

High-Performance GPU-to-CPU Transpilation and Optimization via High-Level Parallel Constructs

William S. Moses*, Ivan R. Ivanov†, Jens Domke‡, Toshio Endo†, Johannes Doerfert¶, and Oleksandr Zinenko§

* Massachusetts Institute of Technology, USA wmoses@mit.edu

† Tokyo Institute of Technology, Japan ivanov.i.aa@m.titech.ac.jp, endo@is.titech.ac.jp

‡ RIKEN Center for Computational Science, Japan jens.domke@riken.jp

¶ Argonne National Laboratory, USA jdoerfert@anl.gov

§ Google, France zinenko@google.com

摘要—*警告：该PDF由GPT-Academic开源项目调用大语言模型+Latex翻译插件一键生成，版权归原文作者所有。翻译内容可靠性无保障，请仔细鉴别并以原文为准。项目Github地址 https://github.com/binary-husky/gpt_academic/。当前大语言模型：gpt-4o-mini，当前语言模型温度设定：0。为了防止大语言模型的意外谬误产生扩散影响，禁止移除或修改此警告。

在并行性仍然是性能的主要来源的同时，体系结构实现和编程模型随着每一代新硬件而变化，这往往导致昂贵的应用重构。大多数性能可移植性工具都需要手动和昂贵的应用移植到另一种编程模型。

我们提出了一种替代方法，它基于Polygeist/MLIR自动将用一种编程模型（CUDA）编写的程序转换为另一种（CPU线程）。我们的方法包括并行构造的表征（representation），允许常规编译器转换透明且无需修改地应用，并支持并行性特定的优化。我们通过为多核CPU转译并优化CUDA Rodinia基准测试套件来评估我们的框架，取得了相较于手写的OpenMP代码76%的几何平均加速。此外，我们展示了来自PyTorch的CUDA内核如何在无用户干预的情况下高效地在仅限CPU的超级计算机Fugaku上运行和扩展。我们的PyTorch兼容层利用转译的 CUDA PyTorch内核的性能超过了PyTorch CPU原生后端2.7×。

I. INTRODUCTION

单核性能扩展的结束意味着并行化和领域特定性现在是效率提升的主要来源。超级计算机架构师在创新上展开竞争，以支持从物理仿真（simulation）到机器学习的计算和内存密集型应用。最新最快的超级计算机Fugaku完全基于A64FX CPU，这与商品CPU不同，提供高带宽内存访问的支持，并且在能效方面可与最近的GPU相媲美 [1]。

然而，如何高效且富有成效地利用这些计算机进行实际应用是一个挑战，因为最近的框架和高性能库是在考虑NVIDIA GPU的情况下开发的。例如，将PyTorch [2]移植到A64FX的尝试遇到了多重挑战。“原生”（native）

默认的CPU PyTorch后端仅为关键内核提供了简单版本，如实现为六层嵌套循环的2D卷积。英特尔的oneDNN [3]在Arm CPU上表现不佳并不令人意外，因为它是为没有高带宽内存的商品CPU量身定制的。富士通定制的oneDNN [4]被更好地调优，但在与GPU的竞争中并不具备普遍的优势。

```
__device__ float sum(float* data, int n) { ... }
__global__
void normalize(float *out, float* in, int n) {
    int tid = blockIdx.x + blockDim.x * threadIdx.x;
    // Optimization: Compute the sum once per block.
    // __shared__ int val;
    // if (threadIdx.x == 0) val = sum(in, n);
    // __syncthreads;
    float val = sum(in, n);
    if (tid < n)
        out[tid] = in[tid] / val;
}
void launch(int *d_out, int* d_in, int n) {
    normalize<<<(n+31)/32, 32>>>(d_out, d_in, n);
}
```

图 1. 一个示例 CUDA 程序 normalize，它对一个向量进行归一化，以及一个调用内核的 CPU 函数 launch。目前，sum 在每个线程中被调用，导致总的工作量为 $O(N^2)$ 。通过使用共享内存（见上面的注释），工作量可以部分减少到 $O(N^2/B)$ ，这使得每个块中 sum 的计算只需执行一次，或者通过在内核之前计算 sum 一次，可以完全减少到 $O(N)$ 。

许多非库方法被提出以实现性能可移植性。这些方法从语言扩展如OpenCL [5]或OpenACC [6]，到并行编程框架如Kokkos [7]，再到特定领域语言如SPIRAL [8]、Halide [9]或Tensor Comprehensions [10]。由于语言或原始程序与目标框架之间的编程模型（model）差异，所有这些方法仍然要求遗留（legacy）应用程序被移植，有时需要完全重

写。

我们探索了一种基于完全自动化编译器的替代方法,该编译器接受一种编程模型 (model) (CUDA) 中的代码并生成目标另一个模型 (model) (CPU线程) 的二进制文件。尽管过去曾探索GPU到CPU的转换 [11]–[13],但很少能生成高效的代码。事实上,针对CPU的优化,甚至是通用编译器变换,如公共子表达式消除或循环不变代码运动,因缺乏可以分析的并行结构的表征 (representation) 而受到阻碍 [14]。以图 1 中的求和为例,可以通过每个块的单个线程完成,或者在所有计算发生在CPU上后,由一个总线程完成。主流编译器内的并行性表征 (representation) 在最近才开始探索 [15]–[19],现有的变换是有限的,通常仅适用于简单的CPU代码。

我们提出了一种针对最常见GPU构造的编译器模型 (model): 多级并行性、级别范围同步和级别局部存储器。这与提供单级并行性、统一内存和对等同步的CPU并行性不同。与在优化管道之前操作的源和AST级方法以及将同步建模为“黑盒”优化障碍的现有编译器方法相反,我们完全从内存语义建模同步。这既允许基于同步的代码与现有优化进行交互,也使得新颖的特定并行性优化成为可能。

我们的模型 (model) 在LLVM编译器基础设施的MLIR层 [20]实现,并利用MLIR的嵌套模块方法来处理GPU代码 [22]。我们扩展了Polygeist [23] C/C++前端以支持CUDA,并生成保留高级并行性和程序结构的MLIR。我们的原型编译器能够将PyTorch CUDA内核以及其他计算密集型基准编译为LLVM支持的任何CPU架构。除了考虑执行模型差异的变换之外,我们还通过OpenMP利用CPU上的并行性。最后,我们的MocCUDA PyTorch集成使用我们的方法在没有GPU的情况下编译和执行CUDA内核,同时替代不受支持的调用。

我们通过编译Rodinia CUDA基准 [24]以及PyTorch CUDA内核来评估我们端到端转换的正确性和效率。当以商品CPU为目标时,我们的OpenMP加速的CUDA代码与Rodinia套件中的参考OpenMP实现性能相当,并且可扩展性有所提高。当使用我们的框架在仅CPU的Fugaku超级计算机上运行PyTorch时,我们在处理每秒图像数量上达到了与基于OneDNN的PyTorch CPU后端相比,约为Resnet-50 [25]的conv2d内核的两倍,同时整体训练的性能可比。

总体而言,我们的论文做出了以下贡献:

- 一种通用的高级和与平台无关的表征 (representation), 基于障碍同步的语义定义,通过内存语义确保正确

性,从而确保现有优化的透明应用。

- 新的并行特定优化可以利用我们的高层并行语义来优化程序。
- 对Polygeist C/C++前端的扩展,旨在将GPU和CPU并行结构直接映射到我们的高级并行原语中。
- 对Rodinia [24]基准套件的一个子集以及在PyTorch [2]中运行Resnet-50所需的内部CUDA内核进行端到端的CUDA到CPU的转译,以便在仅支持CPU的Fugaku超级计算机上运行。

II. BACKGROUND

主流编译器如 Clang 和 GCC 缺乏统一的高层次并行性 (parallelism) 表征。具体而言,在诸如 CUDA、OpenMP 或 SYCL 等框架中编译并行构造时,强制并行区域的主体存在于一个单独的 (closure) 函数中,该函数由相应的运行时调用。线程索引 (thread index) 或同步 (synchronization) 等概念通常会被单独表征,往往通过不透明的内置调用 (intrinsic calls)。由于编译器历史上缺乏关于并行性和所涉及的运行时效果的信息,任何并行构造也无意中成为优化的障碍。尽管近年来已经有一些尝试 [14]–[19], [26] 来改善 CPU 并行构造的表征,但加速器编程带来了额外的挑战。独特的编程模型和复杂的内存层次结构导致主流编译器中 GPU 并行性 (parallelism) 的高层次表征仍然探讨不足。

A. GPU Compilation

考虑图 1 中的 CUDA 程序,该程序对一个向量进行归一化。当使用 LLVM/Clang 编译时,GPU 程序是一个独立的编译单元,如图 II 所示。这阻止了 GPU 内核与 CPU 调用代码之间的任何优化。在图 1 的情况下,传统编译器中的程序总工作量为 $O(N^2)$,因为每个线程都执行 $O(N)$ 的 sum 调用。然而,如果在内核调用之前仅能执行一次 sum 调用,例如,通过执行循环不变代码移动 (LICM) 变换,则工作量将减少到 $O(N)$ 。这种优化的一个效果较差的变体可以通过使用共享内存将工作量减少到 $O(\frac{N^2}{B})$ 。MLIR 为 GPU 程序提供了嵌套模块表征 (representation),支持主机/设备代码移动 [22],但尚未实现并行代码移动。在 GPU 到 CPU 的代码移动中,LICM 从并行循环中移出始终是合法的,因为以前的设备内存存在主机上也可用。

B. MLIR Infrastructure

MLIR 是一个为重用和可扩展性设计的近期编译器基础设施 [20]。其内部表征 (IR) 可以看作是广泛采用的 LLVM IR 的概念性继承者 [21],并且对传统 CPU 模型

```

target triple = "x86_64-unknown-linux-gnu"

define @launch(float* %d_out, float* %d_in, i32 %n) {
  call @__cudaPushCallConfiguration(...)
  call @__cudaLaunchKernel(@normalize_stub, ...)
  ret
}

target triple = "nvptx64"

define @normalize(float* %out, float* %in, i32 %n) {
  %tid = call i32 @llvm.nvvm.ptx.tid.x()
  %sum = call i32 @sum(i32* %in, i32 %n)
  %cmp = icmp slt i32 %tid, %n
  br i1 %cmp, label %body, label %exit

body:
  %gep = getelementptr float* %in, i32 %tid
  %load = load float, float* %gep
  %nrm = fdiv float %load, %sum
  %ptr = getelementptr float* %out, i32 %tid
  store float %nrm, float* %ptr
  br label %exit

exit:
  ret
}

```

图 2. 在与LLVM/Clang编译时，图1中的launch（启动）和normalize（归一化）函数的简化降低。由于这两个函数生成不同的汇编代码，它们被放置在没有调用上下文的独立模块中。

之外的编程模型提供了更好的支持。MLIR 并不是提供一组预定义的指令和类型，而是操作在包含可互操作的用户定义操作、属性和类型集合的方言上。操作是对IR指令的泛化，可以是任意复杂的，特别是可以包含更多IR的区域，从而创建嵌套表征。操作定义并使用遵循单一静态赋值（SSA）[27]的值。例如，MLIR方言可以对整个物理或虚拟指令集建模，例如NVVM（针对Nvidia GPU的虚拟IR）、其他IR（例如LLVM IR）、更高层次的控制流构造（如仿射循环）、并行编程模型（如OpenMP和OpenACC）、机器学习图等。MLIR支持并鼓励在同一编译单元中混合来自不同方言的操作。

MLIR支持GPU，得益于其同名方言，该方言定义了高层次的SIMT编程模型以及主机/设备通信机制，并且一组低层次的特定平台方言：NVVM（CUDA）、ROCDL（ROCm）和SPIR-V。MLIR对GPU编程的特殊方法在于其统一代码表征。得益于IR的灵活性，一个模块可以包含其他模块，例如，“主机”翻译单元可以将“设备”翻译单元嵌入为IR，而不是文件引用或二进制数据。这种方法提供了其他编译器无法获得的主机/设备优化机会，特别是通过在主机和设备之间自由移动代码[22]。

```

// Kernel launch is available within the calling
// function, enabling optimizations across the
// GPU/CPU boundary.
func @launch(%h_out : memref<?xf32>,
             %h_in  : memref<?xf32>, %n : i64) {
  // Parallel for across all blocks in a grid.
  parallel.for (%gx, %gy, %gz) = (0, 0, 0)
    to (grid.x, grid.y, grid.z) {

    // Shared memory = stack allocation in a block.
    %shared_val = memref.alloca : memref<f32>

    // Parallel for across all threads in a block.
    parallel.for (%tx, %ty, %tz) = (0, 0, 0)
      to (blk.x, blk.y, blk.z) {
        // Control-flow is directly preserved.
        if %tx == 0 {
          %sum = func.call @sum(%d_in, %n)
          memref.store %sum, %shared_val[] : memref<f32>
        }
        // Synchronization via explicit operation.
        polygeist.barrier(%tx, %ty, %tz)
        %tid = %gx + grid.x * %tx
        if %tid < %n {
          %res = ...
          store %res, %d_out[%tid] : memref<?xf32>
        }
      }
    }
}

```

图 3. 在Polygeist/MLIR中，图1所示的CUDA launch/normalize代码的共享内存变量的表征（representation）。内核调用在主机代码中直接可用，该主机代码调用它。通过在块和线程之间使用并行for循环，显式地体现了并行性，而共享内存则放置在块内，以表明它可以被同一块中的任何线程访问，但不能被来自不同块的线程访问。

C. Polygeist

Polygeist是基于Clang的MLIR的C和C++前端[23]。它能够将广泛的C++程序翻译成多种MLIR方言的混合，这些方言保留了程序高层结构的元素。特别是，Polygeist将结构化控制流（循环和条件语句）作为MLIR SCF方言保留。它还通过依赖于MLIR的多维内存引用（memref）类型，在可能的情况下保留多维数组构造，从而简化分析。最后，Polygeist能够识别适合多面体优化（polyhedral）的程序部分[28]，并使用仿射（Affine）方言表示它们。

III. APPROACH

我们扩展了Polygeist编译器[23]，以直接从CUDA生成并行MLIR。这利用了统一的CPU/GPU表征（representation），使优化器能够理解主机/设备执行，

```

__global__ f() {
    codeA();
    barrier();
    codeB();
}

```

图 4. 在两个任意指令之间包含一个障碍的程序。

并在内核边界之间实现优化。现有MLIR的第一级并行构造（`scf.parallel`, `affine.parallel`）的使用使我们能够针对现有的CPU和GPU后端。最后，MLIR的可扩展操作集允许我们定义自定义指令，带有相关属性和自定义优化。

我们将GPU内核启动的表征（`representation`）定义如下（如图 3所示）：

- 在网格中的所有块上进行 3D 并行 for 循环。
- 对任何共享内存进行堆栈分配，范围限定为每个块唯一。
- 在一个块中的所有线程上进行3D并行for循环。
- 一个自定义的Polygeist（Polygeist）屏障操作，提供等同于CUDA/ROCm（CUDA/ROCm）同步的语义。

这个过程使我们能够以一种保留所需语义的形式表示任何GPU程序（GPU program）。它被编译器完全理解，因此适合进行编译器优化。此外，通过使用一般的并行性（`parallelism`）、分配（`allocation`）和同步（`synchronization`）构造来表示GPU程序，我们不仅能够优化原始程序，还能将其重新定向到不同的架构（`architecture`）。

A. Barrier Semantics

一个 CUDA 或 ROCm `__syncthreads` 函数保证了在函数调用之前，块中的所有线程都已完成所有指令的执行，并且在调用之后，没有任何线程会执行任何指令。传统上，编译器将此类函数表示为不透明的优化屏障，这些屏障可能访问所有内存，并禁止涉及它们的任何变换。

在我们的系统中，我们选择通过一个新的 `polygeist.barrier` 操作来表示这种线程级同步。与其他方法不同，`polygeist.barrier`（因此简单称为 `barrier`）旨在仅防止改变外部可见行为的转换。我们并不禁止在 `barrier` 之间的任何代码移动，而是通过定义 `barrier` 具有特定的内存属性来成功实现所需的语义，这些内存属性表示为一组内存位置（包括未知）和内存效应类型（读取、写入、分配、释放），这在 MLIR 中是标准的。考虑图 4 中的简单程序。只有当 `codeA` 和 `codeB` 访问相同的内存时，才可以观察到同步的影响。此外，如

果两者仅读取相同的内存位置，则同步也是不必要的。我们可以列举其余的三种情况：

- 1) `codeA` 写入，`codeB` 加载
- 2) `codeA` 加载，`codeB` 写入
- 3) `codeA` 写，`codeB` 写

具有`codeA`写行为的屏障可以确保案例 1 的正确性：在`codeB`中的加载操作不能被提升到屏障之上，因为它将似乎读取到一个不同的（状态）。对称地，具有（写入）行为的屏障将确保案例 2 的正确性。因此，（写入）和（写入）的写行为的联合足以防止加载操作被错误地移动到屏障之外。

然而，这并不阻止写操作被移动。例如，（写入）可以在案例 3 中被复制到屏障之上，并且它似乎会具有相同的最终内存（状态），因为在屏障之前的多余写操作将永远不会被读取。因此，我们还定义屏障具有（写入）和（写入）的读取行为。

这可以扩展以包括在给定屏障之前或之后可能执行的并行循环中的所有操作的内存效应。在具有显式分支的控制流图上，可以通过分别探索前驱或后继中的操作来实现。然而，在具有显式循环和条件操作的（循环）结构控制流级别上进行操作，允许分析被简化。此外，如果在同一块中存在多个屏障，则无需超越它。

考虑到一个足够表达的副作用模型，屏障的内存语义可以进一步扩展。屏障对来自不同线程的同一位置的读取/写入操作进行排序，自然执行顺序在一个线程内是足够的。因此，屏障不需要捕捉操作的内存效应，其中地址是线程标识符的单射函数。在可能的情况下，将内存访问提升到（线性）方言中的线性形式，可以实现精确分析。考虑图 5 中的代码。屏障周围的读取和写入表达式具有仿射访问集 $\mathcal{A}_o = \{A(i) : i = tx\}$ ，其中 tx 是线程标识符。屏障具有仿射访问集 $\mathcal{A}_b = \{A(i) : i \neq tx\}$ 。由于被访问的地址集没有重叠， $\mathcal{A}_o \cap \mathcal{A}_b = \emptyset$ ，因此允许在屏障之间移动代码。从概念上讲，对（写入）的写总是在同一线程内的读取之前发生，因此屏障是多余的。相反，如果对（写入）的加载和存储偏移为（写入），则屏障就是必要的，因为在屏障之后加载的数据将由不同的线程存储。更一般地，通过收集所有相关操作的内存位置和效应类型来定义给定屏障的内存属性时，只需包括由不同线程使用的内存位置。当涉及多个基地址时，必须检查别名保证。


```

__global__ f() {      // 0 <= t.x < blockDim.x
  A[threadIdx.x] = ...; // W A[i]: i == t.x
  barrier();          // RW A[i]: i != t.x
  ... = A[threadIdx.x]; // R A[i]: i == t.x
}

```

图 5. 障碍语义可以被细化为访问所有线程中在其之上/之下的操作所访问的内存地址，除了当前线程。

B. Barrier Lowering

为了使程序能够在上运行，我们必须高效地模拟程序的同步行为。尽管第节中的内存语义使我们能够在优化期间保持屏障的正确性，但本节讨论如何在实现屏障。

架构没有线程块的概念，也没有等待这种概念化线程组的屏障指令。相反，我们使用常规的线程和工作共享在它们之间分配线程块的循环迭代。从概念上讲，这与执行模型不同，后者是多个显式线程每个执行一次迭代。工作共享要求每个线程顺序执行多个迭代，使得在迭代中间进行同步变得不可能，只能在循环结束时同步。

为了应对这个问题，我们为我们的（多级中间表示）表征（）开发了一种新的屏障消除技术。如第节所述，过去已经探索了几种方法，包括循环分裂（）和继续传递（）。我们的方法是对前者的扩展，结合了两种风格的转换：并行循环分割和并行循环互换。

1) Parallel Loop Splitting: 假设一个障碍（）的核函数（）作为其直接父节点。通过将障碍周围的循环拆分成两个并行的循环，分别执行障碍之前和之后的代码，可以消除该障碍。如果障碍之前的代码创建了在其之后使用的价值（），这些必须在第二个并行循环中存储或重新计算。我们采用类似于中的技术来确定需要存储的最小数据量。具体来说，我们可以通过创建所有价值的图来表征（）这个问题。然后，我们将所有无法重新计算的价值定义（例如，从被覆盖的内存中加载）标记为源（），而在障碍之后使用的价值标记为汇（）。通过对该图执行最小分支割（），我们可以推导出需要存储的最小数据量。

```

parallel %i = 0 to 10 {
  %x = load data[%i]
  %y = load data[2 * %i]
  %a = fmul %x, %x
  %b = fmul %y, %y
  %c = fsub %x, y
  barrier
  call @use(%a, %b, %c)
  ...
}

%x_cache = memref<10xf32>
%y_cache = memref<10xf32>
parallel %i = 0 to 10 {
  %x = load data[%i]
  %y = load data[2 * %i]
  store %x, %x_cache[%i]
  store %y, %y_cache[%i]
}
parallel %i = 0 to 10 {
  %x = load %x_cache[%i]
  %y = load %y_cache[%i]
  %a = fmul %x, %y
  %b = fsub %y, %z
  call @use(%a, %b)
  ...
}

```

图 6. 示例：在屏障周围拆分的并行循环。这里，屏障上方的代码（包含2次加载和3次浮点操作）与屏障下方的代码（包含对@use的调用和其他操作）放置在一个单独的并行for循环中。此变换消除了屏障，同时保持了语义。由于值%a、%b和%c在屏障后被使用，因此必须格外小心以确保它们可用。在这里，最小割算法（min-cut algorithm）保留了%x和%y，然后重新计算%a、%b和%c，因为这样会导致保留2个值而不是3个值。

2) Parallel Loop Interchange: 不是所有的屏障操作都有一个并行作为其直接父节点，有些可能嵌套在其他控制流操作中。我们创建了一个模型，指定可以并行运行的指令。除了（屏障）以外，我们的表征不要求程序的任何特定顺序或并发性。因此，添加额外的屏障是合法的（尽管可能会降低并行性）。我们可以利用这个特性来实现控制流的屏障降低。

考虑一个控制流构造，它包含一个屏障并嵌套在一个并行中。在的周围立即添加屏障将导致在的上下立即进行并行循环拆分。因此，位于上方和下方的操作将被分离到各自的并行中，将成为中间循环中唯一的操作。然后，我们可以应用以下技术之一，将与并行进行交换，从而使后者成为（屏障）的直接父节点。

考虑包含屏障的串行循环的情况，见图。这种模式在代码中很常见，例如，实现线程之间的归约。由于必须等待所有线程，因此每个线程必须执行相同数量的。因此，内部循环的迭代次数对所有线程都是相同的，允许循环交换。

虽然语句可以被视为具有零次或一次迭代的循环，但在必要时，直接与周围的并行进行交

```

parallel for %id=0 to N {
  for %j = 5 to 0 {
    if (%id < 2^%j)
      A[%id] += \
        A[%id + 2^%j]
    barrier
  }
}

```

```

for %j = 5 to 0 {
  parallel for %id=0 to N {
    if (%id < 2^%j)
      A[%id] += A[%id + 2^%j]
    barrier
  }
}

```

图 7. **左:** 一个共享内存加法, 包含一个内核调用, 其中有一个带屏障 (barrier) 的 for 循环。 **右:** 通过将并行 for 循环与串行 for 循环进行交换, 现在在并行循环中直接包含屏障 (barrier) 的相同代码。

```

parallel for %i=0 to N {
  do {
    run(%i)
    barrier
  } while(condition())
}

```

```

%helper = alloca memref<i1>
scf.do {
  parallel for %i=0 to N {
    run(%i)
    barrier
    %c = condition()
    if %i == 0 {
      store %c, %helper[]
    }
  }
  %c = load %helper[]
} while(%c)

```

图 8. 在 while 循环周围的并行交换。由于 condition() 函数调用必须在每个线程上执行以保持正确性, 因此使用一个辅助变量来保存第一个线程上调用的价值 (value)。

换更有效。

在 中, 循环必须在循环之前计算其迭代次数。循环支持动态退出条件。例如, 考虑图 中的代码。为了保持正确性, 我们必须在每个线程中执行

C. Usage

转译 () 的设计旨在通过允许作为现有 编译器 (如) 的替代品, 便于使用。具体而言, 扩展了 前端, 因此使用与 相同的标志和语法。然而, 还引入了一些额外的标志, 例如 用于指定 到 的转换, 以及 用于指定用于生成 程序的特定方法和并行优化集 (参见第 节)。

IV. PARALLEL OPTIMIZATION

提供 的并行性和 程序的高级表征 () 使各种优化成为可能。这些优化包括适用于任何并行程序的一般优化, 以及在 转换为 的背景下的特定优化。

A. Barrier Elimination & Motion

由于 样式的屏障必须特别转换以支持 架构, 因此消除或简化任何屏障可能会产生显著的影响。此外, 即使在 上运行 代码时, 屏障的消除也是非常有用的, 因为任何同步都会减少并行性。大部分关于屏障消除和简化的基础设施直接来源于第 节中定义的其内存行为。给定一个屏障, 设 M_{before}^\dagger 为在 之前, 直到另一个屏障或并行区域开始的内存效果的并集, 设在 之后, 直到并行区域结束的内存效果的并集为 M_{after} 。如果在屏障跨越的同一位置没有内存效果 除了读取后读取 () 之外 (即 $(M_{before}^\dagger \cap M_{after}) \setminus \text{RAR} = \emptyset$), 则这个屏障的行为被先前的屏障所涵盖, 可以被消除。对称条件 $M_{before} \cap M_{after}^\dagger \setminus \text{RAR} = \emptyset$ 表示该屏障被后续的屏障所涵盖。可消除屏障的一个特定简单情况是完全没有内存效果的屏障。

```

__global__ void bpnnp_layerforward(...) {
    __shared__ float node[HEIGHT];
    __shared__ float weights[HEIGHT][WIDTH];
    if ( tx == 0 )
        node[ty] = input[index_in] ;
    // Unnecessary Barrier #1
    __syncthreads();
    // Unnecessary Store #1
    weights[ty][tx] = hidden[index];
    __syncthreads();

    // Unnecessary Load #1
    weights[ty][tx] = weights[ty][tx] * node[ty];
    __syncthreads();

    for ( int i = 1 ; i <= log2(HEIGHT) ; i++){
        if( ty % pow(2, i) == 0 )
            weights[ty][tx] += weights[ty+pow(2, i-1)][tx];
        __syncthreads();
    }

    hidden[index] = weights[ty][tx];
    // Unnecessary Barrier #2
    __syncthreads();

    if ( tx == 0 )
        output[by * hid + ty] = weights[tx][ty];
}

```

图 9. 一个来自Rodinia反向传播测试的CUDA内核示例，其中包含不必要的同步和不必要的共享内存使用，而寄存器（register）就足够了。

例如，考虑来自基准测试中的代码。第一个和最后一个指令是不必要的。可以通过上面的基于内存的屏障消除算法证明这一点。对于第一个屏障， M_{before} （一直到开始）仅包含对 $node[ty]$ 的写入和对 $weights[ty][tx]$ 的读取。 M_{after}^{\dagger} （直到第二个屏障）包含对 $weights[ty][tx]$ 的写入和对 $weights[ty+pow(2, i-1)][tx]$ 的读取。如果在调用语境中已知指针不会发生别名，则这些之间没有冲突。因此，消除这个屏障是安全的。

同样的内存分析也可以用于执行屏障移动。只需在拟议移动的位置放置一个虚构的屏障，并检查先前的内存分析是否会推导出当前的屏障是多余的，从而允许屏障移动。

B. Memory-to-register promotion across barriers

从其内存行为定义语义的目标之一是使内存优化能够在包含障碍的代码中正确有效地进行。如在第一节中所述，障碍有上下文代码的内存行为，显著的例外是来自当前线程的访问。这一缺口很重

要，因为它使得内存到寄存器的提升能够在线程本地内存（例如局部变量）上进行。此外，这种优化能够成功地用快速寄存器替换慢速内存读取。例如，再次考虑图中的代码。考虑标记为“无用存储”和“无用加载”的对 $node[ty]$ 的加载和存储，以及两者之间的同步。在那个时刻，唯一可以加载的值是之前由同一读取存储的值。此外，由于在任何其他人可以从 $node[ty]$ 读取之前，该位置已被覆盖，因此只要加载使用包含从 $node[ty]$ 加载的值的寄存器，第一次存储就可以安全地消除。在内存到寄存器优化期间，现在能够成功推导出这一转发属性，因为第一节中描述的内存属性的缺口使其能够推断出障碍操作不会覆盖当前线程的存储。因此，传统的加载和存储转发能够正确地在障碍代码上操作。

C. Parallel loop-invariant code motion

传统的循环不变代码移动优化旨在将指令移动到串行循环之外，从而减少串行的执行次数。如果可能访问内存或具有其他副作用，除了检查的操作数本身是循环不变的之外，编译器还必须检查循环内没有其他代码与执行的内存访问冲突。

在目前的编译器上，虽然可以将循环不变代码移动应用于核心内的串行循环，但无法将其应用于推动指令（到核心调用外。这部分归因于核心与调用它们的代码保存在不同模块中的事实，以及传统编译器对并行性的理解不足（见图）。

反直觉的是，凭借正确的语义，即使我们无法将循环不变代码移动应用于等效的串行循环，我们也可以将其应用于并行循环。我们将依赖于程序的语义，它允许我们在保持屏障所需顺序的前提下，任意交错并行循环的迭代。因此，可以合法地想象在锁步运行程序。也就是说，如果一个并行循环有 n 条指令，则每个线程会在任何线程执行指令之前先执行指令，依此类推。因此，只要其操作数是不变的，并且并行循环中的没有先前指令与之冲突，就可以合法地推动一条指令。换句话说，不需要检查是否与并行循环中的任何后续指令存在冲突即可允许推动。

```

omp.parallel {
  omp.wsloop %i= 1 to 10 {
    codeA(%i)
  }
}
omp.parallel {
  omp.wsloop %i= 1 to 10 {
    codeA(%i)
  }
}

```

```

omp.parallel {
  omp.wsloop %i=1 to 10 {
    codeA(%i)
  }
  omp.barrier
  omp.wsloop %i=1 to 10 {
    codeA(%i)
  }
}

```

图 10. 在 MLIR 上应用的 OpenMP 并行区域融合示例。给定两个相邻的 OpenMP 并行区域，每个区域都用给定的闭包（closure）初始化线程，沿着线程屏障（thread barrier）融合这两个闭包，从而允许线程只初始化一次，而不是两次。

```

for (i=0; i<N; i++) {
  #pragma omp parallel for
  for (j=0; j<10; j++) {
    body(i, j);
  }
}

```

```

#pragma omp parallel
for (i=0; i<N; i++) {
  #pragma omp for
  for (j=0; j<10; j++) {
    body(i, j);
  }
  #pragma omp barrier
}

```

图 11. OpenMP 并行区域提升的示例，如在 C/C++ 中应用。这是图 10 中 OpenMP 并行区域合并的扩展，不同之处在于它应用于整个 for 循环。与其为外层循环的每次迭代 i 创建一个单独的闭包（closure）/线程初始化，不如只创建一次闭包/线程初始化。

D. Block Parallelism Optimizations

开放是我们在上进行并行执行的主要目标。它通过两个构造实现并行“”循环。首先，循环被概述为一个函数，该函数每个线程调用一次，代表了开放的“并行”（）构造。然后，在概述的函数内，迭代空间在线程之间分配，代表了开放的“工作共享循环”（）构造。开放还有一个“屏障”（）构造，但其语义与屏障的语义不同。

当多个并行循环连续执行时，例如，在第节中的屏障降低之后，通过融合相邻的开放“并行”构造而不融合工作共享循环（见图），可以减少线程管理的开销，从而不撤销屏障降低。并行区域融合可以扩展到各种构造，例如将开放并行区域移至图中包围的“”之外。这将使线程初始化只调用一次，而不是 N 次。将其一般应用于控制流构造使得通过对一个块执行并行循环裂变生成的所有并行循环都可以进行开放并行（但不是工作共享循环）的融合。

由于程序往往是在考虑高并行性的情况下编写的，不同块提供的并行性可能已经单独饱和了可用核心的数量。如果没有使用共享内存，块和线程并行性可以折叠为一个单一的开放并行，这将在单个并行区域中均匀划分总的迭代空间。然而，如果存在共享内存，我们的工具将生成嵌套并行区域以表示共享内存的分配。在这种情况下，嵌套的开放并行区域带来的额外开销可能会超过潜在的附加并行性。此外，内循环的并行化可能会导致不利的内存效应（如假共享）进一步影响性能。

因此，我们还支持一种优化，用于序列化任何嵌套的开放并行区域。进行这样的序列化可能会利用内存局部性来提高性能。

V. MOCUDA: INTEGRATION INTO PYTORCH

A. Targeting CPU-only Supercomputers

我们工作的一个目标是对中的执行模型进行建模，以便将为原本编写的高性能代码在仅有超级计算机上执行，尤其是在配备处理器的机器上。作为一个主要例子，我们考虑，它尚未成功移植到架构。使用其“”（）后端支持执行，但对于许多内核只提供了一个（）且表现不佳的实现，例如，对于二维卷积没有内存优化的层嵌套循环。在英特尔上，计算密集型内核的执行被委托给库，而在上表现不佳，因为它是为普通定制的，而上没有可用的高带宽内存。富士通的分支（我们在第节中也进行了比较）需要手动调优以提高性能，而且仍然无法在普遍情况下与竞争。有趣的是，具备高带宽内存的可以受益于与类似的计算组织。因此，我们设计了“”来模拟中的后端，通过将内核转化为支持并行的内核，使用前面章节中描述的机制。

B. Architecture of MocCUDA and Integration of Polygeist

我们实现了一个兼容层，以便在上透明地执行的，后端，目的是避免手动重构和与现有库的互操作性。虽然确实有其自身的后端，但其设计与在上运行的超大规模并行计算机的兼容

性较差。尤其是，性能问题是由同步内核执行、不匹配的内存布局和假设低速内存引起的。

我们并没有尝试对 的 后端进行艰巨而耗时的重新设计，而是决定基于 的 后端设计来构建 ，因为该平台与 有相对的相似性。因此， 必须处理并替换四种类型的操作：

- 1) 在 中调用深度学习特定函数，
- 2) 对辅助 库的调用，并且
- 3) 在 中执行 内核 ()。

只有后者可以被编译，因为其余的功能作为二进制库分发。

为了最小化原型的工作范围，我们分析了 与 后端的所有交互，以确定 () () 对于一个众所周知的深度学习问题，即 神经网络。

从分析中可以看出， 与 的交互主要限于识别已安装 的属性、内存管理以及 流的管理和同步。对于原型，我们将自己限制为每个 节点模拟一个 ，并只管理显式内存分配和数据传输。为了让 访问 设备属性，我们将真实 () 的数据转储到一个文件中，后者稍后由没有 的系统上的 使用。我们通过苹果的

模拟 流，使得模拟 操作与 的管理层之间的异步性得以实现。 模型仅使用一部分闭源 函数（卷积层的前向和后向传递、批量归一化层以及张量加法）。我们重新实现了所有必要的 函数变体，重点是 并行化和 友好的 加上 卷积，以及对 布局（批次、通道、高、宽）的支持。此外，我们识别了对 库的调用，并实现了包装函数来拦截这些调用，并将它们分发到为 设计的 库，如 、 或富士通的 对于 。对于其他链接到的 的库（如 ），采用类似的方法

是可能的，但对 来说并不是必要的。

除了 和其他库调用外，我们观察到自定义内核，即在 中用 实现的 内核，而不是由（二进制）库提供的内核，并识别出调用它们的高阶函数。 所需的函数包括跨步张量内核（加法、乘法等）、聚合操作如 “ ” 等。负对数似然损失内核使用 的 屏障。将这些基于 的 函数移植到 需要大量的改造工作，而 则执行自动翻译。为了演示目的，我们使用 自动将使用的 的 数及其递归调用的函数从 翻译到 ，并将结果集成到 中。

我们将上述所有的包装器和重新实现合并到 中，该工具可与 一起使用，以拦截 调用，用于在 上训练带有 后端的 。

VI. EVALUATION

我们展示了我们方法在两个著名的 基准测试套件上的优势和适用性： 基准测试套件的一个子集 和 神经网络的 实现。这些基准测试的选择是为了 提供我们 在具有手工编码 版本的基准测试套件 () 上的与 编译的粗略性能比较，以及 展示我们的系统成功地集成到一个有用且实际的应用 () 中，在不具有任何 的超级计算机 上。此外，我们还将我们方法的性能与现有的 工具在 矩阵乘法上的性能进行了比较。

对于 ，我们将翻译后的 与 代码与 版本的基准进行比较（在存在的情况下），以及在 上的运行。对于 ，我们与“原生” () 和 后端进行比较。

¹ 是在提交

¹ 相关版本的MocCUDA和Polygeist可以在 <https://anonymous-data.s3.amazonaws.com/MocCUDA-master.zip> 和 <https://anonymous-data.s3.amazonaws.com/Polygeist-main.zip> 获得。

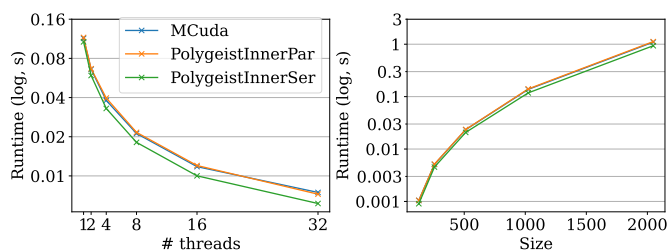


图 12. PolygeistInnerPar 的性能与 MCUDA 相似,而 PolygeistInnerSer 超过了 MCUDA。PolygeistInnerSer 类似于 MCUDA 禁用内循环并行化,而 PolygeistInnerPar 保持了块和线程的并行性。左:平均运行时间随线程数的变化 (对矩阵大小进行平均)。右:平均运行时间随矩阵大小的变化 (对线程数进行平均)。

下针对 编译的。对于 , 我们使用 的 为 , , 以 及 的 库进 行 的 编译。对于基 准 测量,我们使用 预安装 的 ()。 我们在运行 的 实例(双插槽 , 频率为 , 每个插槽有 个核心和 内存)上评 估 和矩阵乘法测试。测量是在第一个插槽 上进行的,同时禁用了超线程和涡轮提升。每个数字是至少 次重复的中位数。

A. Comparison to MCUDA

首先,我们将我们的方法与 中 的先前工作进行比较。 是一个 级 别的工具,它生成新的 作为输 出,并使用类似的循环分裂技术来处理同步。作为一种源到源的工具,它仅处理输入语言的一部分,这使得它无法处理 程序。相反,我们在图 中比较了在不同线程数 () 和 矩阵大小 () 下的矩阵乘法内核的运行 时间。 在不包含内循环序列化 () 的所有 优化下 平均产生的代码与 相差。具体而 言 在 个 线 程下有的减速 在 个线程下有的加速。这种行为是由 在处理嵌套并行构造时的开销

²尽管我们将在没有GPU的系统上运行PyTorch,但我们必须在支持CUDA的系统上编译PyTorch,以确保发出正确的代码。我们还防止了三个Pytorch函数的内联。

引起的。实际上 仅对最外层循环进行 并行化。当 对内循环进行序列化 () 时 它在整体 上实现了相较于 的加速 其中在 个线程下加速 在 个线程下加速。

B. Use case 1: Rodinia Benchmarks

我们基准测试了一 由 个基准 组成,这些基准目前被 支持,并且具 有非平凡的运行时间。 我们通过比较用 编译 并在 上执行的程序输出与通过我们的流程编译并 在 上执行的程序输出来验证我们流程的正确性。我 们在内核和 或包含内核的代码的计算部分中插入了时间 测量,在某些情况下,每个基准多次测量。可能的情况下, 我们计时相同基准的 版本的等效部分。 我们在图 右中比较了编译 为 的 基 准 与 基准版本。尽管基 准之间存在一些变化,但总体而言,我们的方法与手动 编码的基准版本相当,甚至在启用内部串行化优化时净 获得了 的几何平均性能提升。在没有内部串行化的 情况下,我们仍然看到 的几何平均加速。一 些基 准如 和 采 用了用于模板计算的优化技术,通过在线程 之间复制计算以减少同步开销,更好地利用 中可用的并行性。这使得 代码 比 版本复杂得多,导致其性能较差。 由于程序在共享内存中执行额外的工作进行数据缓 存 和 的 版 本速度较慢。 的 本获得了相对加速,因为纯 版的代 码通过单独的 “ ” 循环获得期望的依赖结构,而在 代 码中则用 来 完成。 能够成功优化围绕屏障的代码,从而 实现加速。 的加速部分归因于并行优化

³hybridsort、 kmeans、 leukocyte、 mummergpu、 huffman和heartwall在Polygeist中使用了不被支持的C++或CUDA特 性(虚函数和纹理内存)。lavaMD和dwt2d基准使用了格式不正确 的C++,由于从未初始化内存中读取而导致未定义行为(GPU驱动程序 将共享内存零初始化,但没有这个要求)。 nn和gaussian测试 在≤ 0.005秒内完成。

⁴由于nn的运行时间很短,已经被排除,这两个版本在数据加载方面 有所不同(在运行时动态加载,在OpenMP代码中,以及在CUDA代码中 预加载所有数据),因此也应该排除。

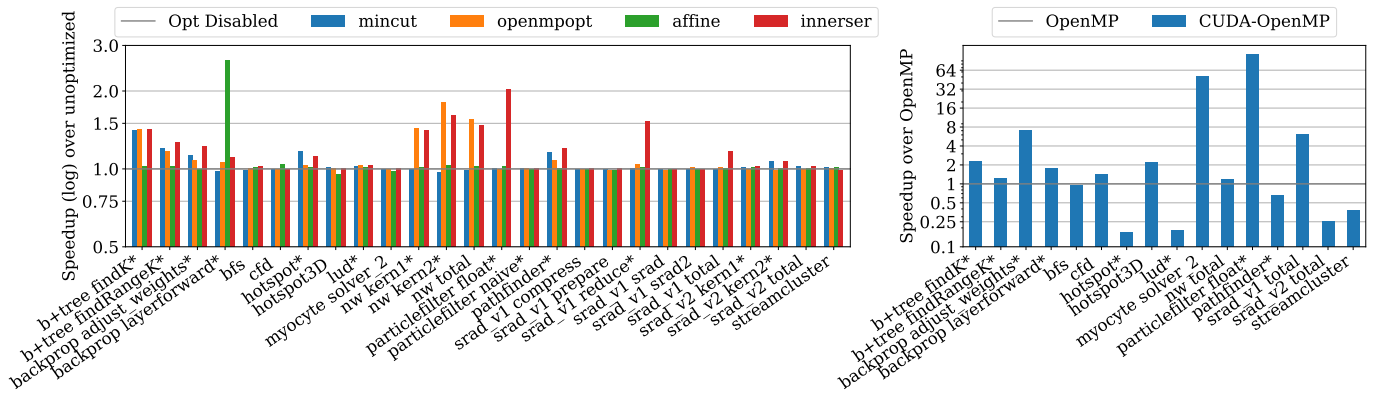


图 13. 左侧：应用各种并行和/或CPU优化的内核的相对加速（越高越好）。右侧：Rodinia CUDA代码编译为OpenMP与原生Rodinia OpenMP代码的加速对比（在可用的情况下）。包含屏障的基准测试用星号标记。

（见图 左），部分原因是 代码采用线性数组实现，这是 的要求，而不是 代码中使用的双指针。 和 均因在并行区域边界内的代码优化以及内部串行化而获得加速。

我们通过在图 中比较转换后的 与原生 内核来测试我们方法的缩放特性。转换后的 代码通常比原生 版本扩展得更好。由于大多数 程序在设计时考虑了成千上万的线程，这表明我们的框架能够在将 特定构造重写为 兼容的等效项时保留这种并行性。在 一个线程且没有内部串行化的情况下，所有测试的转换后 代码具有的几何平均加速。由于并非所有测试都存在 版本的基准，如果我们仅考虑存在 版本的 代码，我们发现其几何平均加速为，而 的加速仅为。对内部循环进行串行化略微降低了可扩展性，但仍然在所有启用内部串行化的测试中获得的几何平均加速，并在具有 版本的代码中获得的加速。

我们进行了一次消融分析，以研究单个优化对性能的影响。图 左 中的 系列显示了我们的方法在第 节中概述的优化下的性能测量，以减少在屏障之间保留的数据量。这仅与包含屏障的基准测试相关（在图中以星号标记）。在适用的基准测试中， 提供了 的几何平均加速比（ ）。图 左 中的 系列展示了 区域合并及类似优化的影响，并带来了 的几何平均加速比。图 左 中的 系列显示了将控制流提升

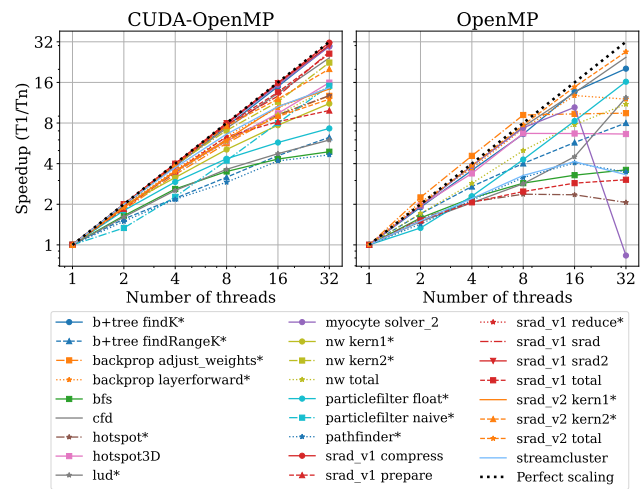


图 14. CUDA Rodinia 内核的可缩放性行为，当在使用 32 个线程的 CPU 上运行并利用 OpenMP 时，以及 OpenMP Rodinia 内核（如有）之间的比较。并非所有的 Rodinia CUDA 内核都有 OpenMP 版本。

到其仿射（ ）变体并启用简单循环优化（例如循环展开）的结果。尽管这在整体上产生了 的几何平均加速比，但对于反向传播层前向测试（

），它产生了的加速，因为它将一个包含同步的循环完全展开并简化为 语句。

C. Use case 2: Pytorch/Resnet50 Test

为了评估 ，我们在 超算的一个 单元上执行了完整的节点并行训练，与原生 后端和优化的 后端进行比较。

我们以数据并行的方式在 的正向和反向传播。我

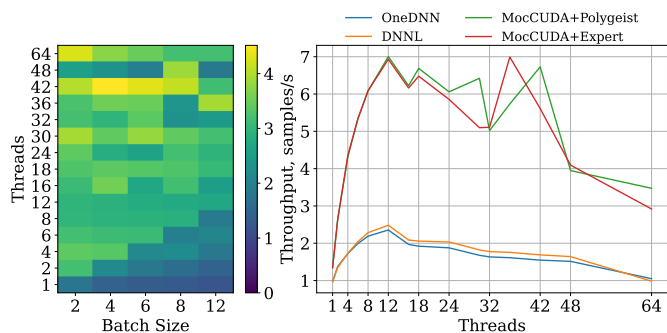


图 15. 在Fugaku节点上对ResNet50进行训练。左:“MocCUDA+Polygeist”相对于富士通-调优版oneDNN的相对吞吐量增加的热图，值越高越好。右：不同批次大小的几何平均吞吐量；“MocCUDA+Expert”使用专家编写的OpenMP内核；“MocCUDA+Polygeist”使用生成的内核。

们使用 的合成基准测试脚本（为 神经网络模型配置）。我们构建了带有 、 和富士通 库的 ，以支持基于 的多节点、分布式深度学习。我们为每个 核心内存组（ ）分配一个 排名，模拟每个节点最多 个 ，并将测试范围从一个节点（ 个排名）扩展到一个节点（ 个排名）在一个 单元中（最小为 环形拓扑），同时将 线程数固定为 ，以容纳每个核心一个线程。我们的方法使用 ，而其他后端依赖于 。性能测量使用 进行，该工具通过 组织设置神经网络，创建图像，使用 执行层，并返回图像每秒的吞吐量指标。我们在批量大小 线程数 的情况下运行，跨轮次进行平均。

我们观察到 在批量大小和线程数量方面系统性地优于 调优的 最高可实现 的吞吐量提升（几何均值 最小值 ）如图 所示。使用我们专家编写的内核的 与使用 由 生成的内核的 相当，如 节所述。

这种改进可以通过第 节中描述的设计和 性能特征相结合来解释。由于 的 未考虑 上可用的 ，因此它使用对缓存友好的直接卷积，而不是基于 的卷积，在存在高带宽内存的 上效率较低。

尽管 调优的 的定制分支提高了 的性能（尽管几何均值仅提高），但仍留有改进性能的空间。

这表明，我们的方法能够从 版本自动推导出深度学习内核（以及潜在的其他应用程序）的高效版本，从而解决缺失或低效内核在具有高带宽内存的 上的局限，而不需要反向工程或重新工程应用程序。

VII. RELATED WORK

A. GPU to CPU Synchronization

直接提供的第一个用于在 上模拟 的工具之一是为了调试目的而创建的，它在 上模拟每个线程，并配备一个独立的 线程。尽管功能上可行，但可用线程数的巨大差距使得仿真效率低下。

（ ）对 代码执行 转换，以生成新的 代码，该代码调用一个与线程无关的并行 例程。 开创了使用“深裂变”（ ）处理同步的使用方法，该方法在同步点拆分并行循环和其他构造，目的是消除它们。这种裂变技术也应用于其他工具： （ ），一个将 汇编语言解析为 并即时编译内核函数的二进制翻译工具； （ ），一个用于 的 编译器传递； （ ），另一个采用裂变的 翻译 转换传递，并显著处理 级原语；甚至还有本研究。虽然裂变方法背后的直觉与此处使用的方法类似，但我们是在一个高级编译器内部应用裂变而不是在源代码或低级 中。如在第 节所示，对结构化程序执行裂变使代码变换更加高效。在源级别执行裂变会错失在裂变之前运行优化（如屏障消除）的机会，而在低级别执行裂变又需尝试重建高级结构，在 中工作使我们能够同时应用优化并保留高级结构。此外，源级工具往往相当脆弱，因为它们必须重新实现解析和目标语言（如 ）的语义，因此仅在输入语言的有限子集上操作，导致需要工程重构以替代不受支持的构造（如指针算术）。

另一种方法使用继续传递 （ ）来处理同步，通过创建所有同步点的状态机（例如“微线程”（ ））。和

() 提出一种在 中采用继续传递的方法, 包括检测和减少控制流图中分歧的算法, 后续工作旨在最小化活跃值以减少内存流量。

() 类似于 已知 并行区域扩展的使用是有益于 的原始虚拟 , 但现在 的 使用 的 并通 过 执 并行区域 行同步。共享内存存在 中作为单个全局变量实现, 并按块数进行扩展。

在低级 上工作的先前研究付出了巨大努力, 以重建高层构造, 例如循环和内核配置, 这些在高效裂变或继续传递中是必需的。例如,

进行各种规范化和循环变换, 以重写控制流图并尝试将其识别为可以处理的几种形式之一。在源 级别工作的先前研究(例如), 除了仍需识别 级概念外, 无法受益于简化代码的优化, 从而导致更易于控制的流。

相比之下, 通过在 的混合抽象上操作, 我们能够同时保留源级结构并执行程序变换, 例如循环展开或 运动, 可以, 如, 删除嵌套同步。

B. Parallel Portability/IR, & OpenMP Optimizations

几个工具在主机语言中定义了新的抽象, 这些抽象自然适合于 或 的执行。示例包括 中的 、 、 或 (仅限于 计算), 以及 中的 , 以及 中的

。这些方法为使用它们编写的任何新代码提供了性能可移植性。然而, 任何现有代码都必须在所述框架中重写(并且可能与其他框架或其他语言编写的代码不兼容)。

几项先前的工作讨论了并行性的中间表征(), 例如, 用于在 中表征 ; 用于在 中表征 , 用于模式树, 以及 方 言; 以及 用于以控制流图的形式直观地表征并发。这些工作主要集中在特定风格并行性的表征() (例如 中的 任务), 而不包括 风格的屏障, 而不是并行特定的转换

() (如屏障消除) 或 优化, 除了一致性 竞争检查或自动并行化

并行区域扩展的使用是有益于 的。 可选择以较弱的形式支持该转换, 即在相同控制级别内合并 并行区域。

C. Barriers

几项先前的研究探讨了屏障或同步指令的语义, 包括与 的相关性。已有研究验证了屏障的正确性。 探索并实验评估了不同 供应商的前进进度 公平性模型。 实现了一种在工作组之间应用的 屏障操作, 而不仅仅是在一个工作组内。 将 内存屏障添加到程序中, 以确保弱一致性和顺序一致性语义。他们发现, 在没有同步和延迟设置分析的情况下, 引入一致性语义的平均减速为 , 而使用这些分析插入更少的同步操作时, 弱一致性和顺序一致性分别可以实现 和 的减速。

VIII. CONCLUSION

通过扩展 , 我们开发了一个端到端系统, 能够表征()、 优化和转译 和 并行程序。能够同时表征和在不同并行框架之间转换是至关重要的, 因为高性能计算() 越来越依赖(异构)() 并行性来处理其工作负载。我们框架的重要组成部分是高阶屏障操作的开发, 这对于表征 程序至关重要, 其语义可以通过其内存行为完全定义。与以前的并行屏障表征不同, 我们的语义使得在优化过程中能够直接集成屏障。为了验证我们方法的有效性, 我们展示了在商品 上对 基准套件的一个子集进行 到 的优化和转译, 并将一个高效的 从 源 转译为可在仅有 的超级计算机 上运行的代码。 尽管由于 和 之间实现的差异 性能在每个案例中有所不同 但 基准套件在转译的 代码与手写的 版本 之间实现了 的几何平均加速。同样 内核在原生 后端 之上也观察到了的加速。

目前 在 上运行的转译 代码保持相同的调度 除了最内层循环的序列化 这是为了提高性能。未

□ □ □ □ □ □ □ □ □ □ □ □ □ □ □ □ □ □ □

□ □ □ □ □

□□□□□□□□ □□ □□□□□□□□
 □□□□□□□□□□ □□□□□□□□□□
 □□□□□□□□

[illegible]

□□□□□□□□ □□□□□□□□

□□□□ □□□□□ □□□□□□□□ □□
□□□□□□□□□□ □□□□□□□□□□
□□□□□ □□□□□□□□□□

□□□□□□□□

□□ □□□□□□□□ □□□ □□□□□□□□□□□□

□□□□□□□□□□

□□□□□
 □□□□□□□□□□□□□□ □□□□□□□□□□
 □□□ □□□□ □□□□□□□□□□□□
 □□□□□□□□□□□ □□□□□□□□□□□□□□

□ □ □ □ □ □ □ □ □ □ □ □ □ □ □ □ □ □ □

□□□□ □□□□ □□□□□□□□□□□□
 □□□□□□□□□□ □□ □□□□□□□□
 □□□□□□□□□□□□□□ □□□□
 □□□□□□□□□□□□ □□□□□□□□□□
 □□□□□□□□□□□□□□□□□□□□



□□□□□□□□□□ □□ □□□ □□□□
 □□□ □□□□□□□ □□□□□□□□□□
 □□ □□□□□□□□□□ □□□□□□□□
 □□□□□□ □□□ □□□□□□□□□□□□□□□□□

□□□□□□□□□□ □□ □□□ □□□□
 □□□ □□□□□□□ □□□□□□□□□ □□
 □□□□□□□□□□ □□□ □□□□□□□□
 □□ □□□□□□□□ □□□□□□□□□□□
 □□□□□□□

□□□□□□□□ □□□ □□□□□□□□□□ □□□□
□□□□□□□□ □□□□□□□□□□

000000000000 00 000
 000000 000000000 00 000 0000
 000000000 0000000000000000 00
 000

0000000000 000 0000000000 000
 0000000000 0000000000 0 0000
 000000000000000 0000000000
 0000 00000 0000 0000 00000
 000 0000 00000000 00000
 00000 00000000 000000000
 0000000

□□□ □□□□□□□□□□□□ □□□
□□□□□ □□□□□□□□ □□□ □□□
□□□ □□ □

□□□□□□□□□□ □□
□□□□□□□ □□□□□□□□

..

□□□□□
□□□□□□□ □□□□□□□□□□□□□□

□□□□□□□□□□ □□ □□□
□□□□□□□□□□ □□□□□□□□
□□□ □□□□ □□□□□□□□□□
□□□□□□□□ □□□□□□□□□□
□□□□□□ □□□ □□□□□□□□

e1f22da320b0c7384cbd2f4df50255c7c2fa6b9d

.. ..

□□ □□□
□□□□□□
□□□□□□□ □□□□□□□□

□□□□□□□□□□ □□□□□□ □□
□□□□□□ □□□□□□□□□□

□□□□□□□□□□ □□□□□□□
□□ □□□□□

□□□□□□□□□□ □□ □□□
□□□ □□□□□ □□□□□□□
□□□□□□□□□□ □□□□□□□□
□□ □□□ □□□□□□□□ □□□
□□□□□□□□□□

□□□□□□□□□□ □□ □□□ □□□□
□□□□□□□□□□ □□□□□□□□ □□
□□□□□□□ □□□□□□□□□□□□

□□□□□□□□□□ □□□□□□□□□
□□ □□□□□□ □□□□□□□□□□

□□□□□□ □□□□□□ □□□□□□

□□□□□□□□□□ □□ □□□ □□□□
□□□□□□□□□□ □□□□□□□□ □□
□□□□□□ □□□□□□□□□□□□ □□
□□□□ □□□□□□□□ □□□□□□
□□□□ □□□□ □□□□

• •

• •

^

000000
0000000000000000 000000000000 00
000000000000 00 000000000000
00 000000 000000 000000000000

□□□□□□ □□□□□□ □□□

□□□□□□□□□□ □□ □□□ □□□□□
□□□ □□□□□□ □□□□□□□□ □□
□□□□□□□□□ □□□ □□□□□□□ □□
□□□□□□□ □□□□□□□□□□

□□□□□□□□□□ □□ □□□ □□□□
□□□ □□□□□□ □□□□□□□□□ □□
□□□□□□□□□□ □□□□□□ □□□□□□
□□□ □□□□□□□□□□□□□□□□

□□□□□□□□□□ □□ □□□ □□□□ □□□
□□□□□□□□□□□□□ □□□□□□□□ □□
□□□□□□□□□ □□ □□□□□□□□□□
□□□□□□□□□□

-

□□□□□□□□□□ □□
□□□ □□□ □□ □□□□□□□□□□
□□□□□□□□□□

,

□□□□□□□□□□
□□ □□□ □□□□ □□□ □□□□□□
□□□□□□□□□□□ □□□□□□□□□
□□ □□□□□□□□□□□□□□□□
□□□□□□□□□□ □□□□□□□□
□□□□□□□□□ □□□□
□□□□□□□□□□ □□□□
□□□□□□□□□□