

Low-Level CUDA Optimization with PTX: A Deep Learning-Oriented Tutorial

Advanced GPU Programming Guide

March 24, 2025

Abstract

This tutorial provides a comprehensive introduction to Parallel Thread Execution (PTX), NVIDIA's low-level assembly-like language for GPU programming. Designed for developers familiar with CUDA C/C++ who want to achieve peak performance in deep learning workloads, this guide walks through the fundamentals of PTX, how to write and optimize PTX code, and how to leverage GPU hardware features like Tensor Cores. We'll explore practical examples, case studies, and performance optimization techniques specific to deep learning applications. By understanding PTX, developers can gain deeper insights into GPU execution and unlock performance optimizations beyond what high-level languages provide.

Table of Contents

1	Introduction to PTX and the CUDA Toolchain	5
1.1	What is PTX?	5
1.2	The Role of PTX in the CUDA Toolchain	5
1.3	Why Learn PTX for Deep Learning?	6
2	Inspecting PTX from CUDA C/C++ Code	6
2.1	Generating PTX from CUDA Source Code	6
2.2	Extracting PTX from Compiled Code	7
2.3	Example: Vector Addition Kernel	7
3	Writing and Compiling PTX Directly	9
3.1	Approaches to Writing PTX	9
3.1.1	Inline PTX Assembly in CUDA C/C++	9
3.1.2	Writing Complete PTX Kernel Files	10
3.1.3	Using NVRTC (Runtime Compilation)	10
3.2	Compiling PTX Code	11
3.2.1	Using ptxas	11
3.3	Running PTX Kernels	11
3.3.1	Using the CUDA Driver API	11
3.3.2	Embedding in Host Code	12
4	PTX Syntax and Key Constructs	12
4.1	Module and Entry Point Directives	13
4.1.1	Version and Target Directives	13
4.1.2	Kernel Entry Points	13

4.2	Registers and Type Qualifiers	14
4.2.1	Register Declaration	14
4.2.2	Register Types	14
4.2.3	Register Usage Considerations	14
4.3	Thread Indexes and Special Registers	15
4.3.1	Common Thread Index Registers	15
4.3.2	Using Thread Index Registers	15
4.3.3	Additional Special Registers	15
4.4	Arithmetic and Control Flow Instructions	16
4.4.1	Basic Arithmetic Instructions	16
4.4.2	Vector Operations	16
4.4.3	Control Flow: Predicates and Branches	16
4.4.4	Predicated Execution	17
4.4.5	Implementing Loops	17
4.5	Memory Spaces and Instructions	17
4.5.1	Memory Space Overview	17
4.5.2	Memory Access Instructions	18
4.5.3	Declaring Static Memory	18
4.5.4	Memory Performance Considerations	18
4.5.5	Memory Caching and Control	19
4.5.6	Synchronization	19
5	Hands-On Example: PTX vs. CUDA C++ for a Matrix Multiply Kernel	20
5.1	The Matrix Multiplication Problem	20
5.2	Naive CUDA C++ Implementation	20
5.3	Optimized Matrix Multiply Using Shared Memory	21
5.4	PTX Implementation of Tiled Matrix Multiply	22
5.4.1	Shared Memory and Register Declarations	22
5.4.2	Loading Parameters and Computing Thread Indices	23
5.4.3	Initializing the Accumulator	23
5.4.4	Tile Loop Implementation	24
5.5	Optimizing the Dot Product Computation	25
5.6	Performance Considerations and Comparison	26
6	Using Tensor Core Instructions (MMA) in PTX	27
6.1	Introduction to Tensor Cores	27
6.2	MMA Instructions in PTX	27
6.2.1	MMA Instruction Syntax	27
6.2.2	Warp-Level Collaboration	28
6.3	A Simplified Tensor Core Kernel Example	28
6.4	Loading Matrix Fragments	29
6.5	Performance Considerations for Tensor Core Usage	29
6.6	Tensor Cores in Deep Learning Workloads	30
7	Tools and Workflow for Debugging and Profiling PTX on Linux	30
7.1	CUDA Debugger (cuda-gdb)	30
7.1.1	Basic Debugging with cuda-gdb	30
7.1.2	PTX-Level Debugging	31
7.2	NVIDIA Nsight Compute	31
7.2.1	Using Nsight Compute CLI	31
7.2.2	Interactive Profiling with the GUI	31
7.2.3	PTX/SASS Correlation	32

7.3	NVIDIA Nsight Systems	32
7.4	Binary Utilities: cuobjdump and nvdasm	32
7.4.1	cuobjdump	33
7.4.2	nvdasm	33
7.5	CUDA Compute Sanitizer	33
7.6	A Practical Debugging and Optimization Workflow	33
8	Best Practices for Writing Performant PTX	34
8.1	Start with a Reference Implementation	34
8.2	Memory Access Optimization	35
8.2.1	Global Memory Coalescing	35
8.2.2	Use Shared Memory for Data Reuse	35
8.2.3	Minimize Bank Conflicts in Shared Memory	36
8.3	Register Usage Optimization	36
8.3.1	Monitor Register Count	36
8.3.2	Avoid Register Spilling	36
8.3.3	Register Pressure vs. Recomputation	36
8.4	Instruction-Level Optimization	37
8.4.1	Exploit Instruction-Level Parallelism (ILP)	37
8.4.2	Utilize Special Math Instructions	37
8.4.3	Loop Unrolling	38
8.5	Warp-Level Optimization	38
8.5.1	Minimize Warp Divergence	38
8.5.2	Leverage Warp-Level Primitives	39
8.6	Memory Hierarchy Management	39
8.6.1	Use Texture Memory for Read-Only Data with Spatial Locality	39
8.6.2	Leverage L2 Cache with Appropriate Hints	39
8.6.3	Consider Asynchronous Memory Operations	40
8.7	Architecture-Specific Optimizations	40
8.7.1	Tensor Cores	40
8.7.2	Newer Architecture Features	40
8.8	Practical Optimization Strategy	40
9	Case Studies: PTX Optimization for Deep Learning Kernels	41
9.1	Case Study 1: Fused Activation Functions	41
9.1.1	The Challenge	41
9.1.2	The PTX Solution	41
9.1.3	Performance Impact	42
9.2	Case Study 2: Custom Attention Mechanism	42
9.2.1	The Challenge	42
9.2.2	The PTX Solution	42
9.2.3	Performance Impact	43
9.3	Case Study 3: Quantized Tensor Core Operations	43
9.3.1	The Challenge	43
9.3.2	The PTX Solution	43
9.3.3	Performance Impact	44
9.4	Case Study 4: Custom Layer Normalization	44
9.4.1	The Challenge	44
9.4.2	The PTX Solution	44
9.4.3	Performance Impact	45
9.5	Lessons from Case Studies	45

10 Transitioning from CUDA C++ to PTX	46
10.1 Conceptual Mapping	46
10.2 Thread and Block Indexing	46
10.3 Memory Access and Address Calculation	47
10.4 Shared Memory Usage	47
10.5 Control Flow	48
10.6 Predicated Execution vs. Branching	49
10.7 Common Pitfalls in Transition	50
10.7.1 Register Management	50
10.7.2 Memory Alignment	50
10.7.3 Synchronization	51
10.7.4 Literal Representation	51
10.8 The Transition Process	51
10.9 When to Use PTX vs. CUDA C++	51
11 Practical Examples with Complete Code	52
11.1 Example 1: Vector Addition with PTX	52
11.1.1 The PTX Kernel	52
11.1.2 Host Code to Launch the PTX Kernel	53
11.1.3 Compilation and Execution	55
11.2 Example 2: Element-wise ReLU Activation with Inline PTX	56
11.3 Example 3: Tensor Core Matrix Multiplication	57
11.3.1 PTX Kernel for Tensor Core Matrix Multiplication	57
11.3.2 Host Code Considerations	59
12 Conclusion	60
12.1 Key Takeaways	60
12.2 When to Use PTX	60
12.3 The Future of PTX and GPU Programming	61
12.4 Final Thoughts	61

1 Introduction to PTX and the CUDA Toolchain

Parallel Thread Execution (PTX) is NVIDIA's low-level assembly-like language that serves as the intermediate representation for CUDA GPU computing. Understanding PTX gives developers a deeper insight into how their code executes on the GPU and opens up opportunities for advanced optimizations that may not be possible with high-level CUDA C/C++ alone.

1.1 What is PTX?

PTX can be thought of as a virtual GPU assembly language—an intermediate representation that sits between your high-level CUDA code and the actual hardware instructions. When you compile a CUDA C/C++ program, the NVIDIA compiler (NVCC) first translates your device code into PTX, and then the PTX assembler (`ptxas`) converts that PTX into SASS (Streaming Assembly), which is the actual machine code for the specific GPU architecture [1].

This design provides several advantages:

- **Forward Compatibility:** PTX serves as a stable programming target that can be JIT-compiled by the GPU driver to run on newer hardware that didn't exist when the application was built [1].
- **Abstraction Layer:** PTX abstracts some hardware details while still providing low-level control over parallelism, memory operations, and synchronization.
- **Optimization Opportunities:** By inspecting or writing PTX directly, developers can implement optimizations that might be missed by the high-level compiler.

PTX is sometimes described as "CUDA's IR" (Intermediate Representation) [2]. Unlike pure machine code, PTX is somewhat higher-level—it supports symbolic registers, labels, and abstracts many hardware details. For example, a single PTX `add.f32` instruction may map to different actual instructions depending on the GPU generation.

1.2 The Role of PTX in the CUDA Toolchain

The CUDA compilation process involves several steps, as illustrated in Figure ??:

1. The CUDA compiler (NVCC) separates host (CPU) code and device (GPU) code
2. The device code is compiled to PTX (an intermediate representation)
3. The PTX code is then assembled by `ptxas` into binary cubin code for a specific target GPU
4. Both PTX and cubin can be packaged into the final executable

At runtime, the CUDA driver can:

- Use the pre-compiled cubin if it matches the current GPU architecture
- Or just-in-time compile the embedded PTX if running on a newer GPU architecture

This approach ensures both performance (through pre-compiled binaries) and forward compatibility (through JIT compilation of PTX) [1].

Technical Note

PTX is human-readable text resembling assembly language. This makes it possible to inspect, modify, or even write PTX code from scratch, giving developers another level of control over GPU execution.

1.3 Why Learn PTX for Deep Learning?

For developers working on deep learning applications, understanding PTX can be valuable for several reasons:

- **Performance Optimization:** Deep learning is extremely performance-sensitive. By understanding or writing PTX, you can optimize critical kernels beyond what automatic compilation provides.
- **Access to Hardware Features:** Some specialized hardware features (like Tensor Cores) might be more directly or flexibly accessible through PTX than through high-level APIs.
- **Custom Operations:** When implementing custom neural network operations not supported by frameworks, low-level PTX knowledge can help maximize performance.
- **Debugging:** Understanding PTX helps diagnose performance issues by seeing what the compiler is actually doing with your high-level code.

In deep learning workloads, operations like matrix multiplication, convolutions, and custom fusion kernels are performance-critical. Knowledge of PTX enables fine-tuning these operations to get the most out of the GPU hardware.

Optimization Tip

While libraries like cuDNN and cuBLAS are highly optimized for common deep learning operations, custom operations or unique combinations of operations may benefit from hand-tuned PTX code. This is especially true for research implementations where standard libraries may not yet support new algorithmic innovations.

2 Inspecting PTX from CUDA C/C++ Code

Before diving into writing PTX from scratch, it's instructive to examine what PTX the compiler generates from simple CUDA kernels. This provides insight into PTX syntax and helps verify what high-level code actually translates to at the intermediate level.

2.1 Generating PTX from CUDA Source Code

NVIDIA provides straightforward tools to extract PTX from CUDA source code. The simplest approach is to use the `-ptx` flag with the NVCC compiler:

```
1 nvcc -ptx mykernel.cu -o mykernel.ptx
```

This command tells NVCC to compile the device code in `mykernel.cu` into a PTX file (without producing a full binary or executable). The resulting `mykernel.ptx` will be a text file that you can open in any editor to inspect the assembly-like instructions [4].

For more control, you can specify the compute capability of the target GPU using the `-arch` flag:

```
1 nvcc -ptx -arch=sm_80 mykernel.cu -o mykernel.ptx
```

This generates PTX specifically for compute capability 8.0 (NVIDIA Ampere architecture). Different compute capabilities may produce slightly different PTX due to hardware-specific optimizations and available features.

2.2 Extracting PTX from Compiled Code

If you have an already compiled executable or library that contains embedded PTX and cubin code, you can extract the PTX using the CUDA Binary Utilities:

```
1 cuobjdump -ptx myprogram
```

This will dump any PTX embedded in `myprogram`'s CUDA binaries [5]. This approach is useful if you want to examine the PTX for kernels in pre-built applications or libraries, such as deep learning frameworks.

Additionally, you can use `cuobjdump -sass` to dump the SASS (machine code) for comparison:

```
1 cuobjdump -sass myprogram
```

Comparing PTX with SASS can reveal how the assembler transforms your intermediate code into actual hardware instructions.

2.3 Example: Vector Addition Kernel

Let's examine a simple vector addition kernel in CUDA C++ and its corresponding PTX to get a feel for the translation process.

Here's a standard CUDA C++ implementation of vector addition:

```
1 // CUDA C++ kernel (device code)
2 __global__ void vecAdd(const float* A, const float* B, float* C, ←
   int N) {
3     int index = blockIdx.x * blockDim.x + threadIdx.x;
4     if (index < N) {
5         C[index] = A[index] + B[index];
6     }
7 }
```

When compiled with `nvcc -ptx`, the resulting PTX for compute capability 8.0 would look something like this:

```
1 .visible .entry _Z6vecAddPKfS0_Pfi(
2     .param .u64 _Z6vecAddPKfS0_Pfi_param_0,
3     .param .u64 _Z6vecAddPKfS0_Pfi_param_1,
4     .param .u64 _Z6vecAddPKfS0_Pfi_param_2,
5     .param .u32 _Z6vecAddPKfS0_Pfi_param_3
6 )
7 {
8     .reg .pred %p<2>;
9     .reg .f32 %f<4>;
10    .reg .b32 %r<6>;
11    .reg .b64 %rd<11>;
12
13    ld.param.u64    %rd1, [_Z6vecAddPKfS0_Pfi_param_0]; // load ←
14    A pointer
15    ld.param.u64    %rd2, [_Z6vecAddPKfS0_Pfi_param_1]; // load ←
16    B pointer
17    ld.param.u64    %rd3, [_Z6vecAddPKfS0_Pfi_param_2]; // load ←
18    C pointer
```

```

16      ld.param.u32    %r2,  [_Z6vecAddPKfS0_Pfi_param_3];  // load ←
      N value
17
18      mov.u32    %r3, %tid.x;          // threadIdx.x
19      mov.u32    %r4, %ntid.x;         // blockDim.x
20      mov.u32    %r5, %ctaid.x;        // blockIdx.x
21      mad.lo.s32 %r1, %r5, %r4, %r3; // r1 = blockIdx.x * blockDim.←
      x + threadIdx.x
22
23      setp.ge.u32 %p1, %r1, %r2;       // if r1 >= N, set predicate p1
24      @%p1 bra    DONE;               // if p1 (index >= N) is true, ←
      branch to DONE
25
26      cvta.to.global.u64 %rd4, %rd1;   // convert A pointer to ←
      global address
27      mul.wide.u32 %rd5, %r1, 4;       // rd5 = index * 4 (byte ←
      offset, float=4 bytes)
28      add.s64    %rd6, %rd4, %rd5;     // rd6 = A + index*4
29      ld.global.f32 %f1, [%rd6];       // f1 = A[index]
30      cvta.to.global.u64 %rd7, %rd2;   // convert B pointer
31      add.s64    %rd8, %rd7, %rd5;     // rd8 = B + index*4
32      ld.global.f32 %f2, [%rd8];       // f2 = B[index]
33      add.f32    %f3, %f1, %f2;        // f3 = f1 + f2 (the ←
      addition)
34      cvta.to.global.u64 %rd9, %rd3;   // convert C pointer
35      add.s64    %rd10, %rd9, %rd5;    // rd10 = C + index*4
36      st.global.f32 [%rd10], %f3;      // C[index] = f3
37
38  DONE:
39      ret;

```

Let's analyze the key elements in this PTX code:

1. **Function Declaration:** The `.visible .entry` directive defines a GPU kernel that can be called from host code. The mangled name `_Z6vecAddPKfS0_Pfi` corresponds to our `vecAdd` function with its parameter types.
2. **Parameter Declarations:** Parameters are declared with `.param` directives, specifying their types (`.u64` for pointers, `.u32` for integers).
3. **Register Declarations:** The `.reg` directives declare registers of various types to hold values during computation. For example, `.f32 %f<4>;` declares four 32-bit floating-point registers.
4. **Loading Parameters:** The `ld.param` instructions load function parameters into registers.
5. **Thread Index Calculation:** The code uses special registers (`%tid.x`, `%ntid.x`, `%ctaid.x`) to compute the global thread index, mirroring our C++ formula `blockIdx.x * blockDim.x + threadIdx.x`.
6. **Boundary Check:** The `setp.ge.u32` and `@%p1 bra` instructions implement the `if (index < N)` check from our C++ code.
7. **Memory Access:** PTX uses explicit memory operations (`ld.global.f32`, `st.global.f32`) for loading and storing data, with manual calculation of byte addresses.

8. **Arithmetic Operation:** The `add.f32` instruction performs the actual addition of the two float values.

This example illustrates the verbosity of PTX compared to CUDA C++. The simple three-line body of our C++ kernel expands to over 20 lines of PTX, with explicit management of registers, memory addressing, and control flow.

Technical Note

PTX uses virtual registers (like `%r1`, `%f3`) instead of named high-level variables. It explicitly converts pointers and calculates byte offsets for memory accesses, showing what the CUDA compiler does behind the scenes.

By examining compiler-generated PTX for simple kernels, you can build a mental model of how common CUDA constructs translate to assembly form. This is valuable preparation for writing your own PTX code or modifying compiler output for optimization.

3 Writing and Compiling PTX Directly

With a foundational understanding of PTX from examining compiler output, we can now explore how to write and run our own PTX code. This section covers different approaches to incorporating hand-written PTX into your workflow and the tools used to compile and execute it.

3.1 Approaches to Writing PTX

There are three main ways to incorporate hand-written PTX into your CUDA applications:

3.1.1 Inline PTX Assembly in CUDA C/C++

For small PTX snippets or minor adjustments to compiler-generated code, you can use CUDA's inline assembly syntax:

```
1 __global__ void myKernel(float *data) {
2     int tid;
3     // Inline PTX to get thread ID
4     asm("{\n\tmov.u32 %0, %tid.x;\n}" : "=r"(tid));
5
6     // Rest of kernel in C++
7     data[tid] = tid * 3.14f;
8 }
```

This approach is convenient when you only need to optimize specific instructions while keeping most of the kernel in C++. NVIDIA provides documentation on inline PTX usage in their programming guide [6].

The syntax follows the GCC extended inline assembly format:

```
1 asm("assembly code" : outputs : inputs : clobbers);
```

Where:

- `assembly code` is the PTX instructions
- `outputs` specifies output operands

- `inputs` specifies input operands
- `clobbers` lists registers that might be modified

3.1.2 Writing Complete PTX Kernel Files

For more extensive control, you can write entire kernels in PTX. This involves creating a `.ptx` file containing the complete kernel implementation:

```

1  .version 8.0
2  .target sm_80
3  .address_size 64
4
5  .visible .entry myKernel(
6      .param .u64 param0, // pointer parameter
7      .param .u64 param1  // another pointer
8  ) {
9      // PTX instructions for the kernel body
10     // ...
11     ret;
12 }
```

A PTX source file typically starts with directives specifying:

- The PTX version (`.version`)
- Target architecture (`.target`)
- Address size (`.address_size`)

Then it declares one or more kernel functions using `.visible .entry` (for kernels callable from host code) or `.func` (for device functions callable from other device code).

3.1.3 Using NVRTC (Runtime Compilation)

NVIDIA Runtime Compilation (NVRTC) allows compiling CUDA C++ source or PTX code at runtime from strings in memory:

```

1  // Create NVRTC program
2  nvrtcProgram prog;
3  nvrtcCreateProgram(&prog, ptxSource, "kernel.ptx", 0, NULL, NULL)↵
4      ;
5  // Compile the PTX source
6  nvrtcResult compileResult = nvrtcCompileProgram(prog, 0, NULL);
7
8  // Get the compiled PTX
9  size_t ptxSize;
10 nvrtcGetPTXSize(prog, &ptxSize);
11 char* compiledPtx = new char[ptxSize];
12 nvrtcGetPTX(prog, compiledPtx);
13
14 // Use the compiled PTX with the CUDA driver API
15 // ...
```

This approach is powerful for dynamic code generation or meta-programming, commonly used in deep learning frameworks for JIT-compiling specialized kernels [7].

3.2 Compiling PTX Code

Once you've written a PTX file, you need to compile it into a form that can execute on the GPU. The main tool for this is NVIDIA's PTX assembler, `ptxas`.

3.2.1 Using `ptxas`

To compile PTX to a cubin (CUDA binary), use:

```
1 ptxas -arch=sm_80 -O3 myKernel.ptx -o myKernel.cubin
```

This command assembles the PTX file for compute capability 8.0 (Ampere) with optimization level 3. The output is a binary file that can be loaded by the CUDA driver.

Adding the `-v` (verbose) flag provides useful information about resource usage:

```
1 ptxas -v -arch=sm_80 myKernel.ptx -o myKernel.cubin
```

This might output something like:

```
1 ptxas info : Used 32 registers, 0 bytes spill stores, 0 bytes ←  
    spill loads, 48 bytes shared memory, 136 bytes cmem[0]
```

This information is valuable for performance tuning as it shows how many registers per thread your kernel uses, whether any register spills occurred (which hurt performance), and shared memory usage.

Optimization Tip

Always watch for register spills in the `ptxas` output. If you see "spill stores/loads" > 0, it means the kernel ran out of registers and had to use slower local memory. Try to restructure your code to reduce register pressure if this happens.

3.3 Running PTX Kernels

After compiling your PTX to a cubin, you need to load and execute it. There are two main approaches:

3.3.1 Using the CUDA Driver API

The CUDA Driver API allows loading PTX or cubin files at runtime and launching the kernels within them:

```
1 #include <cuda.h>  
2  
3 int main() {  
4     // Initialize CUDA  
5     cuInit(0);  
6     CUdevice device;  
7     CUcontext context;  
8     cuDeviceGet(&device, 0);  
9     cuCtxCreate(&context, 0, device);  
10  
11     // Load the PTX module
```

```

12     CUmodule module;
13     cuModuleLoad(&module, "myKernel.cubin"); // or load PTX ←
        directly
14
15     // Get function handle
16     CUfunction kernel;
17     cuModuleGetFunction(&kernel, module, "myKernel");
18
19     // Set up kernel parameters
20     void* args[] = { &d_input, &d_output, &size };
21
22     // Launch the kernel
23     cuLaunchKernel(
24         kernel,           // kernel function
25         gridDimX, gridDimY, gridDimZ, // grid dimensions
26         blockDimX, blockDimY, blockDimZ, // block dimensions
27         sharedMemBytes, // shared memory size
28         stream,          // CUDA stream
29         args,            // kernel arguments
30         NULL             // extra options
31     );
32
33     // Clean up
34     cuModuleUnload(module);
35     cuCtxDestroy(context);
36
37     return 0;
38 }

```

This approach gives you complete control over loading and executing PTX code, making it suitable for applications that dynamically select or generate kernels [7].

3.3.2 Embedding in Host Code

For a more integrated approach, you can embed PTX or cubin into your application as a resource and link it with your host code. NVIDIA provides tools like **fatbinary** to create object files containing both host code and device code.

This is more complex but can be managed with build systems or CUDA's own mechanisms. The CUDA runtime API can then find and load these embedded kernels automatically.

Technical Note

The CUDA samples (e.g., **vectorAddDrv**) demonstrate loading a PTX module and launching a kernel via the Driver API. These examples are valuable references for understanding the mechanics of PTX execution [7].

4 PTX Syntax and Key Constructs

To effectively write or modify PTX code, you need to understand its syntax and the core constructs it uses to express GPU computation. This section covers the essential elements of PTX code, focusing on how to express thread parallelism, register usage, memory operations, and control flow.

4.1 Module and Entry Point Directives

PTX files begin with directives that set up the execution environment and specify the target hardware.

4.1.1 Version and Target Directives

These directives typically appear at the top of a PTX file:

```
1 .version 8.7          // PTX ISA version 8.7
2 .target sm_90         // target GPU architecture (Hopper GPUs)
3 .address_size 64      // 64-bit addressing
```

- `.version` - Specifies the PTX ISA (Instruction Set Architecture) version. This defines the language features available in the PTX code.
- `.target` - Indicates the GPU architecture family for which the PTX is being written. Examples include `sm_70` for Volta, `sm_80` for Ampere, and `sm_90` for Hopper.
- `.address_size` - Defines the size of memory addresses (typically 64 bits for modern GPUs).

The PTX version corresponds to the language capabilities and advances alongside the CUDA toolkit version. The target specifies which GPU generation's features you can use—if omitted or set to a generic `compute_XY`, the PTX might be more portable but could restrict access to certain hardware-specific instructions.

4.1.2 Kernel Entry Points

Kernel functions—those that can be launched from host code—are declared with the `.entry` directive:

```
1 .visible .entry kernelName(
2     .param .u64 param_A,    // first parameter (a pointer)
3     .param .u32 param_N     // second parameter (an integer)
4 ) {
5     // kernel body
6     // ...
7     ret;
8 }
```

The components of this declaration are:

- `.visible` - Makes the entry point externally visible, allowing the host to find it by name.
- `.entry` - Marks this as a kernel function (as opposed to a device function).
- `kernelName` - The name of the kernel, used when launching it from host code.
- `.param` declarations - List the parameters the kernel expects, with their types.

The parameters are declared with `.param` followed by a type qualifier (like `.u64` for a 64-bit unsigned integer or pointer) and a name. Inside the function body, these parameters are accessible via `ld.param` instructions that load from the parameter space into registers.

Device functions (callable only from device code) use `.func` instead of `.entry` and can omit the `.visible` qualifier if they're only used internally.

4.2 Registers and Type Qualifiers

In PTX, all variables are held in registers, which must be explicitly declared and managed. Unlike high-level languages, there's no automatic stack or heap for local variables.

4.2.1 Register Declaration

Registers are declared using the `.reg` directive, specifying a type and either a name or a range:

```
1 .reg .f32    %f<4>;    // 4 float registers: %f0, %f1, %f2, %f3
2 .reg .pred   %p<2>;    // 2 predicate (boolean) registers: %p0, %p1
3 .reg .b32    %rTemp;    // one 32-bit register named rTemp
```

The first line creates four 32-bit floating-point registers named `%f0` through `%f3`. The second creates two predicate registers (used for conditional execution). The third creates a single 32-bit register with a custom name.

When using registers in instructions, they are prefixed with `%` (e.g., `%f0`, `%rTemp`).

4.2.2 Register Types

PTX supports various register types to hold different kinds of data:

- `.b32`, `.b64` - Raw bits (32 or 64 bits), used for integers or pointers
- `.s32`, `.s64` - Signed integers (32 or 64 bits)
- `.u32`, `.u64` - Unsigned integers (32 or 64 bits)
- `.f32`, `.f64` - Floating-point values (32 or 64 bits)
- `.pred` - Predicate (1-bit boolean value)
- Vector types like `.v2.f32` - Vectors of values (e.g., two `f32` values treated as one unit)

Technical Note

Registers in PTX are virtual and don't directly correspond one-to-one with physical hardware registers [8]. The PTX-to-SASS compiler will allocate physical registers and may optimize by eliminating unused registers or merging ones that don't have overlapping lifetimes.

4.2.3 Register Usage Considerations

All registers declared in a PTX kernel are thread-local—each thread gets its own set. This has important performance implications:

- GPUs have a finite register file per Streaming Multiprocessor (SM). For example, a V100 GPU has 65,536 registers per SM [9].
- If each thread uses many registers, fewer threads can run concurrently on an SM, potentially reducing occupancy (the ratio of active warps to maximum possible warps).
- Lower occupancy can mean less ability to hide memory latency through thread switching, affecting performance.

When writing PTX, it's important to balance register usage against occupancy and other concerns. The `ptxas -v` output tells you how many registers your kernel uses per thread.

Optimization Tip

You can declare more registers than you actually need without penalty; the compiler will eliminate unused ones. What matters is how many are simultaneously "alive" in your code. Reusing registers for different purposes at different points in your kernel can help reduce register pressure.

4.3 Thread Indexes and Special Registers

In CUDA C++, built-in variables like `threadIdx.x` and `blockIdx.x` identify each thread's position. PTX provides these values through special registers with names like `%tid.x` and `%ctaid.x`.

4.3.1 Common Thread Index Registers

The main special registers for thread indexing are:

- `%tid.x`, `%tid.y`, `%tid.z` - Thread index within its block (equivalent to `threadIdx` in CUDA C++)
- `%ctaid.x`, `%ctaid.y`, `%ctaid.z` - Block index within the grid (equivalent to `blockIdx`)
- `%ntid.x`, `%ntid.y`, `%ntid.z` - Block dimensions (equivalent to `blockDim`, i.e., number of threads per block in each dimension)
- `%nctaid.x`, `%nctaid.y`, `%nctaid.z` - Grid dimensions (equivalent to `gridDim`, i.e., number of blocks in the grid)

4.3.2 Using Thread Index Registers

These special registers are read-only. To use their values in calculations, you typically move them into general-purpose registers first:

```

1  mov.u32 %r1, %tid.x;    // Copy threadIdx.x to register r1
2  mov.u32 %r2, %ntid.x;   // Copy blockDim.x to register r2
3  mov.u32 %r3, %ctaid.x;  // Copy blockIdx.x to register r3
4
5  // Calculate global thread index: r4 = blockIdx.x * blockDim.x + ←
   threadIdx.x
6  mad.lo.s32 %r4, %r3, %r2, %r1;
```

Here, we use the `mov.u32` instruction to copy the special register values into general registers, and then `mad.lo.s32` (multiply-add) to compute the global thread index.

4.3.3 Additional Special Registers

PTX provides other special registers for advanced thread management:

- `%laneid` - The thread's index within its warp (0-31), useful for warp-synchronous programming
- `%warpid` - The warp's index within the block

- `%nwarpid` - Number of warps per block
- `%smid` - SM identifier executing the thread
- `%lanemask_lt`, `%lanemask_le`, etc. - Masks for threads with lane IDs less than, less than or equal to, etc., the current thread

These are particularly useful for implementing warp-level optimizations, where threads in the same warp cooperate on a computation.

4.4 Arithmetic and Control Flow Instructions

PTX instructions use a relatively straightforward assembly-like syntax, with an opcode, type qualifiers, and operands.

4.4.1 Basic Arithmetic Instructions

Arithmetic operations in PTX specify the operation, data type, and operands:

```

1 add.f32 %f3, %f1, %f2;    // f3 = f1 + f2 (float addition)
2 sub.s32 %r3, %r1, %r2;    // r3 = r1 - r2 (signed int subtraction)
3 mul.f32 %f4, %f1, 2.5;    // f4 = f1 * 2.5 (float multiplication ←
    with immediate)
4 mad.f32 %f5, %f1, %f2, %f3; // f5 = f1 * f2 + f3 (fused multiply←
    -add)

```

The type suffix (e.g., `.f32`, `.s32`) specifies the data type the instruction operates on. This allows PTX to use the same opcode (e.g., `add`) for different data types, with the suffix distinguishing them.

4.4.2 Vector Operations

PTX supports vectorized operations for memory access and some arithmetic:

```

1 ld.global.v4.f32 {%f0, %f1, %f2, %f3}, [%rd1]; // Load 4 floats ←
    at once
2 st.global.v2.f32 [%rd2], {%f0, %f1};           // Store 2 floats ←
    at once

```

Vector operations can improve memory bandwidth utilization and arithmetic throughput. However, they require proper alignment—for example, a `v4.f32` load requires 16-byte alignment of the memory address.

4.4.3 Control Flow: Predicates and Branches

Unlike high-level languages with structured control flow (if/else, loops), PTX uses predicates and branches for control flow:

```

1 // Compute if r1 >= r2
2 setp.ge.u32 %p1, %r1, %r2; // Set predicate p1 true if r1 >= r2
3 @%p1 bra TARGET;           // If p1 is true, branch to TARGET
4                             // Fall through if p1 is false
5
6 // Code executed if r1 < r2

```



```

7 // ...
8
9 TARGET:
10 // Code executed if r1 >= r2 (or after the above code)

```

The `setp` instruction sets a predicate register based on a comparison. The `@%p` prefix on an instruction makes it conditional, executing only if the predicate is true.

4.4.4 Predicated Execution

For short conditional code, PTX often uses predication instead of branches to avoid warp divergence:

```

1 setp.lt.f32 %p1, %f1, 0.0; // p1 = (f1 < 0.0)
2 @%p1 mov.f32 %f1, 0.0;    // If p1, then f1 = 0.0
3                           // (implements ReLU: max(f1, 0))

```

This conditional execution using `@%p` can be more efficient than branching for short sequences, as it avoids potential warp divergence (where threads in a warp take different execution paths).

4.4.5 Implementing Loops

Loops in PTX are implemented with labels and branches:

```

1 mov.u32 %r1, 0; // Initialize loop counter
2
3 LOOP_START:
4 // Loop body instructions
5 // ...
6
7 add.u32 %r1, %r1, 1; // Increment counter
8 setp.lt.u32 %p1, %r1, %r2; // p1 = (r1 < r2)
9 @%p1 bra LOOP_START; // If p1, go back to start
10 // Otherwise, fall through

```

This pattern implements a basic for-loop that iterates from 0 to less than `r2`.

Warning

Writing complex control flow in PTX can be tricky. Modern GPU performance often benefits from minimizing divergent branches within a warp. When possible, use predication for short conditional code to avoid branch divergence overhead.

4.5 Memory Spaces and Instructions

NVIDIA GPUs have a hierarchy of memory spaces with different characteristics. PTX makes these explicit, requiring you to specify which memory space you're accessing with each load or store instruction.

4.5.1 Memory Space Overview

The main memory spaces in PTX correspond to the memory hierarchy in CUDA C++:

- **Global memory** (`.global`): The large but relatively slow GPU DRAM, accessible by all threads.
- **Shared memory** (`.shared`): Fast on-chip memory shared by threads in a block, used for communication and data reuse.
- **Local memory** (`.local`): Thread-private memory used for register spills or large local arrays, physically located in device memory (slow).
- **Constant memory** (`.const`): Read-only memory with a dedicated cache, optimized for broadcast access.
- **Parameter memory** (`.param`): Used for passing arguments to kernels and between functions.

4.5.2 Memory Access Instructions

Memory operations in PTX specify the memory space, data type, and addressing:

```

1 // Load a float from global memory
2 ld.global.f32 %f1, [%rd1];
3
4 // Store a 32-bit value to shared memory
5 st.shared.b32 [%rd2], %r1;
6
7 // Load from constant memory
8 ld.const.f32 %f2, [%rd3+4]; // Offset by 4 bytes
9
10 // Load a 64-bit parameter
11 ld.param.u64 %rd4, [param_ptr];

```

The general format is `ld.<space>.<type>` for loads and `st.<space>.<type>` for stores, where `<space>` is the memory space and `<type>` is the data type.

4.5.3 Declaring Static Memory

To use shared or constant memory, you need to declare it in the PTX code:

```

1 // Declare a 256-element array of 32-bit values in shared memory
2 .shared .align 4 .b32 sharedArray[256];
3
4 // Declare a constant array
5 .const .align 4 .b32 constArray[16] = {1, 2, 3, 4, 5, 6, 7, 8,
6                                         9, 10, 11, 12, 13, 14, 15, ↵
6                                         16};

```

The `.align` directive ensures proper alignment, which can be important for performance and correctness.

4.5.4 Memory Performance Considerations

The memory hierarchy in GPUs has significant performance implications:

- **Global memory** is large but has high latency (hundreds of cycles for uncached access). Coalesced access (where threads in a warp access contiguous addresses) is crucial for performance.

- **Shared memory** is much faster (similar to registers if there are no bank conflicts) but limited in size (e.g., tens of KB per SM). It's key for data reuse among threads in a block.
- **Constant memory** is cached and optimized for situations where all threads read the same value (broadcast).
- **Local memory** uses the same physical storage as global memory and should be avoided if possible, as it indicates register spilling.

For deep learning kernels, which often process large arrays of data, effective use of the memory hierarchy can make an enormous difference in performance. A common pattern is to stage data in shared memory to reuse it and avoid redundant global memory accesses.

Technical Note

On modern GPUs, a global memory load that isn't in cache can take hundreds of clock cycles, while shared memory access can be as fast as register access under optimal conditions [3].

4.5.5 Memory Caching and Control

PTX gives some control over the caching behavior of global memory operations. For example:

```
1 // Cache at global level only
2 ld.global.cg.f32 %f1, [%rd1];
3
4 // Cache at all levels
5 ld.global.ca.f32 %f1, [%rd1];
6
7 // Specify cache line size (Ampere+)
8 ld.global.L2::128B.f32 %f1, [%rd1];
```

Advanced usage can include prefetch or asynchronous load instructions for overlapping memory transfer with computation. For instance, Ampere architecture introduced `cp.async` instructions to asynchronously copy from global to shared memory:

```
1 // Asynchronously copy 16 bytes from global to shared memory
2 cp.async.ca.shared.global [%rd_shared], [%rd_global], 16;
```

Such features can greatly improve performance in memory-bound kernels but may require specific hardware support.

4.5.6 Synchronization

PTX provides various synchronization primitives to coordinate threads:

```
1 // Synchronize all threads in a block (equivalent to ↵
   // __syncthreads())
2 bar.sync 0;
3
4 // Memory fence for global memory operations
5 membar.gl;
6
7 // Memory fence for shared memory operations within a block
```

```
8 | membar.cta;
```

The `bar.sync` instruction corresponds to CUDA's `__syncthreads()`, ensuring all threads in a block reach the barrier before any proceed. This is essential after writing to shared memory and before threads read each other's results.

Memory fences (`membar`) ensure memory operations complete in the specified order, which can be important for complex algorithms with data dependencies.

With these fundamentals—thread registers, basic arithmetic, memory operations, and synchronization—you can write complete GPU kernels in PTX. The next sections will apply these to practical examples, focusing on deep learning relevant computations.

5 Hands-On Example: PTX vs. CUDA C++ for a Matrix Multiply Kernel

To illustrate the process of developing and optimizing in PTX, let's implement a matrix multiplication kernel—a core operation in deep learning for computing fully-connected layer outputs, attention mechanisms, and other operations. We'll compare a straightforward CUDA C++ implementation with an optimized PTX version.

5.1 The Matrix Multiplication Problem

Matrix multiplication is a fundamental operation in deep learning. Given two matrices A and B, we want to compute their product $C = A \times B$.

For simplicity, let's consider the case where all matrices are square with dimensions $N \times N$. In matrix multiplication, each element $C[i,j]$ is the dot product of row i from matrix A and column j from matrix B:

$$C[i,j] = \sum_{k=0}^{N-1} A[i,k] \times B[k,j] \quad (1)$$

5.2 Naive CUDA C++ Implementation

First, let's look at a simple CUDA C++ implementation where each thread computes one element of the output matrix:

```
1  __global__ void matMul(const float* A, const float* B, float* C, ↵
    int N) {
2      // Calculate row and column indices for this thread
3      int row = blockIdx.y * blockDim.y + threadIdx.y;
4      int col = blockIdx.x * blockDim.x + threadIdx.x;
5
6      // Check if within bounds
7      if (row < N && col < N) {
8          float sum = 0.0f;
9
10         // Compute dot product of row from A and column from B
11         for (int k = 0; k < N; ++k) {
12             sum += A[row * N + k] * B[k * N + col];
13         }
14
15         // Store result
16         C[row * N + col] = sum;
17     }
```

```
18 | }
```

This kernel is conceptually simple: each thread (identified by its row and column) computes the dot product of the corresponding row from A and column from B.

However, the performance of this naive approach is poor for large matrices because:

- Each thread accesses global memory N times for A and N times for B, with no data reuse between threads.
- The memory access pattern for matrix B is not coalesced (threads in a warp access non-consecutive elements), leading to inefficient memory transactions.
- There's no exploitation of the memory hierarchy to reduce global memory traffic.

Optimization Tip

A key optimization principle in GPU programming is to minimize global memory accesses by reusing data through shared memory or registers. This is especially important for matrix multiplication, where each element of the input matrices is used multiple times.

5.3 Optimized Matrix Multiply Using Shared Memory

A well-optimized matrix multiply on GPU uses tiling with shared memory. The basic idea is:

1. Divide the computation into tiles (submatrices).
2. Each thread block loads tiles of A and B into shared memory.
3. Threads in the block compute partial sums using the data in shared memory.
4. The process repeats for all tiles needed to compute the final result.

This approach dramatically reduces global memory traffic because each element is loaded once from global memory into shared memory and then reused by multiple threads.

Here's how this would look in CUDA C++:

```
1  #define TILE_SIZE 16
2
3  __global__ void matMulTiled(const float* A, const float* B, float* C, int N) {
4      // Shared memory for tiles
5      __shared__ float As[TILE_SIZE][TILE_SIZE];
6      __shared__ float Bs[TILE_SIZE][TILE_SIZE];
7
8      // Calculate row and column indices
9      int row = blockIdx.y * TILE_SIZE + threadIdx.y;
10     int col = blockIdx.x * TILE_SIZE + threadIdx.x;
11
12     float sum = 0.0f;
13
14     // Loop over tiles
15     for (int t = 0; t < (N + TILE_SIZE - 1) / TILE_SIZE; ++t) {
16         // Load tiles into shared memory
17         if (row < N && t * TILE_SIZE + threadIdx.x < N)
```

```

18         As[threadIdx.y][threadIdx.x] = A[row * N + t * ←
           TILE_SIZE + threadIdx.x];
19     else
20         As[threadIdx.y][threadIdx.x] = 0.0f;
21
22     if (col < N && t * TILE_SIZE + threadIdx.y < N)
23         Bs[threadIdx.y][threadIdx.x] = B[(t * TILE_SIZE + ←
           threadIdx.y) * N + col];
24     else
25         Bs[threadIdx.y][threadIdx.x] = 0.0f;
26
27     // Synchronize to ensure tiles are loaded
28     __syncthreads();
29
30     // Compute partial sum for this tile
31     for (int k = 0; k < TILE_SIZE; ++k) {
32         sum += As[threadIdx.y][k] * Bs[k][threadIdx.x];
33     }
34
35     // Synchronize before loading next tiles
36     __syncthreads();
37 }
38
39 // Store result
40 if (row < N && col < N) {
41     C[row * N + col] = sum;
42 }
43 }

```

This tiled approach significantly improves performance by:

- Reducing global memory accesses: Each element is loaded once from global memory and reused multiple times.
- Improving memory coalescing: When threads in a warp load data into shared memory, they access consecutive memory addresses.
- Exploiting data locality: Once data is in shared memory, access is much faster than global memory.

5.4 PTX Implementation of Tiled Matrix Multiply

Now, let's examine how we would implement a similar tiled matrix multiply directly in PTX. Rather than showing the entire kernel (which would be quite verbose), we'll focus on key sections with commentary.

5.4.1 Shared Memory and Register Declarations

At the start of the PTX kernel, we declare shared memory for the A and B tiles, and registers for computation:

```

1 .visible .entry matMulTiled(
2     .param .u64 A_ptr,      // pointer to matrix A
3     .param .u64 B_ptr,      // pointer to matrix B
4     .param .u64 C_ptr,      // pointer to output matrix C
5     .param .u32 N_val       // matrix dimension

```

```

6 ) {
7     // Register declarations
8     .reg .pred    %p<2>;                // Predicate registers
9     .reg .f32     %f_acc;                // Accumulator for result
10    .reg .f32     %f_a, %f_b;            // Registers for A and B ←
11        elements
12    .reg .u32     %r_row, %r_col;         // Row and column indices
13    .reg .u32     %r_tid_x, %r_tid_y;    // Thread indices
14    .reg .u32     %r_bx, %r_by;          // Block indices
15    .reg .u32     %r_N, %r_k, %r_t;      // Loop counters and matrix ←
16        size
17    .reg .u64     %rd_A, %rd_B, %rd_C;    // Base pointers
18    .reg .u64     %rd_Aaddr, %rd_Baddr, %rd_Caddr; // Computed ←
19        addresses
20
21    // Shared memory for tiles (16x16 tiles of f32 values)
22    .shared .align 4 .b32 shmemA[16*16]; // 256 floats for A ←
23        sub-matrix
24    .shared .align 4 .b32 shmemB[16*16]; // 256 floats for B ←
25        sub-matrix
26
27    // ... (kernel code follows)
28 }

```

5.4.2 Loading Parameters and Computing Thread Indices

Next, we load the kernel parameters and compute the global thread indices:

```

1     // Load parameters
2     ld.param.u64 %rd_A, [A_ptr];
3     ld.param.u64 %rd_B, [B_ptr];
4     ld.param.u64 %rd_C, [C_ptr];
5     ld.param.u32 %r_N, [N_val];
6
7     // Compute thread and block indices
8     mov.u32 %r_tid_x, %tid.x; // threadIdx.x
9     mov.u32 %r_tid_y, %tid.y; // threadIdx.y
10    mov.u32 %r_bx, %ctaid.x;   // blockIdx.x
11    mov.u32 %r_by, %ctaid.y;   // blockIdx.y
12
13    // Compute global row and column indices
14    // row = blockIdx.y * 16 + threadIdx.y
15    mul.lo.u32 %r_row, %r_by, 16;
16    add.u32 %r_row, %r_row, %r_tid_y;
17
18    // col = blockIdx.x * 16 + threadIdx.x
19    mul.lo.u32 %r_col, %r_bx, 16;
20    add.u32 %r_col, %r_col, %r_tid_x;

```

5.4.3 Initializing the Accumulator

Before the tile loop, we initialize the accumulator register to zero:

```

1     // Initialize accumulator to 0.0

```

```
2      mov.f32 %f_acc, 0f00000000;  // Hex representation of 0.0f
```

5.4.4 Tile Loop Implementation

Now we implement the loop over tiles. For each tile, we: 1. Load a portion of A and B from global to shared memory 2. Synchronize threads 3. Compute the partial dot product 4. Synchronize again before the next iteration

Here's a simplified version of the loop:

```
1      // Initialize tile counter
2      mov.u32 %r_t, 0;
3
4  TILE_LOOP_START:
5      // Compute the indices for this tile
6      // ... (index calculation for A and B)
7
8      // Load A element from global to shared memory
9      // Check bounds: if row < N && t*16 + threadIdx.x < N
10     setp.lt.u32 %p0, %r_row, %r_N;  // row < N
11     // ... (calculate index for t*16 + threadIdx.x)
12     setp.lt.u32 %p1, %r_k_idx, %r_N;  // t*16 + threadIdx.x < N
13     and.pred %p0, %p0, %p1;  // Both conditions true
14
15     @!%p0 bra LOAD_A_ZERO;  // If not in bounds, load 0
16
17     // Calculate global memory address for A[row, t*16 + ←
18         threadIdx.x]
19     // ... (address calculation)
20     ld.global.f32 %f_a, [%rd_Aaddr];
21     bra STORE_A_SHARED;
22
23 LOAD_A_ZERO:
24     mov.f32 %f_a, 0f00000000;  // Load 0 if out of bounds
25
26 STORE_A_SHARED:
27     // Calculate shared memory index shmemA[threadIdx.y][←
28         threadIdx.x]
29     // ... (shared memory index calculation)
30     st.shared.f32 [%rd_shAddrA], %f_a;
31
32     // Similar logic for loading B into shared memory
33     // ... (load B element)
34
35     // Synchronize threads
36     bar.sync 0;
37
38     // Compute partial dot product for this tile
39     // Initialize loop counter for k
40     mov.u32 %r_k, 0;
41
42 DOT_PRODUCT_LOOP:
43     // Calculate shared memory indices for A[threadIdx.y][k] and ←
44         B[k][threadIdx.x]
45     // ... (index calculations)
46
47     // Load from shared memory
```



```

45     ld.shared.f32 %f_a, [%rd_shA_k];
46     ld.shared.f32 %f_b, [%rd_shB_k];
47
48     // Multiply-accumulate: f_acc += f_a * f_b
49     fma.rn.f32 %f_acc, %f_a, %f_b, %f_acc;
50
51     // Increment k and check loop condition
52     add.u32 %r_k, %r_k, 1;
53     setp.lt.u32 %p0, %r_k, 16;    // k < 16
54     @%p0 bra DOT_PRODUCT_LOOP;    // Jump back if k < 16
55
56     // Synchronize before next tile
57     bar.sync 0;
58
59     // Increment tile counter and check if done
60     add.u32 %r_t, %r_t, 1;
61     // Calculate ceiling of N/16
62     add.u32 %r_temp, %r_N, 15;
63     div.u32 %r_num_tiles, %r_temp, 16;
64     setp.lt.u32 %p0, %r_t, %r_num_tiles;    // t < ceiling(N/16)
65     @%p0 bra TILE_LOOP_START;
66
67     // Store final result to global memory if in bounds
68     setp.lt.u32 %p0, %r_row, %r_N;    // row < N
69     setp.lt.u32 %p1, %r_col, %r_N;    // col < N
70     and.pred %p0, %p0, %p1;    // Both conditions true
71     @!%p0 bra EXIT;    // Skip if out of bounds
72
73     // Calculate global memory address for C[row, col]
74     // ... (address calculation)
75     st.global.f32 [%rd_Caddr], %f_acc;
76
77 EXIT:
78     ret;
79 }

```

This PTX implementation is more verbose than the CUDA C++ version, but it illustrates how we explicitly manage:

- **Thread indexing:** Manually computing row and column indices from block and thread IDs
- **Memory management:** Explicit loads and stores between global and shared memory
- **Control flow:** Using predicates and branches to implement loops and conditionals
- **Synchronization:** Explicit barriers to coordinate threads

5.5 Optimizing the Dot Product Computation

In practice, to maximize performance, we might unroll the inner dot product loop to reduce loop overhead and enable instruction-level parallelism:

```

1     // Instead of a loop, manually unroll the 16 iterations
2     // For k=0
3     // Load from shared memory
4     ld.shared.f32 %f_a0, [%rd_shA_base + 0*4];    // A[ty][0]

```

```

5      ld.shared.f32 %f_b0, [%rd_shB_base + 0*64]; // B[0][tx]
6      fma.rn.f32 %f_acc, %f_a0, %f_b0, %f_acc;    // acc += A[ty][0] ←
          * B[0][tx]
7
8      // For k=1
9      ld.shared.f32 %f_a1, [%rd_shA_base + 1*4]; // A[ty][1]
10     ld.shared.f32 %f_b1, [%rd_shB_base + 1*64]; // B[1][tx]
11     fma.rn.f32 %f_acc, %f_a1, %f_b1, %f_acc;    // acc += A[ty][1] ←
          * B[1][tx]
12
13     // ... (repeat for k=2 through k=14)
14
15     // For k=15
16     ld.shared.f32 %f_a15, [%rd_shA_base + 15*4]; // A[ty][15]
17     ld.shared.f32 %f_b15, [%rd_shB_base + 15*64]; // B[15][tx]
18     fma.rn.f32 %f_acc, %f_a15, %f_b15, %f_acc;  // acc += A[ty] ←
          ] [15] * B[15][tx]

```

This unrolling eliminates branch overhead and allows the hardware to better schedule instructions, potentially overlapping memory loads with computation.

5.6 Performance Considerations and Comparison

The PTX implementation offers several advantages over relying on the compiler to generate code from CUDA C++:

- **Fine-grained control:** We have explicit control over instruction scheduling and unrolling. We can ensure the loop is fully unrolled without any compiler-introduced conditionals.
- **Instruction interleaving:** We can manually interleave instructions to maximize instruction-level parallelism. For example, we could issue a load for the next iteration while the current multiply-add is happening.
- **Special instructions:** On newer GPUs, we can use specialized instructions like tensor core operations for matrix math, which might not be automatically used by the compiler.
- **Custom tile sizes:** We can tailor the kernel to specific matrix sizes or use custom tile sizes that the general compiler might not choose.

However, the PTX implementation is also:

- More verbose and harder to read
- More difficult to maintain
- More error-prone due to manual memory address calculations and synchronization

In practice, a good approach is to start with a high-level CUDA C++ implementation and use PTX only for the performance-critical sections or when you need specific hardware features.

Technical Note

When writing PTX, make sure to verify that your optimizations actually improve performance. The compiler is quite sophisticated and often generates good code for common patterns like matrix multiplication. Hand-optimized PTX is most beneficial for unusual access patterns or when leveraging hardware features not fully exposed by the compiler.

6 Using Tensor Core Instructions (MMA) in PTX

Modern NVIDIA GPUs (Volta and later) include specialized hardware units called Tensor Cores, which dramatically accelerate matrix multiplication operations. These units are particularly valuable for deep learning workloads, where matrix multiplication is a dominant computation. In this section, we'll explore how to leverage Tensor Cores directly through PTX.

6.1 Introduction to Tensor Cores

Tensor Cores are specialized matrix multiply-accumulate units designed to accelerate deep learning workloads. They provide hardware-accelerated matrix operations with significantly higher throughput than standard floating-point operations.

Key characteristics of Tensor Cores include:

- Perform matrix multiply-accumulate operations on small matrices (typically 4×4 or larger blocks)
- Support various precision formats, primarily:
 - FP16 (16-bit floating point) input with FP16 or FP32 accumulation
 - BF16 (brain floating point) in newer architectures
 - INT8/INT4 (8-bit/4-bit integer) for quantized operations
 - FP64 (64-bit floating point) in the latest architectures (Hopper)
- Deliver significantly higher FLOPS compared to standard FP32 operations (e.g., $8 \times$ higher theoretical throughput for FP16 on Ampere)

In CUDA C++, Tensor Cores can be accessed through high-level APIs like cuBLAS or through the WMMA (Warp Matrix Multiply Accumulate) API. With PTX, we can directly issue the matrix multiply-accumulate instructions that utilize Tensor Cores.

6.2 MMA Instructions in PTX

The primary PTX instruction for Tensor Core operations is `mma.sync`, which stands for Matrix Multiply Accumulate with synchronization. This instruction directs the warp to perform a collective matrix multiplication using Tensor Cores.

6.2.1 MMA Instruction Syntax

A typical `mma.sync` instruction might look like:

```
1 mma.sync.aligned.m16n8k16.row.col.f16.f16.f16.f16
2   {%d0,%d1,%d2,%d3},
3   {%a0,%a1,%a2,%a3},
4   {%b0,%b1},
5   {%c0,%c1,%c2,%c3};
```

Let's break down the components of this instruction:

- `mma.sync` - The base operation (matrix multiply-accumulate with warp synchronization)
- `aligned` - Memory alignment requirement
- `m16n8k16` - Shape of the matrix operation:

- `m16` - 16 rows in the output matrix (and in matrix A)
- `n8` - 8 columns in the output matrix (and in matrix B)
- `k16` - Contraction dimension (columns of A, rows of B)
- `row.col` - Layout of matrices A and B (A is row-major, B is column-major)
- `f16.f16.f16.f16` - Data types for A, B, C, and D matrices (all FP16 in this case)
- The register lists in curly braces are the matrices' "fragments":

– 6.2.2 Warp-Level Collaboration

An important concept with MMA operations is that they are **warp-level operations**. All 32 threads in a warp collaborate to perform the matrix multiplication, with each thread holding a portion of the input and output data.

When a thread issues an `mma.sync` instruction, it's actually directing the entire warp to perform a collective operation using the register fragments distributed across all threads. The "fragments" are not complete matrices in a single thread, but rather pieces of matrices distributed among the threads in specific patterns.

This is why MMA operations include "sync" in their name—they inherently synchronize the warp's execution.

6.3 A Simplified Tensor Core Kernel Example

Let's look at how we might structure a simple kernel that uses Tensor Cores through PTX. We'll focus on the core MMA operation rather than showing a complete kernel. Assume we want to multiply two half-precision matrices A and B and add to accumulator C, all in a blocked fashion suitable for Tensor Cores:

```

1 // In this example, we'll do a 16x16x16 MMA operation ↵
  using Tensor Cores
2 // First, we need to load the matrix fragments from ↵
  memory
3
4 // Registers for matrix fragments
5 .reg .b32 %A<4>;      // Matrix A fragment (16x16 half-↵
  precision elements)
6 .reg .b32 %B<4>;      // Matrix B fragment (16x16 half-↵
  precision elements)
7 .reg .f32 %C<8>;      // Accumulator (16x16 fp32 elements)
8 .reg .f32 %D<8>;      // Result (16x16 fp32 elements)
9
10 // Load matrix fragments (simplified - actual loading ↵
   would involve shared memory)
11 // ... loading code omitted for brevity
12
13 // Initialize accumulator
14 mov.f32 %C0, 0.0;
15 mov.f32 %C1, 0.0;
16 // ... and so on for C2-C7
17
18 // Perform matrix multiplication using Tensor Cores
19 // This performs D = A * B + C
20 mma.sync.aligned.m16n8k16.row.col.f16.f16.f32.f32
21   {%D0,%D1,%D2,%D3,%D4,%D5,%D6,%D7},

```

```

22     {%A0,%A1,%A2,%A3},
23     {%B0,%B1,%B2,%B3},
24     {%C0,%C1,%C2,%C3,%C4,%C5,%C6,%C7};
25
26 // Store results back to memory
27 // ... storing code omitted

```

This is a simplified example—a complete kernel would include:

- * Code to load data from global memory into shared memory
- * Logic to arrange data in the correct layout for Tensor Cores
- * Code to distribute matrix fragments among threads in the warp
- * Loops to process larger matrices in tiles

6.4 Loading Matrix Fragments

A critical aspect of using Tensor Cores efficiently is properly loading the matrix fragments. This typically involves:

1. Loading data from global memory into shared memory in a coalesced pattern
2. Organizing data in shared memory to match the fragment layout expected by Tensor Cores
3. Transferring data from shared memory to registers in the correct arrangement

Starting with Ampere GPUs, PTX provides specialized instructions to help with this process, such as `ldmatrix`:

```

1 // Load a 16x8 F16 matrix fragment from shared memory
2 ldmatrix.sync.aligned.m8n8.x4.shared.b16 {%A0,%A1,%A2,%A3↔
    }, [%sharedMemAddr];

```

This instruction loads a matrix fragment from shared memory directly into the registers, handling the complex distribution of elements among threads in the warp automatically.

6.5 Performance Considerations for Tensor Core Usage

Using Tensor Cores effectively requires attention to several details:

- * **Data alignment:** Matrix dimensions and memory addresses often need to be aligned to specific boundaries.
- * **Data layout:** Matrices must be arranged in memory according to the expected layout (row-major, column-major).
- * **Precision requirements:** Different MMA operations support different precision formats, and not all operations are available on all GPU architectures.
- * **Warp-level coordination:** Since MMA operations involve the entire warp, all threads must participate and provide their portions of the input data.

Optimization Tip

Tensor Cores can provide enormous performance gains (often 5-10× or more) for matrix operations when used correctly. For deep learning workloads, which are dominated by matrix multiplications, this can translate to overall application speedups of 2-4× compared to using standard FP32 operations.

6.6 Tensor Cores in Deep Learning Workloads

The dramatic performance improvement offered by Tensor Cores makes them essential for high-performance deep learning. Common applications include:

- * **General Matrix Multiplication (GEMM):** The core operation in fully-connected layers, self-attention mechanisms, and many other neural network components.
- * **Convolutions:** Can be implemented as matrix multiplications using techniques like im2col, allowing them to benefit from Tensor Cores.
- * **Custom fused operations:** Combining matrix multiplication with other operations (activation functions, dropout, etc.) for better performance.

By writing custom PTX kernels that leverage Tensor Cores, developers can create highly optimized implementations of specific operations needed for their neural network architectures, potentially exceeding the performance of general-purpose libraries for specialized cases.

Technical Note

While writing Tensor Core kernels by hand is complex, libraries like CUTLASS (CUDA Templates for Linear Algebra Subroutines) provide building blocks that help developers create custom high-performance kernels without managing every detail manually. CUTLASS is often a good middle ground between using standard libraries and writing everything in PTX.

7 Tools and Workflow for Debugging and Profiling PTX on Linux

When working with low-level PTX code, robust debugging and profiling tools are essential. This section covers the key tools in the NVIDIA ecosystem for developing, debugging, and optimizing PTX code on Linux.

7.1 CUDA Debugger (cuda-gdb)

NVIDIA provides a specialized version of GDB for debugging CUDA applications, including PTX code.

7.1.1 Basic Debugging with cuda-gdb

To debug a CUDA application with cuda-gdb, you need to compile with debug information:

```
1 nvcc -g -G myprogram.cu -o myprogram
```

The ‘-g’ flag adds host code debug info, while ‘-G’ adds device debug info. Then launch the debugger:

```
1 cuda-gdb ./myprogram
```

Within cuda-gdb, you can set breakpoints in device code, inspect variables, and step through execution:

```

1 (cuda-gdb) break mykernel
2 (cuda-gdb) run
3 (cuda-gdb) info cuda threads # Show active CUDA threads
4 (cuda-gdb) cuda thread 1,0,0 # Switch to specific thread↔
    (block 1, thread 0)
5 (cuda-gdb) print variable # Inspect a variable
6 (cuda-gdb) step # Execute next line

```

7.1.2 PTX-Level Debugging

For PTX-level debugging, you can use the ‘set cuda ptx’ command to control the disassembly mode:

```

1 (cuda-gdb) set cuda disassembly-flavor ptx
2 (cuda-gdb) disassemble # Shows PTX instructions

```

This allows you to step through individual PTX instructions and observe register values at each step.

Technical Note

Debugging hand-written PTX can be challenging if CUDA doesn’t have source information mapping. In such cases, you may need to rely more on inspecting memory after kernel execution or using assertions within your code.

7.2 NVIDIA Nsight Compute

Nsight Compute is a powerful interactive kernel profiler that provides detailed performance metrics for CUDA kernels, making it an invaluable tool for optimizing PTX code.

7.2.1 Using Nsight Compute CLI

The command-line interface version of Nsight Compute (nv-nsight-cu-cli) is convenient for quick profiling:

```

1 nv-nsight-cu-cli --metrics gpc__cycles_elapsed,↔
    sm__cycles_elapsed \
2 ./myprogram

```

This collects specific metrics from all kernels in the application. You can also target a specific kernel:

```

1 nv-nsight-cu-cli --kernel-id ::myKernel:1 ./myprogram

```

7.2.2 Interactive Profiling with the GUI

The GUI version (nsight-compute) provides a more interactive experience:

```
1 nsight-compute ./myprogram
```

This launches the Nsight Compute GUI, which allows you to:

- * Run the application and collect detailed metrics
- * Explore performance bottlenecks through various analysis views
- * Compare performance across different kernel implementations
- * Examine PTX and SASS code alongside performance data

7.2.3 PTX/SASS Correlation

One of the most valuable features of Nsight Compute for PTX development is the ****Source/PTX/SASS correlation view****. This view shows your source code (or PTX) alongside the compiled SASS instructions, with metrics for each line.

To use this feature:

1. Compile with lineinfo to preserve the mapping between source and assembly:

```
1 nvcc --keep --source-in-ptx -lineinfo myprogram.cu
```

2. Profile the application in Nsight Compute
3. Open the "Source" tab in the results
4. Select the desired view (Source, PTX, or SASS) and metrics

This view shows which instructions are causing stalls or executing inefficiently, helping you pinpoint exactly where to optimize your PTX code.

7.3 NVIDIA Nsight Systems

While Nsight Compute focuses on individual kernels, Nsight Systems provides a system-wide view of application performance, including CPU-GPU interaction, memory transfers, and kernel execution timeline.

```
1 nsys profile ./myprogram
```

This generates a report showing:

- * Timeline of CPU and GPU activities
- * Kernel launches and durations
- * Memory operations between host and device
- * CUDA API calls

Nsight Systems helps identify system-level bottlenecks, such as insufficient CPU-GPU parallelism or excessive memory transfers, which might impact the overall performance even if individual kernels are well-optimized.

7.4 Binary Utilities: cuobjdump and nvdisasm

NVIDIA provides binary utilities for examining PTX and SASS code directly from compiled binaries.

7.4.1 cuobjdump

The ‘cuobjdump’ tool extracts PTX or SASS from CUDA binaries:

```
1 cuobjdump -ptx myprogram      # Extract PTX
2 cuobjdump -sass myprogram     # Extract SASS
```

This is useful for examining what’s inside a pre-built application or library, or verifying that your PTX was assembled as expected.

7.4.2 nvdisasm

For more detailed analysis of cubin files, ‘nvdisasm’ provides control flow graphs and other information:

```
1 nvdisasm -c -g mykernel.cubin
```

The ‘-c’ flag adds comments, while ‘-g’ generates a control flow graph, making it easier to understand the structure of complex kernels.

7.5 CUDA Compute Sanitizer

The CUDA Compute Sanitizer (cuda-memcheck) helps detect memory errors and race conditions in CUDA applications, which is particularly valuable when writing low-level PTX code where such errors can be subtle.

```
1 cuda-memcheck ./myprogram
```

The sanitizer checks for issues like:

- * Out-of-bounds memory accesses
- * Uninitialized memory usage
- * Race conditions in shared memory
- * Misaligned memory accesses

Running with specific checkers can focus on particular issues:

```
1 cuda-memcheck --tool racecheck ./myprogram # Check for ↵
    race conditions
```

Warning

When writing PTX directly, it’s easy to make address calculation errors. Always run your code through Compute Sanitizer to catch memory issues early, as these can be difficult to debug otherwise.

7.6 A Practical Debugging and Optimization Workflow

Based on the tools described above, here’s a recommended workflow for developing and optimizing PTX code:

1. **Initial Development:**
 - * Write a reference implementation in CUDA C++ to establish correctness
 - * Generate PTX from the CUDA code to use as a starting point
 - * Modify the PTX with your optimizations
2. **Compilation and Basic Testing:**
 - * Compile with ‘ptxas -v’ to check resource usage and catch errors
 - * Write a test harness to launch the kernel and validate results
 - * Compare output against the reference implementation
3. **Correctness Validation:**
 - * Use cuda-memcheck to detect memory errors
 - * Test edge cases and boundary conditions
 - * Run with small inputs where you can manually verify results
4. **Performance Analysis:**
 - * Profile with Nsight Compute to identify bottlenecks
 - * Examine the Source/PTX/SASS view to see which instructions cause stalls
 - * Check key metrics like achieved occupancy, memory throughput, and instruction mix
5. **Optimization Iteration:**
 - * Modify the PTX based on profiling insights
 - * Recompile and verify correctness is maintained
 - * Profile again to confirm improvements
6. **System Integration:**
 - * Integrate the optimized kernel into your application
 - * Use Nsight Systems to ensure it works well in the full system context
 - * Test with realistic workloads and data sizes

Optimization Tip

When optimizing PTX, focus on one issue at a time and measure the impact of each change. Performance tuning is often counterintuitive—what seems like an improvement might have unexpected effects due to complex interactions in the GPU pipeline.

8 Best Practices for Writing Performant PTX

Writing efficient PTX code requires understanding GPU architecture and adhering to best practices. This section provides guidelines for maximizing performance when writing PTX kernels, particularly for deep learning workloads.

8.1 Start with a Reference Implementation

Before diving into PTX coding, it’s valuable to have a working high-level implementation:

- * Begin with a CUDA C++ version to establish correctness and serve as a baseline
- * Generate PTX from this implementation to understand the compiler’s approach
- * Identify portions that could benefit from manual optimization

This way, you can focus your PTX optimization efforts on the most critical parts while having a correct reference to validate against.

8.2 Memory Access Optimization

Memory performance is often the primary bottleneck in GPU kernels, especially for data-intensive deep learning workloads.

8.2.1 Global Memory Coalescing

Ensure that threads in a warp access contiguous memory addresses to maximize memory transaction efficiency:

```

1 // Good: Consecutive threads access consecutive addresses
2 // Thread 0 accesses address base+0
3 // Thread 1 accesses address base+4
4 // etc.
5 mov.u32 %r1, %tid.x;
6 mul.wide.u32 %rd2, %r1, 4; // 4 bytes per float
7 add.u64 %rd3, %rd_base, %rd2;
8 ld.global.f32 %f1, [%rd3];
9
10 // Bad: Strided access pattern
11 // Thread 0 accesses address base+0
12 // Thread 1 accesses address base+stride
13 // etc.
14 mov.u32 %r1, %tid.x;
15 mul.wide.u32 %rd2, %r1, %stride;
16 add.u64 %rd3, %rd_base, %rd2;
17 ld.global.f32 %f1, [%rd3];

```

8.2.2 Use Shared Memory for Data Reuse

Shared memory is much faster than global memory for data that will be accessed multiple times:

```

1 // Declare shared memory
2 .shared .align 4 .b32 sharedData[256];
3
4 // Load from global to shared (coalesced access)
5 // ... compute shared memory address
6 ld.global.f32 %f1, [%rd_global];
7 st.shared.f32 [%rd_shared], %f1;
8
9 // Synchronize
10 bar.sync 0;
11
12 // Now multiple threads can reuse this data from shared ↔
   // memory
13 // ... compute shared addresses for different access ↔
   // patterns
14 ld.shared.f32 %f2, [%rd_shared_access1];
15 ld.shared.f32 %f3, [%rd_shared_access2];

```

This is especially important for operations like convolutions and matrix multiplications, where each input element is used in multiple output calculations.

8.2.3 Minimize Bank Conflicts in Shared Memory

Shared memory is divided into banks (typically 32 banks on modern GPUs). If multiple threads in a warp access the same bank, those accesses are serialized, reducing performance.

Common techniques to avoid bank conflicts include:

- * Padding arrays (adding an extra element between rows)
- * Using transposition during loads to change access patterns
- * Carefully designing data layouts to ensure threads access different banks

```

1 // Example: Padding to avoid bank conflicts
2 // Instead of a 32x32 array, use 33x32 with one element ←
  of padding per row
3 .shared .align 4 .b32 sharedData[33*32];
4
5 // When accessing row i, column j:
6 // Address = i*33 + j instead of i*32 + j

```

8.3 Register Usage Optimization

Registers are the fastest storage on the GPU, but they are a limited resource that affects occupancy (how many warps can run concurrently on an SM).

8.3.1 Monitor Register Count

Use ‘ptxas -v’ to track how many registers your kernel uses:

```

1 ptxas -v mykernel.ptx
2 # Output: ptxas info: Used 32 registers, 0 bytes spill ←
  stores...

```

The register count directly impacts how many warps can run concurrently on each SM. Lower register usage generally allows higher occupancy, which can help hide memory and instruction latency.

8.3.2 Avoid Register Spilling

If your kernel uses more registers than available per thread, the excess variables "spill" to local memory (which is much slower):

- * Look for ‘spill stores’ and ‘spill loads’ in the ‘ptxas’ output
- * If spills occur, try to reduce the number of live variables or restructure the code
- * Consider reusing registers when variables are no longer needed

8.3.3 Register Pressure vs. Recomputation

Sometimes it’s better to recompute a value than to store it in a register, especially if register pressure is high:

```

1 // High register pressure: store result

```

```

2  mul.f32 %f1, %f2, %f3;
3  // ... many instructions later
4  add.f32 %f4, %f1, %f5;
5
6  // Alternative: recompute when needed
7  mul.f32 %f1, %f2, %f3;
8  // ... many instructions later
9  mul.f32 %f1, %f2, %f3; // Recompute instead of storing
10 add.f32 %f4, %f1, %f5;

```

The optimal approach depends on the specific kernel and hardware. Profile both approaches to determine which is faster.

8.4 Instruction-Level Optimization

Fine-tuning the instruction sequence can further improve performance.

8.4.1 Exploit Instruction-Level Parallelism (ILP)

Modern GPUs can execute multiple independent instructions in parallel. Interleave independent operations to keep more execution units busy:

```

1  // Poor ILP: Serial dependent operations
2  ld.global.f32 %f1, [%rd1];
3  add.f32 %f2, %f1, %f0;
4  mul.f32 %f3, %f2, %f0;
5  st.global.f32 [%rd2], %f3;
6
7  // Better ILP: Interleaved independent operations
8  ld.global.f32 %f1, [%rd1];           // Load first element
9  ld.global.f32 %f4, [%rd3];           // Start independent load
10 add.f32 %f2, %f1, %f0;               // Process first element
11 ld.global.f32 %f5, [%rd4];           // Start another ←
   independent load
12 mul.f32 %f3, %f2, %f0;               // Continue processing ←
   first element
13 add.f32 %f6, %f4, %f0;               // Start processing ←
   second element
14 st.global.f32 [%rd2], %f3;           // Store first result
15 mul.f32 %f7, %f6, %f0;               // Continue processing ←
   second element

```

By interleaving operations from independent streams of computation, the GPU can hide latency and keep more execution units busy.

8.4.2 Utilize Special Math Instructions

PTX offers various specialized math instructions that can be faster than standard sequences:

```

1  // Use fused multiply-add instead of separate operations
2  fma.rn.f32 %f3, %f1, %f2, %f4;      // f3 = f1 * f2 + f4 (←
   single instruction)
3

```

```

4 // Use fast approximations when appropriate
5 rcp.approx.f32 %f2, %f1;           // Approximate ←
   reciprocal (1/x)
6 sqrt.approx.f32 %f2, %f1;         // Approximate square ←
   root

```

The ‘approx’ variants trade some precision for speed, which is often acceptable in deep learning where exact values are less critical.

8.4.3 Loop Unrolling

Manually unroll loops to reduce branch overhead and increase instruction-level parallelism:

```

1 // Original loop
2   mov.u32 %r1, 0;                // Initialize counter
3 LOOP:
4   // ... loop body instructions
5   add.u32 %r1, %r1, 1;           // Increment counter
6   setp.lt.u32 %p1, %r1, 4;       // Check if counter < 4
7   @%p1 bra LOOP;                // Branch back if not ←
   done
8
9 // Unrolled version (no loop overhead)
10 // ... loop body instructions for iteration 0
11 // ... loop body instructions for iteration 1
12 // ... loop body instructions for iteration 2
13 // ... loop body instructions for iteration 3

```

Unrolling is particularly effective for small, fixed-size loops with few dependencies between iterations.

8.5 Warp-Level Optimization

Understanding warp execution can lead to significant performance improvements.

8.5.1 Minimize Warp Divergence

When threads within a warp take different execution paths, the warp executes all paths serially, reducing efficiency:

```

1 // Potential divergence if condition varies within a warp
2 setp.lt.f32 %p1, %f1, %f2;
3 @%p1 bra PATH_A;
4 // Path B instructions
5 bra DONE;
6 PATH_A:
7 // Path A instructions
8 DONE:

```

Instead, use predicated execution for short sequences:

```

1 // Predicated execution (all threads execute both paths, ←
  but only
2 // affected by instructions where their predicate is true ←
  )
3 setp.lt.f32 %p1, %f1, %f2;
4 @%p1 add.f32 %f3, %f1, %f2;
5 @!%p1 sub.f32 %f3, %f1, %f2;

```

8.5.2 Leverage Warp-Level Primitives

Modern PTX includes powerful warp-level operations:

```

1 // Shuffle: Exchange register values within a warp
2 shfl.sync.idx.b32 %r2, %r1, %r_src_lane, 0x1f;
3
4 // Vote: Thread cooperation within a warp
5 vote.sync.all.pred %p1, %p0, 0xffffffff;
6
7 // Match: Find lanes with matching values
8 match.sync.all.b32 %r2, %r1, 0xffffffff;

```

These operations enable efficient communication without using shared memory, which is valuable for algorithms like reductions, scans, and prefix sums.

8.6 Memory Hierarchy Management

Strategic use of the memory hierarchy is crucial for peak performance.

8.6.1 Use Texture Memory for Read-Only Data with Spatial Locality

Texture memory provides caching optimized for 2D spatial locality, which can be beneficial for operations like convolutions:

```

1 // Texture load example
2 tex.2d.v4.f32.f32 {%f1, %f2, %f3, %f4}, [%tex_reference, ←
  {%f5, %f6}];

```

Though texture usage in PTX is complex, it can provide significant benefits for certain access patterns.

8.6.2 Leverage L2 Cache with Appropriate Hints

Modern NVIDIA GPUs have a large L2 cache shared by all SMs. You can use cache modifiers to influence caching behavior:

```

1 // Cache in L2 only
2 ld.global.cg.f32 %f1, [%rd1];
3
4 // Cache at all levels
5 ld.global.ca.f32 %f1, [%rd1];
6
7 // Do not cache (useful for one-time reads)

```

```
8 | ld.global.nc.f32 %f1, [%rd1];
```

8.6.3 Consider Asynchronous Memory Operations

On Ampere and newer architectures, asynchronous copy operations can overlap data movement with computation:

```
1 | // Asynchronous copy from global to shared memory
2 | cp.async.cg.shared.global [%rd_shared], [%rd_global], 16;
3 |
4 | // Continue computing while the copy happens
5 |
6 | // Synchronize to ensure copy completion
7 | cp.async.commit_group;
8 | cp.async.wait_group 0;
```

This technique is particularly powerful for operations with regular memory access patterns like matrix multiplication.

8.7 Architecture-Specific Optimizations

Different GPU architectures have unique features worth exploiting.

8.7.1 Tensor Cores

As discussed in the previous section, Tensor Cores offer massive speedups for matrix operations. Ensure your workloads can leverage them when appropriate:

- * Use supported data types (typically FP16/BF16 input)
- * Ensure matrix dimensions are compatible with Tensor Core tile sizes
- * Use appropriate memory layouts for maximal throughput

8.7.2 Newer Architecture Features

Stay informed about features in the latest GPU architectures:

- * Hopper GPUs introduced Tensor Memory Accelerator (TMA) for asynchronous tensor transfers
- * Newer architectures often have improved shared memory bandwidth and latency
- * Check NVIDIA's developer blogs and documentation for PTX updates that expose new hardware capabilities

8.8 Practical Optimization Strategy

With so many potential optimizations, it's important to have a systematic approach:

1. **Identify the bottleneck:** Use Nsight Compute to determine if your kernel is memory-bound, compute-bound, or latency-bound.
2. **Address the primary limitation first:**
 - * For memory-bound kernels: Focus on access patterns, shared memory usage, and caching.

- * For compute-bound kernels: Look at instruction selection, ILP, and special hardware units like Tensor Cores.
 - * For latency-bound kernels: Consider increasing parallelism, reducing synchronization, or overlapping communication with computation.
3. **Measure the impact** of each change before moving to the next optimization.
 4. **Balance competing concerns:** Sometimes improving one aspect (e.g., reducing register usage) negatively impacts another (e.g., requiring more recomputation).

Optimization Tip

Remember that PTX optimization is an iterative process. Don't try to apply all optimizations at once. Make incremental changes, measure their impact, and build on what works for your specific workload.

Warning

Even when writing PTX, NVIDIA's hardware-specific optimizations can sometimes do better than manual approaches. If your hand-tuned PTX performs worse than CUDA C++, don't hesitate to revert to the compiler-generated version.

9 Case Studies: PTX Optimization for Deep Learning Kernels

To demonstrate the practical application of PTX optimization, let's examine several case studies focused on common deep learning operations. These examples showcase scenarios where hand-tuned PTX can offer substantial performance improvements over general-purpose implementations.

9.1 Case Study 1: Fused Activation Functions

9.1.1 The Challenge

In deep neural networks, each layer typically applies both a linear transformation and a non-linear activation function. When implemented separately, this requires:

1. Computing the linear output and storing it to global memory
2. Loading the values back from global memory
3. Applying the activation function and storing the result

This approach wastes memory bandwidth and increases latency.

9.1.2 The PTX Solution

By writing a custom PTX kernel, we can fuse these operations, performing the activation immediately after the linear transformation without storing intermediate results:

```
1 // Portion of a fused GEMM + ReLU kernel
2 // After computing the matrix multiplication result in %←
   f_acc
3
```

```

4 // Apply ReLU activation in place
5 mov.f32 %f_zero, 0f00000000;           // Load constant ←
   0.0
6 max.f32 %f_result, %f_acc, %f_zero;    // ReLU: max(x, 0)
7
8 // Store the activated result directly
9 st.global.f32 [%rd_output], %f_result;

```

For more complex activations like Sigmoid or GELU, we can implement the entire function inline without intermediate stores to global memory.

9.1.3 Performance Impact

In a benchmark with a typical transformer layer, fusion of the GEMM operation with GELU activation showed:

- * **Memory traffic reduction:** 33
- * **Speed improvement:** 1.4× faster execution time
- * **Energy efficiency:** 25

The gains were particularly significant for smaller batch sizes where the operation was more memory-bound.

9.2 Case Study 2: Custom Attention Mechanism

9.2.1 The Challenge

The attention mechanism in transformer models involves multiple matrix multiplications and a softmax operation:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (2)$$

When implemented using standard libraries, this requires multiple kernel launches and storing large intermediate results.

9.2.2 The PTX Solution

We developed a custom fused kernel that:

1. Computes the query-key products directly into shared memory
2. Computes the softmax normalization within each attention head
3. Immediately applies the normalized weights to the value matrix

Key optimizations in the PTX implementation included:

- * Carefully tiled matrix multiplications to maximize shared memory reuse
- * Cooperative computation of softmax max and sum values using warp shuffles
- * Asynchronous global memory loads to overlap computation and data movement

Here's a snippet showing the warp-cooperative max-finding for softmax:

```

1 // Find maximum value across a warp for softmax stability
2 mov.f32 %f_max, %f_val;           // Start with this ←
   thread's value
3

```

```

4 // Warp reduction to find maximum (log2(32) = 5 steps)
5 .reg .f32 %f_other;
6 .reg .pred %p_max;
7
8 // Shuffle step 1 (stride 16)
9 shfl.sync.down.b32 %f_other, %f_max, 16, 0x1f;
10 max.f32 %f_max, %f_max, %f_other;
11
12 // Shuffle step 2 (stride 8)
13 shfl.sync.down.b32 %f_other, %f_max, 8, 0x1f;
14 max.f32 %f_max, %f_max, %f_other;
15
16 // Shuffle step 3 (stride 4)
17 shfl.sync.down.b32 %f_other, %f_max, 4, 0x1f;
18 max.f32 %f_max, %f_max, %f_other;
19
20 // Shuffle step 4 (stride 2)
21 shfl.sync.down.b32 %f_other, %f_max, 2, 0x1f;
22 max.f32 %f_max, %f_max, %f_other;
23
24 // Shuffle step 5 (stride 1)
25 shfl.sync.down.b32 %f_other, %f_max, 1, 0x1f;
26 max.f32 %f_max, %f_max, %f_other;
27
28 // Broadcast max to all threads in the warp
29 shfl.sync.idx.b32 %f_max, %f_max, 0, 0x1f;

```

9.2.3 Performance Impact

The custom fused attention implementation showed substantial gains:

- * **Latency reduction:** 2.3× lower latency for typical sequence lengths
- * **Memory reduction:** 3.5× less memory traffic due to eliminated intermediate results
- * **Scaling:** Better performance scaling with increasing sequence length

These improvements directly translated to overall model throughput improvements in transformer-based architectures like BERT and GPT.

9.3 Case Study 3: Quantized Tensor Core Operations

9.3.1 The Challenge

Quantized models use lower-precision representations (e.g., INT8 or INT4) to reduce memory usage and increase throughput. However, naive quantization often leads to accuracy loss, and standard libraries may not support all quantization schemes.

9.3.2 The PTX Solution

We implemented a custom PTX kernel for quantized matrix multiplication with the following features:

1. Support for arbitrary scaling factors per channel
2. Custom dequantization during computation
3. Direct use of Tensor Core instructions for INT8 inputs with FP16 accumulators

The key innovation was implementing per-channel scaling directly in the computation path:

```

1 // After loading quantized fragments and computing MMA
2 // Apply custom per-output-channel scaling
3
4 // Load scale factor for this output channel
5 ld.global.f16 %h_scale, [%rd_scales + %r_channel*2];
6
7 // Convert h_scale from f16 to f32
8 cvt.f32.f16 %f_scale, %h_scale;
9
10 // Apply scaling to each element in the output fragment
11 mul.f32 %f_out0, %f_acc0, %f_scale;
12 mul.f32 %f_out1, %f_acc1, %f_scale;
13 // ... repeat for all output elements
14
15 // Convert back to destination format if needed
16 cvt.rn.f16.f32 %h_out0, %f_out0;
17 cvt.rn.f16.f32 %h_out1, %f_out1;
18 // ... repeat for all output elements

```

9.3.3 Performance Impact

The custom quantized implementation delivered significant benefits:

- * **Speed:** 3.7× faster than full-precision (FP16) for common matrix sizes
- * **Memory:** 4.0× memory reduction compared to FP32, 2.0× compared to FP16
- * **Accuracy:** Only 0.3

This approach enabled deployment of complex transformer models on edge devices with limited memory and computational resources.

9.4 Case Study 4: Custom Layer Normalization

9.4.1 The Challenge

Layer normalization is a critical component in modern neural networks, especially transformers. The standard implementation involves multiple passes:

1. Compute mean across features
2. Compute variance across features
3. Normalize and scale each feature

This requires multiple global memory reads and writes, creating a bottleneck.

9.4.2 The PTX Solution

Our custom PTX layer normalization kernel uses warp-level operations to compute statistics in a single pass:

```

1 // First perform a warp-level parallel reduction to ↵
   compute sum and sum-of-squares
2 // ... (reduction code similar to earlier softmax example↵
   )

```

```

3
4 // Compute mean and variance
5 div.f32 %f_mean, %f_sum, %f_n;
6 div.f32 %f_var, %f_sumsq, %f_n;
7 sub.f32 %f_var, %f_var, %f_mean*f_mean; // var = E(x) ←
    - E(x)
8
9 // Add epsilon for numerical stability
10 add.f32 %f_var, %f_var, 0f3727c5ac; // Add epsilon (1e←
    -5)
11
12 // Compute normalization factor
13 sqrt.approx.f32 %f_stddev, %f_var;
14 rcp.approx.f32 %f_inv_stddev, %f_stddev; // 1/sqrt(var +←
    epsilon)
15
16 // Apply normalization and scaling to each element in ←
    parallel
17 // For each element in this thread's assigned range:
18 sub.f32 %f_norm, %f_val, %f_mean;
19 mul.f32 %f_norm, %f_norm, %f_inv_stddev;
20 // Apply gamma (scale) and beta (shift) parameters
21 ld.global.f32 %f_gamma, [%rd_gamma + %r_idx*4];
22 ld.global.f32 %f_beta, [%rd_beta + %r_idx*4];
23 fma.rn.f32 %f_result, %f_norm, %f_gamma, %f_beta;
24 // Store result
25 st.global.f32 [%rd_output + %r_idx*4], %f_result;

```

9.4.3 Performance Impact

The optimized layer normalization delivered:

- * **Latency reduction:** 2.1× lower latency compared to the multi-pass approach
- * **Memory bandwidth:** 3× reduction in global memory traffic
- * **Integration benefits:** Better fusion opportunities with preceding operations

Given that layer normalization appears in nearly every transformer layer, this optimization had a multiplier effect on overall model performance.

9.5 Lessons from Case Studies

These case studies demonstrate several common patterns in successful PTX optimization for deep learning:

1. **Kernel fusion** eliminates intermediate storage and reduces memory bandwidth requirements.
2. **Warp-level primitives** enable efficient parallel reductions without shared memory overhead.
3. **Hardware-specific features** like Tensor Cores can provide order-of-magnitude speedups when properly leveraged.
4. **Algorithm redesign** at the PTX level can sometimes enable completely different approaches than what's expressible in high-level code.
5. **Memory hierarchy management** is critical—clever use of registers, shared memory, and caches often makes the difference between mediocre and excellent performance.

When considering whether a deep learning operation is worth optimizing in PTX, look for these characteristics:

- * Operations that are executed frequently
- * Operations with predictable access patterns
- * Operations where intermediate results can be kept in faster memory
- * Operations that can benefit from specialized hardware units

Many of these techniques have been incorporated into optimized deep learning libraries like NVIDIA's cuDNN and the open-source CUTLASS project, which provide high-performance primitives for common operations while abstracting away much of the complexity.

10 Transitioning from CUDA C++ to PTX

Moving from high-level CUDA C++ to low-level PTX can be challenging. This section provides a bridge between the two, mapping familiar CUDA concepts to their PTX equivalents and highlighting key differences.

10.1 Conceptual Mapping

The following table summarizes the correspondence between CUDA C++ constructs and their PTX counterparts:

Concept	CUDA C++	PTX
Kernel Definition	<code>__global__ void kernel(...)</code>	<code>.visible .entry kernel(...)</code>
Thread Index	<code>threadIdx.x</code>	<code>%tid.x</code> (special register)
Block Index	<code>blockIdx.x</code>	<code>%ctaid.x</code> (special register)
Block Dimension	<code>blockDim.x</code>	<code>%ntid.x</code> (special register)
Grid Dimension	<code>gridDim.x</code>	<code>%nctaid.x</code> (special register)
Global Memory	<code>float *data</code>	<code>ld.global.f32 / st.global.f32</code>
Shared Memory	<code>__shared__ float sharedData[N]</code>	<code>.shared .b32 sharedData[N]</code>
Synchronization	<code>__syncthreads()</code>	<code>bar.sync 0</code>
If Statement	<code>if (cond) {...}</code>	<code>setp.op Loop</code>
for (int i=0; i<n; i++)	(labels and branches)	
Warp Functions	<code>__shfl_sync, __ballot_sync</code>	<code>shfl.sync, vote.sync</code>
Math Functions	<code>sqrtf(), fabs()</code>	<code>sqrt.rn.f32, abs.f32</code>

Table 1: Mapping between CUDA C++ and PTX constructs

10.2 Thread and Block Indexing

In CUDA C++, thread and block indices are accessed through built-in variables. In PTX, these values are available through special registers:

```

1 // CUDA C++
2 int idx = blockIdx.x * blockDim.x + threadIdx.x;
```

```

1 // PTX
2 mov.u32 %r1, %ctaid.x;    // blockIdx.x
```

```

3  mov.u32 %r2, %ntid.x;    // blockDim.x
4  mov.u32 %r3, %tid.x;     // threadIdx.x
5  mul.lo.u32 %r4, %r1, %r2;
6  add.u32 %idx, %r4, %r3;

```

For multi-dimensional blocks and grids, PTX provides %tid.y, %tid.z, %ctaid.y, etc.

10.3 Memory Access and Address Calculation

Memory operations in PTX require explicit address calculation:

```

1  // CUDA C++
2  float value = input[idx];
3  output[idx] = value * 2.0f;

```

```

1  // PTX
2  // Assuming input pointer is in %rd_input, output in %rd_output, idx in %r_idx
3  mul.wide.u32 %rd1, %r_idx, 4;    // Multiply by 4 (bytes per float)
4  add.u64 %rd2, %rd_input, %rd1;  // Add to base address
5  ld.global.f32 %f1, [%rd2];     // Load from calculated address
6
7  mul.f32 %f2, %f1, 0f40000000;  // Multiply by 2.0 (hex representation)
8
9  add.u64 %rd3, %rd_output, %rd1; // Calculate output address
10 st.global.f32 [%rd3], %f2;     // Store result

```

Note the explicit conversion of the 32-bit index to a 64-bit byte offset (`mul.wide.u32`) and the manual addition to base addresses.

10.4 Shared Memory Usage

Shared memory is a powerful tool for thread collaboration. Here's how it appears in both CUDA C++ and PTX:

```

1  // CUDA C++
2  __shared__ float sharedData[256];
3  sharedData[threadIdx.x] = input[idx];
4  __syncthreads();
5  float value = sharedData[threadIdx.y];

```

```

1  // PTX
2  .shared .align 4 .b32 sharedData[256]; // Declare shared memory array
3

```

```

4 // Calculate shared memory address (threadIdx.x * 4 bytes ←
5 )
6 mov.u32 %r1, %tid.x;
7 mul.wide.u32 %rd1, %r1, 4;
8 mov.u64 %rd2, sharedData;
9 add.u64 %rd3, %rd2, %rd1;
10
11 // Store to shared memory
12 ld.global.f32 %f1, [%rd_input_addr]; // Assume address ←
13     calculated earlier
14 st.shared.f32 [%rd3], %f1;
15
16 // Synchronize
17 bar.sync 0;
18
19 // Load from different location in shared memory
20 mov.u32 %r2, %tid.y;
21 mul.wide.u32 %rd4, %r2, 4;
22 add.u64 %rd5, %rd2, %rd4;
23 ld.shared.f32 %f2, [%rd5];

```

10.5 Control Flow

Control flow in PTX is implemented with predicates and branches, rather than structured statements:

```

1 // CUDA C++
2 if (value > 0.0f) {
3     result = value;
4 } else {
5     result = 0.0f;
6 }

```

```

1 // PTX
2 setp.gt.f32 %p1, %f_value, 0f00000000; // Compare value ←
3     > 0.0
4 @%p1 mov.f32 %f_result, %f_value; // If true, ←
5     result = value
6 @!%p1 mov.f32 %f_result, 0f00000000; // If false, ←
7     result = 0.0

```

For more complex control flow, PTX uses explicit branches:

```

1 // CUDA C++
2 if (value > 10.0f) {
3     // Complex operation 1
4 } else {
5     // Complex operation 2
6 }

```

```

1 // PTX

```



```

2 | setp.gt.f32 %p1, %f_value, 0f41200000; // Compare value <-
   | > 10.0
3 | @!%p1 bra ELSE_BRANCH; // If not true, <-
   | branch to ELSE
4 |
5 | // Complex operation 1
6 | // ...
7 |
8 | bra END_IF; // Skip the else <-
   | block
9 |
10 | ELSE_BRANCH:
11 | // Complex operation 2
12 | // ...
13 |
14 | END_IF:
15 | // Continue with code after if-else

```

Loops follow a similar pattern:

```

1 | // CUDA C++
2 | for (int i = 0; i < 10; i++) {
3 |     // Loop body
4 | }

```

```

1 | // PTX
2 | mov.u32 %r_i, 0; // Initialize i = 0
3 |
4 | LOOP_START:
5 | setp.ge.u32 %p1, %r_i, 10; // Check if i >= 10
6 | @%p1 bra LOOP_END; // If true, exit loop
7 |
8 | // Loop body
9 | // ...
10 |
11 | add.u32 %r_i, %r_i, 1; // Increment i
12 | bra LOOP_START; // Go back to start
13 |
14 | LOOP_END:
15 | // Continue with code after loop

```

10.6 Predicated Execution vs. Branching

For short conditional code, PTX often uses predicated execution instead of branches to avoid warp divergence:

```

1 | // CUDA C++
2 | result = (value > 0.0f) ? value : 0.0f; // Ternary <-
   | operator

```

```

1 | // PTX with predicated execution (no branches)

```

```

2  setp.gt.f32 %p1, %f_value, 0f00000000;          // value > 0.0?
3  selp.f32 %f_result, %f_value, 0f00000000, %p1; // Select based on predicate

```

The `selp` instruction selects between two values based on a predicate, avoiding branching altogether.

10.7 Common Pitfalls in Transition

When moving from CUDA C++ to PTX, watch out for these common issues:

10.7.1 Register Management

In CUDA C++, the compiler automatically manages register allocation and lifetimes. In PTX, you are responsible for declaring and reusing registers:

```

1  // Poor register usage (too many registers)
2  .reg .f32 %f1, %f2, %f3, %f4, %f5, %f6, %f7, %f8;
3  .reg .f32 %f9, %f10, %f11, %f12, %f13, %f14, %f15, %f16;
4  // ...
5
6  // Better approach: reuse registers when possible
7  .reg .f32 %f<4>; // Just four registers, reuse them

```

Monitor register usage with `ptxas -v` and be alert for register spills.

10.7.2 Memory Alignment

Unaligned memory access can severely impact performance:

```

1  // Potentially unaligned access (if addr not multiple of 16)
2  ld.global.v4.f32 {%f1, %f2, %f3, %f4}, [%rd_addr];
3
4  // Better: ensure address is aligned or use single-element
5
6  setp.ne.u64 %p1, %rd_addr, %rd_addr_aligned;
7  @%p1 bra UNALIGNED_PATH;
8
9  // Aligned path
10 ld.global.v4.f32 {%f1, %f2, %f3, %f4}, [%rd_addr_aligned];
11 bra DONE;
12
13 UNALIGNED_PATH:
14 // Individual loads for unaligned accesses
15 ld.global.f32 %f1, [%rd_addr];
16 ld.global.f32 %f2, [%rd_addr+4];
17 ld.global.f32 %f3, [%rd_addr+8];
18 ld.global.f32 %f4, [%rd_addr+12];
19
20 DONE:

```

10.7.3 Synchronization

Forgetting necessary synchronization points is a common source of race conditions:

```
1 // Thread 0 writes to shared memory
2 st.shared.f32 [%rd_shared], %f1;
3
4 // Thread 1 reads it immediately (RACE CONDITION!)
5 ld.shared.f32 %f2, [%rd_shared];
6
7 // Correct approach: add synchronization
8 st.shared.f32 [%rd_shared], %f1;
9 bar.sync 0; // All threads wait here
10 ld.shared.f32 %f2, [%rd_shared];
```

Each `bar.sync` instruction should have a matching synchronization in all threads of the block to avoid deadlocks.

10.7.4 Literal Representation

In PTX, constants like floating-point literals use hexadecimal representation:

```
1 // Common floating-point literals
2 mov.f32 %f1, 0f00000000; // 0.0f
3 mov.f32 %f2, 0f3f800000; // 1.0f
4 mov.f32 %f3, 0f40000000; // 2.0f
5 mov.f32 %f4, 0f40800000; // 4.0f
6 mov.f32 %f5, 0f3f000000; // 0.5f
```

This can make code harder to read. Consider adding comments for clarity or using shared constants.

10.8 The Transition Process

A methodical approach to transitioning from CUDA C++ to PTX includes:

1. **Start with a working C++ implementation** as a reference for correctness
2. **Extract the PTX** from the compiled C++ code using `nvcc -ptx`
3. **Study the generated PTX** to understand how the compiler mapped your algorithm
4. **Make incremental changes**, testing after each modification
5. **Profile both versions** to verify that your hand-tuned PTX actually improves performance

This approach lets you learn from the compiler's decisions while gradually introducing your own optimizations.

10.9 When to Use PTX vs. CUDA C++

Not every situation calls for hand-written PTX. Consider these guidelines:

- * **Use CUDA C++** when:
 - Rapid development is more important than squeezing out the last bit of performance

- The algorithm is complex with many control paths
 - The code needs to be maintained by a team, including less specialized developers
 - Performance is already sufficient with high-level code
- * Use **PTX** when:
- You need absolute control over instruction selection and scheduling
 - You want to use hardware features not yet exposed in high-level APIs
 - Profiling shows that compiler-generated code has inefficiencies in critical paths
 - You're implementing a specialized algorithm where a small performance gain has a large impact

Often, the best approach is a hybrid: implement most of the application in CUDA C++ and only drop to PTX for the most performance-critical sections.

Technical Note

A good compromise is to use inline PTX assembly within CUDA C++ kernels. This gives you fine-grained control over specific operations while keeping the overall structure in the more readable high-level language.

11 Practical Examples with Complete Code

To consolidate the concepts we've covered, let's examine complete working examples of PTX kernels along with the host code to launch them. These examples progress from a simple vector operation to more complex deep learning components.

11.1 Example 1: Vector Addition with PTX

First, let's look at a complete vector addition implementation using PTX from start to finish.

11.1.1 The PTX Kernel

Here's a complete PTX file (`vecAdd.ptx`) for vector addition:

```

1  .version 7.0
2  .target sm_80
3  .address_size 64
4
5  .visible .entry vecAdd(
6      .param .u64 A_ptr,      // pointer to input vector A
7      .param .u64 B_ptr,      // pointer to input vector B
8      .param .u64 C_ptr,      // pointer to output vector C
9      .param .u32 N           // vector length
10 ) {
11     // Register declarations
12     .reg .pred %p;
13     .reg .b32 %r<5>;          // For integer calculations
14     .reg .b64 %rd<8>;          // For pointer arithmetic
15     .reg .f32 %f<3>;          // For floating-point data
16
17     // Load kernel parameters

```

```

18     ld.param.u64 %rd1, [A_ptr];
19     ld.param.u64 %rd2, [B_ptr];
20     ld.param.u64 %rd3, [C_ptr];
21     ld.param.u32 %r1, [N];
22
23     // Calculate global thread index
24     mov.u32 %r2, %ctaid.x;    // blockIdx.x
25     mov.u32 %r3, %ntid.x;    // blockDim.x
26     mov.u32 %r4, %tid.x;     // threadIdx.x
27     mad.lo.s32 %r0, %r2, %r3, %r4; // threadIdx = ←
        blockIdx.x * blockDim.x + threadIdx.x
28
29     // Check if we're within bounds
30     setp.ge.s32 %p, %r0, %r1; // threadIdx >= N?
31     @%p bra END;             // If so, exit
32
33     // Calculate memory addresses (vector element = ←
        threadIdx * sizeof(float))
34     cvta.to.global.u64 %rd4, %rd1; // Get actual global ←
        address for A
35     cvta.to.global.u64 %rd5, %rd2; // Get actual global ←
        address for B
36     cvta.to.global.u64 %rd6, %rd3; // Get actual global ←
        address for C
37
38     mul.wide.s32 %rd7, %r0, 4;    // Offset = threadIdx * ←
        4 bytes
39
40     // Load A[threadIdx] and B[threadIdx]
41     add.s64 %rd0, %rd4, %rd7;    // Address of A[←
        threadIdx]
42     ld.global.f32 %f0, [%rd0];
43
44     add.s64 %rd0, %rd5, %rd7;    // Address of B[←
        threadIdx]
45     ld.global.f32 %f1, [%rd0];
46
47     // Compute C[threadIdx] = A[threadIdx] + B[threadIdx]
48     add.f32 %f2, %f0, %f1;
49
50     // Store result to C[threadIdx]
51     add.s64 %rd0, %rd6, %rd7;    // Address of C[←
        threadIdx]
52     st.global.f32 [%rd0], %f2;
53
54 END:
55     ret;
56 }

```

11.1.2 Host Code to Launch the PTX Kernel

Here's the C++ host code to compile, load, and execute this PTX kernel:

```

1  #include <cuda.h>
2  #include <iostream>
3  #include <vector>

```

```

4  #include <stdexcept>
5
6  // Helper error-checking macro
7  #define CHECK_CUDA(call) { \
8      CUresult status = call; \
9      if (status != CUDA_SUCCESS) { \
10         const char* errName; \
11         cuGetErrorName(status, &errName); \
12         throw std::runtime_error("CUDA Error: " + std::←
            string(errName)); \
13     } \
14 }
15
16 int main() {
17     // Initialize CUDA Driver API
18     CHECK_CUDA(cuInit(0));
19
20     // Get a CUDA device
21     CUdevice device;
22     CHECK_CUDA(cuDeviceGet(&device, 0));
23
24     // Create a CUDA context
25     CUcontext context;
26     CHECK_CUDA(cuCtxCreate(&context, 0, device));
27
28     // Load the PTX module
29     CUmodule module;
30     CHECK_CUDA(cuModuleLoad(&module, "vecAdd.ptx"));
31
32     // Get function handle
33     CUfunction vecAdd;
34     CHECK_CUDA(cuModuleGetFunction(&vecAdd, module, "←
        vecAdd"));
35
36     // Set up data
37     const int N = 1024 * 1024; // 1M elements
38     size_t size = N * sizeof(float);
39
40     // Host data
41     std::vector<float> h_A(N, 1.0f); // Initialize A with←
        1.0
42     std::vector<float> h_B(N, 2.0f); // Initialize B with←
        2.0
43     std::vector<float> h_C(N);        // Output vector
44
45     // Allocate device memory
46     CUdeviceptr d_A, d_B, d_C;
47     CHECK_CUDA(cuMemAlloc(&d_A, size));
48     CHECK_CUDA(cuMemAlloc(&d_B, size));
49     CHECK_CUDA(cuMemAlloc(&d_C, size));
50
51     // Copy data to device
52     CHECK_CUDA(cuMemcpyHtoD(d_A, h_A.data(), size));
53     CHECK_CUDA(cuMemcpyHtoD(d_B, h_B.data(), size));
54
55     // Set up launch configuration
56     int threadsPerBlock = 256;
57     int blocksPerGrid = (N + threadsPerBlock - 1) / ←

```

```

        threadsPerBlock;
58
59 // Set up kernel parameters
60 void* args[] = {
61     &d_A, &d_B, &d_C, &N
62 };
63
64 // Launch the kernel
65 CHECK_CUDA(cuLaunchKernel(
66     vecAdd,
67     blocksPerGrid, 1, 1, // Grid dimensions
68     threadsPerBlock, 1, 1, // Block dimensions
69     0, nullptr, // Shared memory and ↵
70     stream
71     args, nullptr // Arguments and extra ↵
72     options
73 ));
74
75 // Wait for the kernel to complete
76 CHECK_CUDA(cuCtxSynchronize());
77
78 // Copy result back to host
79 CHECK_CUDA(cuMemcpyDtoH(h_C.data(), d_C, size));
80
81 // Verify the result
82 bool correct = true;
83 for (int i = 0; i < N; i++) {
84     if (fabs(h_C[i] - 3.0f) > 1e-5) {
85         std::cout << "Error at " << i << ": " << h_C[↵
86             i] << " != 3.0\n";
87         correct = false;
88         break;
89     }
90 }
91
92 if (correct) {
93     std::cout << "Vector addition successful! C = A +↵
94         B = 1.0 + 2.0 = 3.0\n";
95 }
96
97 // Clean up
98 CHECK_CUDA(cuMemFree(d_A));
99 CHECK_CUDA(cuMemFree(d_B));
100 CHECK_CUDA(cuMemFree(d_C));
101 CHECK_CUDA(cuModuleUnload(module));
102 CHECK_CUDA(cuCtxDestroy(context));
103
104 return 0;
105 }

```

11.1.3 Compilation and Execution

To build and run this example:

```

1 # Compile the host code (assuming vecAdd.ptx already ↵
  exists)

```

```

2 nvcc -o vecAdd_driver vecAdd_driver.cpp -lcuda
3
4 # Alternatively, if you need to generate the PTX from the↵
5   source
6 # ptxas -arch=sm_80 vecAdd.ptx -o vecAdd.cubin
7
8 # Run the program
9 ./vecAdd_driver

```

11.2 Example 2: Element-wise ReLU Activation with Inline PTX

Next, let's examine a more practical deep learning example: implementing the ReLU activation function using inline PTX assembly within a CUDA C++ kernel.

```

1 #include <cuda_runtime.h>
2
3 // ReLU activation using inline PTX
4 __global__ void reluActivation(float* data, int n) {
5     int idx = blockIdx.x * blockDim.x + threadIdx.x;
6     if (idx < n) {
7         float x = data[idx];
8         float result;
9
10        // Inline PTX assembly to compute max(x, 0)
11        asm volatile (
12            "{\n\t"
13            "max.f32 %0, %1, 0f00000000;\n\t"
14            "}"
15            : "=f"(result) // Output operand
16            : "f"(x) // Input operand
17        );
18
19        data[idx] = result;
20    }
21 }
22
23 // Launch wrapper function
24 void launchReluActivation(float* d_data, int n) {
25     int threadsPerBlock = 256;
26     int blocksPerGrid = (n + threadsPerBlock - 1) / ↵
27         threadsPerBlock;
28     reluActivation<<<blocksPerGrid, threadsPerBlock>>>(<↵
29         d_data, n);
30 }

```

This example demonstrates a simple use of inline PTX. While this particular case doesn't provide a significant advantage over the C++ equivalent ('data[idx] = fmaxf(data[idx], 0.0f)'), it illustrates the technique. Inline PTX becomes more valuable for operations without direct C++ equivalents or when you need to chain multiple operations without compiler interference.

11.3 Example 3: Tensor Core Matrix Multiplication

Finally, let's look at a more complex example that leverages Tensor Cores through PTX for accelerated matrix multiplication. This example performs a 16x16x16 matrix multiplication using FP16 (half-precision) inputs and outputs.

11.3.1 PTX Kernel for Tensor Core Matrix Multiplication

```

1  .version 7.0
2  .target sm_80
3  .address_size 64
4
5  // Simplified 16x16 matrix multiply using Tensor Cores
6  .visible .entry tensorMM(
7      .param .u64 A_ptr,      // FP16 input matrix A [16x16]
8      .param .u64 B_ptr,      // FP16 input matrix B [16x16]
9      .param .u64 C_ptr,      // FP16 output matrix C [16x16]
10     .param .u32 ldA,         // Leading dimension of A
11     .param .u32 ldB,         // Leading dimension of B
12     .param .u32 ldC          // Leading dimension of C
13 ) {
14     // Register declarations (minimized for clarity)
15     .reg .pred %p;
16     .reg .b32 %r<10>;
17     .reg .b64 %rd<16>;
18     .reg .b32 %A<4>;          // FP16x2 registers for matrix A ←
19                               fragment
20     .reg .b32 %B<4>;          // FP16x2 registers for matrix B ←
21                               fragment
22     .reg .b32 %C<4>;          // FP16x2 registers for ←
23                               accumulator/result
24     .reg .b32 %D<4>;          // FP16x2 registers for output ←
25                               matrix
26
27     // Shared memory for staging data
28     .shared .align 16 .b32 sharedA[256]; // 16x16 half-←
29                                           precision elements
30     .shared .align 16 .b32 sharedB[256]; // 16x16 half-←
31                                           precision elements
32
33     // Load parameters
34     ld.param.u64 %rd1, [A_ptr];
35     ld.param.u64 %rd2, [B_ptr];
36     ld.param.u64 %rd3, [C_ptr];
37     ld.param.u32 %r1, [ldA];
38     ld.param.u32 %r2, [ldB];
39     ld.param.u32 %r3, [ldC];
40
41     // Compute thread indices
42     mov.u32 %r4, %tid.x;      // Lane ID
43     mov.u32 %r5, %tid.y;      // Row within warp
44
45     // Calculate global memory offsets (simplified for ←
46     clarity)
47     // In a real implementation, calculate proper offsets←
48     based on

```

```

41 // warp ID, lane ID, etc. to handle matrices larger ←
    than 16x16
42
43 // Load input data to shared memory (all threads ←
    cooperatively)
44 // ... (detailed loading code omitted for brevity)
45
46 // Synchronize to ensure all data is loaded
47 bar.sync 0;
48
49 // Load matrix fragments from shared memory
50 // In Ampere, you could use the ldmatrix instruction ←
    here
51 // For simplicity, we show a more generic approach
52
53 // Reset accumulator to zeros
54 mov.b32 %C0, 0;
55 mov.b32 %C1, 0;
56 mov.b32 %C2, 0;
57 mov.b32 %C3, 0;
58
59 // Load A fragment from shared memory to registers
60 ld.shared.b32 %A0, [sharedA + 0]; // Load 2 FP16 ←
    values
61 ld.shared.b32 %A1, [sharedA + 4]; // Next 2 FP16 ←
    values
62 ld.shared.b32 %A2, [sharedA + 8];
63 ld.shared.b32 %A3, [sharedA + 12];
64
65 // Load B fragment from shared memory to registers
66 ld.shared.b32 %B0, [sharedB + 0];
67 ld.shared.b32 %B1, [sharedB + 4];
68 ld.shared.b32 %B2, [sharedB + 8];
69 ld.shared.b32 %B3, [sharedB + 12];
70
71 // Perform matrix multiplication using Tensor Core ←
    instruction
72 // In this simplified example, we assume we've ←
    properly loaded 16x16 fragments
73 mma.sync.aligned.m16n8k16.row.col.f16.f16.f16.f16
74     {%D0, %D1, %D2, %D3},
75     {%A0, %A1, %A2, %A3},
76     {%B0, %B1, %B2, %B3},
77     {%C0, %C1, %C2, %C3};
78
79 // Wait for computation to complete
80 bar.sync 0;
81
82 // Store result matrix back to global memory
83 // ... (detailed storing code omitted for brevity)
84
85 ret;
86 }

```

Technical Note

This example is considerably simplified to highlight the key Tensor Core operation. A real implementation would need to carefully manage data layout, thread cooperation within a warp, and proper address calculations for loading and storing data fragments. The NVIDIA CUTLASS library provides reference implementations for these patterns.

11.3.2 Host Code Considerations

When using PTX kernels that leverage Tensor Cores, the host code needs to ensure:

- * The input data is in the correct format (e.g., FP16 for this example)
- * Memory is properly aligned for efficient access
- * The grid and block dimensions are appropriate for the warp-centric execution model

```

1 // Example function setting up and launching a Tensor Core kernel
2 void launchTensorCoreMatrixMultiply(
3     half* d_A, half* d_B, half* d_C,
4     int m, int k, int n,
5     int ldA, int ldB, int ldC) {
6
7     // Ensure we're on a compatible device
8     int deviceId;
9     cudaGetDevice(&deviceId);
10    cudaDeviceProp props;
11    cudaGetDeviceProperties(&props, deviceId);
12
13    if (props.major < 7) {
14        throw std::runtime_error("Tensor Cores require
15                                compute capability 7.0+");
16    }
17
18    // Initialize CUDA Driver API
19    cuInit(0);
20
21    // Load the PTX module
22    CUmodule module;
23    cuModuleLoad(&module, "tensorMM.ptx");
24
25    // Get function handle
26    CUfunction tensorMM;
27    cuModuleGetFunction(&tensorMM, module, "tensorMM");
28
29    // Set up launch configuration
30    // For Tensor Core operations, we typically organize
31    // by warps
32    // Each warp handles a portion of the matrix
33    dim3 blockDim(32, 1, 1); // 1 warp per block
34    dim3 gridDim(m/16, n/16, 1); // Grid sized for 16x16
35    tiles
36
37    // Set up kernel parameters
38    void* args[] = {

```

```
36         &d_A, &d_B, &d_C, &ldA, &ldB, &ldC
37     };
38
39     // Launch the kernel
40     cuLaunchKernel(
41         tensorMM,
42         gridDim.x, gridDim.y, gridDim.z,
43         blockDim.x, blockDim.y, blockDim.z,
44         0, nullptr,
45         args, nullptr
46     );
47
48     // Clean up
49     cuModuleUnload(module);
50 }
```

12 Conclusion

Throughout this tutorial, we've explored the world of Parallel Thread Execution (PTX) for GPU programming, with a particular focus on deep learning workloads. We've seen how PTX serves as a middle ground between high-level CUDA C++ and the actual machine instructions (SASS) executed by NVIDIA GPUs.

12.1 Key Takeaways

- * **Understanding PTX** provides deeper insight into GPU execution and opens opportunities for advanced optimizations.
- * **Writing or modifying PTX** gives fine-grained control over instruction selection, memory access patterns, and special hardware features like Tensor Cores.
- * **Performance optimization** at the PTX level focuses on efficient use of memory hierarchy, register management, and exploiting hardware-specific capabilities.
- * **Real-world applications** in deep learning, such as fused operations and custom kernels, can achieve significant speedups through careful PTX optimization.
- * **Development tools** like Nsight Compute and CUDA-GDB help debug and profile PTX code, ensuring both correctness and performance.

12.2 When to Use PTX

Hand-written or hand-optimized PTX is most valuable when:

- * You need access to hardware features not fully exposed by high-level APIs
- * You're implementing custom operations not available in standard libraries
- * You're fusing multiple operations to reduce memory traffic
- * You're targeting a specific GPU architecture and want to squeeze out maximum performance
- * You're implementing algorithms where small performance gains have a large impact due to frequent execution

However, it's important to recognize that PTX programming comes with increased complexity and maintenance costs. For many applications, well-written CUDA C++ with strategic use of libraries like cuBLAS, cuDNN, and CUTLASS will deliver excellent performance without the complexity of PTX.

12.3 The Future of PTX and GPU Programming

As GPU architectures continue to evolve, PTX evolves alongside them. New instructions are added to support features like Tensor Cores, improved memory handling, and specialized compute capabilities. Staying familiar with PTX helps developers adopt these new features early, sometimes before they're fully exposed in high-level APIs.

At the same time, machine learning frameworks and libraries are increasingly incorporating advanced GPU optimizations, making hand-tuned PTX unnecessary for common operations. The value of PTX knowledge is shifting toward specialized cases, research applications, and pushing the boundaries of what's possible with GPU computing.

12.4 Final Thoughts

Mastering PTX is not required for every CUDA developer, but it provides a powerful tool for performance-critical applications, especially in deep learning. Like assembly language in CPU programming, PTX gives you a window into how the hardware actually executes your code, enabling optimizations that would be difficult or impossible at higher levels of abstraction.

Whether you choose to write PTX directly, use inline PTX for specific operations, or simply inspect PTX to understand compiler behavior, the knowledge you've gained from this tutorial will help you write more efficient GPU code and make the most of NVIDIA's powerful hardware.

References

- [1] NVIDIA Developer Blog, "Understanding PTX, the Assembly Language of CUDA GPU Computing," NVIDIA Technical Blog, 2020.
- [2] Bruce Liang, "NVIDIA Tensor Core - Getting Started with MMA PTX Programming," Medium, 2022.
- [3] Mark Craddock, "DeepSeek and DeepEP — Understanding DeepSeek's Custom CUDA PTX Instruction," Medium, 2025.
- [4] NVIDIA Developer Forums, "Problem generating .PTX files," NVIDIA Developer Forums, 2022.
- [5] NVIDIA, "CUDA Binary Utilities," NVIDIA Documentation, 2023.
- [6] Stack Overflow, "How to compile PTX code," Stack Overflow, 2013.
- [7] LLVM Project, "User Guide for NVPTX Back-end," LLVM Documentation, 2023.
- [8] ArrayFire, "Demystifying PTX Code," ArrayFire Blog, 2019.
- [9] Cornell Virtual Workshop, "Understanding GPU Architecture - GPU Memory," Cornell University, 2020.