

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

GUILLAUME BERTHOLON, ARTHUR CHARGUÉRAUD, THOMAS KÖHLER, BEGATIM BYTYQI, and DAMIEN ROUHLING, Inria & Université de Strasbourg, CNRS, ICube, 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 *diff*s over conventional C code. 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 $50\times$, 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,

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

	Halide/TVM	Elevate+Rise	Exo	Clay/LoopOpt	ATL	Alpinist	Clava+LARA
Generality	○	◐	◑	◒	◓	◔	◕
Expressiveness	●	●	●	◐	◑	◒	◓
Control	◐	◑	◒	◓	◔	◕	◖
Feedback	◐	◑	◒	◓	◔	◕	◖
Composability	○	●	◐	◑	◒	○	●
Extensibility	○	●	●	○	●	●	●
Trustworthiness	◐	◑	◒	◓	◔	◕	○

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

the optimized code is hard to maintain through hardware and software evolutions. Third, the 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?
- **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 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 at least two 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.

99 Halide [Ragan-Kelley et al. 2013] is an industrial-strength domain-specific compiler for image
100 processing, used e.g. to optimize code running in Photoshop and YouTube. Halide popularized the
101 idea of separating an *algorithm* describing what to compute from a *schedule* describing how to
102 optimize the computation. This separation makes it easy to try different schedules. TVM [Chen
103 et al. 2018] is a tool directly inspired by Halide, but tuned for machine learning applications; it
104 is used by most of the major CPU/GPU manufacturers. Other tools inspired by Halide include
105 Fireiron [Hagedorn et al. 2020a], used at Nvidia, as well as PartIR [Alabed et al. 2024], used at
106 Google. All these tools are inherently limited to the domains (DSLs) that they target. They do
107 not support higher-order composition of transformations, and are not extensible [Barham and
108 Isard 2019; Ragan-Kelley 2023]. Moreover, understanding their output is difficult as the applied
109 transformations are not detailed to the user, even though interactive scheduling systems have been
110 proposed to mitigate this difficulty [Ikarashi et al. 2021].

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

118 Exo [Ikarashi et al. 2022] is an imperative DSL embedded in Python, geared towards the de-
119 velopment of high-performance libraries for specialized hardware. The strength of Exo lies in
120 externalizing target-specific code generation to user-level code instead of compilation backends.
121 Exo programs can be optimized by applying a series of source-to-source transformations. These
122 transformations are described in a Python script, with simple string-based patterns for targeting
123 code points. The user can add custom transformations, possibly defined by composition; higher-
124 order composition seems possible but has not yet been demonstrated. A major limitation of Exo is
125 that it is restricted to static control programs with linear integer arithmetic.

126 Clay [Bagnères et al. 2016a] is a framework to assist in the optimization of loop nests that can be
127 described in the *polyhedral model* [Feautrier 1992]. The polyhedral model only covers a specific
128 class of loop transformations, with restriction over the code contained in the loop bodies, however
129 it has proved extremely powerful for optimizing code falling in that fragment. Clay provides a
130 decomposition of polyhedral optimizations as a sequence of basic transformations with integer
131 arguments. The corresponding transformation script can then be customized by the programmer.
132 Clint [Zinenko et al. 2018b] adds visual manipulation of polyhedral schedules through interactive
133 2D diagrams. LoopOpt [Chelini et al. 2021] provides an interactive interface that helps users design
134 optimization sequences (featuring unrolling, tiling, interchange, and reverse of iteration order) that
135 can be bound in a declarative fashion to loop nests satisfying specific patterns.

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

143 Alpinist [Sakar et al. 2022] is a *pragma*-based tool for optimizing GPU-level, array-based code. It
144 is able to apply basic transformations such as loop tiling, loop unrolling, data prefetching, matrix
145 linearization, and kernel fusion. The key characteristic of Alpinist is that it operates over code
146 formally verified using the VerCors framework [Blom et al. 2017]. Concretely, Alpinist transforms
147

148 not only the code but also its formal annotations. If Alpinist were to leverage transformation scripts
149 instead of pragmas, it might be possible to chain and compose transformations; yet, this possibility
150 remains to be demonstrated.

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

159 1.3 Contribution

160 This paper introduces OptiTrust, the first interactive optimization framework that operates at
161 the level of C syntax, and that supports and validates state-of-the-art optimizations. OptiTrust is
162 open-source and available at the URL: <https://github.com/charguer/optitrust>.
163

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

171 A central aspect of OptiTrust is that it guarantees that the code transformations requested by
172 the programmer preserve the semantics of the program. To that end, OptiTrust leverages our *static*
173 *resource analysis*, which concretely takes the form of a type checking algorithm in a type system
174 featuring linear resources. This type system may be thought of as a variant of the Rust type system,
175 or as a scaled down version of Separation Logic [Reynolds 2002].

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

184 The implementation of OptiTrust distinguishes between *basic* transformations and *combined*
185 transformations. On the one hand, a basic transformation applies minimalistic changes to the
186 abstract syntax tree (AST). The validity of a basic transformation is checked by leveraging the
187 resource analysis. On the other hand, a combined transformation is implemented as a composition
188 of basic transformations. Combined transformations aim to implement high-level strategies, that
189 may trigger the execution of dozens of basic transformation. These more complex combined trans-
190 formations need not be accompanied with code for checking validity: their validity is guaranteed
191 by the validity checks performed by the basic transformations. This two-layer approach enables us
192 to minimize the size of the trusted code base (TCB) of OptiTrust.

193 Internally, OptiTrust does not manipulate an abstract syntax tree (AST) for C syntax, but an
194 AST following a presentation closer to an imperative λ -calculus. We have developed a bidirectional
195

translation between the C language, which we parse using Clang, and the OptiTrust internal language. In particular, our translation is able to eliminate—and subsequently reintroduce—mutable variables and operations involving *l*-values. Considering a syntax and semantics simpler than that of C considerably helps to tame the complexity of the design and implementation of typing rules, code transformations, and correctness criteria associated with transformations.

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

- (1) We restrict ourselves to a subset of the C language. As our case studies show, this subset nevertheless suffices to express numerous practical, high-performance programs, in an idiomatic programming style both for the unoptimized and for the optimized code. Moreover, for simplicity, we ignore complications related to arithmetic overflows.
- (2) We have already implemented more than a hundred transformations, among the most standard ones. These transformations suffice to assess the interest of the OptiTrust approach to code optimization. However, for production usage, we believe there remains a couple hundreds additional transformations to implement.
- (3) We restrict ourselves to a subset of Separation Logic. Our resource-based type system is able to describe the ownership of arrays, matrices, or individual cells, however it does not allow specifying properties about the values stored in data structures. Nevertheless, as our case studies show, shape-based resources suffice to justify the correctness of many practical code optimization patterns.

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

In summary, we present a framework that can readily be exploited to optimize certain classes of programs, and acknowledge that future work remains necessary to achieve full generality. Note that we have taken great care in our design to anticipate for the extension 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 five parts. In Section 3, we describe the abstract syntax tree used internally by OptiTrust. In Section 4, we explain the core of a 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 our 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 justify the correctness of these transformations. Finally, we discuss additional related work in Section 7.

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, highly optimized computer vision library. In Section 2.2, we consider a physics simulation program featuring a kernel typical of particle

```

246 void rowSum(const int n, const int cn, const int kn, const T S[n+kn][cn], ST D[n][cn]) {
247     for (int c = 0; c < cn; c++) { // for each channel (e.g., red, green, and blue)
248         for (int i = 0; i < n; i++) { // for each target pixel in the row described by D
249             ST s = 0;
250             for (int j = i; j < i+kn; j++) // for each source pixel nearby to the right
251                 s += (ST) S[j][c];
252             D[i][c] = s;
253         } } }

```

Fig. 1. Unoptimized C code for the OpenCV case study, using multi-dimensional arrays.

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 Case Study

In image processing, *blur* is typically used to remove noise and smoothen images. 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 the *row-based blur* function, as implemented in the state-of-the-art OpenCV library.

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 $n+kn$ pixels. The parameter kn is a parameter of the transformation—the width of the blur. Each pixel is encoded using cn integers. Practical values include $cn=1$ for grayscale, $cn=3$ for RGB, $cn=4$ for RGBA. The input pixels are encoded on cn integers of type T , whereas the output pixels are encoded on cn integers of type ST . Typically, the type ST is represented on more bits than the type S . The output pixel $D[i]$ is computed as the sum of the values of the input pixels in the range from $S[i]$ to $S[i+kn-1]$. This sum is computed independently for each of the cn color channels.

Multi-dimensional vs flat arrays. The code from Fig. 1 is presented using C syntax for multi-dimensional arrays, for the sake of improved readability. However, the optimized code from Fig. 3 and our contract-annotated code from Fig. 2 instead use a flat array representation. The flat representation is idiomatic 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 multi-dimension arrays and their elimination via a source-to-source transformation.

Optimized code. The handwritten OpenCV library includes an implementation of row-sum blur structured like the code shows 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 our OptiTrust. The only notable difference between the OpenCV code and ours is that OpenCV traverses certain arrays by incrementing pointers, whereas we use explicit array indexing everywhere. In general, for this type of code, array indexing gives equivalent or better performance than pointer shifting; we leave

¹https://github.com/opencv/opencv/blob/4.10.0/modules/imgproc/src/box_filter.simd.hpp#L75: The OpenCV code is implemented as class with the types S and ST 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.


```

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

```

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.

it to future work to extend OptiTrust with a transformation able to introduce pointer shifting over a loop.

The optimized implementation is a *multi-versioned* code, with dedicated execution paths for handling specific values of the parameters. The branches `kn == 3` and `kn == 5` correspond to values of `kn`, the width that are commonly used by library users. For these small constant values of `kn`, the inner loop on `j` from Fig. 2 can be unfolded. Otherwise, when the loop on `j` is not unfolded, a standard algorithmic optimization called *sliding window* can be 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 When using Halide, a programmer either needs to manually refine the code to introduce the sliding window before optimizing the code; or needs to exploit other transformation tools specialized for applying sliding window optimizations [Chaurasia et al. 2015; Kanetaka et al. 2024].

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` can be unfolded, then the multiple occurrences of the loop on `i` that result from this unfolding can be fused into a single loop. The final else-branch in the code from Fig. 2 corresponds to the generic implementation.

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*, in

²<https://github.com/halide/Halide/issues/180>

```

344 void rowSum(const int kn, const T* S, ST* D, const int n, const int cn) {
345     __requires("kn >= 0, n >= 1, cn >= 0");
346     __reads("S ~> Matrix2(n+kn, cn)");
347     __modifies("D ~> Matrix2(n, cn)");
348     __ghost_begin("unswap1", swap_groups, "items := fun i, c -> &D[i][c] ~> Cell");
349     for (int c = 0; c < cn; c++) { // for each channel
350         __xmodifies("for i in 0..n -> &D[MINDEX2(n, cn, i, c)] ~> Cell");
351         for (int i = 0; i < n; i++) { // for each pixel
352             __xmodifies("&D[MINDEX2(n, cn, i, c)] ~> Cell");
353             __ghost(assume, "is_subrange(i..i + kn, 0..n + kn)");
354             ST s = 0;
355             for (int j = i; j < i+kn; j++) {
356                 s += (ST) S[MINDEX2(n+kn, cn, j, c)];
357             }
358             D[MINDEX2(n, cn, i, c)] = s;
359         }
360     }
361     __ghost_end("unswap1");
362 }

```

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

order to guide OptiTrust’s resource-based typechecker. Fig. 3 shows the same C code as in Fig. 1 augmented with a number of no-ops that provide the relevant annotations.³

The clause `__requires` contains assumptions about the input parameters. The clause `__reads` asserts that input array `S` can be accessed in read-only mode. The clause `__modifies` asserts that output array `D` can be modified in place. The syntax `__modifies("D")` is a shorthand for `__modifies("D ~> Matrix2(n, cn)")`. The latter explicit form would be useful to provide the dimensions if the argument of the function was declared as `ST* D` instead of `ST D[n][cn]`, as it is the case for example in the optimized code. The clause `__xmodifies` describes a *loop contract*: it indicates not only that the i -th iteration can modify the cell `D[i][c]`, but also that the other iterations do not access this cell. In other words, the i -th iteration has exclusive access to that cell. The “x” prefix in `__xmodifies` stands for “exclusive”.

The lines introduced by `__ghost_begin`, `__ghost_end`, or sometimes just `__ghost` correspond to *ghost instructions*: no-ops whose purpose is to change the view on the resources. The need for ghost instructions is standard in Separation Logic frameworks. The specialized keywords `__ghost_begin` and `__ghost_end` materialize a pair of ghost instructions that are the reciprocal of one another. The `__ghost_begin` pseudo-instruction in Fig. 3 refines the view on the array `D` from $\star_{i \in 0..n} \star_{c \in 0..cn} \&D[i][c] \rightsquigarrow \text{Cell}$ to $\star_{c \in 0..cn} \star_{i \in 0..n} \&D[i][c] \rightsquigarrow \text{Cell}$, where the star symbol corresponds to the *iterated separating conjunction* operator of Separation Logic. Intuitively, it swaps two “logical loops” that iterate over the sets of memory cells covered by the array `D`. The matching `__ghost_end` pseudo-instruction applies the symmetrical operation, reverting the view on the resources back to the original form.

Optimization script. Fig. 4 shows our optimization script for generating the optimized code of Fig. 2 starting from the annotated unoptimized code of Fig. 3. In OptiTrust, optimizations are dictated by means of a script written in the OCaml programming language. For the reader not familiar with OCaml, $f \ x \ y$ denote the call of `f` on the arguments `x` and `y`; the symbol `~` is used

³Alternatively, one could use *attributes* to attach certain annotations to language constructs, e.g., contracts to loops and functions. We believe, however, that representing ghost instructions as no-op instructions is more natural, because a ghost instruction is part of a sequence without being attached to either the previous or the subsequent instruction.


```

393 Specialize.variable_multi ~mark_then:fst ~mark_else:"nokn"
394   ["kn", int 3; "kn", int 5] [cFunBody "rowSum"; cFor "c"];
395 Loop.unroll [nbMulti; cMark "kn"; cFor "c"];
396 Instr.gather_targets [nbMulti; cMark "kn"; cArrayWrite "D"];
397 Loop.sliding_window [cMark "nokn"; cFor "j"];
398 Loop.swap [nbMulti; cMark "nokn"; cFor "c"];
399 Specialize.variable_multi ~mark_then ~mark_else:"nocn"
400   ["cn" int 1; "cn", int 3; "cn", int 4] [cMark "nokn"];
401 Loop.unroll [nbMulti; cMark "cn"; cFor "c"];
402 Loop.collapse [nbMulti; cMark "cn"; cFor "i"];
403 Target.foreach [nbMulti; cMark "cn"] (fun c ->
404   Loop.fusion_targets ~into:FuseIntoLast [nbMulti; c; cFor "i" ~body:[cArrayWrite "D"]];
405   Instr.gather_targets [c; cStrict; cArrayWrite "D"];
406   Loop.fusion_targets ~into:FuseIntoLast [nbMulti; c; cFor ~stop:[cVar "kn"] "i"];
407   Instr.gather_targets [c; cFor "i"; cArrayWrite "D"];
408 );
409 Cleanup.std [];

```

Fig. 4. Optimization script for the OpenCV case study. [TODO: finalize]

to provide optional (or named) arguments; $[x; y; z]$ denotes a list; (x, y, z) denotes a tuple; $s_1 \wedge s_2$ denotes a string concatenation; and `let f x = e1 in e2` introduces a local function f .

A transformation script generally consists of a series of calls to functions from the OptiTrust library. By convention, the last argument of a transformation always denotes a *target*. A target provides a way to concisely and robustly refer to one or 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.

Targets may also be given as arguments to constraints: for example, `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 . Targets may also include special modifiers. The modifier `nbMulti` in a target indicates that the programmer expects to find not one but multiple AST nodes that match this target. For example, the modifier `tBefore`, which appears in the other two case studies, allows targeting the interstice before an instruction. Finally, the special command `Target.set_ctx` registers a partial target to be used as prefix for all subsequent targets. This command enables the subsequent instructions of the script to focus on a particular piece of the program.

[TODO: comment a little bit more on the contents of the script when it is finalized.]

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. All intermediate versions of the code consist of human-readable, executable C code. [TODO: insert an example diff, for example showing loop fusion.] 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. We encourage the reader to navigate the traces associated with our case studies.⁴

⁴The traces associated with our case studies may be found at: <https://www.chargueraud.org/softs/optitrust/traces/index.html>

```

442 void simulate(const vect* fieldAtCorners,
443 const int nbSteps, const double deltaT,
444 const double pCharge, const double pMass,
445 const int nbParticles, particle* particles) {
446   __reads("fieldAtCorners ~> Array(8)");
447   __modifies("particles ~> Array(nbParticles)");
448   for (int idStep = 0; idStep < nbSteps; idStep++) {
449     for (int idPart = 0; idPart < nbParticles; idPart++) {
450       // Each particle is updated at each time step
451       __xmodifies("&particles[idPart] ~> Cell");
452       const particle p = particles[idPart];
453       // Interpolate the field based on the position relative to the corners of the cell
454       double const coeffs[8];
455       compute_corner_interpolation_coeffs(p.pos, coeffs);
456       const vect fieldAtPos = matrix_vect_mul(coeffs, fieldAtCorners);
457       // Compute the acceleration: F = m*a and F = q*E gives a = q/m*E
458       const vect accel = vect_mul(pCharge / pMass, fieldAtPos);
459       // Compute the new speed and position for the particle
460       const vect speed2 = vect_add(p.speed, vect_mul(deltaT, accel));
461       const vect pos2 = vect_add(p.pos, vect_mul(deltaT, speed2));
462       // Update the particle
463       particles[idPart].speed = speed2;
464       particles[idPart].pos = pos2;
465     } } }

```

Fig. 5. Unoptimized code for the particle simulation case study, with resource annotations.

Validity Checks. The evaluation of the transformation script from Fig. 4 triggers the application of a chain of basic transformations, which modify the AST of the program. For each such basic transformation, OptiTrust checks that the 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. All these aspects associated with correctness criteria and preservation of typing will be discussed 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.⁵ In the present case study, we consider a simplified PIC simulation, focusing on the computations associated with one particular cell. Our goal is to illustrate how OptiTrust can be used to derive a certain number of transformations idiomatic of particle simulations—and also ubiquitous in other physics simulation code.

Unoptimized code. Fig. 5 presents the implementation of a kernel of a particle simulation. It is a simplified version of a realistic simulation code. A number of particles, all with the same mass and charge, move inside a large cube. We assume that the particles do not leave the cube during the simulation. The initial position and speed of every particle is given.

For simplicity, we here assume that the particles do not affect each other—we leave it to future work to optimize the code for the *charge deposit* as in the aforementioned publication. Instead, we

⁵Our case study is inspired by the work by Barsamian et al. [2018], who present a PIC implementation featuring state-of-the-art optimizations, targeting the Intel Xeon Platinum 8160 architecture.

```

491 void simulate(const vect* fieldAtCorners,
492             const int nbSteps, const double deltaT,
493             const double pCharge, const double pMass,
494             const int nbParticles, particle* particles) {
495     const double fieldFactor = deltaT * deltaT * pCharge / pMass;
496     vect* const lFieldAtCorners = (vect*) malloc(nbCorners * sizeof(vect));
497     for (int i = 0; i < nbCorners; i++) {
498         lFieldAtCorners[i].x = fieldAtCorners[i].x * fieldFactor;
499         lFieldAtCorners[i].y = fieldAtCorners[i].y * fieldFactor;
500         lFieldAtCorners[i].z = fieldAtCorners[i].z * fieldFactor;
501     }
502     for (int i = 0; i < nbParticles; i++) {
503         particles[i].speed.x *= deltaT;
504         particles[i].speed.y *= deltaT;
505         particles[i].speed.z *= deltaT;
506     }
507     double const coeffs[8];
508     for (int idStep = 0; idStep < nbSteps; idStep++) {
509         for (int idPart = 0; idPart < nbParticles; idPart++) {
510             compute_corner_interpolation_coeffs(particles[idPart].pos, coeffs);
511             double fieldAtPosX = 0.;
512             double fieldAtPosY = 0.;
513             double fieldAtPosZ = 0.;
514             for (int k = 0; k < nbCorners; k++) {
515                 fieldAtPosX += coeffs[k] * lFieldAtCorners[k].x;
516                 fieldAtPosY += coeffs[k] * lFieldAtCorners[k].y;
517                 fieldAtPosZ += coeffs[k] * lFieldAtCorners[k].z;
518             }
519             const double speed2X = particles[idPart].speed.x + fieldAtPosX;
520             const double speed2Y = particles[idPart].speed.y + fieldAtPosY;
521             const double speed2Z = particles[idPart].speed.z + fieldAtPosZ;
522             particles[idPart].pos.x += speed2X;
523             particles[idPart].pos.y += speed2Y;
524             particles[idPart].pos.z += speed2Z;
525             particles[idPart].speed.x = speed2X;
526             particles[idPart].speed.y = speed2Y;
527             particles[idPart].speed.z = speed2Z;
528         }
529     }
530     for (int i = 0; i < nbParticles; i++) {
531         particles[i].speed.x /= deltaT;
532         particles[i].speed.y /= deltaT;
533         particles[i].speed.z /= deltaT;
534     }
535     free(lFieldAtCorners);
536 }

```

Fig. 6. Optimized code for the particle simulation case study.

assume an external electric field that affects the acceleration of the particles. This external electric field, is described by 8 vectors, which correspond to the field at the corners of the cube. The electric field that applies at a given position inside the cube is obtained by means of a linear interpolation with respect to the fields at 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. 5, 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`, not shown, computes the interpolation coefficients associated with the position of the particle.

```

540 Target.set_prefix [cFunBody "simulate"];
541 Loop.move_out [cVarDef "coeffs"];
542 Function.inline [multi cFun ["matrix_vect_mul"; "vect_add"; "vect_mul"]];
543 Variable.def "fieldFactor" (expr "deltaT * deltaT * pCharge / pMass")
544   [tBefore; cVarDef "lFieldAtCorners"];
545 Record.set_explicit [multi cWrite ~typ:["vect"; "particle"]];
546 Record.to_variables [cVarDefs ["fieldAtPos"; "pos2"; "speed2"; "accel"]];
547 Variable.inline [cVarDef "p"];
548 Accesses.scale_inplace ~factor:(var "fieldFactor") [nbMulti; cVarRe "fieldAtPos[XYZ]"];
549 Accesses.scale_inplace ~factor:(var "deltaT")
550   [nbMulti; sExprRe "particles\\.speed\\.\\.xyz"];
551 Accesses.scale_incopy ~into:"lFieldAtCorners" ~factor:(var "deltaT")
552   [nbMulti; cOr [cVarRe "speed2[XYZ]"]; ];
553 Accesses.scale_immut ~factor:(var "deltaT") [nbMulti; cVarRe "speed2[XYZ]"];
554 Variable.inline [cVarDefsRe ["accel[XYZ]"; "pos2[XYZ]"]];
555 Arith.(simpls_rec [expand; gather_rec]) [];
556 Cleanup.std [];

```

Fig. 7. Optimization script for the particle simulation case study [TODO: finalize].

557 *Optimized code.* Fig. 6 shows our optimized code for the function `simulate`. Three key optimizations
558 are applied. First, the vector and matrix operations are replaced with operations over individual fields
559 (named `pos.x`, `pos.y`, `pos.z`, `speed.x`, `speed.y`, and `speed.z`). Moreover, local vector variables are replaced
560 with families of variables (e.g., `fieldAtPosX`, `fieldAtPosY`, and `fieldAtPosZ`). Second, the array `coeffs`,
561 that was used to store the interpolation coefficients computed by `compute_corner_interpolation_coeffs`,
562 is now allocated once and for all outside the loop. Third, and most importantly, a *scaling* trans-
563 formation is applied on the data in order to simplify the arithmetic computations that need to be
564 performed at every time step. We next explain how this standard optimization technique works,
565 focusing on the x-coordinate.

566 Consider a particle. At a given time step, its speed, written v , its the position, written x , are
567 updated, according to the formulae: $a=qE/m$ and $v += a\Delta_t$ and $x += v\Delta_t$ where E denotes the electric
568 field interpolated at the location of this particle. The constants q , m , and Δ_t corresponds to the
569 program variables `pCharge`, `pMass`, and `deltaT`, respectively. The idea of scaling transformation is to
570 store not the values of E and v , but instead the values E' and v' defined as: $E' = qE\Delta_t^2/m$ and $v' = \Delta_tv$.
571 The interest is that the speed and position updates at a given time step are now described using
572 much simpler formulae: $v' += E'$ and $x' += v'$. These formulae involve nothing more than simple
573 additions. To implement this scaling transformation, the components of the field `speed` of the array
574 `particles` are multiplied, in-place, by a factor Δ_t before starting the simulation; symmetrically, at the
575 end of the simulation, the values are divided by Δ_t . For the electric field array, which is read-only,
576 the scaling factor is applied on an auxiliary array named `lFieldAtCorners`, obtained by multiplying
577 the values of `fieldAtCorners` by $q\Delta_t^2/m$. (An in-place update would also have been possible.) By
578 linearity of the interpolation computations, this scaling propagates to the values computed for the
579 electric field at the particle location (`fieldAtPosX`, `fieldAtPosY`, and `fieldAtPosZ`).

581 *Optimization script.* Fig. 7 shows our optimization script. The operator `multi` is a shorthand for de-
582 scribing a union of targets. For example, `multi cFun ["matrix_vect_mul"; "vect_add"; "vect_mul"]`
583 is equivalent to `nbMulti; cOr [[cFun "matrix_vect_mul"]; [cFun "vect_add"]; [cFun "vect_mul"]]`,
584 where `cOr` is the target disjunction combinator. The constraint such as `cVarRe` match variable
585 names using regular expressions. The constraint `sExprRe` matches C expressions according to their
586 string representation, again using regular expressions. The crux of the script is the evaluation of
587 `Arith.(simpls_rec [expand; gather_rec])` in the script. This transformation performs the arithmetic

```

589 void mm(const int m, const int n, const int p,
590         const float A[m][p], const float B[p][n], float C[m][n]) {
591     for (int i = 0; i < m; i++) {
592         for (int j = 0; j < n; j++) {
593             float sum = 0.0f;
594             for (int k = 0; k < p; k++)
595                 sum += A[i][k] * B[k][j];
596             C[i][j] = sum;
597         } } }
598 void mm1024(const float A[1024][1024], const float B[1024][1024], float C[1024][1024]) {
599     mm(1024, 1024, 1024, A, B, C);
600 }

```

Fig. 8. Unoptimized C code for the matrix-multiply case study, using multi-dimensional arrays, and the specialization to input size 1024.

simplifications that make all the scaling factors cancel out, leaving exactly the trivial additions for updating the speed and the position values of each particle.

[TODO: comment once finalized]

Benefits of using OptiTrust. Applied mathematicians commonly write optimized code such as that of Fig. 6 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.

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 code. The TVM tutorial presents an optimization script⁶ (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.⁷ 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, using a transformation script reasonably short compared with the TVM script, we are able to generate code that features the exact same optimization patterns as in the TVM case study.

Unoptimized code. Fig. 8 shows a naive implementation of matrix-multiply. Like for the first case study, to help the reader we use multi-dimensional arrays as opposed to the flat arrays visible in our contract-annotated code (Fig. 10).

Optimized code. TVM output code is expressed not as C code, but directly in the intermediate representation of LLVM, known as LLVM IR. Hence, our first task has been to manually parse the LLVM IR output associated with TVM’s matrix-multiply case study, and to infer the corresponding C code. This C code is essentially equivalent to that shown in Fig. 9, which contains the code that we produce using OptiTrust.

⁶https://tvm.apache.org/docs/how_to/optimize_operators/opt_gemm.html

⁷The 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 whereas the MKL provides optimized implementation for a fixed set of functions, 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.

```

638 void mm1024(const float* A, const float* B, float* C) {
639     float* pB = (float*)malloc(sizeof(float)[32][256][4][32]);
640     #pragma omp parallel for
641     for (int bj = 0; bj < 32; bj++) {
642         for (int bk = 0; bk < 256; bk++) {
643             for (int k = 0; k < 4; k++) {
644                 for (int j = 0; j < 32; j++) {
645                     pB[32768 * bj + 128 * bk + 32 * k + j] = B[1024 * (4 * bk + k) + 32 * bj + j]; }}}
646     #pragma omp parallel for
647     for (int bi = 0; bi < 32; bi++) {
648         for (int bj = 0; bj < 32; bj++) {
649             float* sum = (float*)malloc(sizeof(float)[32][32]);
650             for (int i = 0; i < 32; i++) {
651                 for (int j = 0; j < 32; j++) {
652                     sum[32 * i + j] = 0.; }}
653             for (int bk = 0; bk < 256; bk++) {
654                 for (int i = 0; i < 32; i++) {
655                     float s[32];
656                     memcpy(s, &sum[32 * i], sizeof(float)[32]);
657                     #pragma omp simd
658                     for (int j = 0; j < 32; j++) {
659                         s[j] += A[1024 * (32 * bi + i) + 4 * bk] * pB[32768 * bj + 128 * bk + j]; }
660                     #pragma omp simd
661                     for (int j = 0; j < 32; j++) {
662                         s[j] += A[1 + 1024 * (32 * bi + i) + 4 * bk] * pB[32 + 32768 * bj + 128 * bk + j]; }
663                     #pragma omp simd
664                     for (int j = 0; j < 32; j++) {
665                         s[j] += A[2 + 1024 * (32 * bi + i) + 4 * bk] * pB[64 + 32768 * bj + 128 * bk + j]; }
666                     #pragma omp simd
667                     for (int j = 0; j < 32; j++) {
668                         s[j] += A[3 + 1024 * (32 * bi + i) + 4 * bk] * pB[96 + 32768 * bj + 128 * bk + j]; }
669                     memcpy(&sum[32 * i], s, sizeof(float)[32]); }}
670             for (int i = 0; i < 32; i++) {
671                 for (int j = 0; j < 32; j++) {
672                     C[1024 * (32*bi + i) + 32*bj + j] = sum[32*i + j]; }}
673             free(sum);
674         } }
675     free(pB);
676 }

```

Fig. 9. Our optimized C code for the matrix-multiply case study. This code features the same optimization patterns and achieves similar performance as the reference output of TVM.

Remark: even though the indices for array accesses involved in this optimized code from Fig. 9 contain somewhat obfuscated arithmetic formulae, these accesses can be displayed until the last optimization step of the OptiTrust script, in the form of multi-dimensional accesses, e.g. `pB[MINDEX4(32, 256, 32, 4, bj, bk, 3, j)]`. Moreover, for the purpose of readability of generated programs, OptiTrust offers the option to print multi-dimensional the same array access in the shorter form `pB[bj; bk; 3; j]`. This form, purposely out of the syntax of C, remains non-ambiguous because the size information appears in the description of the resources at hand.

Compared with the naive code from Fig. 8, the optimized code from Fig. 9 integrates numerous optimizations.


```

687 void mm(float* C, float* A, float* B, int m, int n, int p) { // naive matrix-multiply
688   __reads("A ~> Matrix2(m, p), B ~> Matrix2(p, n)");
689   __modifies("C ~> Matrix2(m, n)");
690   for (int i = 0; i < m; i++) {
691     __xmodifies("for j in 0..n -> &C[MINDEX2(m, n, i, j)] ~> Cell");
692     for (int j = 0; j < n; j++) {
693       __xmodifies("&C[MINDEX2(m, n, i, j)] ~> Cell");
694       float sum = 0.0f;
695       for (int k = 0; k < p; k++)
696         sum += A[MINDEX2(m, p, i, k)] * B[MINDEX2(p, n, k, j)];
697       C[MINDEX2(m, n, i, j)] = sum;
698     }
699   }
700 void mm1024(float* C, float* A, float* B) { // specialization to 1024x1024 matrices
701   __reads("A ~> Matrix2(1024, 1024), B ~> Matrix2(1024, 1024)");
702   __modifies("C ~> Matrix2(1024, 1024)");
703   mm(C, A, B, 1024, 1024, 1024);
704 }

```

Fig. 10. Unoptimized C code for the matrix-multiply case study, using flat arrays and augmented with resource annotations.

- The body of the generic matrix multiply function `mm` is specialized to the size 1024.
- 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.
- 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.
- Results are not accumulated directly in the output matrix `c`. Instead, the results are first accumulated into a stack-allocated array named `s` of size 32. This array `s` can be mapped onto a few 256-bit registers.
- The results accumulated in `s` are then transferred into another stack-allocated array, named `sum`, which has the same size as a block but uses a different memory layout. For each row of the block being processed, two `memcpy` operations are used for synchronizing `s` with `sum`. After processing a block in full, the values from the array `sum` are copied into the output array `c`.
- 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.
- The 4 loops tagged as `#pragma omp simd` in Fig. 9 are very similar. However, if we attempt to factorize them into a loop with 4 iterations, then Intel's compiler (ICX) produces slower code. Hence, the loops need to be unfolded as shown.

Please keep in mind that we have not chosen this particular set of optimizations. We have simply inferred these optimizations from the output of the TVM case study. Our goal is to reproduce the exact same set of optimizations.

Annotated Unoptimized Code. Fig. 10 shows the annotations that we need to insert into the naive matrix-multiply implementation from Fig. 8 for typechecking the resources involved. The syntax of annotations has been explained earlier, except for the iteration construct, which we explain next.

```

736 Function.inline_def [cFunDef "mm"];
737 let tile (id, tile_size) =
738   Loop.tile (int tile_size) ~index:("b" ^ id) ~bound:TileDivides [cFor id] in
739   List.iter tile [("i", 32); ("j", 32); ("k", 4)];
740   Loop.reorder_at ~order:["bi"; "bj"; "bk"; "i"; "k"; "j"] [cPlusEq ()];
741   Loop.hoist_expr ~dest:[tBefore; cFor "bi" "pB" ~indep:["bi"; "i"] [cArrayRead "B"];
742   Matrix.stack_copy ~var:"sum" ~copy_var:"s" ~copy_dims:1 [cFor ~body:[cPlusEq ()] "k"];
743   Omp.simd [cFor ~body:[cPlusEq ()] "j"];
744   Omp.parallel_for [cFunBody "mm1024"; cStrict; cFor ""];
745   Loop.unroll [cFor ~body:[cPlusEq ()] "k"];

```

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

```

746
747 k = tvn.reduce_axis((0, P))
748 A = tvn.placeholder((M, P))
749 B = tvn.placeholder((P, N))
750
751 pB = tvn.compute((N / 32, P, 32),
752   lambda bj, k, j:
753     B[k, bj * 32 + j])
754
755 C = tvn.compute((M, N),
756   lambda i, j:
757     sum(A[i, k] * pB[j // 32, k, j % 32],
758       axis=k))
759
760 CC = s.cache_write(C, "global")
761 bi, bj, i, j = s[C].tile(
762   C.op.axis[0], C.op.axis[1], 32, 32)
763 s[CC].compute_at(s[C], bj)
764 i2, j2 = s[CC].op.axis
765 (kaxis,) = s[CC].op.reduce_axis
766 bk, k = s[CC].split(kaxis, factor=4)
767 s[CC].reorder(bk, i2, k, j2)
768 s[CC].vectorize(j2)
769 s[CC].unroll(k)
770 s[C].parallel(bi)
771 bj3, _, j3 = s[pB].op.axis
772 s[pB].vectorize(j3)
773 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) to the right. Both use Python syntax.

The outer loop on i is annotated with a clause of the form `__xmodifies("for j in 0..n -> Hj")`, where the logical `for` construct denotes an iteration over resources. This clause indicates that the i -th iteration requires exclusive access to the union of the resources described by H_j , for all j in the range $0..n$. The resource $C \rightsquigarrow \text{Matrix2}(m, n)$ is equivalent to `"for i in 0..m -> for j in 0..n -> &C[i][j] ~> Cell"`, which covers all the $n \times m$ cells of the matrix C . The resource `"for j in 0..n -> &C[i][j] ~> Cell"` covers only the i -th row of that matrix. Further in the paper, this same resource will be written using LaTeX-style notation as: $\star_{j \in 0..n} (\&C[i][j] \rightsquigarrow \text{Cell})$.

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. These internal basic transformations may be visualized in the execution trace generated by OptiTrust.⁴ 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. 8 into the 2D-array named `sum` in Fig. 9.

Comparison Against TVM. The TVM matrix-multiply case study appears in Fig. 12. We only comment on specific aspects and refer to TVM’s tutorial⁶ 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 be written more simply as:

```

783 k = tvn.reduce_axis((0, P))
784

```

```

785 A = tvm.placeholder((M, P))
786 B = tvm.placeholder((P, N))
787 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 ρB , as a code transformation. The programmer therefore needs to introduce this auxiliary structure manually in the input source code. Likewise, the blocking strategy needs to be hardwired in the source code. In contrast, our input program for matrix multiply (Fig. 8) is expressed using standard C code, and needs not include any optimization. Starting from a totally unoptimized reference code improves readability, trustworthiness, and maintainability.

Besides, note that our although our source code for matrix-multiply is currently expressed using explicit loops, in the future we could alternatively express it using higher-order combinators similar to those available in TVM (using either C++ syntax or the syntax of another language).

The right-hand side of Fig. 12 shows TVM’s optimization script. Our optimization script for that case study (Fig. 11) is not much longer than that of TVM. As mentioned earlier, we have carefully checked that the C code produced using OptiTrust features the same optimizations and delivers the same performance 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.

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 *diff*s for every transformation step—as well as for every substep—involved in the optimization process.

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 from Section 1.1.

Generality. As pointed out in Section 1.3, this first release of OptiTrust has a number of limitations: it applies to a subset of the language, for a simplified version of Separation Logic, and there remains a fair number of 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, may appear reasonably straightforward. However, by chaining such transformations in the desired manner, the OptiTrust user is able to achieve state-of-the-art high-performance code, which resembles code that an expert might have written by hand. Moreover, the many basic transformations involved need not be explicitly invoked by the user: the use of high-level combined transformations allows to achieve expressiveness via concise scripts—recall, e.g., the call to `Loop.reorder` in the matrix-multiply case study. A key feature of OptiTrust is that, at any stage in the optimization process, the user remains fully in control.

Expressiveness also depends on the generality of the correctness criteria associated with every transformation. In practice, there could be situations where the user may want to legitimately apply a basic transformation, yet OptiTrust’s implementation is unable to recognize this application as correct. In the short term, one option is for the user to treat this particular step as “user-trusted”, and to rely on human review of the *diff* associated with that step.⁸

⁸In the long term, by leveraging a full-blown Separation Logic, possibly combined with the use of interactive proofs, users could replace a piece of code with any other piece of code that provably satisfies the same specification.

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

842
843 *Composability.* OptiTrust transformation scripts are expressed as OCaml programs, and each
844 transformation from our library consists of an OCaml function. Because OCaml is a full-featured
845 programming language, OptiTrust users may define additional transformations at will by combining
846 existing transformations. User-defined transformations may query the abstract syntax tree (AST)
847 that describes the C code, allowing to perform analyses before deciding what transformations
848 to apply. Furthermore, because OCaml is a higher-order programming language, transformation
849 can take other transformations as argument. We use this programming pattern for example to
850 customize the arithmetic simplifications to be performed after certain transformations.

851 *Extensibility.* If the user needs a transformation that is not expressible as a combination of
852 transformations from the OptiTrust library, a custom transformation can be devised. Because
853 OptiTrust does not rely on heuristics, adding a new transformation to OptiTrust does not impact
854 in any way the behavior of existing scripts. To define relatively simple custom transformations,
855 OptiTrust provides a term-rewriting facility based on pattern matching. For more complicated
856 transformations, one can follow the patterns employed in the OptiTrust’s library. For all custom
857 transformations, it is the programmer’s responsibility to work out the criteria under which applying
858 the transformation preserves the semantics of the code, and to adapt contracts if necessary in order
859 to produce well-typed code.

860
861 *Trustworthiness.* Compilers are well-known to be incredibly hard to get 100% correct [Yang et al.
862 2011]. Like compilers, optimization tools are highly subject to bugs. OptiTrust mitigates the risks
863 of producing incorrect code in two ways. First, the *diff* of every step can be thoroughly scrutinized.
864 Secondly, as explained in Section 1.3, we have organized the OptiTrust code base so as to isolate the
865 implementation of the *basic* transformations, which consists of transformations that directly modify
866 the AST. Only basic transformations need to be trusted. We have been careful to systematically
867 minimize the complexity of the interface and of the implementation of our basic transformations.
868 All other transformations—the *combined* transformations—are *not* part of the trusted computing
869 base (TCB).

870
871 This completes our high-level presentation of the OptiTrust framework. The remaining sections
872 present the implementation of OptiTrust: its internal AST, its typechecking algorithm, and its
873 transformations.

874 3 OPTITRUST’S INTERNAL AST

875
876 In OptiTrust, input C programs are encoded into an imperative λ -calculus. All code transformations
877 are performed on that imperative λ -calculus. Then, programs are decoded back into C syntax. This
878 approach enables OptiTrust to report the *diff* associated with a transformation in terms of a concise
879 syntax familiar to the programmer.

880 The OptiTrust abstract syntax tree (AST) is represented as an immutable tree data structure. A
881 program transformation takes as input such an immutable AST, and produces as output another
882

883	$r :=$	range ($t_{\text{start}}, t_{\text{stop}}, t_{\text{step}}$)	range for simple for-loops
884	$\pi :=$	par \cdot	optional parallel flag on simple for-loops
885	$t :=$	x res	variables, and the special variable res
886		b n	boolean values, and number values
887		$(t_1; \dots; t_n)$	sequence
888		let $x = t$	variable definition
889		fun (a_1, \dots, a_n) $\mapsto t$	function definition
890		$t_0(t_1, \dots, t_n)$	function call
891		for $^\pi(i \in r) t$	possibly parallel, simple for-loop
892		if t_0 then t_1 else t_2	conditional
893		$\{f_1 = t_1; \dots; f_n = t_n\}$ $[t_1; \dots; t_n]$	structure and array as values
894		$t_1[t_2]$ $t.f$	projection from array/struct values
895		$t_1 \boxplus t_2$ $t \boxminus f$	address computation
896		$t_1 \boxplus t_2$ $t \boxminus f$	address computation

Fig. 13. Grammar of OptiTrust’s internal λ -calculus. Type annotations on binders and operators are not shown.

AST, which may share subtrees with the input AST. There are two major benefits to following a purely functional programming style using immutable trees. First, this approach avoids numerous bugs typically associated with inadvertent sharing of subtrees when modifying data structures in-place. Second, this approach, by enabling sharing, improves the efficiency of the construction of a complete execution trace for reporting to the user all the intermediate ASTs constructed during the evaluation of the user’s transformation script.

The purpose of this section is to present OptiTrust’s internal language, whose constructs appear throughout the rest of the paper, from the statement of typing rules in Section 4 and usage analysis in Section 5, to the description of transformations in Section 6.

In Section 3.1, we present the grammar of our imperative λ -calculus. It is mostly standard, except perhaps for the special variable **res** used to bind return values, for the two primitive operations to implement stack-allocation, for the construct for (possibly parallel) simple for-loops, and for the operators that compute address arithmetic. These ingredients are not ubiquitous in presentations of λ -calculus, yet they have been employed by others in prior work.

In Section 3.2, we describe, at a high-level, the intuition behind OptiTrust’s translation from C into our internal λ -calculus. Again, such a translation is relatively standard: numerous compilers include a phase that eliminates mutable variables and l -values. The specificity of our translation is that it introduces a few annotations on terms to allow implementing the reciprocal translation.

3.1 OptiTrust’s Internal AST

Fig. 13 gives the grammar of OptiTrust’s AST. 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 equip OptiTrust’s language with a standard call-by-value semantics (not shown). We next describe the grammar, starting with the less common features.

Sequences in OptiTrust. A sequence is a term that consists of a list of subterms or let-bindings, to be executed in order. We impose in the OptiTrust AST the invariant that every function body consists of a sequence block, even if the sequence contains a single instruction. Likewise, loop body and branches of conditionals systematically consist of sequences.

932 *The Special Result Variable.* A special variable, named **res**, is used to denote the result value of a
 933 sequence, typically the body of a function. A statement **return** t that appears in terminal position
 934 of a C function is translated as a binding “**let res = t**”. More generally, in OptiTrust, all sequences
 935 with a non-void type are expected to contain exactly one **let res** construct.

936 The variable **res** also typically appears in function contracts, to specify the return value of a
 937 function. The use of a dedicated name such as **res** is common practice in program verification
 938 tools, such as ESC/Java [Flanagan et al. 2002] or Why3 [Filliâtre 2003]. Besides, treating **return**
 939 as an assignment instruction is a technique that appears, e.g., in the Viper program verification
 940 tool [Müller et al. 2017]. A key interest of the “**let res**” construct is that it allows placing instructions
 941 *after* the point at which the return value is computed. Doing so is specifically useful for *ghost*
 942 *instructions* that must appear at the very end of a sequence. Indeed, the purpose of ghost instructions
 943 is to refine the formulation of certain resources, and these resources may depend on the return
 944 value.

945
 946 *Manipulation of Heap and Stack Cells.* To account for heap-allocated data, OptiTrust provides
 947 the following standard primitive functions: **alloc** for allocating one or several uninitialized cells,
 948 **get** for reading a cell, **set** for writing a cell, and **free** for freeing allocated cells. As usual, a read in
 949 an uninitialized memory cell is undefined behavior. Additionally, to account for stack-allocated
 950 variables, OptiTrust includes two special functions. The operation $\text{ref}(t)$ allocates a memory cell
 951 initialized with a specific content t , and whose space is automatically reclaimed at the end of the
 952 surrounding sequence. The operation $\text{ref_uninit}()$ is similar, but it allocates a memory cell without
 953 initializing its contents. These two special operations are meant to occur as part of a let-binding, for
 954 example **let** $x = \text{ref}(3)$, occurring directly within a sequence. The two stack-allocation operators,
 955 apart from their implicit-free behavior, are treated like other primitive functions. Remark: we have
 956 considered the possibility of encoding **ref** and **ref_uninit** using **alloc**, **set** and **free**, but ultimately
 957 concluded that this approach was slightly less practical, leading to additional complexity in code
 958 transformations and in target resolution.

959
 960 *Possibly Parallel, Simple For Loops.* The construct $\text{for}^\pi(i \in \text{range}(t_{\text{start}}, t_{\text{stop}}, t_{\text{step}})) t_{\text{body}}$ describes
 961 a *simple-for-loop*. In such a loop, the immutable variable i denotes the loop index, and the loop
 962 range, which consists of the loop bounds and the per-iteration step are evaluated only once before
 963 starting the loop. Following the convention used by Python and other languages, the index goes
 964 from the *start* value inclusive to the *stop* value exclusive. If the *step* value is negative, the loop index
 965 iterates downwards. Optionally, the loop may be tagged with a *parallel* flag (i.e., setting π to **par**),
 966 thereby asserting that the loop should be treated as a parallel loop by the compiler and the runtime.
 967 This flag corresponds to the directive: **#pragma openmp parallel**. The restrictions imposed by OpenMP
 968 on the ranges of parallel for-loops essentially constraint them to fit the format $\text{range}(t_{\text{start}}, t_{\text{stop}},$
 969 $t_{\text{step}})$, which is the format that we use for our simple-for-loops.

970
 971 *Structured Data.* The constructs $\{f_1 = t_1; \dots; f_n = t_n\}$ and $[t_1; \dots; t_n]$ build records and arrays as
 972 constant values. Mutable record and arrays are allocated by means of a call to the **alloc** function.
 973 OptiTrust features 4 operations to manipulate structured data. If a corresponds to a constant array
 974 value, then the operation $a[i]$ reads the i -th cell of the array a . If, however, a corresponds to the
 975 address of a heap-allocated or a mutable stack-allocated array, then the memory address of i -th cell
 976 of the array a can be computed by the operation $t \boxplus_{\hat{t}} i$, where \hat{t} denotes the type of the elements
 977 of t . This operation corresponds to the C pointer arithmetic operation $t+i$. The contents of that cell
 978 may be retrieved by evaluating $\text{get}(t \boxplus_{\hat{t}} i)$. Likewise, reading the field f of a constant record r is
 979 described by the operation $r.f$, whereas the memory address of the field f of a record r allocated
 980

in memory is described by the operation $r \boxplus_{\hat{\tau}} f$, 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 $\hat{\tau}$.

Other Language Constructs. The other language constructs from Fig. 13 are standard. They include sequence, variable definition, function definitions, 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, and the metavariable n denotes an integer literal.

Other Primitive Operations. Besides the aforementioned primitive operations for manipulating heap and stack cells, OptiTrust provides primitive functions that correspond to the arithmetic and boolean operators of the C languages. For the special operators `&&` and `||`, we encode them using conditionals whenever the second argument is not a simple expression (in particular, if it might fail or perform side effects).

For the pre/post-increment/decrement operators, such as `i++`, we encode them using conventional function calls, taking advantage of the fact that our type system guarantees the correctness of such an encoding. For example, our type system rules out problematic expressions such as `f(i++, i++)`. More precisely, provided that the input program does not exhibit undefined behavior, and provided that the encoding of the input program is well-typed in our system, then an expression `i++` appearing in that C program is equivalent to the expression `get_incr(&i)`, with the definition `int get_incr(int *p){ int r = *p; *p = r + 1; return r; }`.

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

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

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

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

3.2 Bidirectional Translation between C and OptiTrust’s Language

As mentioned earlier, OptiTrust *encodes* input C programs into OptiTrust’s internal AST, for the purpose of facilitating typechecking and transformation of the code. Because we wish to typecheck terms after their encoding, the encoding must not depend on the result of our resource analysis. After an operation on the internal AST, OptiTrust *decodes* the program back into C syntax, for the purpose of proposing feedback in a familiar programming syntax. In particular, the *diff* associated with a transformation requested by the user is computed over the two pieces of C code. The parts of the input C program that are not altered by the requested transformation do not appear in the *diff*.

In the rest of this section, we explain how the initial program is normalized and how annotations in the OptiTrust AST are used in order to ensure that a *round-trip* property holds. Then, we present the key ideas of the encoding scheme by means of example—as pointed out earlier, the elimination of *l*-values is a standard compiler translation.

Initial Normalization. The current implementation of OptiTrust does not keep in the AST information about comments and spacing—we left this to future work. As a result, for the very first encoding-decoding round-trip, which our implementation automatically applies as preliminary step, some amount of semantically-irrelevant information may be lost. Then, after that initial normalization step, for all programs manipulated by OptiTrust, the *round-trip* property holds: encoding and decoding are exactly reciprocal of each other.

Annotations. To deal with the fact that several C expressions might admit the same encoding in the OptiTrust language, our translation attaches annotation on certain OptiTrust terms. For example, $(*r).f$ and $r \rightarrow f$ are both encoded as $\text{get}(r \boxtimes f)$, therefore we instrument the encoding to attach a “dont-use-arrow” annotation on the `get` term when translating $(*r).f$. As another example, the two terms $e1 \ \&\& \ e2$ and `if (e1){ return e2; } else { return true; }` have the same encoding when $e2$ is a nontrivial expression, therefore we attach a “use- $\&\&$ ” on the conditional when translating $e1 \ \&\& \ e2$. If a transformation step modifies the instruction `return true` from the else-branch into something else, then the annotation “use- $\&\&$ ” is ignored by the decoding operation.

Encoding Scheme and Pure Variables. The core of OptiTrust’s encoding consists of eliminating *l*-values. For a heap allocated piece of data, a read operation $*p$ is encoded as the function call $\text{get}(p)$, and an assignment $*p = v$ is encoded as $\text{set}(p, v)$. For stack-allocated C variables, the encoding distinguishes two cases, depending on whether the address of the variable needs to be manipulated. We say that a stack variable x is *pure* if there is no assignment operation on x and the address of x or one of its sub-components is never taken.⁹¹⁰ For such a pure variable, its definition, say `int x = 3`, is encoded simply as `let x = 3`. For a non-pure variable, its encoding involves a stack-allocation. For example, the definition `int x = 3` is encoded as `let x = ref(3)`, the assignment `x = 4` is encoded as `set(x, 4)`, and an occurrence of $\&x$ is encoded as x . Fig. 14 provides additional examples.

4 COMPUTING PROGRAM RESOURCES

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. Triples will be later generalized 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 declarative entailment relation.

⁹For example, the variable x is *not pure* if there exists an occurrence of $\&x$, or of the form $\&x.f$ or $\&x[i]$. On the contrary, a pointer variable x may be *pure* despite occurring in an expression of the form $\&(x \rightarrow f)$. Indeed, if x is pure, then $\&(x \rightarrow f)$ is encoded as $x \boxtimes f$.

¹⁰Note that a variable declared as `const` is not necessarily pure, because its address might be taken. On the contrary, a variable that is pure may always be annotated as `const`.

1079	<code>int x = 3;</code>	\longleftrightarrow	<code>let_{int} x = 3;</code>	where x pure
1080	<code>f(x);</code>	\longleftrightarrow	<code>f(x);</code>	with <code>void f(int)</code>
1081				
1082	<code>int* a = malloc(sizeof(int));</code>	\longleftrightarrow	<code>let_(int*) a = alloc_{int}(1);</code>	where a pure
1083	<code>*a = *a + 2;</code>	\longleftrightarrow	<code>set_{int}(a, get_{int}(a) + 2);</code>	
1084	<code>free(a);</code>	\longleftrightarrow	<code>free(a);</code>	
1085				
1086	<code>int z;</code>	\longleftrightarrow	<code>let_(int*) z = ref_uninit_{int}();</code>	where z not pure
1087	<code>z = 6;</code>	\longleftrightarrow	<code>set_{int}(z, 6);</code>	
1088	<code>int v = z;</code>	\longleftrightarrow	<code>let_{int} v = get_{int}(z);</code>	where v pure
1089				
1090	<code>int y = 5;</code>	\longleftrightarrow	<code>let_(int*) y = ref_{int}(5);</code>	where y not pure
1091	<code>f(y);</code>	\longleftrightarrow	<code>f(get_{int}(y));</code>	
1092	<code>y = y + 2;</code>	\longleftrightarrow	<code>set_{int}(y, get_{int}(y) + 2);</code>	
1093	<code>y += 4;</code>	\longleftrightarrow	<code>set_add_{int}(y, 4);</code>	
1094	<code>y++;</code>	\longleftrightarrow	<code>get_incr_{int}(y);</code>	
1095				
1096	<code>int* p = &y;</code>	\longleftrightarrow	<code>let_(int*) p = y;</code>	where p pure
1097	<code>*p = *p + 2</code>	\longleftrightarrow	<code>set_{int}(p, get_{int}(p) + 2);</code>	
1098				
1099	<code>int* q = &y;</code>	\longleftrightarrow	<code>let_(int**) q = ref_{(int*)(y);}</code>	where q not pure
1100	<code>q = &z;</code>	\longleftrightarrow	<code>set_{int*}(q, z);</code>	
1101	<code>*q = *q + 2;</code>	\longleftrightarrow	<code>set_{int}(get_{int*}(q), get_{int}(get_{int*}(q)) + 2);</code>	
1102				

Fig. 14. Example translations from C code into the OptiTrust’s internal AST, with type annotations made explicit. A variable x is *pure* if there is no assignment operation on x and no occurrence of $\&x$.

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 pure expressions. Finally, Section 4.7 presents our algorithmic typing rules, which define the judgment $\{\Gamma\} t \{\Gamma'\}$.

Throughout the section, we assume a substitution operator for every entity. Concretely, given a map σ associating variable names to values, we write $\text{Subst}\{\sigma\}(X)$ 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. We next describe these two kind of resources.

Pure Resources. The pure part of a typing context contains bindings of the form “ $x : \tau$ ”, where τ corresponds either to a C type or to a *mathematical type*. A C type is denoted by the meta-variable $\hat{\tau}$. A mathematical type (a.k.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 program variables to C types, and a Coq context, which binds ghost variables to Coq types, including propositions and specifications.

Mathematical types include propositions, of type `Prop`: for example “ $p : n > 0$ ” describes a proof p establishing the proposition “ $n > 0$ ”. Mathematical types also include logical functions, whose type is written $\tau_1 \xrightarrow{\text{pure}} \tau_2$. In particular, a predicate over a type τ has type $\tau \xrightarrow{\text{pure}} \text{Prop}$. *Function specifications* are predicates over function values. As we will see in the next section, such a

Syntax in C	Syntax in the theory	Description
$p \rightsquigarrow \text{Cell}$	$p \rightsquigarrow \text{Cell}_{\hat{t}}$	permission to access the cell at address p of type \hat{t}
$p \rightsquigarrow \text{Matrix1}(n)$	$p \rightsquigarrow \text{Matrix1}_{\hat{t}}(n)$	permission on an array of length n
$p \rightsquigarrow \text{Matrix2}(m, n)$	$p \rightsquigarrow \text{Matrix2}_{\hat{t}}(m, n)$	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

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

hypothesis takes the form $\text{Spec}(f, [a_1, \dots, a_n], \gamma)$, asserting that the function f expects arguments named a_i and admits the *function contract* γ . Note that function contracts may appear in contexts, while contexts are involved in the statement function contracts. This form of impredicativity is standard in higher-order Separation Logic (e.g., [Charguéraud 2020b]). Another particular kind of predicates are heap predicates (e.g., $p \rightsquigarrow \text{Cell}$), which admit the type Hprop . In formalizations of Separation Logic, Hprop is typically defined as $\text{state} \xrightarrow{\text{pure}} \text{Prop}$, where state denotes the type of a memory state, however this definition needs not be revealed to the OptiTrust user.

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 $\text{ptr}(A)$ to denote the C type $A\star$.

Linear Resources. The linear part of a typing context contains bindings of the form “ $y : H$ ”, where H is a *heap predicate*, and where y is a name. 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. Fig. 15 summarizes the most common heap predicates, which have already been discussed in Section 2, but for which we here introduce LaTeX-style notation, moreover making the type annotations explicit. In particular, $p \rightsquigarrow \text{Cell}_{\hat{t}}$ is a resource that corresponds to the ownership of a single cell of type \hat{t} , located at address p . The resource $p \rightsquigarrow \text{Matrix1}_{\hat{t}}(n)$ is syntactic sugar for $\star_{i \in 0..n} p[i] \rightsquigarrow \text{Cell}_{\hat{t}}$. 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, $p \rightsquigarrow \text{Matrix2}_{\hat{t}}(n, m)$ denotes $\star_{i \in 0..n} \star_{j \in 0..m} p[i][j] \rightsquigarrow \text{Cell}_{\hat{t}}$. We leave it to future work to provide mechanisms allowing the user to define *representation predicates* for custom data types.

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 $(\alpha + \beta)H$. As a result, if we have αH at hand 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 `CloseFrac`, which appears in our typing rules. The operation `CloseFrac` 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 H in full.

At this stage in the development of OptiTrust, we assign meaning to a fractional permission αH only when H is a heap predicate that denotes either a single cell $p \rightsquigarrow \text{Cell}_\ell$ or a group of cells (typically, an array or matrix predicate). Finding an appropriate generalization of this definition of fractional resources that applies on user-defined container data structures remains, for the most part, an open problem in Separation Logic.

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.

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

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

Following standard practice in proof assistants, variable names that are nowhere mentioned may 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 resources: $n > 0$ and $p \rightsquigarrow \text{Cell}_{\text{int}}$. Internally, though, all context items are identified by a variable name.

1226 *Bindings of the Special Result Variable.* Recall from Section 3.1 that, in OptiTrust, if a term t
 1227 produces a result value, then this result is bound to the special variable **res**. Therefore, if t has a
 1228 non-void type τ then, in the triple $\{\Gamma\} t \{\Gamma'\}$, the output context Γ' binds a variable **res** of type τ .
 1229 Moreover, if in t the binding on **res** takes the “**let res = v**” for a *pure expression* v , then **res** is bound
 1230 in Γ' as an alias for v . In other words, we would have $(\mathbf{res} : \tau := v) \in \Gamma'$. (The grammar of pure
 1231 expressions is formalized further in Section 4.6.)

1232 *Projection of Context Components.* We define two projection functions. For a context $\Gamma = \langle E \mid F \rangle$,
 1233 the projection “ Γ .pure” returns E , and the projection “ Γ .linear” returns F .
 1234

1235 *Syntax for Contexts with One Component.* As syntactic sugar, we define $[x_0 : \tau_0, \dots, x_n : \tau_n]$ as
 1236 $\langle x_0 : \tau_0, \dots, x_n : \tau_n \mid \emptyset \rangle$, for contexts that are entirely pure. Furthermore, we allow ourselves to write
 1237 F to mean $\langle \emptyset \mid F \rangle$, where F denotes a set of linear resources.
 1238

1239 *Separated Conjunction of Two Contexts.* We write $F_1 \star F_2$ the disjoint union of two sets of linear
 1240 resources. Furthermore, for two contexts Γ_1 and Γ_2 with disjoint domains, we define $\Gamma_1 \otimes \Gamma_2$ as
 1241 $\langle \Gamma_1.\text{pure}, \Gamma_2.\text{pure} \mid \Gamma_1.\text{linear} \star \Gamma_2.\text{linear} \rangle$, where the comma indicates list concatenation. Observe that
 1242 $[E] \otimes F = \langle E \mid F \rangle$.

1243 *Pointwise Operators Over Linear Contexts.* Consider a set of resources F of the form $(y_0 : H_0,$
 1244 $\dots, y_n : H_n)$. 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
 1245 conjunction distributes pointwise over the set of linear resources. Similarly, we define αF as
 1246 $(y_0 : \alpha H_0, \dots, y_n : \alpha H_n)$.
 1247

1248 *Filtering on Contexts.* We define a filtering operation, written $G \vdash X$, where G is a set of resources
 1249 (linear or pure) and X is a set of variable names. This operation computes a set of resources where
 1250 only the entries from G whose name belongs to the set X are kept. Filtering also applies to contexts:
 1251 $\langle E \mid F \rangle \vdash X$ is defined as $\langle E \vdash X \mid F \vdash X \rangle$.
 1252

1253 *Specialization of Contexts.* The specialization operation is used for example to specialize the
 1254 contract of a function for a specific call to that functions. The contract is then specialized on
 1255 the arguments, as well as on the ghost arguments, on which the function is applied. In case of a
 1256 polymorphic function, type arguments are specialized as well. The specialization operation takes
 1257 the form $\text{Specialize}_{\Gamma_0} \{\sigma\}(\Gamma)$. The definition of this operation is fairly technical, yet it is a direct
 1258 generalization of the process of typechecking function applications in higher-order logics. Rather
 1259 than showing the technical details, let us illustrate the specialization operation on an example.

1260 Consider a function f whose input is described by a context $\Gamma \equiv \langle A : \text{Type}, C : \text{Type}, n : \text{int},$
 1261 $p : \text{ptr}(A), b : A, c : C \mid p \rightsquigarrow \text{Matrix}1_A(n) \rangle$, where A and C are type arguments, where p and n
 1262 denote physical arguments, and where b and c are ghost arguments. Consider a function call of
 1263 the form $f(7, q)$, where q is a program variable of type $\text{ptr}(\text{int})$ in scope at the call site. This call
 1264 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
 1265 context describing the pure variables bound at the call site. In particular, we have $(q : \text{ptr}(\text{int})) \in \Gamma_0$.
 1266 For the example considered, the specialization operation yields the context: $\langle C : \text{Type}, b : \text{int},$
 1267 $c : C \mid q \rightsquigarrow \text{Matrix}1_A(7) \rangle$. Observe how the types and arguments being specialized (namely A, n
 1268 and p) are eliminated from the pure part of the context, and the corresponding values (namely $\text{int},$
 1269 7 and q) are substituted in the entities that remain.
 1270

1271 *Renaming on Contexts.* A renaming operation is involved when the programmer explicitly spec-
 1272 ifies the names to assign to the ghost variables obtained as part of the result of a function call.
 1273 The operation $\text{Rename}\{\rho\}(\Gamma)$ renames certain keys from Γ . Here, ρ denotes a map that associates
 1274 resource names to other resource names. The keys from ρ may or may not be bound in Γ . The values

1275		$\{\}$	$\text{alloc}_{\hat{\tau}}()$	$\{[\mathbf{res} : \text{ptr}(\hat{\tau})] \otimes \text{Uninit}(\mathbf{res} \rightsquigarrow \text{Cell}_{\hat{\tau}})\}$
1276				
1277	$\{[a : \text{ptr}(\hat{\tau}), \alpha : \text{frac}] \otimes \alpha(a \rightsquigarrow \text{Cell}_{\hat{\tau}})\}$		$\text{get}_{\hat{\tau}}(a)$	$\{[\mathbf{res} : \hat{\tau}] \otimes \alpha(a \rightsquigarrow \text{Cell}_{\hat{\tau}})\}$
1278	$\{[a : \text{ptr}(\hat{\tau}), b : \hat{\tau}] \otimes \text{Uninit}(a \rightsquigarrow \text{Cell}_{\hat{\tau}})\}$		$\text{set}_{\hat{\tau}}(a, b)$	$\{a \rightsquigarrow \text{Cell}_{\hat{\tau}}\}$
1279	$\{[a : \text{ptr}(\hat{\tau})] \otimes \text{Uninit}(a \rightsquigarrow \text{Cell}_{\hat{\tau}})\}$		$\text{free}_{\hat{\tau}}(a)$	$\{\}$
1280		$\{[b : \hat{\tau}]\}$	$\text{ref}_{\hat{\tau}}(b)$	$\{[\mathbf{res} : \text{ptr}(\hat{\tau})] \otimes \mathbf{res} \rightsquigarrow \text{Cell}_{\hat{\tau}}\}$
1281				
1282		$\{\}$	$\text{ref_uninit}_{\hat{\tau}}()$	$\{[\mathbf{res} : \text{ptr}(\hat{\tau})] \otimes \text{Uninit}(\mathbf{res} \rightsquigarrow \text{Cell}_{\hat{\tau}})\}$

Fig. 16. Contracts assigned to key primitive functions; $\hat{\tau}$ denotes a C type; a and b denote program variables.

from ρ must be fresh from Γ . For example, $\text{Rename}\{x := x', y := y'\}(\langle E_1, x : \tau, E_2 \mid F \rangle)$, where y has no occurrence in E_1, E_2 or F , evaluates to $\langle E_1, x' : \tau, \text{Subst}\{x := x'\}(E_2) \mid \text{Subst}\{x := x'\}(F) \rangle$. As another example, $\text{Rename}\{y := y'\}(\langle E \mid F_1, y : H, F_2 \rangle)$ evaluates to $\langle E \mid F_1, y' : H, F_2 \rangle$.

4.3 Grammar of Contracts

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

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

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

Contracts for Primitive Functions. Fig. 16 gives the contracts that we axiomatize for the operations on heap and stack cells—technically, we present not their contracts but the triples derived from their contracts, to improve readability. The C language does not feature polymorphism, hence we describe in the figure a family of contracts, indexed with the type $\hat{\tau}$ of the memory cell at hand. (For polymorphic operators, there would be a single contract, which would include a quantification over the type $\hat{\tau}$ in the input environment.)

These contracts illustrate, in particular, how the bindings on \mathbf{res} appears in output contexts. An allocation produces an uninitialized permission. A write requires an uninitialized permission and returns a full permission. A read requires a read-only permission and returns it. A free operation requires a full permission and does not return it. Recall that a full permission can be split into read-only resources, and that it may be downgraded at any time into an uninitialized permission. The operation on stack cells are similar.

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

Contracts for Ghost Functions. In addition to contracts for primitive heap-manipulating functions, OptiTrust provides contracts for primitive ghost functions. For example, consider the ghost function `swap_groups`, which was involved in the OpenCV case study (Section 2.1) for swapping two iterators (iterated separating conjunctions), and which typically needs to be involved on every loop-swap operation (as explained further in Section 6.5). It is specified shown below, where H is a heap

1324 predicate that depends on the two indices i and j .

$$1325 \{ [r_i : \text{range}, r_j : \text{range}, H : (\text{int}, \text{int}) \xrightarrow{\text{pure}} \text{Hprop}] \otimes \left(\bigstar_{i \in r_i} \bigstar_{j \in r_j} H(i, j) \right) \} \text{swap_groups} \left\{ \bigstar_{j \in r_j} \bigstar_{i \in r_i} H(i, j) \right\}$$

1327 The OptiTrust user can define custom ghost functions by composing existing ghost functions. Doing
 1328 so can be useful to factorize repetitive resource-manipulation patterns. Ghost functions are written
 1329 and typechecked like regular C functions. Their specificity is that they contain only calls to ghost
 1330 functions and control-flow structures such as for-loops.

1332 *Loop Contracts.* A for-loop annotated with a *loop contract* χ takes the form **for** $(i \in r)_\chi \{t\}$. The
 1333 loop contract χ consists of a record structured as follows.

$$1334 \left\{ \begin{array}{ll} \text{vars} = E & \text{Pure variables, scoping over the other contract components} \\ \text{excl} = \begin{cases} \text{pre} = F_{pre} & \text{Resources consumed exclusively by one iteration} \\ \text{post} = F_{post} & \text{Resources produced exclusively by one iteration} \end{cases} \\ \text{shrd} = \begin{cases} \text{reads} = F_{reads} & \text{Read only resources shared between iterations} \\ \text{inv} = F_{inv} & \text{Sequential invariant, threaded through iterations} \end{cases} \end{array} \right.$$

1341 We call E the *loop ghost variables*. The variables from E scope over F_{pre} , F_{post} , F_{reads} and F_{inv} .
 1342 We call F_{pre} the *consumed per-iteration resources* and F_{post} the *produced per-iteration resources*.
 1343 Resources in F_{pre} and in F_{post} may (and typically do) refer to the loop index. We call F_{reads} the
 1344 *shared reads*. In practice, this set of resources consists of read-only resources. Resources in F_{reads}
 1345 cannot refer to the loop index. We call F_{inv} the *sequential invariant*. It corresponds to a standard
 1346 loop invariant in sequential Separation Logic. In this paper, we consider for simplicity that F_{inv}
 1347 does not depend on the loop index.

1348 As we will see later in Section 4.7, when typechecking a loop of the form **for** $(i \in r)_\chi \{t\}$, the
 1349 loop body t is typechecked in a context that binds an index i of type int , a hypothesis of type $i \in r$,
 1350 the variables from E , the resources F_{pre} , (subfractions of) the resources in F_{reads} , and the resources
 1351 in F_{inv} . The loop body needs to produce the resources F_{post} , and it needs to give back the resources
 1352 that it had received from F_{reads} and from F_{inv} .

1354 *Parallel Loop Contracts.* A loop is parallelizable if it can be typechecked with an empty sequential
 1355 invariant F_{inv} . Hence, we say that a loop contract χ is *parallelizable*, and write $\text{parallelizable}(\chi)$,
 1356 when $\chi.\text{shrd}.\text{inv} = \emptyset$.

1357 4.4 Entailment

1359 We next introduce the *entailment* judgment, written $\Gamma \Rightarrow \Gamma'$. The entailment judgment is used in
 1360 particular to assert that the set of resources obtained at a given program point corresponds to
 1361 the set of resources expected at that same point. For example, the set of resources at the end of a
 1362 function body corresponds to the set described by the postcondition of this function.

1363 The literature on Separation Logic includes two types of entailment: *linear* and *affine* entailment
 1364 relations. OptiTrust is based on a linear entailment relation, disallowing resources to be silently
 1365 “dropped”. The benefits of using linear entailment is that it allows checking the absence of memory
 1366 leaks—every piece of heap allocated data must eventually be freed.

1367 The OptiTrust entailment between two contexts:

$$1368 \langle x_0 : \tau_0, \dots, x_n : \tau_n \mid y_0 : H_0, \dots, y_m : H_m \rangle \Rightarrow \langle x'_0 : \tau'_0, \dots, x'_m : \tau'_m \mid y'_0 : H'_0, \dots, y'_m : H'_m \rangle$$

1369 is defined as

$$1370 \forall x_0 : \tau_0, \dots, \forall x_n : \tau_n, \quad H_0 \star \dots \star H_m \xrightarrow{SL} (\exists x'_0 : \tau'_0, \dots, \exists x'_m : \tau'_m, H'_0 \star \dots \star H'_m)$$

1373 where \Rightarrow^{SL} denotes the standard Separation Logic entailment.

1374 For example, the following entailments hold:

- 1375 • $\langle x : \text{loc} \mid x \rightsquigarrow \text{Cell} \rangle \Rightarrow \langle \alpha : \text{frac} \mid \alpha(x \rightsquigarrow \text{Cell}), (1 - \alpha)(x \rightsquigarrow \text{Cell}) \rangle$
- 1376 • $\langle y : \text{loc} \mid y \rightsquigarrow \text{Cell} \rangle \Rightarrow \langle \mid \text{Uninit}(y \rightsquigarrow \text{Cell}) \rangle$
- 1377 • $\langle A : \text{loc}, n : \text{int}, m : \text{int} \mid \star_{i \in 0..n} \star_{j \in 0..m} A[i][j] \rightsquigarrow \text{Cell} \rangle \Rightarrow \langle \mid \star_{j \in 0..m} \star_{i \in 0..n} A[i][j] \rightsquigarrow \text{Cell} \rangle$
- 1378 • $\langle n : \text{int}, n \text{ even} \mid \rangle \Rightarrow \langle m : \text{int}, n = 2m \mid \rangle$

1379 However, the entailment $\langle x : \text{loc} \mid x \rightsquigarrow \text{Cell} \rangle \Rightarrow \langle \mid \rangle$ does not hold because linear resources
1380 cannot be dropped.

1381 As a shorthand, we write $\Gamma \Leftrightarrow \Gamma'$ to assert that entailment holds both ways, that is, to assert that
1382 the conjunction $(\Gamma \Rightarrow \Gamma') \wedge (\Gamma' \Rightarrow \Gamma)$ holds.

1384 4.5 Subtraction

1385 The *subtraction* operation provides a sound yet incomplete algorithmic implementation of the entail-
1386 ment judgment. The subtraction operation not only allows checking the validity of an entailment, it
1387 also enables a certain amount of inference. At a high level, given Γ and Γ' , the subtraction operation
1388 computes the *frame*, written F , which denotes the set of linear resources such that $\Gamma \Rightarrow \Gamma' \star \langle \emptyset \mid F \rangle$.
1389 The subtraction operation also infers the instantiation map σ providing the witnesses for the
1390 instantiations of the variables that are bound (and therefore existentially quantified) in Γ' . Such
1391 a subtraction operator is found in most—if not all—practical verification frameworks based on
1392 Separation Logic.

1393 The typing rules of OptiTrust actually make use of two variants of the subtraction operation.
1394 The *core subtraction operation*, written $\Gamma \boxminus \Gamma'$, is able to convert uninitialized resources into full
1395 resources on-the-fly, however it does not support splitting read-only resources on-the-fly. The
1396 *carving subtraction operation*, written $\Gamma \ominus \Gamma'$, extends the former with the feature of carving out
1397 a fraction of a read-only permission from Γ every time a corresponding read-only permission is
1398 requested in Γ' . (Carving was described in Section 4.1.)

1399 The core subtraction operation $\Gamma \boxminus \Gamma'$ is formally specified as a partial operation. It may fail
1400 (that is, return \perp) if a resource in Γ' cannot be matched against a corresponding resource in Γ .
1401 Otherwise, the operation returns a result of the form (σ, F) . When $\Gamma \boxminus \Gamma' = (\sigma, F)$, then the
1402 entailment $\Gamma \Rightarrow \text{Specialize}_{\Gamma} \{ \sigma \} (\Gamma') \star \langle \emptyset \mid F \rangle$ holds. In particular, the subtraction operation can be
1403 used to prove an entailment $\Gamma \Rightarrow \Gamma'$, by checking that $\Gamma \boxminus \Gamma'$ evaluates to (σ, \emptyset) for some σ .

1404 The subtraction operation is implemented following a standard scheme.

- 1405 (1) The substitution map σ is initialized with bindings that associates each of the pure variables
1406 of Γ' to a fresh unification variable.
- 1407 (2) Each of the linear resources from Γ' are matched against a corresponding resource from Γ .
1408 This process may trigger unifications, resulting in partial or total resolution of certain
1409 unification variables.
- 1410 (3) If Γ' requests a linear resource of the form $\text{Uninit}(H)$, and if Γ contains the resource H ,
1411 then our algorithm applies an on-the-fly weakening from H to $\text{Uninit}(H)$.
- 1412 (4) The items from Γ that remains at the end are assigned to the frame F .

1413 The carving subtraction operation $\Gamma \ominus \Gamma'$ behaves almost like $\Gamma \boxminus \Gamma'$ but outputs a triple $(E_{frac},$
1414 $\sigma, F)$ where σ and F are the same as in core subtraction and E_{frac} is a pure context for generated
1415 fractions containing only bindings of the form $\alpha : \text{frac}$ and initially empty. Compared to the core
1416 subtraction, the carving subtraction refines step (2) as follows. If Γ' requests a fractional resource
1417 αH , if α is an unconstrained unification variable that denotes a fraction, and if Γ contains a fractional
1418 resource $\beta H'$ for some fraction β and where H unifies with H' , then our algorithm applies an
1419 on-the-fly splitting operation to convert $\beta H'$ into the conjunction of $\alpha' H'$ and $(\beta - \alpha') H'$ for a
1420
1421

1422	$v ::= x$	Variable
1423	n b	Integer or boolean literal
1424	P H $\hat{\tau}$	Mathematical proposition, heap predicate, or C type
1425	$\mathbf{fun}(x_1 : \tau_1, \dots, x_n : \tau_n) \mapsto v$	Pure function definition
1426	$v_0(v_1, \dots, v_n)$	Pure function application
1427		
1428	VAR	INT
1429	$\frac{(v : \tau) \in E}{E \vdash v : \tau}$	$\frac{}{E \vdash n : \text{int}}$
1430		$\frac{}{E \vdash b : \text{bool}}$
1431		$\frac{}{E \vdash \text{int} : \text{Type}}$
1432	PROP	HPROP
1433	$\frac{P \text{ is a mathematical proposition with free variables in } E}{E \vdash P : \text{Prop}}$	$\frac{H \text{ is a heap predicate with free variables in } E}{E \vdash H : \text{Hprop}}$
1434		PtrTYPE
1435		$\frac{}{E \vdash \text{ptr}(A) : \text{Type}}$
1436		PUREAPP
1437	PUREFUN	$\frac{E \vdash v_0 : (\tau_1, \dots, \tau_n) \xrightarrow{\text{pure}} \tau}{E \vdash v_0(v_1, \dots, v_n) : \tau}$
1438	$\frac{(E, x_1 : \tau_1, \dots, x_n : \tau_n) \vdash v : \tau}{E \vdash \mathbf{fun}(x_1 : \tau_1, \dots, x_n : \tau_n) \mapsto v : (\tau_1, \dots, \tau_n) \xrightarrow{\text{pure}} \tau}$	$\frac{E \vdash v_0 : (\tau_1, \dots, \tau_n) \xrightarrow{\text{pure}} \tau}{E \vdash v_0(v_1, \dots, v_n) : \tau}$
1439		$\frac{\forall i \in [1..n], E \vdash v_i : \tau_i}{E \vdash v_0(v_1, \dots, v_n) : \tau}$
1440		

Fig. 17. Grammar of pure expressions, written v , and associated typing judgment, written $E \vdash v : \tau$. Arithmetic operations such as $v_1 + v_2$ are viewed as functions calls and are therefore handled by the rule PUREAPP.

fresh α' added to E_{frac} . Then it adds the binding $\alpha := \alpha'$ in σ . Doing so ensures that a fraction of a piece of βH remains available in Γ , allowing to eliminate other resources of the form $\alpha'' H$ that might appear in the remaining elements from Γ' .

4.6 Typechecking of Pure Expressions

A *pure expression* (a.k.a. *logical expression*), written v , is an expression whose evaluation terminates and whose value does not depend on the memory state. Pure expressions, unlike stateful expressions, may appear in specifications and invariants. Pure expressions include program variables (which are always immutable in the OptiTrust AST), constant literals, mathematical propositions, heap predicates, C types, pure functions definitions, and pure function calls. Figure 17 formalizes the grammar of pure expressions, and defines their typing judgment, written $E \vdash v : \tau$, where E is a pure context—that is, the pure part of a context, in the sense of Section 4.2.

An arithmetic expression $v_1 + v_2$ can be considered as a pure expression if its two arguments are pure. The contract for addition is: $\{[a : \text{int}, b : \text{int}]\} (a + b) \{[\mathbf{res} := a \hat{+} b : \text{int}]\}$, 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 pure expressions, simply with an additional precondition. The contract for division is: $\{[a : \text{int}, b : \text{int}, b \neq 0]\} (a/b) \{[\mathbf{res} := a \hat{/} b : \text{int}]\}$, 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 dividing by 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

1471 produces Γ' and Δ as output. The refinement with usage maps will be discussed further in Section 5.3.
 1472 For now, we focus on describing typing rules for the judgment $\{\Gamma\} t \{\Gamma'\}$.

1473 In general, in a valid triple $\{\Gamma\} t \{\Gamma'\}$, variables from the postcondition Γ' may refer to variables
 1474 from the precondition Γ . For the purpose of the algorithmic typechecking, however, we design the
 1475 typing rules in such a way that Γ' is always *closed*, meaning that variable occurrences in Γ' refer to
 1476 variables previously bound in Γ' . The purpose of this design decision is to maximize the amount of
 1477 information that is propagated forward during the typechecking.

1478 In particular, in the algorithmic typechecking, all the logical bindings (ghost variables and pure
 1479 facts) from Γ are reproduced in Γ' . The pure bindings that appear in Γ' but not in Γ correspond
 1480 either (1) to the binding for **res**, which denotes the result value produced by t , as explained in
 1481 Section 4.3; or (2) to logical bindings (ghost variables and pure facts) that correspond to existentially
 1482 quantified variables and pure postconditions.

1483 Unlike pure bindings, the linear bindings of Γ' may be arbitrarily modified compared with those
 1484 in Γ . The modifications reflect the side effects performed by t . Linear resources that are bound with
 1485 the same name in Γ' as in Γ necessarily correspond to resources that have not been modified by t .

1486 Figure 18 presents our typing rules. The typing rule for applications handles the particular case
 1487 where the subterms are program variables (i.e., functions calls in A-normal form)—the processing of
 1488 effectful subterms depends on resource usage, and is explained further in Section 5.5. The soundness
 1489 of these rules stems from the fact that they correspond to an algorithmic reformulation of the
 1490 standard reasoning rules from Separation Logic. We next describe the rules individually.

1491 *Program Values.* Consider a program value v . In its triple, of the form $\{\Gamma\} v \{\Gamma'\}$, the output
 1492 context Γ' is obtained by extending Γ with an alias binding from **res** to v . Alias bindings were defined
 1493 in Section 4.2. The type of v is computed by means of the typing judgment for pure expressions,
 1494 defined in Section 4.6.
 1495

1496 *Let-Bindings.* Consider an instruction of the form **let** $x = t$. Recall from Section 3.1 that such
 1497 instructions only appear in sequences. The subexpression t produces a value, hence the output
 1498 context Γ_1 associated with t binds the special variable **res**. The expression **let** $x = t$ itself does not
 1499 produce a value, hence its output context Γ_2 does not bind **res**. However, the output context Γ_2 is
 1500 extended with a binding on x . Concretely, Γ_2 is obtained by replacing in Γ_1 the bound name **res**
 1501 with the bound name x .

1502 *Sequence of Instructions.* We decompose the treatment of sequences in two rules: a first rule
 1503 named SEQ for handling the sequence of instructions per se, and a second rule named BLOCK for
 1504 handling the disposal of stack-allocated variables. The rest of this paragraph describes the SEQ rule.

1505 Consider a sequence of instructions t_i . Starting from an input context Γ_0 , each subterm t_i makes
 1506 the context evolves from Γ_{i-1} to Γ_i . There are two particular cases to handle. First, if a t_i is a non-void
 1507 expression, then the output context of t_i features a **res** binding. This binding is materialized as
 1508 a binding over a fresh ghost variable x_i in the final output context Γ_r . Second, if one of the t_i
 1509 instruction is of the form **let res** = t'_i , then the final context Γ_r is patched to include a corresponding
 1510 **res** binding, instead of a binding on a ghost variable x_i .
 1511

1512 *Scope Blocks.* The typing rule BLOCK is responsible for collecting the resources that corresponds to
 1513 stack-allocated variables, when reaching the end of a sequence, that is, the end of their scope. Recall
 1514 from Section 3.1 that stack allocation takes the form **let** $x = \text{ref}(T)$ or **let** $x = \text{ref_uninit}()$, with
 1515 such instructions occurring directly within a sequence. The auxiliary function StackAllocCells(t_1 ,
 1516 ..., t_n) synthesizes, based on the syntax of the terms t_i that appear in the sequence at hand, a
 1517 conjunction of resources, each of the form $\text{Uninit}(p \rightsquigarrow \text{Cell}_\tau)$. These resources are subtracted from
 1518 the set of resources available at the end of the sequence. Crucially, the subtraction operation checks
 1519

$$\begin{array}{c}
1520 \\
1521 \\
1522 \\
1523 \\
1524 \\
1525 \\
1526 \\
1527 \\
1528 \\
1529 \\
1530 \\
1531 \\
1532 \\
1533 \\
1534 \\
1535 \\
1536 \\
1537 \\
1538 \\
1539 \\
1540 \\
1541 \\
1542 \\
1543 \\
1544 \\
1545 \\
1546 \\
1547 \\
1548 \\
1549 \\
1550 \\
1551 \\
1552 \\
1553 \\
1554 \\
1555 \\
1556 \\
1557 \\
1558 \\
1559 \\
1560 \\
1561 \\
1562 \\
1563 \\
1564 \\
1565 \\
1566 \\
1567 \\
1568
\end{array}$$

$$\begin{array}{c}
\frac{\Gamma.\text{pure} \vdash v : \tau}{\{\Gamma\} v \{\Gamma \otimes [\mathbf{res} : \tau := v]\}} \text{VAL} \qquad \frac{\{\Gamma_0\} t \{\Gamma_1\} \quad \Gamma_2 = \text{Rename}\{\mathbf{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 \wedge \quad \Gamma_i = \text{Rename}\{\mathbf{res} := x_i\}(\Gamma'_i)}{\Gamma_r = \begin{cases} \text{Rename}\{x_i := \mathbf{res}\}(\Gamma_n) & \text{if } t_i \text{ is of the form } \mathbf{let} \mathbf{res} = t'_i \\ \Gamma_n & \text{otherwise} \end{cases}} \text{SEQ} \\
\frac{\{\Gamma_0\} (t_1; \dots; t_n) \{\Gamma_r\}}{\{\Gamma_0\} \{t_1; \dots; t_n\} \{\langle \Gamma_r.\text{pure} \mid F \rangle\}} \text{BLOCK} \\
\\
\frac{\begin{array}{c} \{[\Gamma_0.\text{pure}] \otimes \gamma.\text{pre}\} t \{\Gamma_1\} \quad (_, \emptyset) = \Gamma_1 \boxplus \gamma.\text{post} \\ (\mathbf{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 [\mathbf{res} : \hat{t}_f, \text{Spec}(\mathbf{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\}(\text{Subst}\{\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.pre}) \otimes \chi.\text{shrd.reads} \otimes \chi.\text{shrd.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.pre} \otimes \frac{1}{r.\text{len}} \chi.\text{shrd.reads} \otimes \chi.\text{shrd.inv} \\ \{\{\Gamma_0.\text{pure}\} \otimes \Gamma'_p\} t \{\Gamma'_q\} \\ (_, \emptyset) = \Gamma'_q \boxplus \chi.\text{excl.post} \otimes \frac{1}{r.\text{len}} \chi.\text{shrd.reads} \otimes \chi.\text{shrd.inv} \\ \Gamma_q = \text{Subst}\{\sigma'\}((\star_{i \in r} \chi.\text{excl.post}) \otimes \chi.\text{shrd.reads} \otimes \chi.\text{shrd.inv}) \\ \Gamma_r = \text{CloseFrac}([\Gamma_0.\text{pure}, E_{\text{frac}}] \otimes F \otimes \Gamma_q) \\ \text{parallelizable}(\chi) \vee \pi \neq \text{parallel} \end{array}}{\{\Gamma_0\} \mathbf{for}^\pi (i \in r)_\chi t \{\Gamma_r\}} \text{FOR} \\
\\
\frac{\{\Gamma_0\} t_0 \{\Gamma'_0\} \quad \{\text{Learn}\{\mathbf{res} = \text{true}\}(\Gamma'_0)\} t_1 \{\Gamma_1\} \quad \{\text{Learn}\{\mathbf{res} = \text{false}\}(\Gamma'_0)\} t_2 \{\Gamma_2\}}{\begin{array}{c} (_, \emptyset) = \Gamma_1 \boxplus \Gamma_3 \quad (_, \emptyset) = \Gamma_2 \boxplus \Gamma_3 \\ \{\Gamma_0\} \mathbf{if} t_0 \mathbf{then} t_1 \mathbf{else} t_2 \{\Gamma_3\} \end{array}} \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.

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.

Function Definition. Consider a function definition $\mathbf{fun}(a_1 : \hat{t}_1, \dots, a_n : \hat{t}_n)_\gamma \mapsto t$, with arguments a_i of type \hat{t}_i , with body t , and with contract γ . Recall from Section 4.3 that the function contract consists of a precondition $\gamma.\text{pre}$ and a postcondition $\gamma.\text{post}$, both described as contexts. The function closure being built may capture free variables from the current context. In the rule, the pure variables from the current context are described as $\Gamma_0.\text{pure}$. Note, however, that a function closure is not allowed to capture linear resources. Hence, the body of the function is typechecked in an

environment that consists of the conjunction of $\Gamma_0.\text{pure}$ and $\gamma.\text{pre}$. Ultimately, the body of the function must produce a set of resources Γ_r that entails the postcondition $\gamma.\text{post}$. The postcondition includes on **res**, which admits a function type, and a hypothesis $\text{Spec}(\mathbf{res}, [a_1, \dots, a_n], \gamma)$ that describes the specification of the function closure being created. As explained earlier in Section 4.3, this hypothesis captures $\{\gamma.\text{pre}\} \mathbf{res}(a_1, \dots, a_n) \{\gamma.\text{post}\}$, which is indeed the triple intended for the closure named **res**.

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

Two additional technicalities are involved in the statement of the *APP* rule. They correspond to the handling of optional user-provided annotations, name σ and ρ , that may guide the type-checking of an application. Such annotations are commonly found both in proof assistants and in program verification frameworks. The map σ corresponds to a for instantiating a subset of the ghost arguments. Indeed, there could be situations where the subtraction operation would fail to infer a unique possible instantiation, by the only means of the unification process. Hence, user annotations are required to resolve the instantiation. The map ρ corresponds to a renaming map. Its purpose is to allow the user to specify custom names for the ghost bindings and for the linear resources that are produced by the postcondition. In all our case studies, we did not need ρ , and σ was only used on ghost calls, however we anticipate for the general case. Typically, a ghost result produced by a function call is explicitly named using the ρ annotation, in order to provide this same ghost result as part of a σ annotation for a subsequent function call.

Simple for-loops. Consider a possibly parallel, simple for-loop of the form $\mathbf{for}^{\pi}(i \in r)_{\chi} t$. The typechecking of such a loop, driven by the loop contract annotation χ , was already discussed at a high level at the end of Section 4.3. Let us comment on the three remaining aspects. First, the shared-read resources, described by $\chi.\text{shrd}.\text{reads}$, are split into $\frac{1}{r.\text{len}}$ subfractions, where $r.\text{len}$ denotes the number of iterations associated with the range r . Note that, when typechecking the body of the loop for a particular iteration $i \in r$, the denominator $r.\text{len}$ can be assumed to be nonzero—indeed, $i \in r$ is equivalent to $0 \leq i < r.\text{len}$. Second, like for function calls, the instantiation of the contract using the resources from the input environment Γ_0 is computed using a subtraction, involving a frame F as well as an instantiation map σ' . Also, like for function calls, the output context is obtained by invoking the *CloseFrac*s operation. Third, loops, like functions calls, feature optional annotations σ and ρ , which we have omitted from the statement of the rule, for simplicity. The map σ guides how the contract is instantiated in the input environment Γ_0 . The map ρ can be used to explicit the 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_0 **then** t_1 **else** t_2 . The condition t_0 is evaluated in the input context Γ_0 and produces a context Γ'_0 . Then, both branches t_1 and t_2 need to typecheck in the context Γ'_0 . This context needs to be patched to reflect the knowledge that t_0 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_0 is a test of the form $x == y$, and consider the evaluation of $\text{Learn}\{\mathbf{res} = \text{true}\}(\Gamma'_0)$. The output context of t_0 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_1 produces an output context Γ_1 , and likewise the else-branch t_2 produce an output context Γ_2 . The challenge, typical in program logics, is to find a *join* context Γ_3 that both Γ_1 and Γ_2 entail. In the general case, this context Γ_3 cannot be automatically inferred. It needs to be provided, at least in part, via a *join contract*. In practice, numerous heuristics can be devised to synthesis a join contract, to save the need for users to provide explicit join contract. In particular, whenever the if-statement is the last instruction of a function body or a loop body, the join contract is provided by the expected postcondition coming from the function contract or the loop contract. All our case studies followed this pattern. We leave it to future work to experiment with heuristics for join contracts.

5 COMPUTING PROGRAM RESOURCES

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

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

5.1 Grammar of Usage Maps

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

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

- 1667 • “ $y : \text{uninit}$ ” is similar to “ $y : \text{full}$ ” but moreover captures the information that t needs not
1668 read the original contents of the memory cells associated with the resource named y . In
1669 particular, if t performs a write operation in a cell y before any read operation on y , then
1670 the usage of y is uninit .
- 1671 • “ $y : \text{splittedFrac}$ ” can arise when Γ contains a splittable linear resource “ $y : H$ ”, for some
1672 predicate H . The usage “ $y : \text{splittedFrac}$ ” means that t uses an unspecified subfraction of
1673 this resource. In such a situation, the name y is bound both in Γ and in Γ' . It may be the
1674 case, however, that the resource named y carries different fractions in Γ and Γ' .
- 1675 • “ $y : \text{joinedFrac}$ ” can arise when Γ contains a linear resource of the form “ $y : (\alpha - \beta_1 - \dots - \beta_n)H$ ”.
1676 The usage “ $y : \text{joinedFrac}$ ” means that: (1) the linear resource named y is not used by t , and
1677 (2) t produced a resource of the form $(\beta_i - \gamma_1 - \dots - \gamma_m)H$, and (3) these two resources are
1678 merged and the result appears in Γ' under the name y . If a single merge operation is applied,
1679 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
1680 Section 4.1)
- 1681 • “ $y : \text{produced}$ ” means that the linear resource y has been produced by t . In this case, y is
1682 the name of a linear resource in Γ' , and does not occur in Γ .
- 1683 • If a resource name is bound in Γ but not in Δ , then its absence indicates that the corre-
1684 sponding resource is not touched by t . Such a resource is bound under the same name in Γ
1685 and Γ' .

1687 5.2 Operations on Usage Maps

1688 *Projections of Usage Maps.* We define $\Delta.\text{full}$ as the set of names y such that “ $y : \text{full}$ ” appears in Δ .
1689 Likewise, we define $\Delta.\text{required}$, $\Delta.\text{ensured}$, $\Delta.\text{uninit}$, $\Delta.\text{splittedFrac}$, $\Delta.\text{joinedFrac}$ and $\Delta.\text{produced}$.
1690 In addition, we define the following operations.

$$\begin{aligned}
 1691 \quad \Delta.\text{consumed} &= \Delta.\text{full} \cup \Delta.\text{uninit} \\
 1692 \quad \Delta.\text{usedRO} &= \Delta.\text{splittedFrac} \cup \Delta.\text{joinedFrac} \\
 1693 \quad \Delta.\text{notRO} &= \text{dom}(\Delta) \setminus \Delta.\text{usedRO}
 \end{aligned}$$

1694
1695 *Intersection and Filtering.* We define:

$$\begin{aligned}
 1696 \quad \Delta_1 \cap \Delta_2 &= \text{dom}(\Delta_1) \cap \text{dom}(\Delta_2) \\
 1697 \quad \Gamma \vdash \Delta &= \Gamma \vdash \text{dom}(\Delta)
 \end{aligned}$$

1699 *Sequential Composition of Usage Maps.* This section defines the *usage composition operator*, written
1700 $\Delta_1; \Delta_2$. This operator plays a central role in computing the usage of a sequence of terms. Let us
1701 begin with an example.

1702 Consider the sequence “ $(\text{let } r = \text{alloc}(v))^{\Delta_1}; (\text{let } k = \text{get}(r))^{\Delta_2}; \text{free}(r)^{\Delta_3}$ ”. In Δ_1 , we have a
1703 binding “ $y : \text{produced}$ ” because the first instruction produces the resource “ $y : r \rightsquigarrow \text{Cell}$ ”. In Δ_2 ,
1704 we have a binding “ $y : \text{splittedFrac}$ ” because the instruction only reads with y (thus it accepts
1705 any subfraction). In Δ_3 , we have a binding “ $y : \text{uninit}$ ” because the third instruction destroys the
1706 resource y without caring about the value of the Cell.

1707 Let us give three examples of compositions. First, the usage map $\Delta_1; \Delta_2$ contains a binding
1708 “ $y : \text{produced}$ ” because, taken together, the sequential composition of the those two instructions
1709 still creates the resource y . Second, the usage map $\Delta_2; \Delta_3$ contains a binding $y : \text{full}$ because, taken
1710 together, the second and the third instruction consume the Cell, and they read the value that was
1711 contained inside. Third, the usage map $\Delta_1; \Delta_2; \Delta_3$ contains no binding for y because the Cell cannot
1712 be seen from outside the sequence of instruction.

1713 Formally, the usage composition operation $\Delta_1; \Delta_2$ is defined by merging the two usage maps
1714 pointwise by resource name, using the table shown below to compute the “combined usage” in case a
1715

same resource name is bound both in Δ_1 and Δ_2 . For example, Δ_1 contains a binding $y : \text{splittedFrac}$ and Δ_2 contains a binding $y : \text{full}$, then $\Delta_1; \Delta_2$ contains a binding $y : \text{full}$. The \emptyset and \perp entries in the table are explained next.

$\Delta_1; \Delta_2$	\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

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 abstract 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 $y : \text{full}$.

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 $y_2 : (\beta_1 - \gamma_1 - \dots - \gamma_m)H$ 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 (and it is not an uninitialized usage). Therefore, we have to put the usage $t_2 : \text{full}$ 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 `VAL` and `LET` described earlier:

$$\frac{\Gamma, \text{pure} \vdash v : \tau \quad \Delta = \{\mathbf{res} : \text{ensured}\}}{\{\Gamma\} v^\Delta \{\Gamma \otimes [\mathbf{res} : \tau := v]\}} \text{VAL}$$

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

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 by x' in Δ . This renaming is applied on the usage map of the body to follow the renaming in the context.

Since it would be too verbose to explicitly put usage map in each typing rule, the rest of this section will explain how these usage maps are computed in practice without rewriting the rules. 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 , we build Δ'_i to follow renamings in Γ_r .

$$\Delta'_i = \begin{cases} \Delta_i & \text{if } t_i \text{ is of the form } \mathbf{let} \mathbf{res} := t'_i \\ \text{Rename}\{\mathbf{res} := x_i\}(\Delta_i) & \text{otherwise} \end{cases}$$

Then the usage map of the whole sequence is $(\Delta'_1; \dots; \Delta'_n)$.

- 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 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 , the exact usage depends on how the oracle finds a join point between the branch, and if this is made in a way that preserves the invariants of usage map. If the oracle is fully opaque, it is always sound yet imprecise to combine the usage map Δ_0 of the condition expression t_0 to another usage map Δ'_0 that gives a full usage to each linear resource in Γ'_0 and a usage map Δ_3 that contains a produced usage for each linear resource of Γ_3 . For the usage of pure resources, we name Δ_1 (resp. Δ_2) the required usage from t_1 (resp. t_2). Then we take all the pure facts from Γ_3 that are not in Γ'_0 as ensured in a usage map Δ'_3 . This gives $(\Delta_0; \Delta'_0; \Delta_1; \Delta_2; \Delta_3; \Delta'_3)$ as a potential usage map for the if instruction. We leave better usage computation of conditionals for future work.

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 such operation is a quadruplet $(\hat{E}, \hat{F}, \hat{F}', \bar{F})$ where intuitively:

- \hat{E} is the *minimized linear precondition*: a linear context containing resources from $\Gamma.\text{linear}$ 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.
- \bar{F} is the *maximal frame*: a linear context of resources from $\Gamma.\text{linear}$ that were useless in the typechecking of t . It means resources in \bar{F} can be framed during the typechecking of t . Since these resources are not touched by t , they must also occur in $\Gamma'.\text{linear}$.
- \hat{E} 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 $\Gamma.\text{linear}$ in \hat{F} when such a fractional resource suffices to typecheck t .

From the perspective of soundness, the quadruplet $(\hat{E}, \hat{F}, \hat{F}', \bar{F})$ must satisfy the following properties:

- $\{\langle \Gamma.\text{pure}, \hat{E} \mid \hat{F} \rangle\} t \{\langle \Gamma'.\text{pure}, \hat{E} \mid \hat{F}' \rangle\}$
- $\Gamma \Leftrightarrow \langle \Gamma.\text{pure}, \hat{E} \mid \hat{F} \star \bar{F} \rangle$
- $\Gamma' \Leftrightarrow \langle \Gamma'.\text{pure}, \hat{E} \mid \hat{F}' \star \bar{F} \rangle$

From the perspective of completeness, it is desirable for the frame \bar{F} to somehow corresponds to the “maximal” frame. We have not found a way to formally capture the *maximality* aspect of \bar{F} —or, symmetrically, the *minimality* of \hat{F} . It is not clear to us whether these notions could be defined in the general case. What appear clearly is a definition of Minimize ought to at least satisfy the following properties, which capture a partial form of completeness.

- If t can typecheck without a linear resource H , then H should be put in the frame \bar{F} .
- 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 \hat{E} , the resource αH should be placed in \hat{F} , and $(1 - \alpha)H$ should remain in \bar{F} .

1863 Our implementation of Minimize, whose behavior is guided by the entries in the usage map Δ
 1864 computed when typechecking t , satisfy all the aforementioned soundness and partial-completeness
 1865 properties.

1866 5.5 Typechecking of Order-Irrelevant Subexpressions

1868 We next explain how to leverage the minimization procedure for typechecking functions calls that
 1869 are not in A-normal form, but possibly include effectful subexpressions. In C, the arguments of
 1870 a function call may be evaluated in an arbitrary order. Our goal is to ensure that, for well-typed
 1871 OptiTrust programs, the order of evaluation is irrelevant. To that end, we consider a sufficient
 1872 condition: that the arguments can be evaluated in parallel. This condition provides additional
 1873 flexibility for optimizations.

1874 There do exist valid C programs that fail to typecheck in OptiTrust because our condition is
 1875 slightly too restrictive. However, such programs may be easily rewritten by binding variables to
 1876 arguments before the function call. We believe that such a situation would be relatively rare in
 1877 idiomatic C code.

1878 The rule SUBEXPR reduces the typechecking of a term with possibly effectful subexpressions to
 1879 the typechecking of a term whose subexpressions are program variables. In particular, the rule may
 1880 be used to compute the output context associated with a call of the form $f(e_1, \dots, e_n)$ in an input
 1881 context Γ_0 , by reducing the problem to the typechecking of a call of the form $f(x_1, \dots, x_n)$, in an
 1882 input context Γ_p that binds the fresh variables x_i .

1883 The rule SUBEXPR, shown below, applies to a term of the form $e[t_0, \dots, t_n]$, where e denotes a
 1884 multi-evaluation context and where the t_i variables denote the subterms in evaluation position. The
 1885 goal of the rule is to distribute the linear resources from the input context Γ_0 between the subterms
 1886 t_i . If several subterms read the same resource, then this resource needs to be split. If one subterm
 1887 reads a resource and another subterm modifies that same resource, the rule must fail to apply.
 1888 The key idea is to typecheck the subterms one after the other, taking advantage of the Minimize
 1889 operation to remove the minimal amount of resources from the input context, thereby leaving as
 1890 many resources as possible for the remaining subterms.

$$\begin{array}{c}
 1891 \\
 1892 \\
 1893 \\
 1894 \\
 1895 \\
 1896 \\
 1897 \\
 1898 \\
 1899 \\
 1900 \\
 1901 \\
 1902 \\
 1903 \\
 1904 \\
 1905 \\
 1906 \\
 1907 \\
 1908 \\
 1909 \\
 1910 \\
 1911
 \end{array}
 \frac{
 \begin{array}{l}
 \forall i \in [0, n]. \quad \{\Gamma_i\} t_i^{\Delta_i} \{\Gamma'_i\} \wedge (\hat{E}_i, \hat{F}_i, \hat{F}'_i, \bar{F}_i) = \text{Minimize}(\Gamma_i, \Gamma'_i, \Delta_i) \wedge x_i \text{ fresh} \\
 \forall i \in [0, n]. \quad \Gamma_{i+1} = \langle \Gamma_i.\text{pure}, \hat{E}_i \mid \bar{F}_i \rangle \wedge \hat{\Gamma}'_i = \langle \Gamma'_i.\text{pure} \vdash \Delta_i.\text{ensured} \mid \hat{F}'_i \rangle \\
 \Gamma_p = \text{CloseFrac}^{\Delta_p}(\Gamma_{n+1} \otimes \star_{i \in [0, n]} \text{Rename}\{\mathbf{res} := x_i\}(\hat{\Gamma}'_i)) \\
 \{\Gamma_p\} e[x_0, \dots, x_n]^{\Delta_q} \{\Gamma_q\} \\
 \Delta = \text{Rename}\{\mathbf{res} := x_0\}(\Delta_0); \dots; \text{Rename}\{\mathbf{res} := x_n\}(\Delta_n); \Delta_p; \Delta_q
 \end{array}
 }{
 \{\Gamma_0\} e[t_0, \dots, t_n]^\Delta \{\Gamma_q\}
 } \text{SUBEXPR}_\ell$$

1907 *Example application of the rule.* Figure 19 shows the application of the SUBEXPR on an example.
 1908 The columns describe the 4 steps corresponding to the 4 subterms. The rows explain how the
 1909 metavariables from the rule SUBEXPR are instantiated. In particular, observe how the two subexpres-
 1910 sions $\text{get}(p)$ both have read-only access to the same resource H . As the details in the Figure show,
 1911 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 $\text{get}(p)$.

1906 6 JUSTIFYING TRANSFORMATION CORRECTNESS

1907 In this section, we explain how OptiTrust leverages resource typing information to check the
 1908 correctness of the transformations requested by the programmer. The aim of this section is not
 1909 to cover all the transformations implemented in OptiTrust, but to present a representative subset
 1910 thereof. We focus in particular on transformations that leverage the resource information in an
 1911

i	0	1	2	3
Γ_i .pure	$p, q, c : \text{ptr}(\text{int})$	$p, q, c : \text{ptr}(\text{int})$	$p, q, c : \text{ptr}(\text{int})$	$p, q, c : \text{ptr}(\text{int})$ $\alpha : \text{frac}$
Γ_i .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	$\text{get_incr}(c)$ $c : \text{required}$	$\text{get}(p)$	$\text{get}(p)$
Δ_i	$q : \text{required}$ res : ensured	$Hc : \text{full}$ $Hc' : \text{produced}$ res : ensured	$p : \text{required}$ $Hp : \text{splittedFrac}$ res : ensured	$p : \text{required}$ $Hp' : \text{splittedFrac}$ res : ensured
\tilde{E}_i	\emptyset	\emptyset	$\alpha : \text{frac}$	$\beta : \text{frac}$
\tilde{F}_i	\emptyset	$Hc : c \rightsquigarrow \text{Cell}$	$\alpha(p \rightsquigarrow \text{Cell})$	$\beta(p \rightsquigarrow \text{Cell})$
\tilde{F}'_i	\emptyset	$Hc' : c \rightsquigarrow \text{Cell}$	$\alpha(p \rightsquigarrow \text{Cell})$	$\beta(p \rightsquigarrow \text{Cell})$
\tilde{F}_i	$Hp : p \rightsquigarrow \text{Cell}$ $Hq : q \rightsquigarrow \text{Cell}$ $Hc : c \rightsquigarrow \text{Cell}$	$Hp : p \rightsquigarrow \text{Cell}$ $Hq : q \rightsquigarrow \text{Cell}$	$Hp : (1 - \alpha)(p \rightsquigarrow \text{Cell})$ $Hq : q \rightsquigarrow \text{Cell}$	$Hp :$ $(1 - \alpha - \beta)(p \rightsquigarrow \text{Cell})$ $Hq : q \rightsquigarrow \text{Cell}$
$\tilde{\Gamma}'_i$.pure	res := $q : \text{ptr}(\text{int})$	res : int	res : int	res : int
$\tilde{\Gamma}_p$.pure	$p, q, c : \text{ptr}(\text{int}), x_0 := q : \text{ptr}(\text{int}), x_1, x_2, x_3 : \text{int}$			
$\tilde{\Gamma}_p$.linear	$Hp : p \rightsquigarrow \text{Cell}, Hq : q \rightsquigarrow \text{Cell}, Hc' : c \rightsquigarrow \text{Cell}$			

Fig. 19. Example of an application of the SUBEXPR rule on an expression $e[q, \text{get_incr}(c), \text{get}(p), \text{get}(p)]$, in a context $\langle p, q, c : \text{ptr}(\text{int}) \mid Hp : p \rightsquigarrow \text{Cell}, Hq : q \rightsquigarrow \text{Cell}, Hc : c \rightsquigarrow \text{Cell} \rangle$.

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 every basic transformation, we present a generally applicable, *sufficient condition* for the transformation to be correct.

On the one hand, a number of *basic* transformation might seem “trivial” 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.

On the other hand, certain loop transformations might seem “complex” to the reader. We believe this complexity is inherent to advanced loop transformations. In fact, there are even situations where the sufficient conditions that we have considered are *simplified* conditions, that 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 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 maps 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 $\{\Gamma_1\} t^\Delta \{\Gamma_2\}$.

1961 *Groups of Instructions.* Certain transformations operate on groups of consecutive instructions.
 1962 We let the meta-variable T range over a (possibly empty) group of instructions. We generalize our
 1963 alternative syntax by writing $\boxed{\Gamma_1 T; \Delta \Gamma_2}$, where Γ_1 denotes the initial set of resource, Γ_2 denotes
 1964 the final set of resources, and Δ denotes the *composition* of the usage of the instructions from the
 1965 group, as defined in Section 5.3. Formally, if Δ_i denotes the usage of the term t_i , then:

$$1966 \quad \boxed{\Gamma_0 T; \Delta \Gamma_n} \equiv \boxed{\Gamma_0 t_1; \Delta_1 \Gamma_1 t_2; \Delta_2 \Gamma_2 \dots t_n; \Delta_n \Gamma_n} \quad \begin{array}{l} \text{where } T \equiv t_1; t_2; \dots; t_n \\ \text{and } \Delta \equiv \Delta_1; \Delta_2; \dots; \Delta_n \end{array}$$

1967 *Program Contexts.* Transformations generally apply to a program subterm, that is, apply under
 1970 a *program context*. We let the meta-variable \mathcal{E} range over such contexts. For example, evaluating
 1971 a subexpression $1 + 1$ that appears in a context \mathcal{E} is described as the transition from $\mathcal{E}[1 + 1]$ to
 1972 $\mathcal{E}[2]$. We also allow program contexts to denote a hole in the middle of a sequence. For example,
 1973 swapping two instructions that appear inside a sequence is described as the transition from $\mathcal{E}[t_1; t_2]$
 1974 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}'
 1975 denotes the program context associated with the outer sequence that contains $t_1; t_2$. For the first
 1976 few transformations, we will mention a context \mathcal{E} explicitly, however afterwards we will leave it
 1977 implicit.
 1978

1979 *Simple Program Contexts.* Certain transformations operate on subexpressions that appear inside
 1980 an instruction. For those, we may need to restrict the form of the contexts in which the subexpression
 1981 may appear (e.g., to avoid nontrivial control-flow arising from the $\&\&$ operator or from the C ternary
 1982 operator). We therefore introduce the notion of *simple program context*, written $\hat{\mathcal{E}}$, to denote a
 1983 program context that solely consists of possibly-nested functions calls. An example simple context
 1984 is $f(g(\square, 2), g(3, a + 4))$, where \square denotes the hole. One key property for any simple context $\hat{\mathcal{E}}$ is
 1985 that the following rewrite is correct: $\boxed{\hat{\mathcal{E}}[e]_w} \mapsto \boxed{\mathbf{let } x = e; \hat{\mathcal{E}}[x]}$.

1986 Note that the validity of this binding rule, and more generally the interest of simple contexts,
 1987 crucially relies on the hypothesis that the input code typechecks against our typing rules. Indeed, our
 1988 rules ensure that, if a function has multiple arguments, then the available resources are distributed
 1989 across the arguments—only read-only resources can be distributed onto several arguments. For
 1990 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
 1991 term is well-typed, then the effects of $g_2()$ do commute with the effects of $g_1()$ and $g_3()$.
 1992

1993 *Notation for Introducing Ghost Calls.* Recall that a call to a ghost function is an instruction
 1994 that semantically behaves as a no-op, yet updates the set of resources available. In the output of
 1995 transformations, we write $\mathbf{ghost}(\Gamma \longrightarrow \Gamma')$ to mean the insertion of an appropriate ghost function
 1996 call $g()$, such that g admits Γ as precondition and Γ' as postcondition. Concretely, the effect of
 1997 $\mathbf{ghost}(\Gamma \longrightarrow \Gamma')$ is to consume the resources Γ then to produce the resources Γ' .
 1998

1999 We are now ready to present transformations. We begin with transformations on instructions
 2000 and variable bindings, then move on to transformations on storage, and transformations on loops.

2001 6.1 Transformations on Sequences of Instructions

2002 *Moving Instructions.* The basic transformation `Instr.move` allows to move a group of instructions
 2003 to a given destination within the same sequence. Doing so amounts to swapping a group of
 2004 instructions T_1 with an adjacent group of instructions T_2 . The *move* transformation turns a program
 2005 of the form $\mathcal{E}[T_1; T_2]$ into $\mathcal{E}[T_2; T_1]$, where \mathcal{E} denotes a context. The transformation is formalized as
 2006 shown below. The variables Δ_1 and Δ_2 denote the usage associated with T_1 and T_2 . The correctness
 2007 criteria, stated on the right-hand-side, is explained next.
 2008
 2009

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

The expression $\Delta_1.\text{notRO}$ essentially denotes the resources that T_1 modifies. (Technically, this expression denotes the resources that T_1 does not *only* read.) The property $\Delta_1.\text{notRO} \cap \Delta_2 = \emptyset$ captures the idea that if a resource is modified by T_1 , then T_2 must not use it, 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 $\Delta_2.\text{notRO} \cap \Delta_1 = \emptyset$, captures the symmetrical property: if a resource is modified by T_2 , then T_1 must not use it. When both conditions are met, the only resources that both T_1 and T_2 depend on are accessed in read-only mode. In such a situation, the groups of instructions T_1 and T_2 may be safely swapped, without impact on the result of their evaluation.

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 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 context \mathcal{E} .

Intuitively, the deletion operation preserves the semantics of the program if the contents of the resources modified by T are not observed by the rest of the program. If T has been typechecked as “ $\Gamma T; \Delta$ ”, then the set of resources modified by T , written Γ_m , is computed by the filtering operation $\Gamma \vdash \Delta.\text{notRO}$. (Filtering was defined in Section 5.2.) To test the hypothesis that Γ_m is not observed by the rest of the program, we build an auxiliary program, written $\mathcal{E}[G]$, in which the instructions T are replaced with a ghost instruction G that casts the resources from Γ_m into their corresponding “uninitialized form”, written $\text{IntoUninit}(\Gamma_m)$. In other words, if a resource H is modified by T , then the ghost operation G consumes H and produces $\text{Uninit}(H)$.

The transformation can therefore be formalized as follows.

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

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 criteria 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 criteria, 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 $\mathcal{E}[T]$.

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

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.¹¹ A term can be considered

¹¹The 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.

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.

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

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 expression e is that the following program equivalence holds: $\boxed{\text{let } x = e; \hat{\mathcal{E}}[e]}$ \leftrightarrow $\boxed{\text{let } x = e; \hat{\mathcal{E}}[x]}$.

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 an expression e is idempotent, then after the instruction `let x = e`, an instruction `let y = e` 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 context surrounding the code snippet, previously written \mathcal{E} .

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

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.

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

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

$$\boxed{\text{let } x = e; \text{let } y = e; \mathcal{E}[x]} \mapsto \boxed{\text{let } x = e; \text{let } y = e; \mathcal{E}[y]} \quad \text{where } e \text{ is idempotent.}$$

6.3 Transformations on Bindings

Inlining/Binding for Pure Expressions. The basic transformation `Variable.inline_pure` eliminates a binding of the form `let x = v`, where v is a pure 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 pure 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 = e` in programs where x has exactly one occurrence, and this occurrence is contained in the immediately succeeding instruction, under a simple 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.

$$\boxed{\begin{array}{l} \text{let } x = e; \\ \hat{\mathcal{E}}[x]; \end{array}} \mapsto \boxed{\hat{\mathcal{E}}[e];}$$

correct if $\hat{\mathcal{E}}$ is a simple context, with 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 simple 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.

$$\boxed{\begin{array}{l} \text{let } x = e; \\ \hat{\mathcal{E}}[x]; \end{array}} \mapsto \boxed{\begin{array}{l} \text{let } x = e; \\ \text{let } y = e; \\ \hat{\mathcal{E}}[x]; \end{array}} \mapsto \boxed{\begin{array}{l} \text{let } x = e; \\ \text{let } y = e; \\ \hat{\mathcal{E}}[y]; \end{array}} \mapsto \boxed{\begin{array}{l} \text{let } x = e; \\ \hat{\mathcal{E}}[e]; \end{array}}$$

(Instr.dup_let) (Eq.idempotent) (Variable.inline_one)

Inlining a Binding in the Scope of a Sequence. The combined transformation `Variable.inline` eliminates a binding `let x = e` in the general case. If e is a pure 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 = e` without affecting the semantics.

- If x has no occurrences, the effects of e need to be irrelevant to the rest of the program.
- If x has exactly one occurrence, then the effects of e 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, e 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

2157 decomposition of `Variable.inline`, where e is assumed to be idempotent.

2158 $\text{let } x = e; g(); \text{set}(a, x); \text{set}(b, x);$
 2159 $\mapsto g(); \text{let } x = e; \text{set}(a, x); \text{set}(b, x);$ (Instr.swap)
 2160 $\mapsto g(); \text{let } x = e; \text{set}(a, e); \text{set}(b, x);$ (Variable.inline_dup)
 2161 $\mapsto g(); \text{set}(a, e); \text{let } x = e; \text{set}(b, x);$ (Instr.swap)
 2162 $\mapsto g(); \text{set}(a, e); \text{set}(b, e);$ (Variable.inline_one)

2165

2166 We leave to future work the support, in a combined transformation, of more complex patterns
 2167 where occurrences of a non-pure binding appear in depth under control flow constructs.

2168

2169 *Binding Introduction.* The basic transformation `Variable.bind_one` is essentially the reciprocal
 2170 of `Variable.inline_one`. Given an instruction of the form $\hat{\mathcal{E}}[e]$, the transformation introduces a
 2171 binding `let x = e` and turns the instruction into $\hat{\mathcal{E}}[x]$.

2172

2173 $\boxed{\hat{\mathcal{E}}[e];}$ \mapsto $\boxed{\text{let } x = e;$
 2174 $\hat{\mathcal{E}}[x];}$ correct if $\hat{\mathcal{E}}$ is a simple context.

2175

2176 *Folding for Additional Occurrences.* The combined transformation `Variable.bind_dup` is essentially
 2177 the reciprocal of `Variable.inline_dup`. In the scope of a binding `let x = e`, this transformation turns
 2178 $\hat{\mathcal{E}}[e]$ into $\hat{\mathcal{E}}[x]$.

2179

2180 *Common Subexpression Elimination.* The combined transformation `Variable.bind` is essentially
 2181 the reciprocal of `Variable.inline`. Internally, it exploits the transformations `Variable.bind_one` and
 2182 `Variable.bind_dup` to introduce a binding that factorizes the evaluation of common subexpressions.
 2183 For example, if e is idempotent and commutes with $g()$, the program “ $g(); \text{set}(a, e); \text{set}(b, e)$ ” can
 2184 be transformed into “`let x = e; g(); set(a, x); set(b, x)`”.

2184

2185 6.4 Transformations on Storage

2186

2187 The purpose of this section is to present transformations for introducing, eliminating, and converting
 2188 between various forms of storage. We present transformations operating on single cells, and omit
 2189 from the discussion the generalizations to arrays and N -dimensional matrices.

2189

2190 Recall from Section 3 that a pure program variable `const int x = 3` is represented in the OptiTrust
 2191 AST as `let x = 3`, that a non-pure stack-allocated variable `int x = 3` is represented as `let x = ref(3)`,
 2192 and that an uninitialized variable `int x` is represented as `let x = ref_uninit()`. For stack-allocated
 2193 data, the resources produced by `ref` are automatically reclaimed at the end of the scope. For heap-
 2194 allocated data, the resources produced by `alloc` are consumed by the matching call to `free`.

2194

2195 *Separating Declaration from Initialization.* For a stack-allocated variable, the basic transforma-
 2196 tion `Variable.init_detach` separates its declaration from its initialization. This transformation
 2197 is useful as a preliminary step for the combined transformation that hoists a variable declara-
 2198 tion appearing inside a loop into an array allocated outside that loop. The basic transformation
 2199 `Variable.init_attach` applies the reciprocal operation.

2200

2200 $\boxed{\text{let } x = \text{ref}(e);}$ \leftrightarrow $\boxed{\text{let } x = \text{ref_uninit}(); \text{set}(x, e);}$

2201

2202 *Converting between Stack and Heap Allocation.* The basic transformation `Variable.to_heap` trans-
 2203 forms an uninitialized stack-allocated storage into a corresponding heap-allocated storage. The
 2204 transformation takes as optional argument the target at which the free instruction should be

2205

2206 inserted; by default, it is placed at the end of the scope. The reciprocal transformation is named
 2207 `Variable.to_stack`. In the statement below, n denotes the size of the type of x .

$$2208 \quad \boxed{\{T_1; \text{let } x = \text{ref_uninit}(); T_2;\}} \leftrightarrow \boxed{\{T_1; \text{let } x = \text{alloc}(n); T_2; \text{free}(x);\}}$$

2210 *Removal of Unused Storage.* If a stack-allocated storage is never used, it may be removed by means
 2211 of the operation `Instr.delete`. Concretely, the instruction `let x = ref_uninit()` may be deleted if
 2212 x has no occurrences, and the instruction `let x = ref(e)` may be deleted if moreover the effects
 2213 performed by e are not observed by the rest of the program.

2214 If a heap-allocated space is never used, then it may also be removed. To that end, one needs
 2215 to delete both the `alloc` and the corresponding `free` instructions. Neither of them can be removed
 2216 independently, because both depend on each other. However, if we move using `Instr.move` the `alloc`
 2217 instruction next to the `free` instruction, or vice-versa, then the group made of the two instructions
 2218 may be removed at once by means of `Instr.delete`. The combined transformation `Variable.delete`,
 2219 described below, performs this task.

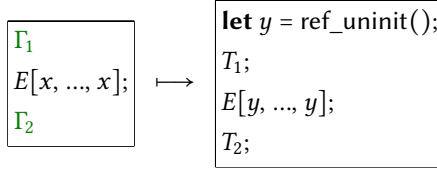
$$2221 \quad \boxed{\begin{array}{l} \text{let } x = \text{alloc}(e); \\ T; \\ \text{free}(x); \end{array}} \xrightarrow{\text{(Instr.move)}} \boxed{\begin{array}{l} \text{let } x = \text{alloc}(e); \\ \text{free}(x); \\ T; \end{array}} \xrightarrow{\text{(Instr.delete)}} \boxed{T}$$

2222 *Temporary Alternative Storage.* The transformation `Variable.local_name` is the most complex that
 2223 we have implemented in terms of operations on plain sequences of instructions. The transformation
 2224 `Variable.local_name` operates over a specified group of instructions, say T , for a specified storage,
 2225 say x . Over this scope, a fresh storage, call it y , is allocated. Just before executing T , the contents
 2226 of x are copied into y . All instructions from T are updated to use y instead of x . Just after these
 2227 instructions, the possibly-updated contents of y is copied into x . Depending on the situation, the
 2228 initial copy from x to y , or the final copy from y into x might be unnecessary—and even ill-typed.
 2229 Such unnecessary copy operations are omitted.

2230 The variable x may be allocated either on the stack or on the heap. The user may choose to
 2231 allocate y on the stack or on the heap. Moreover, our implementation supports the general case
 2232 where x is not just a variable but an N -dimensional matrix. In case where x is a matrix, y may
 2233 correspond to only a subset (i.e., a tile) of the matrix. The interest of the `local_name` transformation
 2234 is to enable the program to operate on a local piece of data. Crucially, the memory layout of this
 2235 data may be refined by subsequent transformations, for example to store the transposed of a matrix
 2236 in a cache friendly way (as in Section 2.3), or to enable vectorization.

2237 The transformation is described in Figure 20. There, the group of instructions T is represented as
 2238 $E[x, \dots, x]$, i.e., as a context with multiple occurrences of x . The context Γ_1 describes the resources
 2239 available before T , and Γ_2 the resources available after T . This typing information is used not only
 2240 for checking the correctness criterion, but also for guiding the generation of the output code.

2241 The correctness criterion appear at the bottom of Figure 20. An essential aspect of this criterion
 2242 is to check that, during the execution of T , the resource R corresponding to the full permission on
 2243 x is *frozen* (i.e., made unavailable) in order to ensure that no operation may be performed on x via
 2244 potential aliases of this pointer. The first ghost call uses a standard technique for enforcing such a
 2245 *freeze* in Separation Logic: introducing a *magic wand* operator (\multimap), guarded by a token named H
 2246 in the postcondition of the ghost operation, and bound as H' in the rest of the sequence. The heap
 2247 predicate H' admits the type `Hprop`, which is the type of all heap predicates in Separation Logic.
 2248 This heap predicate H' serves the role of a *key* for unfreezing R at the desired point—here, the end
 2249 of the scope on which y is used in place of x , where the second ghost call is placed. As far as the
 2250
 2251
 2252
 2253
 2254



2262 where E is a multi-hole context with one hole per occurrence of x , and where:

2263

2264

2265

2266

2267

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

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

2268 correct if the program to the

2269 right typechecks successfully,

2270 where R is one of:

- 2271
- 2272
- 2273
- 2274
- 2275
- 2276
- $(x \rightsquigarrow \text{Cell})$
 - $\alpha(x \rightsquigarrow \text{Cell})$
 - $\text{Uninit}(x \rightsquigarrow \text{Cell})$
- 2277 according to how " $x \rightsquigarrow \text{Cell}$ "
- 2278 appears in Γ_1 .

2268

2269

2270

2271

2272

2273

2274

2275

2276

2277

2278

2279

2280

2281

2282

2283

2284

2285

2286

2287

2288

2289

2290

2291

2292

2293

2294

2295

2296

2297

2298

2299

2300

2301

2302

2303

```

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

```

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

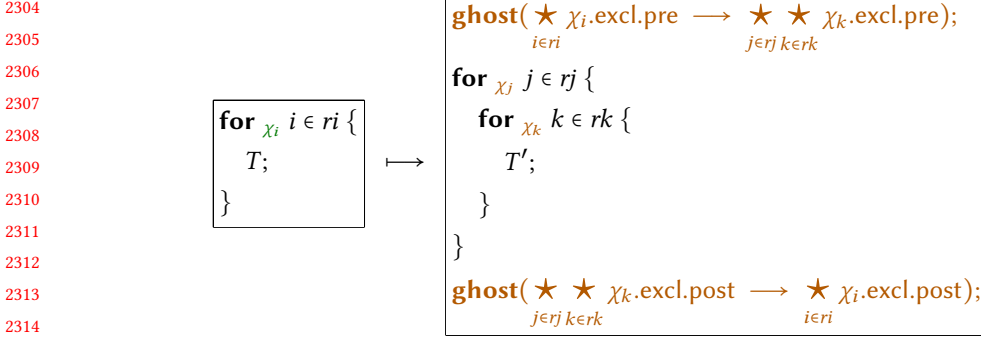
present paper is concerned, the magic wand operator can be viewed as a binary operator on heap predicates whose definition needs not be revealed to the user.

6.5 Transformations on Loops

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

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 ri , rj and rk , 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: $\star_{i \in ri} \chi_i.\text{excl.pre}$ whereas the loop on j consumes instead: $\star_{j \in rj} \star_{k \in rk} \chi_k.\text{excl.pre}$.



2315 where:

2316 $T' \equiv \text{Subst}\{i := \text{RecoverIndex}(j, k)\}(T)$

2317 $\chi_k \equiv \text{Subst}\{i := \text{RecoverIndex}(j, k)\}(\chi_i)$

2318

2319
$$\chi_j \equiv \begin{cases} \text{vars} \equiv \chi_i.\text{vars} \\ \text{shrd} \equiv \text{Subst}\{i := \text{RecoverIndex}(j, rk.\text{start})\}(\chi_i.\text{shrd}) \\ \text{excl} \equiv \{\text{pre} \equiv \star_{k \in rk} \chi_k.\text{excl.pre}; \text{post} \equiv \star_{k \in rk} \chi_k.\text{excl.post}\} \end{cases}$$

2320

2321

2322

2323 with the following possible instantiations for the ranges:

2324

Variant	Range ri	Range rj	Range rk	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

2325 Fig. 21. Description of the 4 variants of the basic transformation `Loop`. tile.

2326

2327 *Loop Interchange.* The basic transformation `Loop.swap` allows interchanging (i.e. swapping) two

2328 loops. It is described at the top of Figure 22. There exists a general criterion capturing when two

2329 loops may be swapped, however this criterion requires reasoning about the resources required by

2330 specific iterations, e.g. $T_{i,j}$ and $T_{i',j'}$ with $i' > i$ and $j > j'$. Instead, we focus on two conditions that

2331 are simpler nevertheless sufficient for many practical situations: if at least one of the outer loop or

2332 the inner loop is parallelizable, then swapping the two loops is correct. Figure 22 describes the case

2333 where the outer loop is parallelizable. The case where the inner loop is parallelizable, not shown, is

2334 treated with just a few changes.

2335 The first step is to partition the resources from the inner loop contract depending on where they

2336 come from relative to the resources from the outer loop. We name partitions by using the first letter

2337 to denote its inner loop origin, and the second letter to denote its outer loop origin. We use I for

2338 invariant, R for shared reads, P for exclusive precondition and Q for exclusive postcondition. For

2339 example, the inner shared reads are partitioned into RP_i that comes from the outer precondition,

2340 and RR that comes from the outer shared reads. Internally, we rely on the subtraction operation for

2341 computing the partitions (additional details may be found in Section ??).

2342 Then, we appropriately place the resources obtained from the partitioning in the contracts χ'_i

2343 and χ'_j associated with the swapped loops. Compared with χ_j , the new contract χ'_j essentially adds

2344 a \star_i operator to certain components. Compared with χ_i , the new contract χ'_i removes occurrences

2345 of the \star_j operators. Note that the loop on i remains parallelizable. Around the new loop nest, a

2346

2347

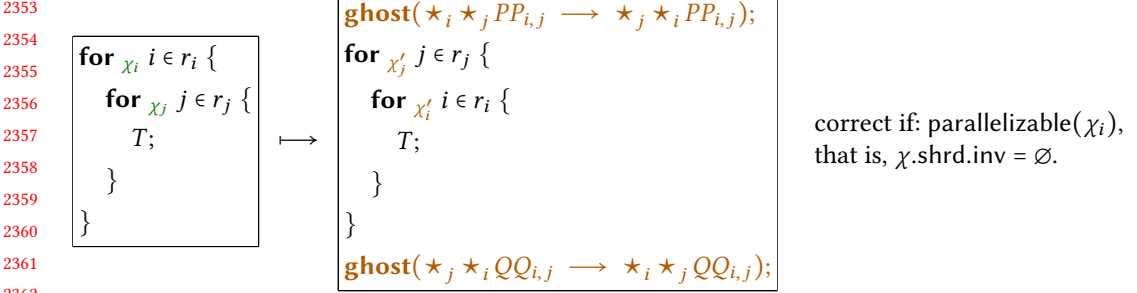
2348

2349

2350

2351

2352



2363 The contracts from the input code are decomposed as follows:

2364 $\chi_j.\text{shrd} = \{\text{inv} = (IP_{i,j} \star IR_j), \text{reads} = (RP_i \star RR)\}$
 2365 $\chi_j.\text{excl} = \{\text{pre} = (PP_{i,j} \star PR_j), \text{post} = (QQ_{i,j} \star PR_j)\}$
 2366 $\chi_i.\text{shrd} = \{\text{inv} = \emptyset, \text{reads} = (\star_j PR_j \star IR_{r_j.\text{start}} \star RR)\}$
 2367 $\chi_i.\text{excl} = \{\text{pre} = (\star_j PP_{i,j} \star IP_{i,r_j.\text{start}} \star RP_i), \text{post} = (\star_j QQ_{i,j} \star IP_{i,r_j.\text{end}} \star RP_i)\}$

2369 The contracts for the output code are built as follows:

2370
2371 $\chi'_i \equiv \left\{ \begin{array}{l} \text{vars} \equiv \chi_j.\text{vars} \\ \text{shrd} \equiv \{\text{inv} \equiv \emptyset, \text{reads} \equiv (PR_j \star IR_j \star RR)\} \\ \text{excl} \equiv \{\text{pre} \equiv (PP_{i,j} \star IP_{i,j} \star RP_i), \text{post} \equiv (QQ_{i,j} \star IP_{i,r_j.\text{next}(j)} \star RP_i)\} \end{array} \right.$
 2372
2373
2374
2375 $\chi'_j \equiv \left\{ \begin{array}{l} \text{vars} \equiv \chi_j.\text{vars} \\ \text{shrd} \equiv \{\text{inv} \equiv (\star_i IP_{i,j} \star IR_j), \text{reads} \equiv (\star_i RP_i \star RR)\} \\ \text{excl} \equiv \{\text{pre} \equiv (\star_i PP_{i,j} \star PR_j), \text{post} \equiv (\star_i QQ_{i,j} \star PR_j)\} \end{array} \right.$
 2376
2377
2378

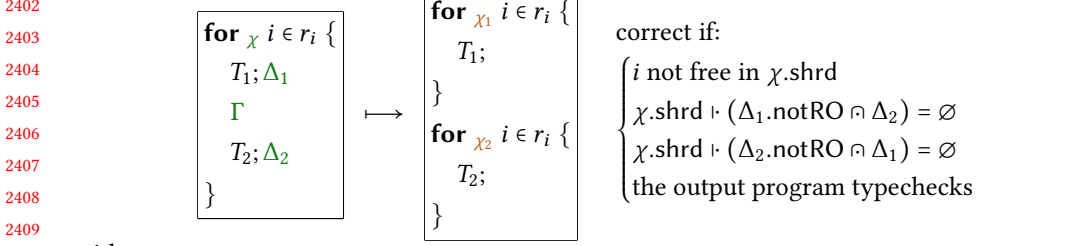
2379 Fig. 22. The basic transformation `Loop.swap`, in the particular case where the outer loop is parallelizable.

2380 pair of ghost operations is inserted for swapping groups of resources—a necessary step to match
 2381 the resources required by the new loop nest.
 2382

2383 *Loop Fission.* The basic transformation `Loop.fission` breaks a loop with body $T_1; T_2$ into two loops
 2384 over the same range, a first loop with body T_1 , and a second loop with body T_2 . The transformation
 2385 is described in Figure 23. As for loop swapping, there exists a general correctness criteria expressed
 2386 using inequalities on indices, but for now we focus on a simpler yet practical criteria.
 2387

2388 Our criterion asserts that loop fission is correct if the resources modified by T_1 at any iteration i
 2389 do not interfere with the resources modified by T_2 at any other iteration $i' \neq i$. To implement this
 2390 check, we inspect the usage maps Δ_1 and Δ_2 associated with T_1 and T_2 , respectively. If T_1 modifies
 2391 one resource from $\chi.\text{shrd}$, then T_2 must not use this same resources; symmetrically, if T_2 modifies a
 2392 resource, then T_1 must not use it. Note, however, that T_1 and T_2 are allowed to both read the same
 2393 resource; moreover, the resources exclusively modified or produced by T_1 at the i -th iteration of
 2394 the first loop may be consumed by T_2 at the i -th iteration of the second loop.
 2395

2396 There remains to explain how to synthesize the contracts χ_1 and χ_2 , associated with the two
 2397 generated loops, from the original contract χ . For `shrd` resources, we project the subsets of $\chi.\text{shrd}$
 2398 resources used by T_1 and T_2 . For `excl` resources, we need to synthesize the invariant at the cut point,
 2399 written R . The first loop takes the exclusive resources from $\chi.\text{excl.pre}$ to R , whereas the second
 2400 loop takes the exclusive resources from R to $\chi.\text{excl.post}$. At a high level, this set of resources R
 2401 is computed by subtracting the shared resources as well as the local allocations from T_1 , described by



2410 with:

2411 $_ , F \equiv \Gamma \boxplus \chi.\text{shrd}.\text{inv}$ $_ , R \equiv F \boxplus \text{StackAllocCells}(T_1)$

2412

2413 $\chi_1 \equiv \left\{ \begin{array}{l} \text{vars} \equiv \chi.\text{vars} \\ \text{shrd} \equiv \chi.\text{shrd} \vdash \Delta_1 \\ \text{excl} \equiv \{\text{pre} \equiv \chi.\text{excl}.\text{pre}, \text{post} \equiv R\} \end{array} \right.$ $\chi_2 \equiv \left\{ \begin{array}{l} \text{vars} \equiv \chi.\text{vars} \\ \text{shrd} \equiv \chi.\text{shrd} \vdash \Delta_2 \\ \text{excl} \equiv \{\text{pre} \equiv R, \text{post} \equiv \chi.\text{excl}.\text{post}\} \end{array} \right.$

2414 where $\chi.\text{shrd} \vdash X$ is a shorthand for $\{\text{inv} = (\chi.\text{shrd}.\text{inv} \vdash X), \text{reads} = (\chi.\text{shrd}.\text{reads} \vdash X)\}$.

2415 Fig. 23. The basic transformation `Loop.fission`.

2416

2417

2418

2419

2420

2421

2422

2423

2424

2425

2426

2427

2428

2429

2430

2431

2432

2433

2434

2435

2436

2437

2438

2439

2440

2441

2442

2443

2444

2445

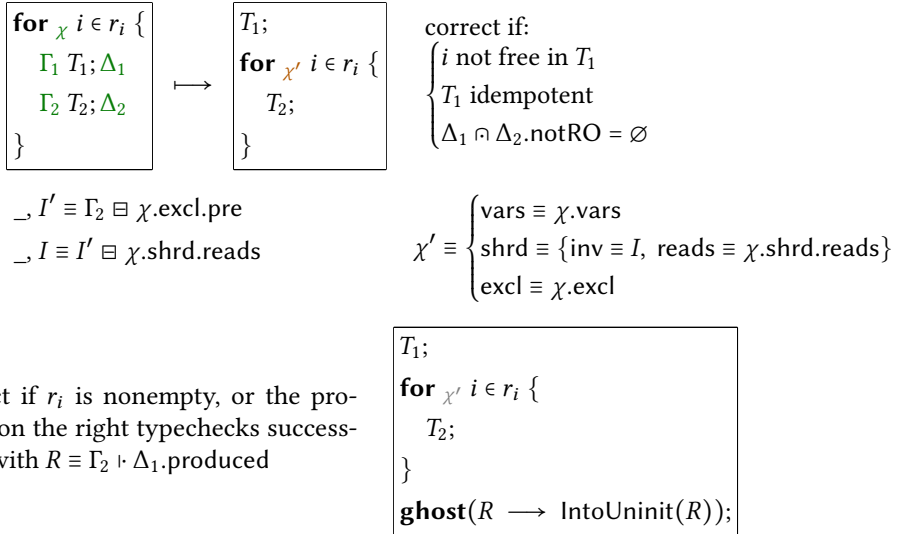
2446

2447

2448

2449

2450



2436 Fig. 24. The basic transformation `Loop.move_out`.

2437 $\chi.\text{shrd}$ and $\text{StackAllocCells}(T_1)$, from the context Γ computed by our typechecker at the location just between T_1 and T_2 .

2438 Observe that the loop contracts χ_1 and χ_2 generated by the loop fission transformation may contain a larger set of resources than strictly necessary. We describe further on, in Section 6.6, a procedure for minimizing loop contracts.

2443 *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.

2444 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

2451 the remaining instructions of the loop, described by T_2 (that is, the condition $\Delta_1 \sqcap \Delta_2.\text{notRO} = \emptyset$).
 2452 Note that, contrarily to the `Instr.move` criteria, it is safe for T_2 to read resources modified by T_1 .
 2453 Additionally, we describe how to check whether T_1 is deletable in the case of a possibly empty loop
 2454 as typechecking an auxiliary program, as we have done before.

2455 *Other Loop Transformations.* There are other important loop transformations that we support.

- 2457 • `Loop.fusion` (reciprocal of `Loop.fission`): it fuses two consecutive loops into a single one.
- 2458 • `Loop.collapse` (reciprocal of `Loop.tile`): it collapses two nested loops into a single one.
- 2459 • `Loop.hoist_alloc`: hoist a variable allocated inside a loop into an array allocated outside
 2460 the loop; more generally, it hoists a matrix of dimension N allocated inside a loop into a
 2461 matrix of dimension $N + 1$ allocated outside the loop.
- 2462 • `Loop.shift`: reindex a loop by applying a positive or negative offset to its values.
- 2463 • `Loop.scale`: reindex a loop using an index that takes either smaller or larger steps.
- 2464 • `Loop.extend_range`: that extends the range of a loop by wrapping its body in a conditional.
- 2465 • `Loop.unroll`: unrolls a loop whose range is statically known.
- 2466 • `Loop.parallel`: to set or unset a parallel flag on a loop using our parallelizable criteria.

2468 6.6 Transformations on Annotations

2469 *Modification to Contracts and Ghost Code.* The semantics of a program is fully determined by
 2470 its actual C code: it does not depend in any way on the *ghost* code nor on the function and loop
 2471 contracts. Therefore, contracts may be freely modified, and ghost instructions may be freely inserted,
 2472 deleted, or modified. The requirement is to reach, after one or several updates, a set of annotations
 2473 for which the typechecking of the same C code succeeds. These modifications can be applied
 2474 performed either directly by the programmer, or during the evaluation of transformations—without
 2475 need for specific correctness criteria.

2476 *Minimization of Loop Contracts.* The aforementioned loop transformations produce correct re-
 2477 source annotations, yet these annotations might be suboptimal for later transformations. Typically,
 2478 the generated loop contracts would include clauses covering a set of resources possibly larger than
 2479 strictly necessary. For example, after the basic loop fission transformation, the contract of the first
 2480 loop would typically mention resources that are in fact only used by the instructions from the sec-
 2481 ond loop. Mentioning unnecessary resources in a contract may impede the applicability of further
 2482 transformations. OptiTrust therefore includes a procedure, implemented as a basic transformation,
 2483 to *minimize* loop contracts. OptiTrust’s combined transformations for loops systematically include
 2484 a call to this procedure.

2485 The loop contract minimization procedure takes as input a loop with contract χ , and updates this
 2486 contract to χ' , without modifying the code. The procedure depends on the usage map Δ computed
 2487 for the instructions T that constitute the loop body.

$$2488 \quad \boxed{\text{for } \overset{\pi}{\chi} i \in r_i \{T; \Delta\}} \quad \mapsto \quad \boxed{\text{for } \overset{\pi}{\chi'} i \in r_i \{T\}}$$

2489 Intuitively, the contract χ' is obtained by filtering out and by weakening resources from χ ,
 2490 depending on their usage in Δ . First, if a resource is unused by T and thus is absent from Δ or has
 2491 usage `joinedFrac`, then it is excluded from χ' . As a result, certain variables that were quantified
 2492 in χ might no longer have occurrence in χ' , hence they can be removed as well. Second, if a
 2493 resource appears with fraction 1 in χ , yet this resource is marked as `splittedFrac` in Δ , then this
 2494 resource is replaced with a read-only version of it. Technically, an additional fraction variable
 2495 must be quantified in χ' , and this fraction variable is used for describing the resource as read-only.

2499

2500 Internally, the implementation of contract minimization reuses our *minimization of triple* procedure
2501 (Section 5.4 and ??). Details may be found in ??.

2502 *Moving and Cancellation of Ghost Instructions.* OptiTrust includes a transformation that attempts
2503 to remove pairs of ghost transformations that cancel each other. Indeed, the sequence $\mathbf{ghost}(H \rightarrow$
2504 $H')$; $\mathbf{ghost}(H' \rightarrow H)$ is equivalent to a no-op. More generally, the user as well as combined
2505 transformations may request a ghost instruction to be moved so as to be (logically) executed as
2506 early as possible in the program; or, symmetrically, to be executed as late as possible. Moving
2507 ghost instructions in such a way may lead to the apparition of pairs of cancellable pairs of ghost
2508 instructions; and, even when ghost instructions do not disappear, moving them away from, e.g., a
2509 loop kernel, may unlock certain transformations.
2510

2511 7 RELATED WORK

2512 The most closely related frameworks were discussed in the introduction. In this section, we comment
2513 on the remaining related work, focusing in turn on each of the ingredients that constitute OptiTrust.
2514

2515 *General-Purpose Compilers.* General purpose compilers such as GCC or ICC are able to apply
2516 a large class of program optimizations, from the classic ones such as inlining, dead code elim-
2517 ination, move of instructions to more advanced ones such as loop fission, loop fusion, or loop
2518 reordering. The same transformations are available in OptiTrust, yet with three major differences.
2519 First, general-purpose compilers apply these transformations on an intermediate representation.
2520 In contrast, OptiTrust applies it at the source level, allowing to produce human-readable feed-
2521 back. Second, general-purpose compilers relies on fully-automated procedures, often guided by
2522 heuristics, to determine what transformations to apply. In contrast, OptiTrust transformations
2523 are fully controlled by the programmer, either directly via basic transformations, or indirectly
2524 via combined transformations. Third, general-purpose compilers rely on static analysis applied
2525 to plain C code to determine whether certain transformations are applicable, and as a result may
2526 lack information to trigger a transformation. In contrast, OptiTrust leverages expressive resource
2527 typing information to justify the correctness of transformations, significantly enlarging the set of
2528 applicable transformations.

2529 *Guidance in General-Purpose Compilers.* To introduce human guidance in general-purpose com-
2530 pilers, a common approach is to insert *pragmas* into the code. For example, Scout [Krzikalla et al.
2531 2011] is a pragma-based tool for guiding source-to-source transformations that introduce vector
2532 instructions. As another example, Radtke and Weinzierl [2024] makes use of C++ attributes for
2533 switching between array-of-structures and structures-of-arrays over the scope of specific computa-
2534 tion kernels; the compiler automatically inserts instructions for copying the data before and after
2535 the loop. A similar approach could be expressed as an OptiTrust transformation, by composing the
2536 `local_name` transformation for arrays (discussed in Section 6.4) with the `aos_to_soa` transformation
2537 (not discussed in the paper). The main limitation of *pragma*-based approaches is that they are
2538 ill-suited for describing sequences of optimizations. Indeed, there is no easy way to attach a pragma
2539 to a line of code that is generated by a first optimization. Kruse and Finkel [2018] suggest the
2540 possibility to stack up pragmas, by providing labels as additional pragma arguments: a second
2541 pragma may refer to the labels introduced by the transformation described in a first pragma. Yet,
2542 this approach does not scale up well beyond a handful of successive transformations. OptiTrust, in
2543 contrast, supports chains of dozens of transformations.
2544

2545 *Domain-Specific Compilers.* Another possible approach to overcome the limitations of general-
2546 purpose compilers is to leverage *domain specific languages* (DSL), such as Halide [Ragan-Kelley et al.
2547 2013], TVM [Chen et al. 2018], Fireiron [Hagedorn et al. 2020a] (used at Nvidia), PartIR [Alabed
2548

2549 et al. 2024] (used at DeepMind). Specialized compilers can benefit from carefully tuned heuristics.
2550 Yet, even for programs expressed in a specific DSL, the optimization search space remains vast,
2551 hence programmer guidance is key to achieve good performance. Halide and its descendants makes
2552 use of a script, called a *schedule*, for guiding the compilation strategy.

2553 For DSLs, the language restriction is also their Achilles' heel: as soon as the user's application
2554 requires a single feature that falls outside of what the DSL can express, the programmer loses
2555 most if not all of the benefits of the DSL. In practice, DSLs typically support the possibility to
2556 include foreign functions (or, inlined general-purpose code), however these foreign functions must
2557 be treated as black box by the DSL compiler, preventing the applications of any domain-specific
2558 optimization across this black box.

2559 In contrast to DSLs, OptiTrust sticks to a standard, general-purpose language. At the same time,
2560 OptiTrust retains the ability to manipulate domain-specific operations and to exploit transforma-
2561 tions that are specific to these operations, as illustrated with the *reduce* function in our OpenCV
2562 case study. At any point in the transformation script, an occurrence of a domain-specific operation
2563 may be lowered into standard C code, thereby enabling further lower-level optimizations.

2564 *Code Transformations via Rewrite Rules.* A rewrite rule maps a code pattern to another code
2565 pattern. A number of tools exploit rewrite rules to perform source-to-source transformations. For
2566 example, TXL [Cordy 2006] is a multi-language rewrite system, whose patterns are expressed at
2567 the level of syntax, using grammars. Coccinelle [Lawall and Muller 2018] allows the programmer
2568 to describe *semantic patches* in C code. CodeBoost [Bagge et al. 2003] applies the Stratego program
2569 transformation language [Bravenboer et al. 2008] to C++ code. CodeBoost can be used to turn
2570 high-level operations on matrices and vectors into typical high-performance source code.

2571 OptiTrust provides a much more expressive language for describing transformations, going far
2572 beyond rewrite rules. Although many transformations *can* be encoded as rewrite rules, the encoding
2573 involved can be cumbersome or inefficient. For example, reconstructing a for-loop for a series of
2574 similar blocks of instructions can be encoded via rewrite rules, yet the blocks must be merged into
2575 the for-loop one by one. Other transformations, especially those involving contracts, would be
2576 challenging to express as rewrite rules. For example, *loop contract minimization* (Section 6.6) would
2577 require the rewrite rule to depend on side-conditions and meta-operations that involve resources
2578 and usage maps.

2580 *Intermediate Languages.* The use of an intermediate language with simpler semantics is common-
2581 place, both in the domain of compilation and in the domain of program verification. Let us cite
2582 a few examples. The Common Intermediate Language (CIL) serves as intermediate compilation
2583 language for the whole .NET ecosystem [Gough and Gough 2001]. Why3 [Filliâtre and Paskevich
2584 2013] serves as intermediate verification language for C, Java, and Ada programs. Viper [Müller
2585 et al. 2017] serves as intermediate verification language for Java, Rust, Go, OpenCL, etc. Although
2586 intermediate languages are commonplace, We are not aware of any framework that leverages a
2587 translation into an intermediate language *and* provides a reciprocal translation back to the source
2588 language, with a round-trip property such as that provided by OptiTrust.

2590 *Source Code Manipulation Frameworks.* Frameworks that offer more expressiveness than rewrite
2591 rules generally give access to the abstract syntax tree (AST) of the source code. Traditional compil-
2592 ers employ an AST, but they are not designed for synthesizing pieces of AST at the source level.
2593 Moreover, traditional compilers operate on intermediate representations, and lose the structure of
2594 the original code. These two limitations of general-purpose compilers have motivated the devel-
2595 opment of frameworks that are specifically designed to support code transformations (and code
2596 analyses) at the level of C code. ROSE [Quinlan 2000; Quinlan and Liao 2011] and Cetus [Bae et al.

2597

2598 2013; Dave et al. 2009] are two such frameworks that provide facilities for manipulating C ASTs.
2599 Source-to-source transformation frameworks have also been employed to produce code target-
2600 ing GPUs [Amini 2012; Konstantinidis 2013; Lebras 2019]. These frameworks implement generic
2601 optimization strategies, in a similar fashion as general-purpose compilers. In contrast, OptiTrust
2602 leverages transformation scripts to guide the optimization of a specific program. Moreover, the
2603 OptiTrust infrastructure supports resource typing, which provides much more precise information
2604 than the classic static code analyses implemented in the frameworks such as ROSE and Cetus.

2605 *Transformation Scripts.* Expressing a series of source-level transformations for a specific program
2606 can be done by means of a transformation script. Such scripts have appeared in particular in the
2607 context of polyhedral transformations [Bagnères et al. 2016b; Bondhugula et al. 2008], for example
2608 in Loopy [Namjoshi and Singhanian 2016] and in work by Zinenko et al. [2018a]. CHILL [Chen et al.
2609 2008; Rudy et al. 2011] includes transformations that go beyond the polyhedral model. It has been
2610 applied to generate finely tuned CUDA code from high-level linear algebra kernels. POET [Yi and
2611 Qasem 2008; Yi et al. 2014] is a scripting language for performing program transformations, for
2612 C/C++ as well as other languages. POET has been employed to generate optimized code for linear
2613 algebra kernels, including semi-automated exploration of a search space of possible optimizations.

2614 Several pieces of work already discussed in the introduction exploit transformation scripts.
2615 Halide [Ragan-Kelley et al. 2013], TVM [Chen et al. 2018] feature schedules that can be viewed as
2616 transformation scripts. Elevate [Hagedorn et al. 2020b] expresses the transformation script in the
2617 form of a composition of functions. ATL [Liu et al. 2022] leverages “tactic”-based proof scripts as
2618 support for expressing transformations scripts. LARA consists of a transformation script featuring
2619 declarative queries as well as arbitrary JavaScript instructions.

2620 All this related work demonstrates a strong interest in leveraging transformation scripts for
2621 putting control of optimizations in the hand of the programmer. Systems differ in what language
2622 they targeted, and what transformations they support. None of the aforementioned systems support
2623 in their transformation scripts a system for targeting program points with the expressiveness and
2624 conciseness offered by OptiTrust targets. Moreover, as far as we know, LARA [Silvano et al. 2019]
2625 and OptiTrust are the only two frameworks making use of transformation scripts for applying
2626 general-purpose transformations at the level of C code. OptiTrust is the first to demonstrate the
2627 use of transformation scripts to produce high-performance code for state-of-the-art benchmarks.
2628 Moreover, unlike LARA, OptiTrust checks that the transformations requested by the programmer
2629 preserve the semantics of the code.

2630 *Proof-Transforming Compilation.* The notion of *Proof Carrying Code* [Necula 1998] refers to the
2631 idea that compilers could be instrumented to carry invariants from high-level source code down
2632 to low-level code. The original line of work on Proof Carrying Code did not aim at full functional
2633 correctness properties, but rather focused on simpler invariants capturing safety properties, such
2634 as the absence of out-of-bound accesses.

2635 Subsequent work introduced the notion of *Proof-Transforming Compilation* to refer to a compiler
2636 that takes as input a formally-verified program and produces as output compiled code accompanied
2637 by a formal proof (i.e., a proof tree in a program logic) that the compiled code satisfies the same
2638 functional correctness specification as the input program. In particular, the PhD work of César
2639 Kunz [Barthe et al. 2009; Kunz 2009] shows how to realize proof-transforming compilation for
2640 standard compiler optimizations, applied at the level of the RTL intermediate language. More
2641 recently, the work on Alpinist [Sakar et al. 2022] demonstrates the feasibility, for a small number
2642 of GPU-oriented optimizations, of transforming GPU code while preserving logical invariants.

2643 Our results obtained so far with OptiTrust demonstrates the feasibility, for a fair number of
2644 general-purpose code optimizations, of transforming C code while preserving resource-based
2645

2646

2647 invariants. In future work, we look forward to extending OptiTrust in order to handle richer logical
2648 invariants and to produce optimized programs accompanied by formal proofs of correctness.

2649

2650 *Separation Logic.* OptiTrust leverages a standard Separation Logic. The most closely related
2651 program logics are VST [Cao et al. 2018], a program verification tool for C, and RefinedC [Sammler
2652 et al. 2021], a very expressive type system for C. Both these systems are grounded on the Iris
2653 framework [Jung et al. 2018a,b], at this day the most advanced formalization of Concurrent Sep-
2654 aration Logic. Other tools, such as Alpinist [Sakar et al. 2022] leverage Viper’s *dynamic frames*
2655 technique [Müller et al. 2017], a cousin of Separation Logic.

2656 Fractional resources [Boyland 2003] are nowadays considered a standard ingredient of Separation
2657 Logic [Jung et al. 2018a]. Following common practice, OptiTrust leverages the notion of fractional
2658 resources to describe read-only resources. The technique of making fractions essentially transparent
2659 to the end-user is directly inspired by the work by Heule et al. [2013] implemented in the Chalice
2660 verification tool.

2661 OptiTrust is, as far as we know, the first framework to compute usage information, describing
2662 precisely how the Separation Logic resources available during the evaluation of a subterm are
2663 actually used by that subterm; and the first framework to exploit Separation Logic information
2664 for guiding the behavior of programmer-requested code transformations (as, e.g., in `local_name`,
2665 described in Section 6.4).

2666

2667 *Contract Inference.* OptiTrust currently requires the programmer to annotate the input program
2668 with function and loop contracts. One may wonder the extent to which such contracts could be
2669 automatically inferred, at least for reasonably simple programs. Journault and Miné [2018] show
2670 that, by leveraging abstract interpretation, for functions such as matrix-multiplication or similar
2671 linear algebra operations, full functional correctness specifications can be automatically computed.
2672 Spies et al. [2024] leverage *bi-abduction* in Separation Logic (a technique at the heart of the *infer*
2673 automated program analysis tool [Calcagno et al. 2019]) to demonstrate automatic inference of
2674 specifications in array- and pointer-manipulating functions of up to a dozen lines of C code. We
2675 leave it to future work to integrate ideas from specification synthesis into OptiTrust.

2676

2677 REFERENCES

- 2678 2022. Unsequenced functions. <https://www.open-std.org/jtc1/sc22/wg14/www/docs/n2956.htm> ISO/IEC JCT1/SC22/WG14
2679 document N2956.
- 2680 Sami Alabed, Daniel Belov, Bart Chrzaszcz, Juliana Franco, Dominik Grewe, Dougal Maclaurin, James Molloy, Tom Natan,
2681 Tamara Norman, Xiaoyue Pan, Adam Paszke, Norman A. Rink, Michael Schaarschmidt, Timur Sitdikov, Agnieszka
2682 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>
- 2683 Vasco Amaral, Beatriz Norberto, Miguel Goulão, Marco Aldinucci, Siegfried Benkner, Andrea Bracciali, Paulo Carreira, Edgars
2684 Celms, Luís Correia, Clemens Grelck, Helen Karatza, Christoph Kessler, Peter Kilpatrick, Hugo Martiniano, Ilias Mavridis,
2685 Sabri Pllana, Ana Respício, José Simão, Luís Veiga, and Ari Visa. 2020. Programming languages for data-intensive HPC
2686 applications: A systematic mapping study. *Parallel Comput.* 91 (2020), 102584. <https://doi.org/10.1016/j.parco.2019.102584>
- 2687 Mehdi Amini. 2012. *Source-to-source automatic program transformations for GPU-like hardware accelerators*. Ph.D. Disserta-
2688 tion. Ecole Nationale Supérieure des Mines de Paris.
- 2689 Hansang Bae, Dheya Mustafa, Jae-Woo Lee, Aurangzeb, Hao Lin, Chirag Dave, Rudolf Eigenmann, and Samuel P. Midkiff.
2013. The Cetus Source-to-Source Compiler Infrastructure: Overview and Evaluation. *Int. J. Parallel Program.* 41, 6 (2013),
2690 753–767. <https://doi.org/10.1007/S10766-012-0211-Z>
- 2691 O.S. Bagge, K.T. Kalleberg, M. Haverlaen, and E. Visser. 2003. Design of the CodeBoost transformation system for domain-
2692 specific optimisation of C++ programs. In *Proceedings Third IEEE International Workshop on Source Code Analysis and*
Manipulation. 65–74. <https://doi.org/10.1109/SCAM.2003.1238032>
- 2693 Lénaïc Bagnères, Oleksandr Zinenko, Stéphane Huot, and Cédric Bastoul. 2016a. Opening Polyhedral Compiler’s Black Box.
2694 In *IEEE/ACM International Symposium on Code Generation and Optimization (CGO)*.

2695

- 2696 Lénaïc Bagnères, Aleksandr Zinenko, Stéphane Huot, and Cédric Bastoul. 2016b. Opening Polyhedral Compiler’s Black Box.
2697 In *IEEE/ACM International Symp. on Code Generation and Optimization*.
- 2698 Paul Barham and Michael Isard. 2019. Machine Learning Systems are Stuck in a Rut. In *Proceedings of the Workshop on Hot
2699 Topics in Operating Systems (Bertinoro, Italy) (HotOS ’19)*. Association for Computing Machinery, New York, NY, USA,
177–183. <https://doi.org/10.1145/3317550.3321441>
- 2700 Y. Barsamian, A. Charguéraud, S. A. Hirstoaga, and M. Mehrenberger. 2018. Efficient Strict-Binning Particle-in-Cell Algorithm
2701 for Multi-Core SIMD Processors. In *24th International Conference on Parallel and Distributed Computing (Euro-Par) (Lecture
2702 Notes in Computer Science, Vol. 11014)*. Springer, Cham, 749–763. https://doi.org/10.1007/978-3-319-96983-1_53
- 2703 Gilles Barthe, Benjamin Grégoire, César Kunz, and Tamara Rezk. 2009. Certificate Translation for Optimizing Compilers.
2704 *ACM Trans. Program. Lang. Syst.* 31, 5, Article 18 (jul 2009), 45 pages. <https://doi.org/10.1145/1538917.1538919>
- 2705 João Bispo and João MP Cardoso. 2020. Clava: C/C++ source-to-source compilation using LARA. *SoftwareX* 12 (2020),
100565. <https://www.sciencedirect.com/science/article/pii/S2352711019302122/pdf>
- 2706 Stefan Blom, Saeed Darabi, Marieke Huisman, and Wytse Oortwijn. 2017. The VerCors Tool Set: Verification of Parallel and
2707 Concurrent Software. In *Integrated Formal Methods*, Nadia Polikarpova and Steve Schneider (Eds.). Springer International
2708 Publishing, Cham, 102–110.
- 2709 Uday Bondhugula, Albert Hartono, J. Ramanujam, and P. Sadayappan. 2008. A practical automatic polyhedral parallelizer
2710 and locality optimizer. In *PLDI’08 ACM Conf. on Programming language design and implementation*.
- 2711 John Boyland. 2003. Checking Interference with Fractional Permissions. In *Static Analysis, 10th International Symposium,
2712 SAS 2003, San Diego, CA, USA, June 11-13, 2003, Proceedings (Lecture Notes in Computer Science, Vol. 2694)*, Radhia Cousot
(Ed.). Springer, 55–72. https://doi.org/10.1007/3-540-44898-5_4
- 2713 Martin Bravenboer, Karl Trygve Kalleberg, Rob Vermaas, and Eelco Visser. 2008. Stratego/XT 0.17. A Language and Toolset
2714 for Program Transformation. *Sci. Comput. Program.* 72, 1–2 (jun 2008), 52–70. <https://doi.org/10.1016/j.scico.2007.11.003>
- 2715 Cristiano Calcagno, Dino Distefano, Peter O’Hearn, and Hongseok. Yang. 2019. Go Huge or Go Home: POPL’19 Most
2716 Influential Paper Retrospective. <https://blog.sigplan.org/2020/03/03/go-huge-or-go-home-popl19-most-influentialpaper-retrospective/>
- 2717 Qinxiang Cao, Lennart Beringer, Samuel Gruetter, Josiah Dodds, and Andrew W. Appel. 2018. VST-Floyd: A Separation
2718 Logic Tool to Verify Correctness of C Programs. *J. Autom. Reason.* 61, 1-4 (2018), 367–422. <https://doi.org/10.1007/S10817-018-9457-5>
- 2719 Arthur Charguéraud. 2020a. Separation logic for sequential programs (functional pearl). *Proc. ACM Program. Lang.* 4, ICFP,
2720 Article 116 (aug 2020), 34 pages. <https://doi.org/10.1145/3408998>
- 2721 Arthur Charguéraud. 2020b. Separation Logic for Sequential Programs (Functional Pearl). *Proc. ACM Program. Lang.* 4,
2722 ICFP, Article 116 (aug 2020), 34 pages. <https://doi.org/10.1145/3408998>
- 2723 Gaurav Chaurasia, Jonathan Ragan-Kelley, Sylvain Paris, George Drettakis, and Frédo Durand. 2015. Compiling high
2724 performance recursive filters. In *Proceedings of the 7th Conference on High-Performance Graphics, HPG 2015, Los Angeles,
2725 California, USA, August 7-9, 2015*, Michael C. Doggett, Steven E. Molnar, Kayvon Fatahalian, Jacob Munkberg, Elmar
Eisemann, Petrik Clarberg, and Stephen N. Spencer (Eds.). ACM, 85–94. <https://doi.org/10.1145/2790060.2790063>
- 2726 Lorenzo Chelini, Martin Kong, Tobias Grosser, and Henk Corporaal. 2021. LoopOpt: Declarative Transformations Made
2727 Easy. In *Proceedings of the 24th International Workshop on Software and Compilers for Embedded Systems (Eindhoven,
2728 Netherlands) (SCOPES ’21)*. Association for Computing Machinery, New York, NY, USA, 11–16. <https://doi.org/10.1145/3493229.3493301>
- 2729 Chun Chen, Jacqueline Chame, and Mary W. Hall. 2008. *CHILL: A Framework for Composing High-Level Loop Transformations*.
2730 Technical Report 08-897. University of Southern California.
- 2731 Tianqi Chen, Thierry Moreau, Ziheng Jiang, Lianmin Zheng, Eddie Yan, Haichen Shen, Meghan Cowan, Leyuan Wang,
2732 Yuwei Hu, Luis Ceze, Carlos Guestrin, and Arvind Krishnamurthy. 2018. TVM: An Automated End-to-End Optimizing
2733 Compiler for Deep Learning. In *OSDI*. USENIX Association. <https://www.usenix.org/system/files/osdi18-chen.pdf>
- 2734 James R Cordy. 2006. The TXL source transformation language. *Science of Computer Programming* 61, 3 (2006), 190–210.
- 2735 Chirag Dave, Hansang Bae, Seung-Jai Min, Seyong Lee, Rudolf Eigenmann, and Samuel Midkiff. 2009. Cetus: A Source-to-
2736 Source Compiler Infrastructure for Multicores. *Computer* 42, 12 (2009), 36–42. <https://doi.org/10.1109/MC.2009.385>
- 2737 Thomas M Evans, Andrew Siegel, Erik W Draeger, Jack Deslippe, Marianne M Francois, Timothy C Germann, William E
2738 Hart, and Daniel F Martin. 2022. A survey of software implementations used by application codes in the Exascale
2739 Computing Project. *The International Journal of High Performance Computing Applications* 36, 1 (2022), 5–12. <https://doi.org/10.1177/10943420211028940>
- 2740 Paul Feautrier. 1992. Some efficient solutions to the affine scheduling problem: one dimensional time. *Intl. Journal of Parallel
2741 Programming* 21, 5 (october 1992), 313–348.
- 2742 Jean-Christophe Filliâtre and Andrei Paskevich. 2013. Why3—Where Programs Meet Provers. In *European Symposium on
2743 Programming (ESOP) (Lecture Notes in Computer Science, Vol. 7792)*. Springer, 125–128. <http://hal.inria.fr/hal-00789533>

- 2745 Jean-Christophe Filliâtre. 2003. *Why: a multi-language multi-prover verification tool*. Research Report 1366. LRI, Université
2746 Paris Sud. <http://www.lri.fr/~filliatr/ftp/publis/why-tool.ps.gz>
- 2747 Cormac Flanagan, K. Rustan M. Leino, Mark Lillibridge, Greg Nelson, James B. Saxe, and Raymie Stata. 2002. Extended
2748 Static Checking for Java. In *Programming Language Design and Implementation (PLDI)*. 234–245. <http://www.soe.ucsc.edu/~cormac/papers/pldi02.ps>
- 2749 John John Gough and K John Gough. 2001. *Compiling for the .Net Common Language Runtime*. Prentice Hall PTR.
- 2750 Bastian Hagedorn, Archibald Samuel Elliott, Henrik Barthels, Rastislav Bodik, and Vinod Grover. 2020a. Fireiron: a data-
2751 movement-aware scheduling language for GPUs. In *Proceedings of the ACM International Conf. on Parallel Architectures
2752 and Compilation Techniques*. 71–82.
- 2753 Bastian Hagedorn, Johannes Lenfers, Thomas Kundefinedhler, Xueying Qin, Sergei Gorbaltch, and Michel Steuwer. 2020b.
2754 Achieving high-performance the functional way: a functional pearl on expressing high-performance optimizations as
2755 rewrite strategies. *Proc. ACM Program. Lang.* 4, ICFP, Article 92 (aug 2020), 29 pages. <https://doi.org/10.1145/3408974>
- 2756 Stefan Heule, K. Rustan M. Leino, Peter Müller, and Alexander J. Summers. 2013. Abstract Read Permissions: Fractional
2757 Permissions without the Fractions. In *Verification, Model Checking, and Abstract Interpretation*, Roberto Giacobazzi, Josh
2758 Berdine, and Isabella Mastroeni (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 315–334.
- 2759 Yuka Ikarashi, Gilbert Louis Bernstein, Alex Reinking, Hasan Genc, and Jonathan Ragan-Kelley. 2022. Exocompilation for
2760 productive programming of hardware accelerators. In *Proceedings of the 43rd ACM SIGPLAN International Conference
2761 on Programming Language Design and Implementation (San Diego, CA, USA) (PLDI 2022)*. Association for Computing
2762 Machinery, New York, NY, USA, 703–718. <https://doi.org/10.1145/3519939.3523446>
- 2763 Yuka Ikarashi, Jonathan Ragan-Kelley, Tsukasa Fukusato, Jun Kato, and Takeo Igarashi. 2021. Guided Optimization for
2764 Image Processing Pipelines. In *2021 IEEE Symposium on Visual Languages and Human-Centric Computing (VL/HCC)*. 1–5.
2765 <https://doi.org/10.1109/VL/HCC51201.2021.9576341>
- 2766 Matthieu Journault and Antoine Miné. 2018. Inferring functional properties of matrix manipulating programs by abstract
2767 interpretation. *Form. Methods Syst. Des.* 53, 2 (oct 2018), 221–258. <https://doi.org/10.1007/s10703-017-0311-x>
- 2768 Ralf Jung, Robbert Krebbers, Jacques-Henri Jourdan, Ales Bizjak, Lars Birkedal, and Derek Dreyer. 2018a. Iris from the
2769 ground up: A modular foundation for higher-order concurrent separation logic. *J. Funct. Program.* 28 (2018), e20.
2770 <https://doi.org/10.1017/S0956796818000151>
- 2771 Ralf Jung, Robbert Krebbers, Jacques-Henri Jourdan, Aleš Bizjak, Lars Birkedal, and Derek Dreyer. 2018b. Iris from the
2772 ground up: A modular foundation for higher-order concurrent separation logic. *Journal of Functional Programming* 28
2773 (2018), e20. <https://people.mpi-sws.org/~dreyer/papers/iris-ground-up/paper.pdf>
- 2774 Yamato Kanetaka, Hiroyasu Takagi, Yoshihiro Maeda, and Norishige Fukushima. 2024. SlidingConv: Domain-Specific
2775 Description of Sliding Discrete Cosine Transform Convolution for Halide. *IEEE Access* 12 (2024), 7563–7583. <https://doi.org/10.1109/ACCESS.2023.3345660>
- 2776 Vasilios Kelefouras and Georgios Keramidas. 2022. Design and Implementation of 2D Convolution on x86/x64 Processors.
2777 *IEEE Transactions on Parallel and Distributed Systems* 33, 12 (2022), 3800–3815. <https://doi.org/10.1109/TPDS.2022.3171471>
- 2778 Athanasios Konstantinidis. 2013. *Source-to-source compilation of loop programs for manycore processors*. Ph.D. Dissertation.
2779 Imperial College London.
- 2780 Michael Kruse and Hal Finkel. 2018. A Proposal for Loop-Transformation Pragmas. *CoRR* abs/1805.03374 (2018).
2781 arXiv:1805.03374 <http://arxiv.org/abs/1805.03374>
- 2782 Olaf Krzikalla, Kim Feldhoff, Ralph Müller-Pfefferkorn, and Wolfgang E. Nagel. 2011. Scout: A Source-to-Source Transform-
2783 ator for SIMD-Optimizations. In *Euro-Par Workshops (2) (LNCS, Vol. 7156)*. Springer.
- 2784 César Kunz. 2009. *Proof preservation and program compilation*. Ph.D. Dissertation. École Nationale Supérieure des Mines de
2785 Paris. <https://pastel.archives-ouvertes.fr/pastel-00004940/file/thesis-ckunz.pdf>
- 2786 Julia Lawall and Gilles Muller. 2018. Coccinelle: 10 Years of Automated Evolution in the Linux Kernel. In *USENIX Conference
2787 on Usenix Annual Technical Conference (USENIX ATC '18)*. USENIX Association, 13 pages.
- 2788 Youenn Lebras. 2019. *Code optimization based on source to source transformations using profile guided metrics*. Ph.D.
2789 Dissertation. Université Paris-Saclay (ComUE). <https://www.theses.fr/2019SACL037.pdf>
- 2790 Amanda Liu, Gilbert Louis Bernstein, Adam Chlipala, and Jonathan Ragan-Kelley. 2022. Verified Tensor-Program Optimiza-
2791 tion via High-Level Scheduling Rewrites. 6, POPL, Article 55 (jan 2022), 28 pages. <https://doi.org/10.1145/3498717>
- 2792 Peter Müller, Malte Schwerhoff, and Alexander J. Summers. 2017. Viper: A Verification Infrastructure for Permission-Based
2793 Reasoning. In *Dependable Software Systems Engineering*. IOS Press, 104–125. <https://doi.org/10.3233/978-1-61499-810-5-104>
- 2794 Kedar S. Namjoshi and Nimit Singhania. 2016. Loopy: Programmable and Formally Verified Loop Transformations. In *Static
2795 Analysis - 23rd International Symposium, SAS (LNCS, Vol. 9837)*. Springer.
- 2796 George Ciprian Necula. 1998. *Compiling with proofs*. Ph.D. Dissertation. Carnegie Mellon University.
- 2797 Peter W. O’Hearn. 2019. Separation logic. *Commun. ACM* 62, 2 (2019), 86–95. <https://doi.org/10.1145/3211968>
- 2798
- 2799
- 2800

- 2794 Pedro Pinto, Joao Bispo, Joao Cardoso, Jorge Gomes Barbosa, Davide Gadioli, Gianluca Palermo, Jan Martinovic, Martin
2795 Golasowski, Katerina Slaninova, Radim Cmar, et al. 2020. Pegasus: Performance Engineering for Software Applications
2796 Targeting HPC Systems. *IEEE Transactions on Software Engineering* (2020). <https://repositorio-aberto.up.pt/bitstream/10216/127756/2/405707.pdf>
- 2797 Dan Quinlan. 2000. ROSE: Compiler support for object-oriented frameworks. *Parallel processing letters* 10, 02n03 (2000),
2798 215–226. https://digital.library.unt.edu/ark:/67531/metadc741175/m2/1/high_res_d/793936.pdf
- 2799 Dan Quinlan and Chunhua Liao. 2011. The ROSE source-to-source compiler infrastructure. In *Cetus users and compiler
2800 infrastructure workshop, in conjunction with PACT*, Vol. 2011. 1.
- 2801 Pawel K. Radtke and Tobias Weinzierl. 2024. Compiler support for semi-manual AoS-to-SoA conversions with data views.
2802 arXiv:2405.12507 [cs.PL] <https://arxiv.org/abs/2405.12507>
- 2803 Jonathan Ragan-Kelley. 2023. Technical Perspective: Reconsidering the Design of User-Schedulable Languages. *Commun.
2804 ACM* 66, 3 (feb 2023), 88. <https://doi.org/10.1145/3580370>
- 2805 Jonathan Ragan-Kelley, Connelly Barnes, Andrew Adams, Sylvain Paris, Frédo Durand, and Saman Amarasinghe. 2013.
2806 Halide: A Language and Compiler for Optimizing Parallelism, Locality, and Recomputation in Image Processing Pipelines.
2807 In *Conference on Programming Language Design and Implementation*. 12 pages. <https://doi.org/10.1145/2491956.2462176>
- 2808 John C. Reynolds. 2002. Separation Logic: A Logic for Shared Mutable Data Structures. In *17th IEEE Symposium on Logic
2809 in Computer Science (LICS 2002), 22-25 July 2002, Copenhagen, Denmark, Proceedings*. IEEE Computer Society, 55–74.
2810 <https://doi.org/10.1109/LICS.2002.1029817>
- 2811 Gabe Rudy, Malik Murtaza Khan, Mary Hall, Chun Chen, and Jacqueline Chame. 2011. A Programming Language Interface
2812 to Describe Transformations and Code Generation. In *Languages and Compilers for Parallel Computing*. Springer Berlin
2813 Heidelberg.
- 2814 Ömer Sakar, Mohsen Safari, Marieke Huisman, and Anton Wijs. 2022. Alpinist: An Annotation-Aware GPU Program
2815 Optimizer. In *Tools and Algorithms for the Construction and Analysis of Systems - 28th International Conference, TACAS
2816 2022, Held as Part of the European Joint Conferences on Theory and Practice of Software, ETAPS 2022, Munich, Germany,
2817 April 2-7, 2022, Proceedings, Part II (Lecture Notes in Computer Science, Vol. 13244)*, Dana Fisman and Grigore Rosu (Eds.).
2818 Springer, 332–352. https://doi.org/10.1007/978-3-030-99527-0_18
- 2819 Michael Sammler, Rodolphe Lepigre, Robbert Krebbers, Kayvan Memarian, Derek Dreyer, and Deepak Garg. 2021. RefinedC:
2820 automating the foundational verification of C code with refined ownership types. In *PLDI '21: 42nd ACM SIGPLAN
2821 International Conference on Programming Language Design and Implementation, Virtual Event, Canada, June 20-25, 2021*,
2822 Stephen N. Freund and Eran Yahav (Eds.). ACM, 158–174. <https://doi.org/10.1145/3453483.3454036>
- 2823 Cristina Silvano, Giovanni Agosta, Andrea Bartolini, Andrea R. Beccari, Luca Benini, Loïc Besnard, João Bispo, Radim
2824 Cmar, João M.P. Cardoso, Carlo Cavazzoni, Daniele Cesarini, Stefano Cherubin, Federico Ficarelli, Davide Gadioli,
2825 Martin Golasowski, Antonio Libri, Jan Martinović, Gianluca Palermo, Pedro Pinto, Erven Rohou, Kateřina Slaninová, and
2826 Emanuele Vitali. 2019. The ANTAREX domain specific language for high performance computing. *Microprocessors and
2827 Microsystems* 68 (2019), 58–73. <https://doi.org/10.1016/j.micpro.2019.05.005>
- 2828 Simon Spies, Lennard Gäher, Michael Sammler, and Derek Dreyer. 2024. Quiver: Guided Abductive Inference of Separation
2829 Logic Specifications in Coq. *Proc. ACM Program. Lang.* 8, PLDI, Article 183 (jun 2024), 25 pages. <https://doi.org/10.1145/3656413>
- 2830 Manish Vachharajani, Neil Vachharajani, David I. August, and Spyridon Triantafyllis. 2003. Compiler Optimization-Space
2831 Exploration. In *Proceedings of the 2003 IEEE/ACM International Symposium on Code Generation and Optimization (CGO)*.
2832 IEEE Computer Society, Los Alamitos, CA, USA, 204. <https://doi.org/10.1109/CGO.2003.1191546>
- 2833 Xuejun Yang, Yang Chen, Eric Eide, and John Regehr. 2011. Finding and Understanding Bugs in C Compilers. In *Conference
2834 on Programming Language Design and Implementation* (San Jose, California, USA). Association for Computing Machinery,
2835 12 pages. <https://doi.org/10.1145/1993498.1993532>
- 2836 Qing Yi and Apan Qasem. 2008. Exploring the Optimization Space of Dense Linear Algebra Kernels. In *LCPC*.
- 2837 Qing Yi, Qian Wang, and Huimin Cui. 2014. Specializing Compiler Optimizations through Programmable Composition for
2838 Dense Matrix Computations. In *Proceedings of the 47th Annual IEEE/ACM International Symposium on Microarchitecture*
2839 (Cambridge, United Kingdom) (*MICRO-47*). IEEE Computer Society, USA, 596–608. <https://doi.org/10.1109/MICRO.2014.14>
- 2840 Oleksandr Zinenko, Lorenzo Chelini, and Tobias Grosser. 2018a. *Declarative Transformations in the Polyhedral Model*.
2841 Research Report RR-9243. <https://hal.inria.fr/hal-01965599>
- 2842 Oleksandr Zinenko, Stéphane Huot, and Cédric Bastoul. 2018b. Visual Program Manipulation in the Polyhedral Model.
2843 *ACM Trans. Archit. Code Optim.* 15, 1, Article 16 (mar 2018), 25 pages. <https://doi.org/10.1145/3177961>