

MIT Open Access Articles

Helium: lifting high-performance stencil kernels from stripped x86 binaries to halide DSL code

The MIT Faculty has made this article openly available. **Please share** how this access benefits you. Your story matters.

Citation: Charith Mendis, Jeffrey Bosboom, Kevin Wu, Shoaib Kamil, Jonathan Ragan-Kelley, Sylvain Paris, Qin Zhao, and Saman Amarasinghe. 2015. Helium: lifting high-performance stencil kernels from stripped x86 binaries to halide DSL code. In Proceedings of the 36th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI 2015). ACM, New York, NY, USA, 391-402.

As Published: <http://dx.doi.org/10.1145/2737924.2737974>

Publisher: Association for Computing Machinery (ACM)

Persistent URL: <http://hdl.handle.net/1721.1/99696>

Version: Author's final manuscript: final author's manuscript post peer review, without publisher's formatting or copy editing

Terms of use: Creative Commons Attribution-Noncommercial-Share Alike



Helium: Lifting High-Performance Stencil Kernels from Stripped x86 Binaries to Halide DSL Code

Charith Mendis[†] Jeffrey Bosboom[†] Kevin Wu[†] Shoaib Kamil[†] Jonathan Ragan-Kelley[‡]
Sylvain Paris[◊] Qin Zhao^{*} Saman Amarasinghe[†]

[†]MIT CSAIL, Cambridge, MA, USA

[‡]Stanford University, Palo Alto, CA, USA

[◊]Adobe, Cambridge, MA, USA ^{*}Google, Cambridge, MA, USA

{charithm,jbosboom,kevinwu,skamil,saman}@csail.mit.edu
jr@cs.stanford.edu sparis@adobe.com zhaoqin@google.com

Abstract

Highly optimized programs are prone to bit rot, where performance quickly becomes suboptimal in the face of new hardware and compiler techniques. In this paper we show how to automatically lift performance-critical stencil kernels from a stripped x86 binary and generate the corresponding code in the high-level domain-specific language Halide. Using Halide’s state-of-the-art optimizations targeting current hardware, we show that new optimized versions of these kernels can replace the originals to rejuvenate the application for newer hardware.

The original optimized code for kernels in stripped binaries is nearly impossible to analyze statically. Instead, we rely on dynamic traces to regenerate the kernels. We perform buffer structure reconstruction to identify input, intermediate and output buffer shapes. We abstract from a forest of concrete dependency trees which contain absolute memory addresses to symbolic trees suitable for high-level code generation. This is done by canonicalizing trees, clustering them based on structure, inferring higher-dimensional buffer accesses and finally by solving a set of linear equations based on buffer accesses to lift them up to simple, high-level expressions.

Helium can handle highly optimized, complex stencil kernels with input-dependent conditionals. We lift seven kernels from Adobe Photoshop giving a 75% performance improvement, four kernels from IrfanView, leading to $4.97\times$ performance, and one stencil from the miniGMG multigrid benchmark netting a $4.25\times$ improvement in performance. We manually rejuvenated Photoshop by replacing eleven of Photoshop’s filters with our lifted implementations, giving $1.12\times$ speedup without affecting the user experience.

Categories and Subject Descriptors D.1.2 [Programming Techniques]: Automatic Programming—Program transformation; D.2.7 [Software Engineering]: Distribution, Maintenance and Enhancement—Restructuring, reverse engineering, and reengineering

Keywords Helium; dynamic analysis; reverse engineering; x86 binary instrumentation; autotuning; image processing; stencil computation

1. Introduction

While lowering a high-level algorithm into an optimized binary executable is well understood, going in the reverse direction—lifting an optimized binary into the high-level algorithm it implements—remains nearly impossible. This is not surprising: lowering eliminates information about data types, control structures, and programmer intent. Inverting this process is far more challenging because stripped binaries lack high-level information about the program. Because of the lack of high-level information, lifting is only possible given constraints, such as a specific domain or limited degree of abstraction to be reintroduced. Still, lifting from a binary program can help reverse engineer a program, identify security vulnerabilities, or translate from one binary format to another.

In this paper, we lift algorithms from existing binaries for the sake of program *rejuvenation*. Highly optimized programs are especially prone to bit rot. While a program binary often executes correctly years after its creation, its performance is likely suboptimal on newer hardware due to changes in hardware and the advancement of compiler technology since its creation. Re-optimizing production-quality kernels by hand is extremely labor-intensive, requiring many engineer-months even for relatively simple parts of the code [21]. Our goal is to take an existing legacy binary, lift the performance-critical components with sufficient accuracy to a high-level representation, re-optimize them with modern tools, and replace the bit-rotted component with the optimized version. To automatically achieve best performance for the algorithm, we lift the program to an even higher level than the original source code, into a high-level domain-specific language (DSL). At this level, we express the original programmer intent instead of obscuring it with performance-related transformations, letting us apply domain knowledge to exploit modern architectural features without sacrificing performance portability.

Though this is an ambitious goal, aspects of the problem make this attainable. Program rejuvenation only requires transforming performance-critical parts of the program, which often apply relatively simple computations repeatedly to large amounts of data. Even though this high-level algorithm may be simple, the generated code is complicated due to compiler and programmer optimizations such as tiling, vectorization, loop specialization, and unrolling. In this paper, we introduce dynamic, data-driven techniques to abstract away optimization complexities and get to the underlying simplicity of the high-level intent.

We focus on stencil kernels, mainly in the domain of image-processing programs. Stencils, prevalent in image processing kernels used in important applications such as Adobe Photoshop, Mi-

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

PLDI’15, June 13–17, 2015, Portland, OR, USA.

Copyright is held by the owner/author(s).

ACM 978-1-4503-3468-6/15/06.

<http://dx.doi.org/10.1145/nnnnnnn.nnnnnn>

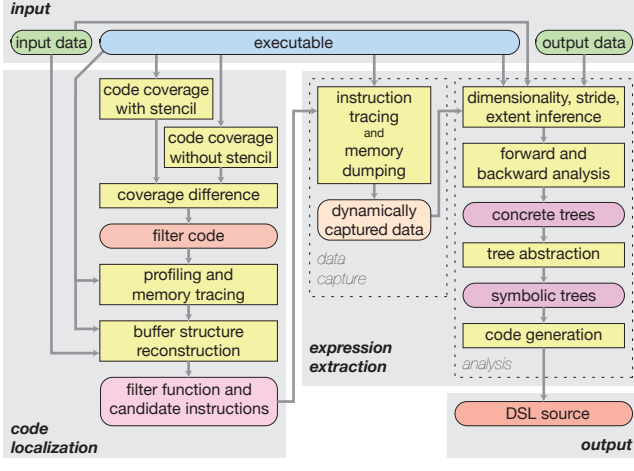


Figure 1. Helium workflow.

crosoft PowerPoint and Google Picasa, use enormous amounts of computation and/or memory bandwidth. As these kernels mostly perform simple data-parallel operations on large image data, they can leverage modern hardware capabilities such as vectorization, parallelization, and graphics processing units (GPUs).

Furthermore, recent programming language and compiler breakthroughs have dramatically improved the performance of stencil algorithms [17, 22, 25, 26]; for example, Halide has demonstrated that writing an image processing kernel in a high level DSL and autotuning it to a specific architecture can lead to 2–10× performance gains compared to hand-tuned code by expert programmers. In addition, only a few image-processing kernels in Photoshop and other applications are hand-optimized for the latest architectures; many are optimized for older architectures, and some have little optimization for any architecture. Reformulating these kernels as Halide programs makes it possible to rejuvenate these applications to continuously provide state-of-the-art performance by using the Halide compiler and autotuner to optimize for current hardware and replacing the old code with the newly-generated optimized implementation. Overall, we lift seven filters and portions of four more from Photoshop, four filters from IrfanView, and the smooth stencil from the miniGMG [31] high-performance computing benchmark into Halide code. We then autotune the Halide schedules and compare performance against the original implementations, delivering an average speedup of 1.75 on Photoshop, 4.97 on IrfanView, and 4.25 on miniGMG. The entire process of lifting and regeneration is completely automated. We also manually replace Photoshop kernels with our rejuvenated versions, obtaining a speedup of 1.12× even while constrained by optimization decisions (such as tile size) made by the Photoshop developers.

In addition, lifting can provide opportunities for further optimization. For example, power users of image processing applications create pipelines of kernels for batch processing of images. Hand-optimizing kernel pipelines does not scale due to the combinatorial explosion of possible pipelines. We demonstrate how our techniques apply to a pipeline kernels by creating pipelines of lifted Photoshop and IrfanView kernels and generating optimized code in Halide, obtaining 2.91× and 5.17× faster performance than the original unfused pipelines.

2. Overview, Challenges & Contributions

Helium lifts stencils from stripped binaries to high-level code. Helium is fully automated, only prompting the user to perform any GUI interactions required to run the program under analysis with

and without the target stencil. In total, the user runs the program five times for each stencil lifted. Figure 1 shows the workflow Helium follows to implement the translation. Overall, the flow is divided into two stages: code localization, described in Section 3, and expression extraction, covered in Section 4. Figure 2 shows the flow through the system for a blur kernel. In this section, we give a high-level view of the challenges and how we address them in Helium.

While static analysis is the only sound way to lift a computation, doing so on a stripped x86 binary is extremely difficult, if not impossible. In x86 binaries, code and data are not necessarily separated, and determining separation in stripped binaries is known to be equivalent to the halting problem [16]. Statically, it is difficult to even find which kernels execute as they are located in different dynamic linked libraries (DLLs) loaded at runtime. Therefore, we use a dynamic data flow analysis built on top of DynamoRIO [8], a dynamic binary instrumentation framework.

Isolating performance critical kernels We find that profiling information alone is unable to identify performance-critical kernels. For example, a highly vectorized kernel of an element-wise operation, such as the invert image filter, may be invoked far fewer iterations than the number of data items. On the other hand, we find that kernels touch all data in input and intermediate buffers to produce a new buffer or the final output. Thus, by using a data-driven approach (described in Section 3.1) and analyzing the extent of memory regions touched by static instructions, we identify kernel code blocks more accurately than through profiling.

Extracting optimized kernels While these stencil kernels may perform logically simple computations, optimized kernel code found in binaries is far from simple. In many applications, programmers expend considerable effort in speeding up performance-critical kernels; such optimizations often interfere with determining the program’s purpose. For example, many kernels do not iterate over the image in a simple linear pattern but use smaller tiles for better locality. In fact, Photoshop kernels use a common driver that provides the image as a set of tiles to the kernel. However, we avoid control-flow complexities due to iteration order optimization by only focusing on data flow. For each data item in the output buffer, we compute an expression tree with input and intermediate buffer locations and constants as leaves.

Handling complex control flow A dynamic trace can capture only a single path through the maze of complex control flow in a program. Thus, extracting full control-flow using dynamic analysis is challenging. However, high performance kernels repeatedly execute the same computations on millions of data items. By creating a forest of expression trees, each tree calculating a single output value, we use *expression forest reconstruction* to find a corresponding tree for all the input-dependent control-flow paths. The forest of expression trees shown in Figure 2(b) is extracted from execution traces of Photoshop’s 2D blur filter code in Figure 2(a).

Identifying input-dependent control flow Some computations such as image threshold filters update each pixel differently depending on properties of that pixel. As we create our expression trees by only considering data flow, we will obtain a forest of trees that form multiple clusters without any pattern to identify cluster membership. The complex control flow of these conditional updates is interleaved with the control flow of the iteration ordering, and is thus difficult to disentangle. We solve this problem, as described in Section 4.6, by first doing a forward propagation of input data values to identify instructions that are input-dependent and building expression trees for the input conditions. Then, if a node in our output expression tree has a control flow dependency on the input, we can predicate that tree with the corresponding input condition. During this forward analysis, we also mark address calculations that depend on input values, allowing us to identify lookup tables during backward analysis.

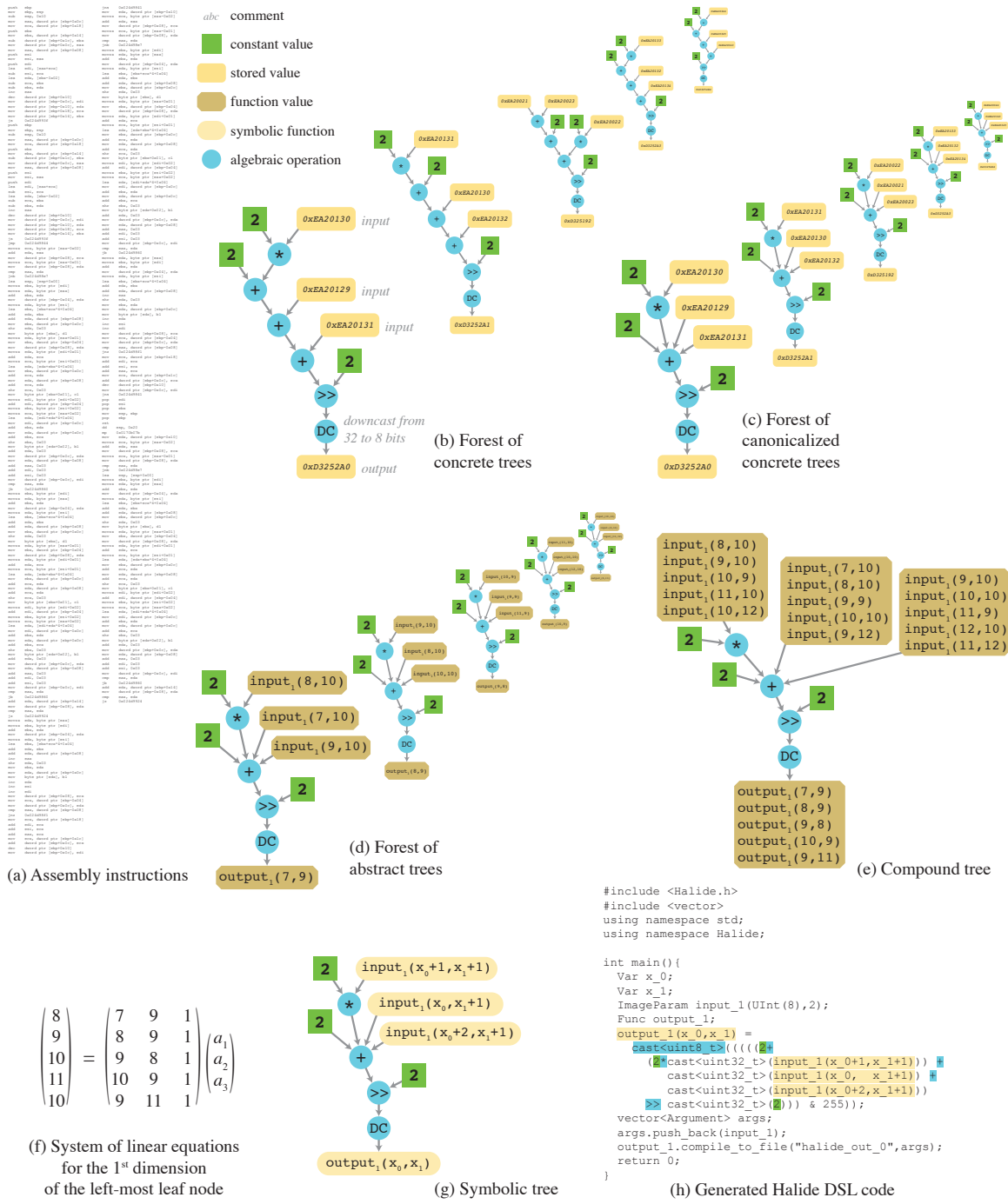


Figure 2. Stages of expression extraction for Photoshop’s 2D blur filter, reduced to 1D in this figure for brevity. We instrument assembly instructions (a) to recover a forest of concrete trees (b), which we then canonicalize (c). We use buffer structure reconstruction to obtain abstract trees (d). Merging the forest of abstract trees into compound trees (e) gives us linear systems (f) to solve to obtain symbolic trees (g) suitable for generating Halide code (h).

Handling code duplication Many optimized kernels have inner loops unrolled or some iterations peeled off to help optimize the common case. Thus, not all data items are processed by the same assembly instructions. Furthermore, different code paths may compute the same output value using different combinations of operations. We handle this situation by canonicalizing the trees and clustering trees representing the same canonical expression during expression forest reconstruction, as shown in Figure 2(c).

Identifying boundary conditions Some stencil kernels perform different calculations at boundaries. Such programs often include loop peeling and complex control flow, making them difficult to handle. In Helium these boundary conditions lead to trees that are different from the rest. By clustering trees (described in Section 4.8), we separate the common stencil operations from the boundary conditions.

Determining buffer dimensions and sizes Accurately extracting stencil computations requires determining dimensionality and the strides of each dimension of the input, intermediate and output buffers. However, at the binary level, multi-dimensional arrays appear to be allocated as one linear block. We introduce *buffer structure reconstruction*, a method which creates multiple levels of coalesced memory regions for inferring dimensions and strides by analyzing data access patterns (Section 3.2). Many stencil computations have ghost regions or padding between dimensions for alignment or graceful handling of boundary conditions. We leverage these regions in our analysis.

Recreating index expressions & generating Halide code Recreating stencil computations requires reconstructing logical index expressions for the multi-dimensional input, intermediate and output buffers. We use access vectors from a randomly selected set of expression trees to create a linear system of equations that can be solved to create the algebraic index expressions, as in Figure 2(f). Our method is detailed in Section 4.10. These algebraic index expressions can be directly transformed into a Halide function, as shown in Figure 2(g)-(h).

3. Code Localization

Helium’s first step is to find the code that implements the kernel we want to lift, which we term *code localization*. While the code performing the kernel computation should be frequently executed, Helium cannot simply assume the most frequently executed region of code (which is often just `memcpy`) is the stencil kernel. More detailed profiling is required.

However, performing detailed instrumentation on the entirety of a large application such as Photoshop is impractical, due to both large instrumentation overheads and the sheer volume of the resulting data. Photoshop loads more than 160 binary modules, most of which are unrelated to the filter we wish to extract. Thus the code localization stage consists of a *coverage difference phase* to quickly screen out unrelated code, followed by more invasive profiling to determine the kernel function and the instructions reading and writing the input and output buffers. The kernel function and set of instructions are then used for even more detailed profiling in the expression extraction stage (in Section 4).

3.1 Screening Using Coverage Difference

To obtain a first approximation of the kernel code location, our tool gathers code coverage (at basic block granularity) from two executions of the program that are as similar as possible except that one execution runs the kernel and the other does not. The difference between these executions consists of basic blocks that only execute when the kernel executes. This technique assumes the kernel code is not executed in other parts of the application (e.g., to draw small preview images), and data-reorganization or UI code

specific to the kernel will still be captured, but it works well in practice to quickly screen out most of the program code (such as general UI or file parsing code). For Photoshop’s blur filter, the coverage difference contains only 3,850 basic blocks out of 500,850 total blocks executed.

Helium then asks the user to run the program again (including the kernel), instrumenting only those basic blocks in the coverage difference. The tool collects basic block execution counts, predecessor blocks and call targets, which will be used to build a dynamic control-flow graph in the next step. Helium also collects a dynamic memory trace by instrumenting all memory accesses performed in those basic blocks. The trace contains the instruction address, the absolute memory address, the access width and whether the access is a read or a write. The result of this instrumentation step enables Helium to analyze memory access patterns and detect the filter function.

3.2 Buffer Structure Reconstruction

Helium proceeds by first processing the memory trace to recover the memory layout of the program. Using the memory layout, the tool determines instructions that are likely accessing input and output buffers. Helium then uses the dynamic control-flow graph to select a function containing the highest number of such instructions.

We represent the memory layout as address regions (lists of ranges) annotated with the set of static instructions that access them. For each static instruction, Helium first coalesces any immediately-adjacent memory accesses and removes duplicate addresses, then sorts the resulting regions. The tool then merges regions of different instructions to correctly detect unrolled loops accessing the input data, where a single instruction may only access part of the input data but the loop body as a whole covers the data. Next, Helium links any group of three or more regions separated by a constant stride to form a single larger region. This proceeds recursively, building larger regions until no regions can be coalesced (see Figure 3). Recursive coalescing may occur if e.g. an image filter accesses the R channel of an interleaved RGB image with padding to align each scanline on a 16-byte boundary; the channel stride is 3 and the scanline stride is the image width rounded up to a multiple of 16.

Helium detects element size based on access width. Some accesses are logically greater than the machine word size, such as a 64-bit addition using an `add/adc` instruction pair. If a buffer is accessed at multiple widths, the tool uses the most common width, allowing it to differentiate between stencil code operating on individual elements and `memcpy`-like code treating the buffer as a block of bits.

Helium selects all regions of size comparable to or larger than the input and output data sizes and records the associated *candidate instructions* that potentially access the input and output buffers in memory.

3.3 Filter Function Selection

Helium maps each basic block containing candidate instructions to its containing function using a dynamic control-flow graph built from the profile, predecessor, and call target information collected during screening. The tool considers the function containing the most candidate static instructions to be the kernel. Tail call optimization may fool Helium into selecting a parent function, but this still covers the kernel code; we just instrument more code than necessary.

The chosen function does not always contain the most frequently executed basic block, as one might naïvely assume. For example, Photoshop’s invert filter processes four image bytes per loop iteration, so other basic blocks that execute once per pixel execute more often.

Helium selects a filter function for further analysis, rather than a single basic block or a set of functions, as a tradeoff between

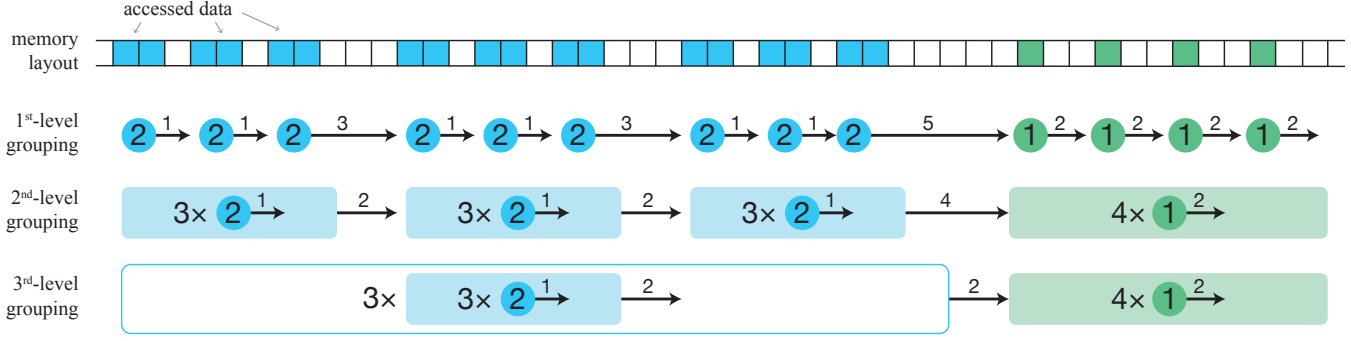


Figure 3. During buffer structure reconstruction, Helium groups the absolute addresses from the memory trace into regions, recursively combining regions of the same size separated by constant stride.

capturing all the kernel code and limiting the instrumentation during the expression extraction phase to a manageable amount of code. Instrumenting smaller regions risks not capturing all kernel code, but instrumenting larger regions generates more data that must be analyzed during expression extraction and also increases the likelihood that expression extraction will extract code that does not belong to the kernel (false data dependencies). Empirically, function granularity strikes a good balance. Helium localizes Photoshop’s blur to 328 static instructions in 14 basic blocks in the filter function and functions it calls, a manageable number for detailed dynamic instrumentation during expression extraction.

4. Expression Extraction

In this phase, we recover the stencil computation from the filter function found during code localization. Stencils can be represented as relatively simple data-parallel operations with few input-dependent conditionals. Thus, instead of attempting to understand all control flow, we focus on data flow from the input to the output, plus a small set of input-dependent conditionals which affect computation, to extract only the actual computation being performed.

For example, we are able to go from the complex unrolled static disassembly listing in Figure 2 (a) for a 1D blur stencil to the simple representation of the filter in Figure 2 (g) and finally to DSL code in Figure 2 (h).

During expression extraction, Helium performs detailed instrumentation of the filter function, using the captured data for buffer structure reconstruction and dimensionality inference, and then applies expression forest reconstruction to build expression trees suitable for DSL code generation.

4.1 Instruction Trace Capture and Memory Dump

During code localization, Helium determines the entry point of the filter function. The tool now prompts the user to run the program again, applying the kernel to known input data (if available), and collects a trace of all dynamic instructions executed from that function’s entry to its exit, along with the absolute addresses of all memory accesses performed by the instructions in the trace. For instructions with indirect memory operands, our tool records the address expression (some or all of $base + scale \times index + disp$). Helium also collects a page-granularity memory dump of all memory accessed by candidate instructions found in Section 3. Read pages are dumped immediately, but written pages are dumped at the filter function’s exit to ensure all output has been written before dumping. The filter function may execute many times; both the instruction trace and memory dump include all such executions.

4.2 Buffer Structure Reconstruction

Because the user ran the program again during instruction trace capture, we cannot assume buffers have the same location as during code localization. Using the memory addresses recorded as part of the instruction trace, Helium repeats buffer structure reconstruction (Section 3.2) to find memory regions with contiguous memory accesses which are likely the input and output buffers.

4.3 Dimensionality, Stride and Extent Inference

Buffer structure reconstruction finds buffer locations in memory, but to accurately recover the stencil, Helium must infer the buffers’ dimensionality, and for each dimension, the stride and extent. For image processing filters (or any domain where the user can provide input and output data), Helium can use the memory dump to recover this information. Otherwise, the tool falls back to generic inference that does not require the input and output data.

Inference using input and output data Helium searches the memory dump for the known input and output data and records the starting and ending locations of the corresponding memory buffers. It detects alignment padding by comparing against the given input and output data. For example, when Photoshop blurs a 32×32 image, it pads each edge by one pixel, then rounds each scanline up to 48 bytes for 16-byte alignment. Photoshop stores the R, G and B planes of a color image separately, so Helium infers three input buffers and three output buffers with two dimensions. All three planes are the same size, so the tool infers each dimension’s stride to be 48 (the distance between scanlines) and the extent to be 32. Our other example image processing application, IrfanView, stores the RGB values interleaved, so Helium automatically infers that IrfanView’s single input and output buffers have three dimensions.

Generic inference If we do not have input and output data (as in the miniGMG benchmark, which generates simulated input at runtime), or the data cannot be recognized in the memory dump, Helium falls back to generic inference based on buffer structure reconstruction. The dimensionality is equal to the number of levels of recursion needed to coalesce memory regions. Helium can infer buffers of arbitrary dimensionality so long padding exists between dimensions. For the dimension with least stride, the extent is equal to the number of adjacent memory locations accessed in one grouping and the stride is equal to the memory access width of the instructions affecting this region. For all other dimensions, the stride is the difference between the starting addresses of two adjacent memory regions in the same level of coalescing and the extent is equal to the number of independent memory regions present at each level.

If there are no gaps in the reconstructed memory regions, this inference will treat the memory buffer as single-dimensional, regardless of the actual dimensionality.

Inference by forward analysis We have yet to encounter a stencil for which we lack input and output data and for which the generic inference fails, but in that case the application must be handling boundary conditions on its own. In this case, Helium could infer dimensionality and stride by looking at different tree clusters (Section 4.8) and calculating the stride between each tree in a cluster containing the boundary conditions.

When inference is unnecessary If generic inference fails but the application does not handle boundary conditions on its own, the stencil is pointwise (uses only a single input point for each output point). The dimensionality is irrelevant to the computation, so Helium can assume the buffer is linear with a stride of 1 and extent equal to the memory region’s size.

4.4 Input/Output Buffer Selection

Helium considers buffers that are read, not written, and not accessed using indices derived from other buffer values to be input buffers. If output data is available, Helium identifies output buffers by locating the output data in the memory dump. Otherwise (or if the output data cannot be found), Helium assumes buffers that are written to with values derived from the input buffers to be output buffers, even if they do not live beyond the function (e.g., temporary buffers).

4.5 Instruction Trace Preprocessing

Before analyzing the instruction trace, Helium preprocesses it by renaming the x87 floating-point register stack using a technique similar to that used in [13]. More specifically, we recreate the floating point stack from the dynamic instruction trace to find the top of the floating point stack, which is necessary to recover non-relative floating-point register locations. Helium also maps registers into memory so the analysis can treat them identically; this is particularly helpful to handle dependencies between partial register reads and writes (e.g., writing to `eax` then reading from `ah`).

4.6 Forward Analysis for Input-Dependent Conditionals

While we focus on recovering the stencil computation, we cannot ignore control flow completely because some branches may be part of the computation. Helium must distinguish these *input-dependent conditionals* that affect *what* the stencil computes from the control flow arising from optimized loops controlling *when* the stencil computes.

To capture these conditionals, the tool first identifies which instructions read the input directly using the reconstructed memory layout. Next, Helium does a forward pass through the instruction trace identifying instructions which are affected by the input data, either directly (through data) or through the flags register (control dependencies). The input-dependent conditionals are the input-dependent instructions reading the flag registers (conditional jumps plus a few math instructions such as `adc` and `sbb`).

Then for each static instruction in the filter function, Helium records the sequence of taken/not-taken branches of the input-dependent conditionals required to reach that instruction from the filter function entry point. The result of the forward analysis is a mapping from each static instruction to the input-dependent conditionals (if any) that must be taken or not taken for that instruction to be executed. This mapping is used during backward analysis to build predicate trees (see Figure 5).

During the forward analysis, Helium flags instructions which access buffers using indices derived from other buffers (indirect access). These flags are used to track index calculation dependencies during backward analysis.

4.7 Backward Analysis for Data-Dependency Trees

In this step, the tool builds data-dependency trees to capture the exact computation of a given output location. Helium walks backwards

through the instruction trace, starting from instructions which write output buffer locations (identified during buffer structure reconstruction). We build a data-dependency tree for each output location by maintaining a frontier of nodes on the leaves of the tree. When the tool finds an instruction that computes the value of a leaf in the frontier, Helium adds the corresponding operation node to the tree, removes the leaf from the frontier and adds the instruction’s sources to the frontier if not already present.

We call these *concrete trees* because they contain absolute memory addresses. Figure 2 (b) shows a forest of concrete trees for a 1D blur stencil.

Indirect buffer access Table lookups give rise to indirect buffer accesses, in which a buffer is indexed using values read from another buffer (`buffer_1(input(x,y))`). If one of the instructions flagged during forward analysis as performing indirect buffer access computes the value of a leaf in the frontier, Helium adds additional operation nodes to the tree describing the address calculation expression (see Figure 4). The sources of these additional nodes are added to the frontier along with the other source operands of the instruction to ensure we capture both data and address calculation dependencies.

Recursive trees If Helium adds a node to the data-dependency tree describing a location from the same buffer as the root node, the tree is recursive. To avoid expanding the tree, Helium does not insert that node in the frontier. Instead, the tool builds an additional non-recursive data-dependency tree for the initial write to that output location to capture the base case of the recursion (see Figure 4). If all writes to that output location are recursively defined, Helium assumes that the buffer has been initialized outside the function.

Known library calls When Helium adds the return value of a call to a known external library function (e.g., `sqrt`, `floor`) to the tree, instead of continuing to expand the tree through that function, it adds an external call node that depends on the call arguments. Handling known calls specially allows Helium to emit corresponding Halide intrinsics instead of presenting the Halide optimizer with the library’s optimized implementation (which is often not vectorizable without heroic effort). Helium recognizes these external calls by their symbol, which is present even in stripped binaries because it is required for dynamic linking.

Canonicalization Helium canonicalizes the trees during construction to cope with the vagaries of instruction selection and ordering. For example, if the compiler unrolls a loop, it may commute some but not all of the resulting instructions in the loop body; Helium sorts the operands of commutative operations so it can recognize these trees as similar in the next step. It also applies simplification rules to these trees to account for effects of fix-up loops inserted by the compiler to handle leftover iterations of the unrolled loop. Figure 2 (c) shows the forest of canonicalized concrete trees.

Data types As Helium builds concrete trees, it records the sizes and kinds (signed/unsigned integer or floating-point) of registers and memory to emit the correct operation during Halide code generation (Section 4.11). Narrowing operations are represented as downcast nodes and overlapping dependencies are represented with full or partial overlap nodes.

Predication Each time Helium adds an instruction to the tree, if that instruction is annotated with one or more input-dependent conditionals identified during the forward analysis, it records the tree as predicated on those conditionals. Once it finishes constructing the tree for the computation of the output location, Helium builds similar concrete trees for the dependencies of the predicates the tree is predicated on (that is, the data dependencies that control whether the branches are taken or not taken). At the end of the backward analysis, Helium has built a concrete *computational tree* for each output location (or two trees if that location is updated recursively), each with zero or more concrete *predicate trees* attached. During

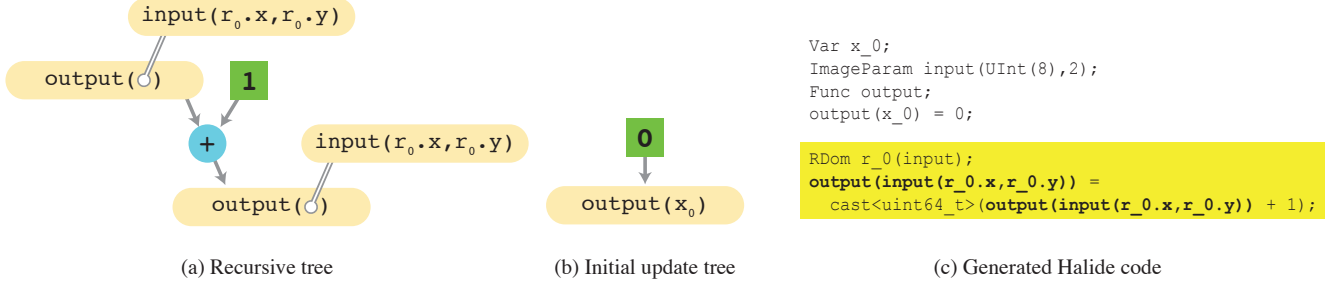


Figure 4. The trees and lifted Halide code for the histogram computation in Photoshop’s histogram equalization filter. The initial update tree (b) initializes the histogram counts to 0. The recursive tree (a) increments the histogram bins using indirect access based on the input image values. The Halide code generated from the recursive tree is highlighted and indirect accesses are in bold.

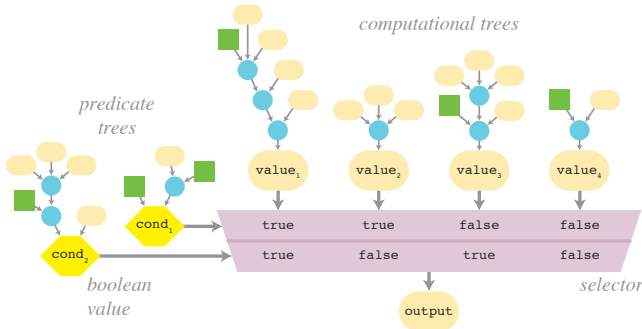


Figure 5. Each computational tree (four right trees) has zero or more predicate trees (two left trees) controlling its execution. Code generated for the predicate trees controls the execution of the code generated for the computational trees, like a multiplexer.

code generation, Helium uses predicate trees to generate code that selects which computational tree code to execute (see Figure 5).

4.8 Tree Clustering and Buffer Inference

Helium groups the concrete computational trees into clusters, where two trees are placed in the same cluster if they are the same, including all predicate trees they depend on, modulo constants and memory addresses in the leaves of the trees. (Recall that registers were mapped to special memory locations during preprocessing.) The number of clusters depends on the control dependency paths taken during execution for each output location. Each control dependency path will have its own cluster of computational trees. Most kernels have very few input-dependent conditionals relative to the input size, so there will usually be a small number of clusters each containing many trees. Figure 5 shows an example of clustering (only one computational tree is shown for brevity). For the 1D blur example, there is only one cluster as the computation is uniform across the padded image. For the threshold filter in Photoshop, we get two clusters.

Next, our goal is to abstract these trees. Using the dimensions, strides and extents inferred in 4.3, Helium can convert memory addresses to concrete indices (e.g., memory address 0xD3252A0 to `output_1(7,9)`). We call this *buffer inference*.

At this stage the tool also detects function parameters, assuming that any register or memory location that is not in a buffer is a parameter. After performing buffer inference on the concrete computational trees and attached predicate trees, we obtain a set of *abstract computational and predicate trees* for each cluster. Figure

2 (d) shows the forest of abstract trees for the 1D blur stencil. The leaves of these trees are buffers, constants or parameters.

4.9 Reduction Domain Inference

If a cluster contains recursive trees, Helium must infer a *reduction domain* specifying the range in each dimension for which the reduction is to be performed. If the root nodes of the recursive trees are indirectly accessed using the values of another buffer, then the reduction domain is the bounds of that other buffer. If the initial update tree depends on values originating outside the function, Helium assumes the reduction domain is the bounds of that input buffer.

Otherwise, the tool records the minimum and maximum buffer indices observed in trees in the cluster for each dimension as the bounds of the reduction domain. Helium abstracts these concrete indices using the assumption that the bounds are a linear combination of the buffer extents or constants. This heuristic has been sufficient for our applications, but a more precise determination could be made by applying the analysis to multiple sets of input data with varying dimensions and solving the resulting set of linear equations.

4.10 Symbolic Tree Generation

At this stage, the abstract trees contain many relations between different buffer locations (e.g., `output(3,4)` depends on `input(4,4)`, `input(3,3)`, and `input(3,4)`). To convert these dependencies between specific index values into symbolic dependencies between buffer coordinate locations, Helium assumes an affine relationship between indices and solves a linear system. The rest of this section details the procedure that Helium applies to the abstract computational and predicate trees to convert them into *symbolic trees*.

We represent a stencil with the following generic formulation. For the sake of brevity, we use a simple addition as our example and conditionals are omitted.

```

for  $x_1 = \dots$ 
  ...
  for  $x_D = \dots$ 
    output[ $x_1 \dots x_D$ ] =
      buffer1[ $f_{1,1}(x_1, \dots, x_D)$ ] ... [ $f_{1,k}(x_1, \dots, x_D)$ ] +
      ... + buffern[ $f_{n,1}(x_1, \dots, x_D)$ ] ... [ $f_{n,k}(x_1, \dots, x_D)$ ]
  
```

where `buffer` refers to the buffers that appear in leaf nodes in an abstract tree and `output` is the root of that tree. The functions $f_{1,1}, \dots, f_{n,k}$ are *index functions* that describe the relationship between the buffer indices and the output indices. Each index function is specific to a given leaf node and to a given dimension. In our work, we consider only affine index functions, which covers many practical scenarios. We also define the *access vector* $\vec{x} = (x_1, \dots, x_D)$.

For a D -dimensional output buffer at the root of the tree with access vector \vec{x} , a general affine index function for the leaf node ℓ and dimension d is $f_{\ell, \vec{a}}(\vec{x}) = [\vec{x}; 1] \cdot \vec{a}$ where \vec{a} is the $(D + 1)$ -dimensional vector of the affine coefficients that we seek to estimate. For a single abstract tree, this equation is underconstrained but since all the abstract trees in a cluster share the same index functions for each leaf node and dimension, Helium can accumulate constraints and make the problem well-posed. In each cluster, for each leaf node and dimension, Helium formulates a set of linear equations with \vec{a} as unknown.

In practice, for data-intensive applications, there are always at least $D + 1$ trees in each cluster, which guarantees that our tool can solve for \vec{a} . To prevent slowdown from a hugely overconstrained problem, Helium randomly selects a few trees to form the system. $D + 1$ trees would be enough to solve the system, but we use more to detect cases where the index function is not affine. Helium checks that the rank of the system is $D + 1$ and generates an error if it is not. In theory, $D + 2$ random trees would be sufficient to detect such cases with some probability; in our experiments, our tool uses $2D + 1$ random trees to increase detection probability.

Helium solves similar sets of linear equations to derive affine relationships between the output buffer indices and constant values in leaf nodes.

As a special case, if a particular cluster’s trees have an index in any dimension which does not change for all trees in that cluster, Helium assumes that dimension is fixed to that particular value instead of solving a linear system.

Once the tool selects a random set of abstract trees, a naïve solution to form the systems of equations corresponding to each leaf node and each dimension would be to go through all the trees each time. However, all the trees in each cluster have the same structure. This allows Helium to merge the randomly selected trees into a single *compound tree* with the same structure but with extended leaf and root nodes that contain all relevant buffers (Fig. 2(e)). With this tree, generating the systems of equations amounts to a single traversal of its leaf nodes.

At the end of this process, for each cluster, we now have a symbolic computational tree possibly associated with symbolic predicate trees as illustrated in Figure 2(g).

4.11 Halide Code Generation

The symbolic trees are a high-level representation of the algorithm. Helium extracts only the necessary set of predicate trees, ignoring control flow arising from loops. Our symbolic trees of data dependencies between buffers match Halide’s functional style, so code generation is straightforward.

Helium maps the computational tree in each cluster to a Halide function predicated on the predicate trees associated with it. Kernels without data-dependent control flow have just one computational tree, which maps directly to a Halide function. Kernels with data-dependent control flow have multiple computational trees that reference different predicate trees. The predicate trees are lifted to the top of the Halide function and mapped to a chain of `select` expressions (the Halide equivalent of C’s `?:` operator) that select the computational tree code to execute.

If a recursive tree’s base case is known, it is defined as a Halide function; otherwise, Helium assumes an initialized input buffer will be passed to the generated code. The inferred reduction domain is used to create a Halide `RDom` whose variables are used to define the recursive case of the tree as a second Halide function. Recursive trees can still be predicated as above.

5. Limitations

Lifting stencils with Helium is not a sound transformation. In practice, Helium’s lifted stencils can be compared against the

original program on a test suite – validation equivalent to release criteria commonly used in software development. Even if Helium were sound, most stripped binary programs do not come with proofs of correctness, so testing would still be required.

Some of Helium’s simplifying assumptions cannot always hold. The current system can only lift stencil computations with few input-dependent conditionals, table lookups and simple repeated updates. Helium cannot lift filters with non-stencil or more complex computation patterns. Because high performance kernels repeatedly apply the same computation to large amounts of data, Helium assumes the program input will exercise both branches of all input-dependent conditionals. For those few stencils with complex input-dependent control flow, the user must craft an input to cover all branches for Helium to successfully lift the stencil.

Helium is only able to find symbolic trees for stencils whose tree shape is constant. For trees whose shape varies based on a parameter (for example, box blur), Helium can extract code for individual values of the parameter, but the resulting code is not generic across parameters.

Helium assumes all index functions are affine, so kernels with more complex access functions such as radial indexing cannot be recognized by Helium.

By design, Helium only captures computations derived from the input data. Some stencils compute weights or lookup tables from parameters; Helium will capture the application of those tables to the input, but will not capture table computation.

6. Evaluation

6.1 Extraction Results

We used Helium to lift seven filters and portions of four more from Photoshop CS 6 Extended, four filters from IrfanView 4.38, and the smooth stencil from the miniGMG high-performance computing benchmark into Halide code. We do not have access to Photoshop or IrfanView source code; miniGMG is open source.

Photoshop We lifted Photoshop’s blur, blur more, sharpen, sharpen more, invert, threshold and box blur (for radius 1 only) filters. The blur and sharpen filters are 5-point stencils; blur more, sharpen more and box blur are 9-point stencils. Invert is a pointwise operation that simply flips all the pixel bits. Threshold is a pointwise operation containing an input-dependent conditional: if the input pixel’s brightness (a weighted sum of its R, G and B values) is greater than the threshold, the output pixel is set to white, and otherwise it is set to black.

We lifted portions of Photoshop’s sharpen edges, despeckle, histogram equalization and brightness filters. Sharpen edges alternates between repeatedly updating an image-sized side buffer and updating the image; we lifted the side buffer computation. Despeckle is a composition of blur more and sharpen edges. When run on despeckle, Helium extracts the blur more portion. From histogram equalization, we lifted the histogram calculation, but cannot track the histogram through the equalization stage because equalization does not depend on the input or output images. Brightness builds a 256-entry lookup table from its parameter, which we cannot capture because it does not depend on the images, but we do lift the application of the filter to the input image.

Photoshop contains multiple variants of its filters optimized for different x86 instruction sets (SSE, AVX etc.). Our instrumentation tools intercept the `cpuid` instruction (which tests CPU capabilities) and report to Photoshop that no vector instruction sets are supported; Photoshop falls back to general-purpose x86 instructions. We do this for engineering reasons, to reduce the number of opcodes our backward analysis must understand; this is not a fundamental limitation. The performance comparisons later in this section do not intercept `cpuid` and thus use optimized code paths in Photoshop.

Filter	total BB	diff BB	filter func BB	static ins. count	mem dump	dynamic ins. count	tree size
Invert	490663	3401	11	70	32 MB	5520	3
Blur	500850	3850	14	328	32 MB	64644	13
Blur More	499247	2825	16	189	38 MB	111664	62
Sharpen	492433	3027	30	351	36 MB	79369	31
Sharpen More	493608	3054	27	426	37 MB	105374	55
Threshold	491651	2728	60	363	36 MB	45861	8/6/19
Box Blur (radius 1)	500297	3306	94	534	28 MB	125254	253
Sharpen Edges	499086	2490	11	63	46 MB	80628	33
Despeckle	499247	2825	16	189	38 MB	111664	62
Equalize	501669	2771	47	198	8 MB	38243	6
Brightness	499292	3012	10	54	32 MB	21645	3

Figure 6. Code localization and extraction statistics for Photoshop filters, showing the total static basic blocks executed, the static basic blocks surviving screening (Section 3.1), the static basic blocks in the filter function selected at the end of localization (Section 3.3), the number of static instructions in the filter function, the memory dump size, the number of dynamic instructions captured in the instruction trace (Section 4.1), and the number of nodes per concrete tree. Threshold has two computational trees with 8 and 6 nodes and one predicate tree with 19 nodes. The filters below the line were not entirely extracted; the extracted portion of despeckle is the same as blur more. The total number of basic blocks executed varies due to unknown background code in Photoshop.

Figure 6 shows statistics for code localization, demonstrating that our progressive narrowing strategy allows our dynamic analysis to scale to large applications.

All but one of our lifted filters give bit-identical results to Photoshop’s filters on a suite of photographic images, each consisting of 100 megapixels. The lifted implementation of box blur, the only filter we lifted from Photoshop that uses floating-point, differs in the low-order bits of some pixel values due to reassociation.

IrfanView We lifted the blur, sharpen, invert and solarize filters from IrfanView, a batch image converter. IrfanView’s blur and sharpen are 9-point stencils. Unlike Photoshop, IrfanView loads the image data into floating-point registers, computes the stencil in floating-point, and rounds the result back to integer. IrfanView has been compiled for maximal processor compatibility, which results in unusual code making heavy use of partial register reads and writes.

Our lifted filters produce visually identical results to IrfanView’s filters. The minor differences in the low-order bits are because we assume floating-point addition and multiplication are associative and commutative when canonicalizing trees.

miniGMG To demonstrate the applicability of our tool beyond image processing, we lifted the Jacobi smooth stencil from the miniGMG high-performance computing benchmark. We added a command-line option to skip running the stencil to enable coverage differencing during code localization. Because we do not have input and output image data for this benchmark, we manually specified an estimate of the data size for finding candidate instructions during code localization and we used the generic inference described in section 4.3 during expression extraction. We set `OMP_NUM_THREADS=1` to limit miniGMG to one thread during analysis, but run using full parallelism during evaluation.

Because miniGMG is open source, we were able to check that our lifted stencil is equivalent to the original code using the SymPy¹ symbolic algebra system. We also checked output values for small data sizes.

6.2 Experimental Methodology

We ran our image filter experiments on an Intel Core i7 990X with 6 cores (hyperthreading disabled) running at 3.47GHz with 8 GB RAM and running 64-bit Windows 7. We used the Halide release built from git commit 80015c.

Helium We compiled our lifted Halide code into standalone executables that load an image, time repeated applications of the filter,

and save the image for verification. We ran 10 warmup iterations followed by 30 timing iterations. We tuned the schedules for our generated Halide code for six hours each using the OpenTuner-based Halide tuner [4]. We tuned using a 11267 by 8813 24-bit truecolor image and evaluated with a 11959 by 8135 24-bit image.

We cannot usefully compare the performance of the four Photoshop filters that we did not entirely lift this way; we describe their evaluation separately in Section 6.5.

Photoshop We timed Photoshop using the ExtendScript API² to programmatically start Photoshop, load the image and invoke filters. While we extracted filters from non-optimized fallback code for old processors, times reported in this section are using Photoshop’s choice of code path. In Photoshop’s performance preferences, we set the tile size to 1028K (the largest), history states to 1, and cache tiles to 1. This amounts to optimizing Photoshop for batch processing, improving performance by up to 48% over default settings while dramatically reducing measurement variance. We ran 10 warmup iterations and 30 evaluation iterations.

IrfanView IrfanView does not have a scripting interface allowing for timing, so we timed IrfanView running from the command line using PowerShell’s `Measure-Command`. We timed 30 executions of IrfanView running each filter and another 30 executions that read and wrote the image without operating on it, taking the difference as the filter execution time.

miniGMG miniGMG is open source, so we compared unmodified miniGMG performance against a version with the loop in the `smooth` stencil function from the OpenMP-based Jacobi smoother replaced with a call to our Halide-compiled lifted stencil. We used a generic Halide schedule template that parallelizes the outer dimension and, when possible, vectorizes the inner dimension. We compared performance against miniGMG using OpenMP on a 2.4GHz Intel Xeon E5-2695v2 machine running Linux with two sockets, 12 cores per socket and 128GB RAM.

6.3 Lifted Filter Performance Results

Photoshop Figure 7 compares Photoshop’s filters against our standalone executable running our lifted Halide code. We obtain an average speedup of 1.75 on the individual filters (1.90 excluding box blur).

Profiling using Intel VTune shows that Photoshop’s blur filter is not vectorized. Photoshop does parallelize across all the machine’s

¹<http://sympy.org>

²<https://www.adobe.com/devnet/photoshop/scripting.html>

Filter	Photoshop	Helium	speedup
Invert	102.23 ± 1.65	58.74 ± .52	1.74x
Blur	245.87 ± 5.30	93.74 ± .78	2.62x
Blur More	317.97 ± 2.76	283.92 ± 2.52	1.12x
Sharpen	270.40 ± 5.80	110.07 ± .69	2.46x
Sharpen More	305.50 ± 4.13	147.01 ± 2.29	2.08x
Threshold	169.83 ± 1.37	119.34 ± 8.06	1.42x
Box Blur	273.87 ± 2.42	343.02 ± .59	.80x

Filter	IrfanView	Helium	speedup
Invert	215.23 ± 37.98	105.94 ± .78	2.03x
Solarize	220.51 ± 46.96	102.21 ± .55	2.16x
Blur	3129.68 ± 17.39	359.84 ± 3.96	8.70x
Sharpen	3419.67 ± 52.56	489.84 ± 7.78	6.98x

Figure 7. Timing comparison (in milliseconds) between Photoshop and IrfanView filters and our lifted Halide-implemented filters on a 11959 by 8135 24-bit truecolor image.

hardware threads, but each thread only achieves 10-30% utilization. Our lifted filter provides better performance by blocking and vectorizing in addition to parallelizing.

Photoshop implements box blur with a sliding window, adding one pixel entering the window and subtracting the pixel leaving the window. Our lifted implementation of box blur is slower because Helium cancels these additions and subtractions when canonicalizing the tree, undoing the sliding window optimization.

IrfanView Figure 7 compares IrfanView’s filters against our Halide applications. We obtain an average speedup of 4.97.

miniGMG For miniGMG, we measure the total time spent in the stencil we translate across all iterations of the multigrid invocation. Unmodified miniGMG spends 28.5 seconds in the kernel, while miniGMG modified to use our lifted Halide smooth stencil finishes in 6.7 seconds for a speedup of 4.25.

6.4 Performance of Filter Pipelines

We also compare filter pipelines to demonstrate how lifting to a very high-level representation enables additional performance improvements through stencil composition. For Photoshop, our pipeline consists of blur, invert and sharpen more applied consecutively, while for IrfanView we ran a pipeline of sharpen, solarize and blur.

We obtain a speedup of 2.91 for the Photoshop pipeline. Photoshop blurs the entire image, then inverts it, then sharpens more, which has poor locality. Halide inlines blur and invert inside the loops for sharpen more, improving locality while maintaining vectorization and parallelism.

We obtain a speedup of 5.17 for the IrfanView pipeline. IrfanView improves when running the filters as a pipeline, apparently by amortizing the cost of a one-time preparation step, but our Halide code improves further by fusing the actual filters.

6.5 In Situ Replacement Photoshop Performance

To evaluate the performance impact of the filters we partially extracted from Photoshop, we replaced Photoshop’s implementation with our automatically-generated Halide code using manually-implemented binary patches. We compiled all our Halide code into a DLL that patches specific addresses in Photoshop’s code with calls to our Halide code. Other than improved performance, these patches are entirely transparent to the user. The disadvantage of this approach is that the patched kernels are constrained by optimization decisions made in Photoshop, such as the granularity of tiling, which restricts our ability to fully optimize the kernels.

When timing the replacements for filters we entirely lift, we disabled Photoshop’s parallelism by removing `MULTIPROCESSOR`

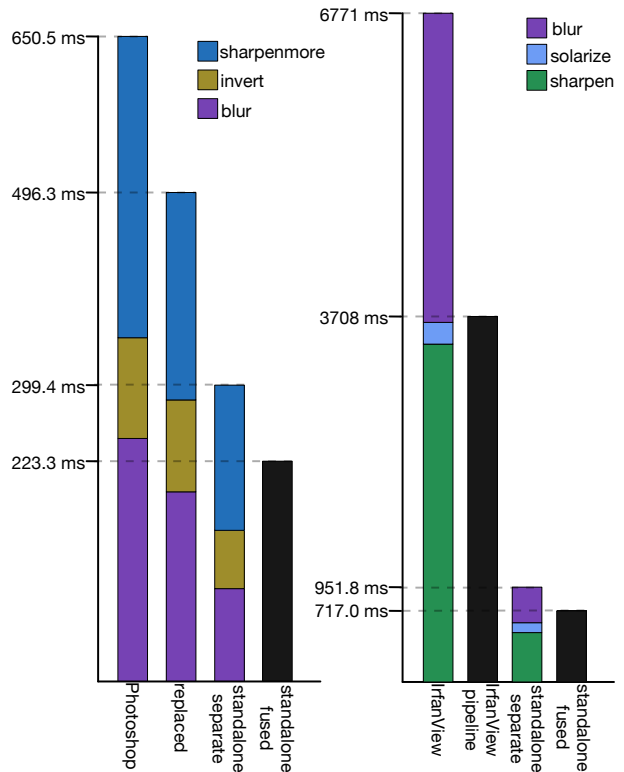


Figure 8. Performance comparison of Photoshop and IrfanView pipelines. Left-to-right, the left graph shows Photoshop running the filters in sequence, Photoshop hosting our lifted implementations (Section 6.5), our standalone Halide executable running the filters in sequence, and our Halide executable running the fused pipeline. The right graph shows IrfanView running the filters in sequence, IrfanView running the filters as a pipeline (in one IrfanView instance), our Halide executable running the filters in sequence, and our Halide executable running the fused pipeline.

`SUPPORT_8BX` from the Photoshop installation, allowing our Halide code to control parallelism subject to the granularity limit imposed by Photoshop’s tile size. When timing the filters we only partially lift, we removed parallelism from the Halide schedule and allow Photoshop to parallelize around our Halide code. While we would prefer to control parallelism ourselves, enough Photoshop code is executing outside the regions we replaced to make disabling Photoshop’s parallelism a large performance hit.

Figure 9 compares unmodified Photoshop (same numbers as in the previous section) with Photoshop after in situ replacement. For the fully-lifted filters, we are still able to improve performance even while not fully in control of the environment. Our replacement for box blur is still slower for the reason described in Section 6.3. The portions of histogram equalization and brightness we lift are too simple to improve, but the replaced sharpen edges is slightly faster, demonstrating that even when Helium cannot lift the entire computation, it can still lift a performance-relevant portion.

7. Related Work

Binary static analysis Phoenix [23], BitBlaze [24], BAP [9], and other tools construct their own low-level IR (*e.g.*, register transfer language (RTL)) from binaries. These low-level IRs allow only limited analysis or low-level transformations.

Filter	Photoshop	replacement	speedup
Invert	102.23 ± 1.65	93.20 ± .71	1.10x
Blur	245.87 ± 5.30	191.83 ± 1.12	1.28x
Blur More	317.97 ± 2.76	310.70 ± .88	1.02x
Sharpen	270.40 ± 5.80	194.80 ± .66	1.39x
Sharpen More	305.50 ± 4.13	210.20 ± .71	1.45x
Threshold	169.83 ± 1.37	124.10 ± .76	1.37x
Box Blur	273.87 ± 2.42	395.40 ± .72	.69x
Sharpen Edges	798.43 ± 1.45	728.63 ± 1.85	1.10x
Despeckle	763.87 ± 1.59	756.40 ± 1.59	1.01x
Equalize	405.50 ± 1.45	433.87 ± .90	.93x
Brightness	498.00 ± 1.31	503.47 ± 1.17	.99x

Figure 9. Timing comparison (in milliseconds) between Photoshop filters and in situ replacement with our lifted Halide-implemented filters on a 11959 by 8135 24-bit truecolor image.

Other static analysis aims for high-level representations of binaries. Value set analysis [6] is a static analysis that tracks the possible values of pointers and indices to analyze memory access in stripped binaries; instead of a complicated static analysis, Helium recovers buffer structure from actual program behavior captured in traces. SecondWrite [3], [13] decompiles x86 binaries to LLVM IR; while the resulting IR can be optimized and recompiled, this IR is too low-level to get more than minor speedup over existing optimized binaries. [18] uses SecondWrite for automatic parallelization of affine loops, but must analyze existing loop structure and resolve aliasing, while we lift to Halide code expressing only the algorithm. McSema [2] also decompiles x86 to LLVM IR for analysis. The Hex-Rays decompiler [1] decompiles to a C-like pseudocode which cannot be recompiled to binaries. SmartDec [14] is a binary to C++ decompiler that can extract class hierarchies and try/catch blocks; we extract high-level algorithms independent of particular language constructs.

Dynamic translation and instrumentation Binary translation systems like QEMU [7] translate machine code between architectures using RISC-like IR; RevNIC [11] and S2E [12] translate programs from x86 to LLVM IR by running them in QEMU. Dynamic instrumentation systems like Valgrind [20] present a similar RISC-like IR for analysis and instrumentation, then generate machine code for execution. These IRs retain details of the original binary and do not provide enough abstraction for high-level transformations.

Microarchitecture-level dynamic binary optimization Some systems improve the performance of existing binaries through microarchitecture-level optimizations that do not require building IR. Dynamo [5] improves code locality and applies simple optimizations on frequently-executed code. Ubiquitous memory introspection [33] detects frequently-stalling loads and adds prefetch instructions. [19] translates x86 binaries to x86-64, using the additional registers to promote stack variables. We perform much higher-level optimizations on our lifted stencils.

Automatic parallelization Many automatic parallelization systems use dynamic analysis to track data flow to analyze communication to detect parallelization opportunities, but these systems require source code access (often with manual annotations). [27] uses dynamic analysis to track communication across programmer-annotated pipeline boundaries to extract coarse-grained pipeline parallelism. Parallax [28] performs semi-automatic parallelization, using dynamic dependency tracking to suggest programmer annotations (e.g., that a variable is killed). HELIX [10] uses a dynamic loop nesting graph to select a set of loops to parallelize. [29] uses dynamic analysis of control and data dependencies as input to a trained predictor to autoparallelize loops, relying on the user to check correctness.

Pointer and shape analysis Pointer analyses have been written for assembly programs [15]. Shape analyses [30] analyze programs statically to determine properties of heap structures. [32] uses dynamic analysis to identify pointer-chasing that sometimes exhibits strides to aid in placing prefetch instructions. Because we analyze concrete memory traces for stencils, our buffer structure reconstruction and stride inference is indifferent to aliasing and finds regular access patterns.

8. Conclusion

Most legacy high-performance applications exhibit bit rot during the useful lifetime of the application. We can no longer rely on Moore’s Law to provide transparent performance improvements from clock speed scaling, but at the same time modern hardware provides ample opportunities to substantially improve performance of legacy programs. To rejuvenate these programs, we need high-level, easily-optimizable representations of their algorithms. However, high-performance kernels in these applications have been heavily optimized for a bygone era, resulting in complex source code and executables, even though the underlying algorithms are mostly very simple. Current state-of-the-art techniques are not capable of extracting the simple algorithms from these highly optimized programs. We believe that fully dynamic techniques, introduced in Helium, are a promising direction for lifting important computations into higher-level representations and rejuvenating legacy applications.

Helium source code is available at <http://projects.csail.mit.edu/helium>.

Acknowledgments

We would like to thank Vladimir Kiriansky, Derek Bruening, and the DynamoRIO user group for invaluable help debugging DynamoRIO clients and Sarah Kong, Chris Cox, Joseph Hsieh, Alan Erickson, and Jeff Chien of the Photoshop team for their helpful input. This material is based upon work supported by DOE awards DE-SC0005288 and DE-SC0008923 and DARPA agreement FA8750-14-2-0009. Charith Mendis was supported by a MITEI fellowship.

References

- [1] Idapro, hexrays. URL <http://www.hex-rays.com/idapro/>.
- [2] Mcsema: Static translation of x86 into llvm. 2014.
- [3] K. Anand, M. Smithson, K. Elwazeer, A. Kotha, J. Gruen, N. Giles, and R. Barua. A compiler-level intermediate representation based binary analysis and rewriting system. In *Proceedings of the 8th ACM European Conference on Computer Systems, EuroSys ’13*, pages 295–308, New York, NY, USA, 2013. ACM. ISBN 978-1-4503-1994-2. URL <http://doi.acm.org/10.1145/2465351.2465380>.
- [4] J. Ansel, S. Kamil, K. Veeramachaneni, J. Ragan-Kelley, J. Bosboom, U.-M. O’Reilly, and S. Amarasinghe. Opentuner: An extensible framework for program autotuning. In *International Conference on Parallel Architectures and Compilation Techniques*, Edmonton, Canada, August 2014.
- [5] V. Bala, E. Duesterwald, and S. Banerjia. Dynamo: A transparent dynamic optimization system. In *Proceedings of the ACM SIGPLAN 2000 Conference on Programming Language Design and Implementation, PLDI ’00*, pages 1–12, New York, NY, USA, 2000. ACM. ISBN 1-58113-199-2. URL <http://doi.acm.org/10.1145/349299.349303>.
- [6] G. Balakrishnan and T. Reps. Analyzing memory accesses in x86 executables. In E. Duesterwald, editor, *Compiler Construction*, volume 2985 of *Lecture Notes in Computer Science*, pages 5–23. Springer Berlin Heidelberg, 2004. ISBN 978-3-540-21297-3.
- [7] F. Bellard. QEMU, a fast and portable dynamic translator. In *Proceedings of the Annual Conference on USENIX Annual Technical Conference, ATEC ’05*, pages 41–41, Berkeley, CA, USA, 2005. USENIX Association. URL www.qemu.org.

- [8] D. Bruening, Q. Zhao, and S. Amarasinghe. Transparent dynamic instrumentation. In *Proceedings of the 8th ACM SIGPLAN/SIGOPS Conference on Virtual Execution Environments*, VEE '12, pages 133–144, New York, NY, USA, 2012. ACM. ISBN 978-1-4503-1176-2. . URL <http://doi.acm.org/10.1145/2151024.2151043>.
- [9] D. Brumley, I. Jager, T. Avgerinos, and E. J. Schwartz. BAP: A binary analysis platform. In *Proceedings of the 23rd International Conference on Computer Aided Verification*, CAV'11, pages 463–469, Berlin, Heidelberg, 2011. Springer-Verlag. ISBN 978-3-642-22109-5. URL <http://dl.acm.org/citation.cfm?id=2032305.2032342>.
- [10] S. Campanoni, T. Jones, G. Holloway, V. J. Reddi, G.-Y. Wei, and D. Brooks. HELIX: Automatic parallelization of irregular programs for chip multiprocessing. In *Proceedings of the Tenth International Symposium on Code Generation and Optimization*, CGO '12, pages 84–93, New York, NY, USA, 2012. ACM. ISBN 978-1-4503-1206-6. . URL <http://doi.acm.org/10.1145/2259016.2259028>.
- [11] V. Chipounov and G. Candea. Reverse engineering of binary device drivers with RevNIC. In *Proceedings of the 5th European Conference on Computer Systems*, EuroSys '10, pages 167–180, New York, NY, USA, 2010. ACM. ISBN 978-1-60558-577-2. . URL <http://doi.acm.org/10.1145/1755913.1755932>.
- [12] V. Chipounov, V. Kuznetsov, and G. Candea. S2E: A platform for in-vivo multi-path analysis of software systems. In *Proceedings of the Sixteenth International Conference on Architectural Support for Programming Languages and Operating Systems*, ASPLOS XVI, pages 265–278, New York, NY, USA, 2011. ACM. ISBN 978-1-4503-0266-1. . URL <http://doi.acm.org/10.1145/1950365.1950396>.
- [13] K. ElWazeer, K. Anand, A. Kotha, M. Smithson, and R. Barua. Scalable variable and data type detection in a binary rewriter. In *Proceedings of the 34th ACM SIGPLAN Conference on Programming Language Design and Implementation*, PLDI '13, pages 51–60, New York, NY, USA, 2013. ACM. ISBN 978-1-4503-2014-6. . URL <http://doi.acm.org/10.1145/2491956.2462165>.
- [14] A. Fokin, E. Derevenec, A. Chernov, and K. Troshina. Smartdec: Approaching c++ decompilation. In *Proceedings of the 2011 18th Working Conference on Reverse Engineering*, WCRE '11, pages 347–356, Washington, DC, USA, 2011. IEEE Computer Society. ISBN 978-0-7695-4582-0. . URL <http://dx.doi.org/10.1109/WCRE.2011.49>.
- [15] B. Guo, M. Bridges, S. Triantafyllis, G. Ottoni, E. Raman, and D. August. Practical and accurate low-level pointer analysis. In *Code Generation and Optimization*, 2005. CGO 2005. *International Symposium on*, pages 291–302, March 2005. .
- [16] R. N. Horspool and N. Marovac. An approach to the problem of detranslation of computer programs. *The Computer Journal*, 23(3): 223–229, 1980.
- [17] S. Kamil, C. Chan, L. Oliker, J. Shalf, and S. Williams. An auto-tuning framework for parallel multicore stencil computations. In *Parallel Distributed Processing (IPDPS), 2010 IEEE International Symposium on*, pages 1–12, April 2010. .
- [18] A. Kotha, K. Anand, M. Smithson, G. Yellareddy, and R. Barua. Automatic parallelization in a binary rewriter. In *Proceedings of the 2010 43rd Annual IEEE/ACM International Symposium on Microarchitecture*, MICRO '10, pages 547–557, Washington, DC, USA, 2010. IEEE Computer Society. ISBN 978-0-7695-4299-7. . URL <http://dx.doi.org/10.1109/MICRO.2010.27>.
- [19] J. Li, C. Wu, and W.-C. Hsu. Dynamic register promotion of stack variables. In *Proceedings of the 9th Annual IEEE/ACM International Symposium on Code Generation and Optimization*, CGO '11, pages 21–31, Washington, DC, USA, 2011. IEEE Computer Society. ISBN 978-1-61284-356-8. URL <http://dl.acm.org/citation.cfm?id=2190025.2190050>.
- [20] N. Nethercote and J. Seward. Valgrind: A framework for heavyweight dynamic binary instrumentation. In *Proceedings of the 2007 ACM SIGPLAN Conference on Programming Language Design and Implementation*, PLDI '07, pages 89–100, New York, NY, USA, 2007. ACM. ISBN 978-1-59593-633-2. . URL <http://doi.acm.org/10.1145/1250734.1250746>.
- [21] S. Paris. Adobe systems. personal communication, 2014.
- [22] J. Ragan-Kelley, C. Barnes, A. Adams, S. Paris, F. Durand, and S. Amarasinghe. Halide: A language and compiler for optimizing parallelism, locality, and recomputation in image processing pipelines. In *Proceedings of the 34th ACM SIGPLAN Conference on Programming Language Design and Implementation*, PLDI '13, 2013. . URL <http://doi.acm.org/10.1145/2491956.2462176>.
- [23] M. Research. Phoenix compiler and shared source common language infrastructure. URL <http://www.research.microsoft.com/phoenix>.
- [24] D. Song, D. Brumley, H. Yin, J. Caballero, I. Jager, M. G. Kang, Z. Liang, J. Newsome, P. Poosankam, and P. Saxena. BitBlaze: A new approach to computer security via binary analysis. In *Proceedings of the 4th International Conference on Information Systems Security. Keynote invited paper*, Hyderabad, India, Dec. 2008.
- [25] K. Stock, M. Kong, T. Grosser, L.-N. Pouchet, F. Rastello, J. Ramanujam, and P. Sadayappan. A framework for enhancing data reuse via associative reordering. In *Proceedings of the 35th ACM SIGPLAN Conference on Programming Language Design and Implementation*, PLDI '14, pages 65–76, New York, NY, USA, 2014. ACM. ISBN 978-1-4503-2784-8. . URL <http://doi.acm.org/10.1145/2594291.2594342>.
- [26] Y. Tang, R. A. Chowdhury, B. C. Kuszmaul, C.-K. Luk, and C. E. Leiserson. The pochoir stencil compiler. In *Proceedings of the Twenty-third Annual ACM Symposium on Parallelism in Algorithms and Architectures*, SPAA '11, pages 117–128, New York, NY, USA, 2011. ACM. ISBN 978-1-4503-0743-7. . URL <http://doi.acm.org/10.1145/1989493.1989508>.
- [27] W. Thies, V. Chandrasekhar, and S. Amarasinghe. A practical approach to exploiting coarse-grained pipeline parallelism in c programs. In *Proceedings of the 40th Annual IEEE/ACM International Symposium on Microarchitecture*, MICRO 40, pages 356–369, Washington, DC, USA, 2007. IEEE Computer Society. ISBN 0-7695-3047-8. . URL <http://dx.doi.org/10.1109/MICRO.2007.7>.
- [28] H. Vandierendonck, S. Rul, and K. De Bosschere. The paralax infrastructure: Automatic parallelization with a helping hand. In *Proceedings of the 19th International Conference on Parallel Architectures and Compilation Techniques*, PACT '10, pages 389–400, New York, NY, USA, 2010. ACM. ISBN 978-1-4503-0178-7. . URL <http://doi.acm.org/10.1145/1854273.1854322>.
- [29] Z. Wang, G. Tournavitis, B. Franke, and M. F. P. O'boyle. Integrating profile-driven parallelism detection and machine-learning-based mapping. *ACM Trans. Archit. Code Optim.*, 11(1):2:1–2:26, Feb. 2014. ISSN 1544-3566. . URL <http://doi.acm.org/10.1145/2579561>.
- [30] R. Wilhelm, M. Sagiv, and T. Reps. Shape analysis. In D. Watt, editor, *Compiler Construction*, volume 1781 of *Lecture Notes in Computer Science*, pages 1–17. Springer Berlin Heidelberg, 2000. ISBN 978-3-540-67263-0. . URL http://dx.doi.org/10.1007/3-540-46423-9_1.
- [31] S. Williams, D. D. Kalamkar, A. Singh, A. M. Deshpande, B. Van Straalen, M. Smelyanskiy, A. Almgren, P. Dubey, J. Shalf, and L. Oliker. Optimization of geometric multigrid for emerging multi- and manycore processors. In *Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis*, page 96. IEEE Computer Society Press, 2012.
- [32] Y. Wu. Efficient discovery of regular stride patterns in irregular programs and its use in compiler prefetching. In *Proceedings of the ACM SIGPLAN 2002 Conference on Programming Language Design and Implementation*, PLDI '02, pages 210–221, New York, NY, USA, 2002. ACM. ISBN 1-58113-463-0. . URL <http://doi.acm.org/10.1145/512529.512555>.
- [33] Q. Zhao, R. Rabbah, S. Amarasinghe, L. Rudolph, and W.-F. Wong. Ubiquitous memory introspection. In *International Symposium on Code Generation and Optimization*, San Jose, CA, Mar 2007. URL <http://groups.csail.mit.edu/commit/papers/07/zhao-cgo07-umi.pdf>.