck in a Rut. In *Proceedings of the Workshop on Hot*
3402 *Topics in Operating Systems (Bertinoro, Italy) (HotOS ’19)*. Association for Computing Machinery, New York, NY, USA,
3403 177–183. <https://doi.org/10.1145/3317550.3321441>
- 3404 Y. Barsamian, A. Charguéraud, S. A. Hirstoaga, and M. Mehrenberger. 2018. Efficient Strict-Binning Particle-in-Cell Algorithm
3405 for Multi-Core SIMD Processors. In *24th International Conference on Parallel and Distributed Computing (Euro-Par) (Lecture*
3406 *Notes in Computer Science, Vol. 11014)*. Springer, Cham, 749–763. https://doi.org/10.1007/978-3-319-96983-1_53
- 3407 Gilles Barthe, Benjamin Grégoire, César Kunz, and Tamara Rezk. 2009. Certificate Translation for Optimizing Compilers.
3408 *ACM Trans. Program. Lang. Syst.* 31, 5, Article 18 (jul 2009), 45 pages. <https://doi.org/10.1145/1538917.1538919>
- 3409 Guillaume Bertholon and Arthur Charguéraud. 2025. Bidirectional Translation between a C-like Language and an Imperative
3410 Lambda-calculus. In *36es Journées Francophones des Langages Applicatifs (JFLA 2025)*. Roiffé, France. <https://inria.hal.science/hal-04859522>
- 3411 João Bispo and João MP Cardoso. 2020. Clava: C/C++ source-to-source compilation using LARA. *SoftwareX* 12 (2020),
3412 100565. <https://www.sciencedirect.com/science/article/pii/S2352711019302122/pdf>
- 3413 Stefan Blom, Saeed Darabi, Marieke Huisman, and Wytse Oortwijn. 2017. The VerCors Tool Set: Verification of Parallel and
3414 Concurrent Software. In *Integrated Formal Methods, Nadia Polikarpova and Steve Schneider (Eds.)*. Springer International
3415 Publishing, Cham, 102–110.
- 3416 Uday Bondhugula, Albert Hartono, J. Ramanujam, and P. Sadayappan. 2008. A practical automatic polyhedral parallelizer
3417 and locality optimizer. In *PLDI’08 ACM Conf. on Programming language design and implementation*.
- 3418 John Boyland. 2003. Checking Interference with Fractional Permissions. In *Static Analysis, 10th International Symposium,*
3419 *SAS 2003, San Diego, CA, USA, June 11-13, 2003, Proceedings (Lecture Notes in Computer Science, Vol. 2694)*, Radhia Cousot
3420 (Ed.). Springer, 55–72. https://doi.org/10.1007/3-540-44898-5_4
- 3421 Gary Bradschi, Adrian Kaehler, et al. 2000. OpenCV. *Dr. Dobb’s journal of software tools* 3, 2 (2000). <https://opencv.org/>.
- 3422 Martin Bravenboer, Karl Trygve Kalleberg, Rob Vermaas, and Eelco Visser. 2008. Stratego/XT 0.17. A Language and Toolset
3423 for Program Transformation. *Sci. Comput. Program.* 72, 1–2 (jun 2008), 52–70. <https://doi.org/10.1016/j.scico.2007.11.003>
- 3424 Cristiano Calcagno, Dino Distefano, Peter O’Hearn, and Hongseok. Yang. 2019. Go Huge or Go Home: POPL’19 Most
3425 Influential Paper Retrospective. [https://blog.sigplan.org/2020/03/03/go-huge-or-go-home-popl19-most-influentialpaper-](https://blog.sigplan.org/2020/03/03/go-huge-or-go-home-popl19-most-influentialpaper-retrospective/)
3426 [retrospective/](https://blog.sigplan.org/2020/03/03/go-huge-or-go-home-popl19-most-influentialpaper-retrospective/)
- 3427 Qinxiang Cao, Lennart Beringer, Samuel Gruetter, Josiah Dodds, and Andrew W. Appel. 2018. VST-Floyd: A Separation
3428 Logic Tool to Verify Correctness of C Programs. *J. Autom. Reason.* 61, 1-4 (2018), 367–422. [https://doi.org/10.1007/S10817-](https://doi.org/10.1007/S10817-018-9457-5)
3429 [018-9457-5](https://doi.org/10.1007/S10817-018-9457-5)
- 3430 Arthur Charguéraud. 2020a. Separation logic for sequential programs (functional pearl). *Proc. ACM Program. Lang.* 4, ICFP,
3431 Article 116 (aug 2020), 34 pages. <https://doi.org/10.1145/3408998>
- 3432 Arthur Charguéraud. 2020b. Separation Logic for Sequential Programs (Functional Pearl). *Proc. ACM Program. Lang.* 4,
3433 ICFP, Article 116 (aug 2020), 34 pages. <https://doi.org/10.1145/3408998>
- 3434 Arthur Charguéraud, Adam Chlipala, Andres Erbsen, and Samuel Gruetter. 2022. Omnisemantics: Smooth Handling of
3435 Nondeterminism. (Sept. 2022). <https://hal.inria.fr/hal-03255472> To appear in ACM Transactions on Programming
3436 Languages and Systems (TOPLAS).

- 3431 Arthur Charguéraud, Adam Chlipala, Andres Erbsen, and Samuel Gruetter. 2023. Omnisemantics: Smooth Handling of
3432 Nondeterminism. *ACM Trans. Program. Lang. Syst.* 45, 1, Article 5 (March 2023), 43 pages. <https://doi.org/10.1145/3579834>
- 3433 Gaurav Chaurasia, Jonathan Ragan-Kelley, Sylvain Paris, George Drettakis, and Frédo Durand. 2015. Compiling high
3434 performance recursive filters. In *Proceedings of the 7th Conference on High-Performance Graphics, HPG 2015, Los Angeles,*
3435 *California, USA, August 7-9, 2015*, Michael C. Doggett, Steven E. Molnar, Kayvon Fatahalian, Jacob Munkberg, Elmar
3436 Eisemann, Petrik Clarberg, and Stephen N. Spencer (Eds.). ACM, 85–94. <https://doi.org/10.1145/2790060.2790063>
- 3437 Lorenzo Chelini, Martin Kong, Tobias Grosser, and Henk Corporaal. 2021. LoopOpt: Declarative Transformations Made Easy.
3438 In *Proceedings of the 24th International Workshop on Software and Compilers for Embedded Systems* (Eindhoven, Netherlands)
3439 (*SCOPES '21*). Association for Computing Machinery, New York, NY, USA, 11–16. <https://doi.org/10.1145/3493229.3493301>
- 3440 Chun Chen, Jacqueline Chame, and Mary W. Hall. 2008. *CHiLL: A Framework for Composing High-Level Loop Transformations*.
3441 Technical Report 08-897. University of Southern California.
- 3442 Tianqi Chen, Thierry Moreau, Ziheng Jiang, Lianmin Zheng, Eddie Yan, Haichen Shen, Meghan Cowan, Leyuan Wang,
3443 Yuwei Hu, Luis Ceze, Carlos Guestrin, and Arvind Krishnamurthy. 2018. TVM: An Automated End-to-End Optimizing
3444 Compiler for Deep Learning. In *OSDI*. USENIX Association. <https://www.usenix.org/system/files/osdi18-chen.pdf>
- 3445 James R Cordy. 2006. The TXL source transformation language. *Science of Computer Programming* 61, 3 (2006), 190–210.
- 3446 Chirag Dave, Hansang Bae, Seung-Jai Min, Seyong Lee, Rudolf Eigenmann, and Samuel Midkiff. 2009. Cetus: A Source-to-
3447 Source Compiler Infrastructure for Multicores. *Computer* 42, 12 (2009), 36–42. <https://doi.org/10.1109/MC.2009.385>
- 3448 Thomas M Evans, Andrew Siegel, Erik W Draeger, Jack Deslippe, Marianne M Francois, Timothy C Germann, William E
3449 Hart, and Daniel F Martin. 2022. A survey of software implementations used by application codes in the Exascale
3450 Computing Project. *The International Journal of High Performance Computing Applications* 36, 1 (2022), 5–12. <https://doi.org/10.1177/10943420211028940>
- 3451 Paul Feautrier. 1992. Some efficient solutions to the affine scheduling problem: one dimensional time. *Intl. Journal of Parallel*
3452 *Programming* 21, 5 (october 1992), 313–348.
- 3453 Jean-Christophe Filliâtre and Andrei Paskevich. 2013. Why3—Where Programs Meet Provers. In *European Symposium on*
3454 *Programming (ESOP) (Lecture Notes in Computer Science, Vol. 7792)*. Springer, 125–128. <http://hal.inria.fr/hal-00789533>
- 3455 Jean-Christophe Filliâtre. 2003. *Why: a multi-language multi-prover verification tool*. Research Report 1366. LRI, Université
3456 Paris Sud. <http://www.lri.fr/~filliatr/ftp/publis/why-tool.ps.gz>
- 3457 Cormac Flanagan, K. Rustan M. Leino, Mark Lillibridge, Greg Nelson, James B. Saxe, and Raymie Stata. 2002. Extended
3458 Static Checking for Java. In *Programming Language Design and Implementation (PLDI)*. 234–245. <http://www.soe.ucsc.edu/~cormac/papers/pldi02.ps>
- 3459 John John Gough and K John Gough. 2001. *Compiling for the .Net Common Language Runtime*. Prentice Hall PTR.
- 3460 Bastian Hagedorn, Archibald Samuel Elliott, Henrik Barthels, Rastislav Bodik, and Vinod Grover. 2020a. Fireiron: a data-
3461 movement-aware scheduling language for GPUs. In *Proceedings of the ACM International Conf. on Parallel Architectures*
3462 *and Compilation Techniques*. 71–82.
- 3463 Bastian Hagedorn, Johannes Lenfers, Thomas Kundenedhler, Xueying Qin, Sergei Gorlatch, and Michel Steuwer. 2020b.
3464 Achieving high-performance the functional way: a functional pearl on expressing high-performance optimizations as
3465 rewrite strategies. *Proc. ACM Program. Lang.* 4, ICFP, Article 92 (aug 2020), 29 pages. <https://doi.org/10.1145/3408974>
- 3466 Stefan Heule, K. Rustan M. Leino, Peter Müller, and Alexander J. Summers. 2013. Abstract Read Permissions: Fractional
3467 Permissions without the Fractions. In *Verification, Model Checking, and Abstract Interpretation*, Roberto Giacobazzi, Josh
3468 Berdine, and Isabella Mastroeni (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 315–334.
- 3469 Yuka Ikarashi, Gilbert Louis Bernstein, Alex Reinking, Hasan Genc, and Jonathan Ragan-Kelley. 2022. Exocompilation for
3470 productive programming of hardware accelerators. In *Proceedings of the 43rd ACM SIGPLAN International Conference*
3471 *on Programming Language Design and Implementation* (San Diego, CA, USA) (*PLDI 2022*). Association for Computing
3472 Machinery, New York, NY, USA, 703–718. <https://doi.org/10.1145/3519939.3523446>
- 3473 Yuka Ikarashi, Jonathan Ragan-Kelley, Tsukasa Fukusato, Jun Kato, and Takeo Igarashi. 2021. Guided Optimization for
3474 Image Processing Pipelines. In *2021 IEEE Symposium on Visual Languages and Human-Centric Computing (VL/HCC)*. 1–5.
3475 <https://doi.org/10.1109/VL/HCC51201.2021.9576341>
- 3476 Matthieu Journault and Antoine Miné. 2018. Inferring functional properties of matrix manipulating programs by abstract
3477 interpretation. *Form. Methods Syst. Des.* 53, 2 (oct 2018), 221–258. <https://doi.org/10.1007/s10703-017-0311-x>
- 3478 Ralf Jung, Robbert Krebbers, Jacques-Henri Jourdan, Ales Bizjak, Lars Birkedal, and Derek Dreyer. 2018a. Iris from the
3479 ground up: A modular foundation for higher-order concurrent separation logic. *J. Funct. Program.* 28 (2018), e20.
3480 <https://doi.org/10.1017/S0956796818000151>
- 3481 Ralf Jung, Robbert Krebbers, Jacques-Henri Jourdan, Aleš Bizjak, Lars Birkedal, and Derek Dreyer. 2018b. Iris from the
3482 ground up: A modular foundation for higher-order concurrent separation logic. *Journal of Functional Programming* 28
3483 (2018), e20. <https://people.mpi-sws.org/~dreyer/papers/iris-ground-up/paper.pdf>
- 3484 Yamato Kanetaka, Hiroyasu Takagi, Yoshihiro Maeda, and Norishige Fukushima. 2024. SlidingConv: Domain-Specific
3485 Description of Sliding Discrete Cosine Transform Convolution for Halide. *IEEE Access* 12 (2024), 7563–7583. <https://doi.org/10.1109/ACCESS.2024.3388888>
- 3486

- 3480 //doi.org/10.1109/ACCESS.2023.3345660
- 3481 Vasilios Kelefouras and Georgios Keramidas. 2022. Design and Implementation of 2D Convolution on x86/x64 Processors. *IEEE Transactions on Parallel and Distributed Systems* 33, 12 (2022), 3800–3815. <https://doi.org/10.1109/TPDS.2022.3171471>
- 3482 Athanasios Konstantinidis. 2013. *Source-to-source compilation of loop programs for manycore processors*. Ph. D. Dissertation. Imperial College London.
- 3483 Michael Kruse and Hal Finkel. 2018. A Proposal for Loop-Transformation Pragmas. *CoRR* abs/1805.03374 (2018). arXiv:1805.03374 <http://arxiv.org/abs/1805.03374>
- 3484 Olaf Krzikalla, Kim Feldhoff, Ralph Müller-Pfefferkorn, and Wolfgang E. Nagel. 2011. Scout: A Source-to-Source Transformer for SIMD-Optimizations. In *Euro-Par Workshops (2) (LNCS, Vol. 7156)*. Springer.
- 3485 César Kunz. 2009. *Proof preservation and program compilation*. Ph. D. Dissertation. École Nationale Supérieure des Mines de Paris. <https://pastel.archives-ouvertes.fr/pastel-00004940/file/thesis-ckunz.pdf>
- 3486 Chris Lattner, Mehdi Amini, Uday Bondhugula, Albert Cohen, Andy Davis, Jacques Pienaar, River Riddle, Tatiana Shpeisman, Nicolas Vasilache, and Aleksandr Zinenko. 2021. MLIR: Scaling Compiler Infrastructure for Domain Specific Computation. In *2021 IEEE/ACM International Symposium on Code Generation and Optimization (CGO)*, 2–14. <https://doi.org/10.1109/CGO51591.2021.9370308>
- 3487 Julia Lawall and Gilles Muller. 2018. Coccinelle: 10 Years of Automated Evolution in the Linux Kernel. In *USENIX Conference on Usenix Annual Technical Conference (USENIX ATC '18)*. USENIX Association, 13 pages.
- 3488 Youenn Lebras. 2019. *Code optimization based on source to source transformations using profile guided metrics*. Ph. D. Dissertation. Université Paris-Saclay (ComUE). <https://www.theses.fr/2019SACLV037.pdf>
- 3489 Amanda Liu, Gilbert Louis Bernstein, Adam Chlipala, and Jonathan Ragan-Kelley. 2022. Verified Tensor-Program Optimization via High-Level Scheduling Rewrites. 6, *POPL*, Article 55 (jan 2022), 28 pages. <https://doi.org/10.1145/3498717>
- 3490 Martin Paul Lücke, Oleksandr Zinenko, William S. Moses, Michel Steuwer, and Albert Cohen. 2024. The MLIR Transform Dialect. Your compiler is more powerful than you think. *CoRR* abs/2409.03864 (2024). <https://doi.org/10.48550/ARXIV.2409.03864> arXiv:2409.03864
- 3491 Peter Müller, Malte Schwerhoff, and Alexander J. Summers. 2017. Viper: A Verification Infrastructure for Permission-Based Reasoning. In *Dependable Software Systems Engineering*. IOS Press, 104–125. <https://doi.org/10.3233/978-1-61499-810-5-104>
- 3500 Kedar S. Namjoshi and Nimit Singhania. 2016. Loopy: Programmable and Formally Verified Loop Transformations. In *Static Analysis - 23rd International Symposium, SAS (LNCS, Vol. 9837)*. Springer.
- 3501 George Ciprian Necula. 1998. *Compiling with proofs*. Ph. D. Dissertation. Carnegie Mellon University.
- 3502 Peter W. O’Hearn. 2019. Separation logic. *Commun. ACM* 62, 2 (2019), 86–95. <https://doi.org/10.1145/3211968>
- 3503 Pedro Pinto, Joao Bispo, Joao Cardoso, Jorge Gomes Barbosa, Davide Gadioli, Gianluca Palermo, Jan Martinovic, Martin Golasowski, Katerina Slaninova, Radim Cmar, et al. 2020. Pegasus: Performance Engineering for Software Applications Targeting HPC Systems. *IEEE Transactions on Software Engineering* (2020). <https://repositorio-aberto.up.pt/bitstream/10216/127756/2/405707.pdf>
- 3504 Dan Quinlan. 2000. ROSE: Compiler support for object-oriented frameworks. *Parallel processing letters* 10, 02n03 (2000), 215–226. https://digital.library.unt.edu/ark:/67531/metadc741175/m2/1/high_res_d/793936.pdf
- 3505 Dan Quinlan and Chunhua Liao. 2011. The ROSE source-to-source compiler infrastructure. In *Cetus users and compiler infrastructure workshop, in conjunction with PACT*, Vol. 2011. 1.
- 3506 Pawel K. Radtke and Tobias Weinzierl. 2024. Compiler support for semi-manual AoS-to-SoA conversions with data views. arXiv:2405.12507 [cs.PL] <https://arxiv.org/abs/2405.12507>
- 3507 Jonathan Ragan-Kelley. 2023. Technical Perspective: Reconsidering the Design of User-Schedulable Languages. *Commun. ACM* 66, 3 (feb 2023), 88. <https://doi.org/10.1145/3580370>
- 3508 Jonathan Ragan-Kelley, Connelly Barnes, Andrew Adams, Sylvain Paris, Frédo Durand, and Saman Amarasinghe. 2013. Halide: A Language and Compiler for Optimizing Parallelism, Locality, and Recomputation in Image Processing Pipelines. In *Conference on Programming Language Design and Implementation*. 12 pages. <https://doi.org/10.1145/2491956.2462176>
- 3509 John C. Reynolds. 2002. Separation Logic: A Logic for Shared Mutable Data Structures. In *17th IEEE Symposium on Logic in Computer Science (LICS 2002), 22-25 July 2002, Copenhagen, Denmark, Proceedings*. IEEE Computer Society, 55–74. <https://doi.org/10.1109/LICS.2002.1029817>
- 3510 Gabe Rudy, Malik Murtaza Khan, Mary Hall, Chun Chen, and Jacqueline Chame. 2011. A Programming Language Interface to Describe Transformations and Code Generation. In *Languages and Compilers for Parallel Computing*. Springer Berlin Heidelberg.
- 3511 Ömer Sakar, Mohsen Safari, Marieke Huisman, and Anton Wijs. 2022. Alpinist: An Annotation-Aware GPU Program Optimizer. In *Tools and Algorithms for the Construction and Analysis of Systems - 28th International Conference, TACAS 2022, Held as Part of the European Joint Conferences on Theory and Practice of Software, ETAPS 2022, Munich, Germany, April 2-7, 2022, Proceedings, Part II (Lecture Notes in Computer Science, Vol. 13244)*, Dana Fisman and Grigore Rosu (Eds.). Springer, 332–352. https://doi.org/10.1007/978-3-030-99527-0_18
- 3528

- 3529 Michael Sammler, Rodolphe Lepigre, Robbert Krebbers, Kayvan Memarian, Derek Dreyer, and Deepak Garg. 2021. RefinedC:
3530 automating the foundational verification of C code with refined ownership types. In *PLDI '21: 42nd ACM SIGPLAN*
3531 *International Conference on Programming Language Design and Implementation, Virtual Event, Canada, June 20-25, 2021*,
3532 Stephen N. Freund and Eran Yahav (Eds.). ACM, 158–174. <https://doi.org/10.1145/3453483.3454036>
- 3533 Cristina Silvano, Giovanni Agosta, Andrea Bartolini, Andrea R. Beccari, Luca Benini, Loïc Besnard, João Bispo, Radim
3534 Cmar, João M.P. Cardoso, Carlo Cavazzoni, Daniele Cesarini, Stefano Cherubin, Federico Ficarelli, Davide Gadioli,
3535 Martin Golasowski, Antonio Libri, Jan Martinovič, Gianluca Palermo, Pedro Pinto, Erven Rohou, Kateřina Slaninová, and
3536 Emanuele Vitali. 2019. The ANTAREX domain specific language for high performance computing. *Microprocessors and*
Microsystems 68 (2019), 58–73. <https://doi.org/10.1016/j.micpro.2019.05.005>
- 3537 Simon Spies, Lennard Gäher, Michael Sammler, and Derek Dreyer. 2024. Quiver: Guided Abductive Inference of Separation
3538 Logic Specifications in Coq. *Proc. ACM Program. Lang.* 8, PLDI, Article 183 (jun 2024), 25 pages. <https://doi.org/10.1145/3656413>
- 3539 Manish Vachharajani, Neil Vachharajani, David I. August, and Spyridon Triantafyllis. 2003. Compiler Optimization-Space
3540 Exploration. In *Proceedings of the 2013 IEEE/ACM International Symposium on Code Generation and Optimization (CGO)*.
3541 IEEE Computer Society, Los Alamitos, CA, USA, 204. <https://doi.org/10.1109/CGO.2003.1191546>
- 3542 Xuejun Yang, Yang Chen, Eric Eide, and John Regehr. 2011. Finding and Understanding Bugs in C Compilers. In *Conference*
3543 *on Programming Language Design and Implementation* (San Jose, California, USA). Association for Computing Machinery,
12 pages. <https://doi.org/10.1145/1993498.1993532>
- 3544 Qing Yi and Apan Qasem. 2008. Exploring the Optimization Space of Dense Linear Algebra Kernels. In *LCPC*.
- 3545 Qing Yi, Qian Wang, and Huimin Cui. 2014. Specializing Compiler Optimizations through Programmable Composition for
3546 Dense Matrix Computations. In *Proceedings of the 47th Annual IEEE/ACM International Symposium on Microarchitecture*
3547 *(Cambridge, United Kingdom) (MICRO-47)*. IEEE Computer Society, USA, 596–608. [https://doi.org/10.1109/MICRO.2014.](https://doi.org/10.1109/MICRO.2014.14)
14
- 3548 Oleksandr Zinenko, Lorenzo Chelini, and Tobias Grosser. 2018a. *Declarative Transformations in the Polyhedral Model*.
3549 Research Report RR-9243. <https://hal.inria.fr/hal-01965599>
- 3550 Oleksandr Zinenko, Stéphane Huot, and Cédric Bastoul. 2018b. Visual Program Manipulation in the Polyhedral Model.
3551 *ACM Trans. Archit. Code Optim.* 15, 1, Article 16 (mar 2018), 25 pages. <https://doi.org/10.1145/3177961>
- 3552
- 3553
- 3554
- 3555
- 3556
- 3557
- 3558
- 3559
- 3560
- 3561
- 3562
- 3563
- 3564
- 3565
- 3566
- 3567
- 3568
- 3569
- 3570
- 3571
- 3572
- 3573
- 3574
- 3575
- 3576
- 3577