ML-Triton, A Multi-Level Compilation and Language Extension to Triton GPU Programming

Dewei Wang

Intel Corporation
Shanghai, China
dewei.wang@intel.com

Ettore Tiotto
Intel Corporation
Toronto, Canada
ettore.tiotto@intel.com

Wei Zhu
Intel Corporation
Shanghai, China
wei2.zhu@intel.com

Quintin Wang
Intel Corporation
Shanghai, China
quintin.wang@intel.com

Liyang Ling
Intel Corporation
Shanghai, China
liyang.ling@intel.com

Whitney Tsang
Intel Corporation
Toronto, Canada
whitney.tsang@intel.com

Julian Oppermann

Codeplay Software

London, United Kingdom julian.oppermann@codeplay.com

Abstract

In the era of Large Language Models (LLMs), dense operations such as General Matrix Multiplication (GEMM) and Multi-Head Attention (MHA) are critical components. These operations are well-suited for parallel execution using a tile-based approach. While traditional GPU programming often relies on low level interfaces like CUDA or SYCL, Triton [1] has emerged as a domain-specific language (DSL) that offers a more user-friendly and portable alternative by programming at a higher level.

The current Triton starts at the workgroup (aka threadblock) level, and directly lowers to per-thread level. And then attempt to coalesce and amend through a series of passes, promoting information from low-level representation. We believe this is pre-mature lowering based on the below observations.

 GPU has a hierarchical structure both physically and logically. Modern GPUs often feature SIMD units capable of directly operating on tiles on a warp or warpgroup basis, such as blocked load and blocked matrix multiply-accumulate (MMA). Jacky Deng
Intel Corporation
Shanghai, China
jacky.deng@intel.com

- Multi-level gradual lowering can make compiler decoupled and clean by separating considerations inter and intra a logical layer.
- 3. Kernel developers often need fine control to get good performance on the latest hardware. FlashAttention2 [2] advocates explicit data partition between warps to make a performance boost.

In this context, we propose ML-Triton which features multi-level compilation flow and programming interface. Our approach begins at the workgroup level and progressively lowers to the warp and intrinsic level, implementing a multi-level lowering align with the hierarchical nature of GPU. Additionally, we extend triton language to support user-set compiler hint and warp level programming, enabling researchers to get good out-of-the box performance without awaiting compiler updates.

Experimental results demonstrate that our approach achieves performance above 95% of expert-written kernels on Intel GPU, as measured by the geometric mean.

CCS Concepts: • Software and its engineering \rightarrow Compilers.

Keywords: Triton, MLIR, AI Compiler, GPU, Code Generation, Parallel Computing

ACM Reference Format:

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Table 1. Hardware Specification for PVC max 1550

Hardware	Logical	Memory	Capacity
Level	Level	Hierarchy	
Chip	Grid	HBM	128 GB / GPU
XeCore	CTA	L1 Cache	512 KB / XeCore
	(Workgroup)	SLM	128 KB / XeCore
EU Lane	Warp Thread	Register File	512 KB / XeCore

1 Background

1.1 Intel GPU

Released in 2022, the Intel Ponte Vecchio GPU (PVC) [3] is architected with a strong emphasis on AI inference and training workloads. This GPU features a modular design, consisting of two tiles, each containing 64 XeCores. Each XeCore houses 8 Execution Units (EUs), with each EU capable of supporting 8 hardware contexts, where each context operates using native SIMD16 instructions. This architecture positions PVC as a formidable General-Purpose GPU (GPGPU), delivering robust performance and efficiency across a broad range of applications.

PVC boasts an advanced memory hierarchy that significantly enhances efficiency and speed. The global memory offers extensive storage accessible by all XeCores within the GPU. The L2 cache acts as a high-speed intermediary, bridging global memory and core processing units to reduce latency and optimize overall performance. Additionally, the Shared Local Memory (SLM) and L1 cache, which share the same physical space, facilitate rapid data communication between warps, further optimizing performance. In terms of logical hierarchy, Intel employs the concepts of *Workgroup*, *Subgroup*, *Workitem*, which are the counterpart of Nvidia *CTA/ThreadBlock*, *Warp*, *Thread* respectively. The specification of PVC (max 1550) is summarized in Table 1.

Intel complements the hardware with a comprehensive instruction set that maximizes the GPU's potential. Recognizing the importance of memory bandwidth, PVC features 1D/2D blocked load/store/prefetch instructions, which enhances data transfer efficiency from HBM. These instructions support cache hint, address calculation offload, hardware data padding and hardware transpose. To accelerate compute-intensive tasks, PVC also has a DPAS (Dot Product Accumulation with Systolic Array) instruction which accelerates GEMM on blocked data.

PVC supports both SIMT and SIMD programming models, offering flexibility to developers. For SIMT, Intel promotes SYCL, a cross-platform C++ programming model with a thread-parallel execution approach closely aligned with CUDA. For SIMD, there is a SYCL extension called eSIMD [4], enables developers to write explicitly vectorized code.

1.2 Triton

A high-performance kernel is essential for maximizing the computational power of a GPU, thereby accelerating both training and inference processes. Traditionally, handwritten vendor libraries such as CUTLASS [5] have been the go-to solutions for achieving optimal performance. However, these approaches often demand a deep understanding of hardware intricacies and lack the flexibility required by researchers who frequently experiment with novel ideas and seek solutions that deliver good performance out-of-the-box.

Triton is emerging as a new programming language tailored specifically for GPU kernel development. With its Python-like syntax and workgroup level programming interface, Triton makes GPU programming more accessible to AI researchers and engineers. It allows researchers to focus on algorithmic innovation and high-level optimizations without being bogged down by low-level hardware details.

Currently, Triton is the default backend for TorchInductor [6], enabling PyTorch ATen operators to be dispatched to pre-written Triton template kernels. TorchInductor also supports Just-In-Time (JIT) generation of element-wise and reduction operations in Triton, which can then be fused with Triton template kernels. This approach provides greater flexibility compared to traditional pre-defined operator fusion patterns, allowing for more dynamic and efficient execution.

Moreover, an increasing number of frameworks and tools are adopting Triton for kernel development, including vLLM [7], Mamba [8] and DeepSpeed [9].

1.2.1 Triton Dialect.

The latest Triton compiler is built on top of MLIR [10], leveraging its extensive set of built-in utilities. Besides reusing existing dialects such as arith, math, and scf for computation and control flow, Triton introduces its own dialect, known as the Triton dialect (abbreviated as tt), specifically designed to express block-level operations on tensors.

In Triton, tensor represents an N-dimensional array of either values or pointers. By default, Triton uses tensor of pointers as the primary mechanism for memory access. This means that each element in the tensor is a pointer, representing a block of pointers. Later Triton introduced pointer to a block tensor (block pointer) to represent a contiguous block of data. However, in the default compilation pipeline, all block pointers are eventually rewritten into tensors of pointers. While this approach is general enough to handle sparse operations, it necessitates heavy memory analysis to determine data contiguity. For dense operations, we argue that using block pointer is a more efficient approach because it explicitly conveys contiguity information.

Table 2 outlines the primary operations associated with the Triton dialect that will be covered in the following chapters. For more detailed information, please refer to the Triton Dialect definition [11].

Table 2. Triton Dialect

Operations	Description
get_program_id	get ID of the current program/workgroup
load	load a tensor from pointer
store	store a tensor to pointer
dot	matrix multiplication
reduce	reduce along tensor's specified axis
make_tensor_ptr	returns a pointer to a block in a tensor
advance	advance offsets of the tensor pointer

1.2.2 Layout Encoding.

A tensor's RankedTensorType includes a default *encoding* [12] attribute that can provide additional information to the tensor. Triton takes advantage of this and introduces *layout encoding* with careful design. The layout encoding indicates how data should be partitioned across threads [13]. Listing 1 shows the layout encoding that will be used.

BlockedEncoding represents a contiguous portion of a tensor. The parameters for this encoding include:

- sizePerThread: Specifies the block size that each thread operates on.
- *threadsPerWarp*: Defines the arrangement of threads within a warp
- warpsPerCTA: Defines the arrangement of warps within a CTA.
- *order*: Determines the memory access order, with the fastest-changing axis first.

Figure 1 provides a visual representation of a typical BlockedEncoding in Triton. In this example, each thread processes a 2x2 block, 4 threads in a row and 8 threads in a column form a warp, 4 warps in a row and 2 warps in a column form a CTA. As a result, each warp handles a 16x8 block, each CTA works on a 32x32 block.

DotOperandEncoding is used for operands in a dot operation. Take $d = tt.dot \ a$, b, c for example, both c and d share the same layout encoding, a and b have DotOperandEncoding with their parent being c's layout encoding, a's opIdx is 0, b's opIdx is 1, indicating their respective positions in the dot operation.

SliceEncoding indicates its layout is squeezed along the *dim* dimension of the *parent* layout encoding. Take *dst* = *tt.reduce src*, *dim* for example, *dst* has a SliceEncoding with its *parent* being *src*'s layout encoding and *dim* being the dimension to be reduced.

1.2.3 Compilation Flow.

In Triton, kernel functions are decorated with *triton.jit*. Triton compiler will first walk the Abstract Syntax Tree (AST) of the kernel function to generate Triton IR on-the-fly using a standard SSA construction algorithm [14]. Later, Triton IR is converted to Triton GPU IR by adding a naive layout

Listing 1. Triton layout encoding

Tensor: 32×32 , sizePerThread = [2, 2], threadsPerWarp = [8, 4], warpsPerCTA = [2, 4]

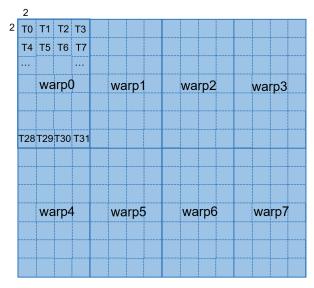


Figure 1. Triton BlockedEncoding

encoding to each tensor type. The Triton GPU IR then undergoes a series of middle-end optimizations aimed at analyzing and simplifying the code. These optimizations include memory coalescing, dot-product-specific enhancements, software pipelining, etc. Finally, the Triton GPU IR is converted to LLVM IR which can then be passed to GPU backend compiler to generate binary for execution. Figure 2 illustrates the compilation flow.

2 Compilation Flow

Our proposed compilation flow, illustrated in Figure 3, implements a multi-level lowering that reflects the GPU hierarchy. This approach decouples considerations at different layers, allowing for more efficient and targeted optimizations.

Initially, Triton IR operates at the **workgroup level**, then we convert it to TritonGPU IR by adding appropriate layout encoding to specify its data distribution between warps. The following *distribute-to-warps* pass will transform the kernel workload to **warp level** i.e. what each warp should work on.

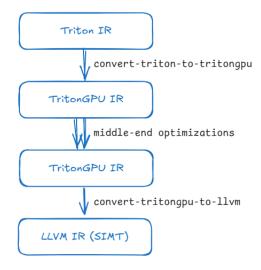


Figure 2. Triton compilation flow

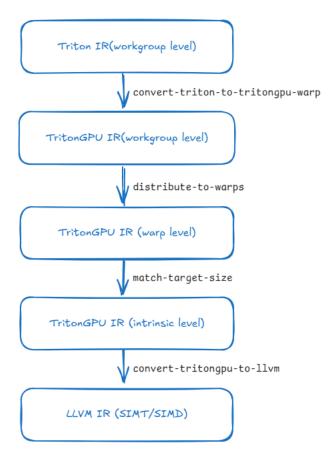


Figure 3. proposed compilation flow

The *match-target-size* pass further split operations to match the LLVM intrinsic size that vendor target can support which we refer to as the **intrinsic level**. Finally, the TritonGPU IR is converted to LLVM IR with either SIMT or SIMD style. Next, we will use a GEMM example [15] to illustrate the compilation flow, as GEMM is a typical AI workload people are most familiar with.

Given $A \in \mathbb{R}^{m \times k}$ $B \in \mathbb{R}^{k \times n}$ $C \in \mathbb{R}^{m \times n}$, the GEMM is C += A * B. User can configure the triton kernel to process a workload of 256x256 += 256x32 * 32x256 in the loop body, with the number of warps (numWarps) set to 32. Listing 2 presents the Triton IR after parsing the AST.

2.1 Convert-triton-to-tritongpu-warp

This pass begins by analyzing the kernel's workload pattern (e.g., element-wise, reduction, gemm, attention) and then figures out the optimal layout encoding for root operation such as *tt.dot* operation. Subsequently, we get all other value's layout encoding through def-use chain propagation.

We make the propagation rules straight-forward: apart from the rules for *tt.dot* and *tt.reduce* introduced earlier, other operations- including *tt.load*, *tt.store*, *tt.advance* and *arith/math.unary/binary* operations require that all source operands and results share the same layout encoding.

There are three major differences from triton upstream.

- 1. Workload-Aware. A dot operation may require different partition strategies depending on the workload to achieve optimal performance. For instance, a square partition is most beneficial for typical GEMM, while FlashAttention-2 [2] prefers a partition along the row dimension to minimize inter-warp communication and achieve peak performance.
- One-Off layout encoding: Our approach determines the layout encoding in a single step, whereas the upstream triton initially assigns a naive layout encoding and refines it in subsequent passes.
- 3. **Focus on sizePerWarp**: We aim to get *sizePerWarp* (the block size per warp works on) rather than *sizePerThread*. It would be a pre-mature lowering to get what each thread works on at the beginning.

So, for the GEMM example, firstly we need to figure out the layout encoding for the root operation - $c += tt.dot \ a, \ b$. Given that

```
c's workgroupSize = [256, 256], numWarps = 32
```

By applying square partitioning between warps , we get c's BlockedEncoding:

```
warpsPerCTA = [8, 4]
```

sizePerWarp = [workgroupSize / warpsPerCTA] = [32, 64]

Then a and b have DotOperandEncoding respectively. Finally by propagation rule described above, all tensor types are annotated with a layout encoding as shown in Listing 2. Note that our generated Triton GPU IR retains the same operations as Triton IR. 1

 $^{^1\}mathrm{Upstream}$ triton would introduce many "convert-layout" operations to help the lowering work.

Listing 2. GEMM Triton IR (w/o highlighted layout encoding) GEMM TritonGPU IR (w/ highlighted layout encoding)

2.2 Distribute-to-warps

This pass distributes the workload of a workgroup across warps according to the corresponding layout encoding. After the pass, we get what each warp works on. Previously we modify BlockedEncoding to include <code>sizePerWarp</code> and <code>WarpsPerCTA</code>—these parameters determine how the workload is distributed.

So, the first step is to get the equivalent BlockedEncoding for every layout encoding. For DotOperandEncoding and SliceEncoding, we derive from its parent layout encoding. The mapping rules are detailed in Table 3.

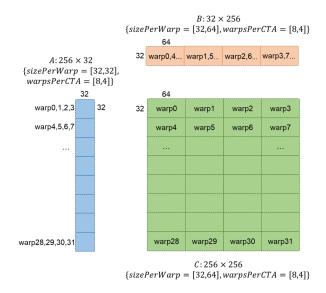


Figure 4. GEMM warp distribution

For the GEMM example, Figure 4 illustrates the data distribution between warps. Matrix C is evenly distributed, with

each warp processing a 32x64 block. For matrix A, the second dimension of *sizePerWarp* already matches *workgroupSize*, yet we have 4 warps in a row to arrange, so warp 0-3 work on the same 32x32 sub-block of A. Similarly, warps 0,4,8..28 work on the same 32x64 sub-block of B.

After the pass, as shown in Listing 3, *tt.dot* is transformed from 256x256 = 256x32 * 32x256 to 32x64 = 32x32 * 32x64, *offsets* of *tt.make_tensor_ptr* are adjusted from a function of *tt.program_id* to a function of *tt.program_id* and *gpu.subgroup_id* (aka warp_id).

2.3 Match-target-size

This pass splits operations into multiple smaller operations to match the target LLVM intrinsic size. All values sharing the same layout encoding are split consistently unless a specific operation requires a different size, in which case a *tt.extract*² operation is introduced to extract a sub-block from the input tensor. Users can specify options such as maximum load size and maximum dot size.

For the GEMM example, PVC's max load size is 32x32, max dot size is 8x16 = 8x16 * 16x16. The data partitioning is illustrated in Figure 5.

After the pass, as shown in Listing 4, the load for A with size 32x32 already matches the target load size, so it remains unchanged. However, the load for B with size 32x64 is spilt into 2 load operations. The dot operation, originally 32x64 = 32x32 * 32x64, is split into 32 smaller dot operations of size 8x16 = 8x16 * 16x16. All values in the def-use chain of the load are spilt into operations working on 32x32 block. Since tt.dot requires different block size, tt.extract is added to extract a 8x16 sub-block from A and a 16x16 sub-block from B, the sub-blocks are then fed to the dot operation. tt.extract

²In practice, tt.extract is moved to triton intel gpu dialect

Original Layout Encoding	Equivalent BlockedEncoding	
#blocked: sizePerWarp = pS, warpsPerCTA = pW ^a	NA	
#dot0 = #triton_gpu.dot_op<{opIdx = 0, parent = #blocked}>	$sizePerWarp = [pS[0], shape^b, warpsPerCTA = pW$	
<pre>#dot1 = #triton_gpu.dot_op<{opIdx = 1, parent = #blocked}></pre>	sizePerWarp = [shape[0], pS[1]], warpsPerCTA = pW	
#slice = #triton_gpu.slice<{dim = i, parent = #blocked}>	sizePerWarp = pS.erase(i), warpsPerCTA = pW.erase(i)	

 $^{^{}a}$ pS and pW are Integer Arrays

```
tt.func public @matmul_kernel_with_block_pointers(%arg0: !tt.ptr<f16>, %arg1: !tt.ptr<f16>, %arg2: !tt.ptr<f32>) {
    %cst = arith.constant dense<0.000000e+00> : tensor<32x64xf32, #blocked>
    %warp_id = gpu.subgroup_id : index
    ...
    %a_ptr = tt.make_tensor_ptr %arg0, [%c4096, %c4096], [%c4096, %c1], [%offsetY, %c0] : <tensor<32x32xf16, #dot0>>
    %b_ptr = tt.make_tensor_ptr %arg1, [%c4096, %c4096], [%c4096, %c1], [%c0, %offsetX] : <tensor<32x64xf16, #dot1>>
    %loop:3 = scf.for %arg3 = %c0 to %c4096 step %c32 iter_args(%c = %cst, %arg5 = %a_ptr, %arg6 = %b_ptr) ... {
    %a = tt.load %arg5 : !tt.ptr<tensor<32x32xf16, #dot0>>
    %b = tt.load %arg6 : !tt.ptr<tensor<32x32xf16, #dot1>>
    %d = tt.dot %a, %b, %c : tensor<32x32xf16, #dot0> * tensor<32x64xf16, #dot1> -> tensor<32x64xf32, #blocked>
}
%d_ptr = tt.make_tensor_ptr %arg2, [%c4096, %c4096], [%c4096, %c1], [%offsetY, %offsetX] : <tensor<32x64xf32, #blocked>>
tt.store %d_ptr, %loop#0 : !tt.ptr<tensor<32x64xf32, #blocked>>
}
```

Listing 3. GEMM TritonGPU IR after distribute-to-warps

```
tt.func public @matmul_kernel_with_block_pointers(%arg0: !tt.ptr<f16>, %arg1: !tt.ptr<f16>, %arg2: !tt.ptr<f32>) {
    %cst = arith.constant dense<0.000000e+00> : tensor<8x16xf32>
    %a_ptr = tt.make_tensor_ptr %arg0, [%c4096, %c4096], [%c4096, %c1], [%offsetY, %c0] : <tensor<32x32xf16>>
    %b_ptr0 = tt.make_tensor_ptr %arg1, [%c4096, %c4096], [%c4096, %c1], [%c0, %offsetX] : <tensor<32x32xf16>>
    %b_ptr1 = tt.make_tensor_ptr %arg1, [%c4096, %c4096], [%c4096, %c1], [%c0, %offsetX + %c32] : <tensor<32x32xf16>>
    %loop:4 = scf.for %arg3 = %c0 to %c4096 step %c32 iter_args(%c_sub = %cst, %arg5 = %a_ptr, %arg6 = %b_ptr0, %arg7 = %b_ptr1) ... {
    %a = tt.load %arg5 : !tt.ptr<tensor<32x32xf16>>
    %b0 = tt.load %arg6 : !tt.ptr<tensor<32x32xf16>>
    %b1 = tt.load %arg7 : !tt.ptr<tensor<32x32xf16>>
    %a_sub0 = tt.extract %a[0] : tensor<32x32xf16> -> tensor<8x16xf16>
    %b_sub0 = tt.extract %b0[0] : tensor<32x32xf16> -> tensor<16x16xf16>
    %accumulate = tt.dot %a_sub0, %b_sub0, %c_sub : tensor<8x16xf16> * tensor<16x16xf16> -> tensor<8x16xf32>
    ... // 32 tt.dot in all }
%d_ptr0 = tt.make_tensor_ptr %arg2, [%c4096, %c4096], [%c4096, %c1], [%offsetY, %offsetX] : <tensor<8x16xf32>>
    tt.store %d_ptr0, %loop#0 : !tt.ptr<tensor<8x16xf32>>
}
```

Listing 4. GEMM TritonGPU IR after match-target-size

will be lowered to sub-register access in the assembly code, without introducing any register moves.

2.4 Convert-tritongpu-to-llvm

This pass converts all operations to LLVM IR. It reuses upstream MLIR conversions for arith, math and scf operations. For triton operations, separate conversion patterns are used to map them to LLVM intrinsic. PVC GPU backend compiler

^b shape is the static size of the RankedTensorType to which this layout encoding is attached

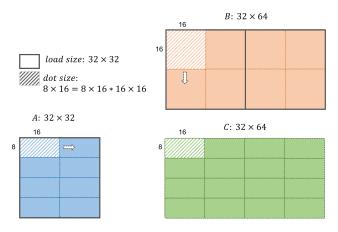


Figure 5. data partition to match target intrinsic size

provides two sets of intrinsic: VectorCompute-Intrinsic [16] for SIMD programming and GenISA-Intrinsic [17] for SIMT programming.

Basically, the conversion is a 1 to 1 mechanical mapping since we already have the operations match the target supported intrinsic size. Only that for SIMT conversion, the data is evenly distributed to each thread lane, meaning the vector size in the intrinsic need to be divided by *threadsPerWarp*.

Table 4. Triton to LLVM

Triton Ops	LLVM Ops - SIMT	LLVM Ops - SIMD
tt.load A	2DBlockRead.v64i16	load2d.stateless.v512i32
tt.load B	2DBlockRead.v32i32	load2d.stateless.v512i32
tt.dot	dpas.v8f32.v8i16.v8i32	dpas2.v128f32.v128i32.v64i32
tt.store	2DBlockWrite.v8i32	store2d.stateless.v128i32

^{*}due to intrinsic constraint, i16/i32 are used instead of f16/f32

Table 4 shows the conversion. Take *tt.load A* for example, the workload size *32x32xf16* is flattened to *v512i32*, when divided by PVC's threadsPerWarp(16), results in *v64i16*. Listing 5 shows the LLVM IR after the conversion.

2.5 FlashAttention-2

As demonstrated in the compilation flow, each value's layout encoding is the key. The layout encoding dictates how work is distributed among warps and serves as a guide for how each operation should be split to match the target intrinsic size.

Once each value is annotated with the correct layout encoding, the subsequent passes can be applied effectively. And once we figure out the root operation's layout encoding, the encoding for all other values can be inferred by tracing the def-use chain. This approach naturally facilitates pre- and post-operation fusion, as they can be seamlessly expanded from the root.

Let's take Flash attention-2 [18] as another example. It can be seen as a fused kernel of two back-to-back GEMMs with an online softmax[19] in between, as outlined in Algorithm 1.

We adopt the work partitioning from the original paper [2] which distributes output matrix O along the row dimension across all warps and K/V accessed by all warps. Balancing the data shared among warps and register pressure, we arrive at the following kernel configuration for PVC:

By horizontally partitioning between warps, we get *O's BlockedEncoding*:

Subsequently, the layout encoding for all other tensors is inferred, as summarized in Table 5. Figure 6 illustrates their relationships.

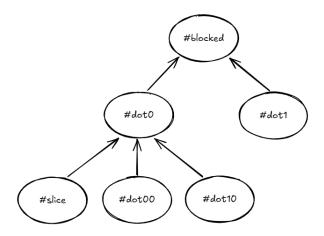


Figure 6. relation of FlashAttention-2 layout encoding

3 Language Extension

Triton, as a DSL for accelerated computing, is inherently extensible. During our development process, we found that by incorporating simple compiler hint and warp level programming, Triton significantly simplifies GPU programming and accelerates the journey to peak performance. The tradeoff is that what is the algorithm workload specific or compiler hackable optimization, we choose to let user have control.

3.1 Compiler Hint

In the code generation pipeline, we resort to a workload-aware pass to detect specific patterns and determine the root operation's tiling partition - layout encoding. Actually this setting is derived from the best-known practice of expert's kernel tuning experience. Also, there are instances where researchers may want to manually define the tiling partition. For example, FlashAttention-2 [2] explicitly proposes how

```
;SIMD style
%a = call <512 x i32> @llvm.genx.lsc.load2d.stateless.v512i32.i1.i64(..., i64 %a_ptr, ...)
%b0 = call <512 x i32> @llvm.genx.lsc.load2d.stateless.v512i32.i1.i64(..., i64 %b_ptr0, ...)
%b1 = call <512 x i32> @llvm.genx.lsc.load2d.stateless.v512i32.i1.i64(..., i64 %b_ptr1, ...)
%a_sub0 = shufflevector <512 x i32> %a, <512 x i32> poison, <64 x i32> <i32 0..63>
%b_sub0 = shufflevector <512 x i32> %b0, <512 x i32> poison, <128 x i32> <i32 0..127>
%d = call <128 x float> @llvm.genx.dpas2.v128f32.v128i32.v64i32(%c_sub, <128 x i32> %b_sub0, <64 x i32> %a_sub0, ...)
;SIMT style
%a = call <64 x i16> @llvm.genx.GenISA.LSC2DBlockRead.v64i16(i64 %a_ptr, ...)
%b0 = call <32 x i32> @llvm.genx.GenISA.LSC2DBlockRead.v32i32(i64 %b_ptr0, ...)
%b1 = call <32 x i32> @llvm.genx.GenISA.LSC2DBlockRead.v32i32(i64 %b_ptr0, ...)
%a_sub0 = shufflevector <64 x i16> %a, <64 x i16> undef, <8 x i32> <i32 0, i32 1, i32 2, i32 3, i32 4, i32 5, i32 6, i32 7>
%b_sub0 = shufflevector <32 x i32> %b0, <32 x i32> undef, <8 x i32> <i32 0, i32 1, i32 2, i32 3, i32 4, i32 5, i32 6, i32 7>
%d = call <8 x float> @llvm.genx.GenISA.sub.group.dpas.v8f32.v8i16.v8i32(%c_sub, <8 x i16> %a_sub0, <8 x i32> %b_sub0, ...)
```

Listing 5. GEMM LLVM IR

```
Algorithm 1: FlashAttention-2 forwardinput: Q, K, V \in \mathbb{R}^{N \times D}, N is the sequence length,<br/>D is the head dimension1 J \leftarrow N/BlockSize2 Load Q from HBM3 for j in (0, J) do4 Load K_j from HBM5 Compute QK = Q * K_j6 Compute P = online\_softmax QK<br/>(intermediate result: m = rowmax, l = rowsum)7 Load V_j from HBM8 Compute O + = P * V_j9 end10 Store O back to HBM.
```

Table 5. FlashAttention-2 Layout Encoding

Symol	Layout Encoding
0	#blocked : sizePerWarp = [16, 64], warpsPerCTA = [8, 1]
V	<pre>#dot1 = #triton_gpu.dot_op<{opIdx = 1, parent = #blocked}></pre>
QK, P	#dot0 = #triton_gpu.dot_op<{opIdx = 0, parent = #blocked}>
K	#dot10 = #triton_gpu.dot_op<{opIdx = 1, parent = #dot0}>
Q	#dot00 = #triton_gpu.dot_op<{opIdx = 0, parent = #dot0}>
m, l	<pre>#slice = #triton_gpu.slice<{dim = 1, parent = #dot0}></pre>

to partition work between different warps to get optimal performance.

Hence, we provide a compiler hint that allows users to specify the root operation's tiling partition between warps. Below are the tiling options available for 2D tensor.

Horizontal: Evenly tiles along the first(row) dimension. **Vertical**: Evenly tiles along the last(column) dimension.

Square: Tiles to form square sub-blocks.

For flash attention-2, merely setting the second dot's *tiling* to *horizontal* is sufficient, no other source code changes are needed. The compiler can then figure out all values' layout encoding accordingly.

```
o = tl.dot(p, v, o, tiling="horizontal")
```

3.2 Warp Level API

Writing kernels at the workgroup level reduces the burden on developers, but performance is highly dependent on compiler, which need time to evolve. Rather than relying solely on compiler-specific optimizations, we believe it is more effective to give developers fine-grained control over their code. For instance, FlashAttention-3 [20] proposes better warp level management to leverage the latest hardware capabilities. Similarly, many kernel libraries like CUTLASS [5] offer warp level C++ templates.

Thereby, we introduce a warp level language extension. The key elements are as follows:

warp_level: Metadata indicating this a warp level kernel.
tl.warp_id(): Returns linear ID of the current warp within
the workgroup.

tl.alloc(shape, data type): Allocates and returns a pointer to a block in the SLM with the specified *shape* and *data type*.

tl.reduce(..., cross_warp, dst_warps): Adds keyword parameters for reduction operations (e.g., max, sum). When cross_warp is set to true, it is a reduction across all warps, otherwise it is a reduction within the current warp. The dst_warps parameter allows the reduction result to be broadcast only to the specified destination warps. If not set, the result will be broadcast to all warps.

LLM inference often involves long sequences for key-value (KV) pairs. To boost throughput, Flash Decoding [21] proposes splitting the keys and values. For efficient memory management, Paged Attention [7] divides the request's KV cache into blocks.

Listing 6 shows how a paged attention triton kernel could be implemented. While the core algorithm remains similar to FlashAttention-2, the warp distribution differs significantly because the sequence length of the query is typically 1. This makes writing a paged attention kernel relatively easy, but achieving optimal performance out-of-the-box is challenging. The compiler must perform specific analyses and optimizations to enhance performance.

Listing 6. paged attention triton kernel - workgroup level

```
# warp 0 load Q from HBM and store it to SLM
    slm_block_ptr = tl.alloc(shape=(1, D), dtype=tl.float16)
2
    if tl.warp_id() == 0:
        q = tl.load(Q_block_ptr)
        tl.store(slm_block_ptr, q)
    tl.barrier()
    q = tl.load(slm_block_ptr)
    k = tl.load(Ki_block_ptr)
    qk = tl.dot(q, k)
    m_i = tl.max(qk, axis=1)
10
11
    m_i = tl.max(m_i, cross_warp = True) # sync partial max
12
    p = tl.exp((qk - m_i[:, None]))
    l_i = tl.sum(p, axis = 1)
13
    l_i = tl.sum(l_i, cross_warp = True) # sync partial sum
14
    p /= l_i[:, None]
15
    v = tl.load(Vi_block_ptr)
16
    o = tl.dot(p.to(tl.float16), v)
17
    # reduce the Output to warp 0 and store it back to HBM
18
    o = tl.sum(o, cross_warp = True, dst_warps=(0))
19
    if tl.warp_id() == 0:
20
        tl.store(0_block_ptr, o)
```

Listing 7. paged attention triton kernel - warp level

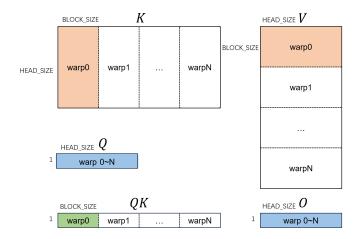


Figure 7. paged attention warp distribution

To ensure sufficient parallelism on the GPU, one approach is to further partition the KV cache between warps, as depicted in Figure 7. Each warp works on separate KV cache blocks, and reduction across warps is needed to synchronize each warp's partial result.

However, by programming Triton at warp level, users can easily express the above decomposition. Listing 7 shows what the warp level kernel would be.

4 Experimental Results

In this section, we aim to demonstrate the effectiveness of our design by evaluating the performance of several popular AI workload kernels.

The experiments were conducted on Intel's PVC max 1550 using OneAPI 2024.1. Performance was measured based on the kernel's GPU execution time recorded by SYCL profiling event [22]. For comparison, we benchmark Triton against Intel's XeTLA [23], Xe template-based linear algebra library optimized as a peak-performance reference, similar to NVIDIA's CUTLASS. To ensure a fair comparison, we used identical configurations for both XeTLA and Triton, including parameters such as tile size, minimizing any performance discrepancies due to these settings.

4.1 **GEMM**

GEMM is a fundamental operation in the AI domain, constituting a significant portion of the computational workload. We used the GEMM kernel [15] from the Triton tutorial for our tests.

We evaluated two types of GEMM operations: memorybound and compute-bound. All matrix shapes were derived from LLM models such as LLama-2 and LLama-3.

Compute-bound GEMM is relevant for both LLM training and inference. As researchers focus on long context length on a single GPU, we tested matrix sizes ranging from m = 1k to 16k, large enough to fully utilize the GPU and achieve peak

hardware throughput. Figure 8 shows that Triton achieves a geometric mean of 96% of XeTLA's performance.

Memory-bound GEMM is a common scenario in LLM inference, particularly during the next-token prediction stage. We evaluated this on cases with large m, large k, and large n to demonstrate Triton's robustness. As shown in Figure 9, Triton's performance is comparable to XeTLA, with a 94% geometric mean.



Figure 8. compute-bound GEMM performance

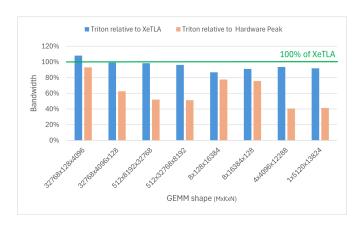


Figure 9. memory-bound GEMM performance

4.2 FlashAttention-2

FlashAttention-2 is widely used for MHA, playing a critical role in modern transformer models. We used the kernel in triton tutorial [18] for our tests.

We evaluated the forward pass with a total 32k tokens and sequence length ranging from 1k to 32k, aligned with the context length of most popular LLMs. The hidden dimension was set to 2048, with head dimension to be either 64 or 128 (i.e.,32 heads or 16 heads). The benchmark results in Figure 10 and Figure 11 show less than a 5% performance gap, demonstrating the high quality of our code generation.

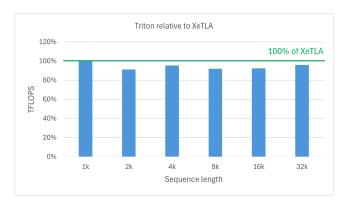


Figure 10. FlashAttention-2 forward performance head dimension = 64

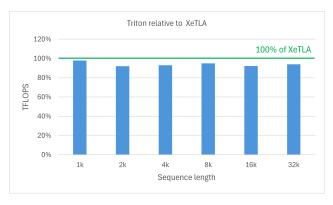


Figure 11. FlashAttention-2 forward performance head dimension = 128

4.3 Paged Attention

Paged attention is widely used in LLM inference engines. Unlike flash attention, the key/value pairs in paged attention are not stored contiguously and must be accessed through a block table mapping, which increases the strain on memory access.

Compared to traditional triton workgroup level implementation, our warp level kernel directly express the distribution between warps, requiring only a few additional lines of code.

As shown in Figure 12, Triton's performance is above 95% of XeTLA's, demonstrating its capability to handle complex kernels effectively.

5 Conclusion and Future work

In this paper, we presented ML-Triton which features multilevel lowering and programming interface. The multi-level compilation flow is closely aligned with the GPU's layered hierarchy. By progressively lowering operations from the workgroup level to the warp level and finally to the intrinsic level, we decompose high-level operations step by step guided by the layout encoding in an innovative straightforward way.

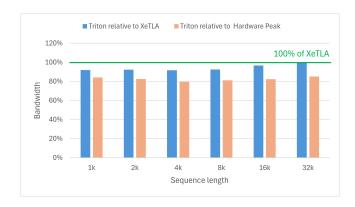


Figure 12. PagedAttention performance

Additionally, we extend Triton language by introducing user-defined compiler hints and warp level programming. These enhancements provide researchers with fine-grained control over their code, reducing dependency on compiler specific optimizations, leading to better out-of-the-box performance.

We thoroughly evaluated three popular kernels—GEMM, FlashAttention-2, and Paged Attention—based on our approach, achieving a performance gap of less than 5% compared to expert tuned implementation.

Overall, building on top of Triton, our proposal further bridges the gap between ease of use and high performance in GPU programming.

Looking ahead, we plan to polish our design as we encounter more use cases from the rapidly evolving AI land-scape. We also anticipate that this programming and compilation paradigm could be extended beyond GPUs to other many-core architectures.

References

- [1] Philippe Tillet, H. T. Kung, and David Cox. 2019. Triton: an intermediate language and compiler for tiled neural network computations. In Proceedings of the 3rd ACM SIGPLAN International Workshop on Machine Learning and Programming Languages (Phoenix, AZ, USA) (MAPL 2019). Association for Computing Machinery, New York, NY, USA, 10–19. https://doi.org/10.1145/3315508.3329973
- [2] Tri Dao. 2023. FlashAttention-2: Faster Attention with Better Parallelism and Work Partitioning. arXiv:2307.08691 [cs.LG] https://arxiv.org/abs/2307.08691
- [3] H. Jiang. 2022. Intel's Ponte Vecchio GPU: Architecture, Systems & Software. In 2022 IEEE Hot Chips 34 Symposium (HCS). IEEE Computer Society, Los Alamitos, CA, USA, 1–29. https://doi.org/10.1109/HCS55958.2022.9895631
- [4] Guei-Yuan Lueh, Kaiyu Chen, Gang Chen, Joel Fuentes, Wei-Yu Chen, Fangwen Fu, Hong Jiang, Hongzheng Li, and Daniel Rhee. 2021. C-formetal: High performance simd programming on intel gpus. In 2021 IEEE/ACM International Symposium on Code Generation and Optimization (CGO). IEEE, 289–300.
- [5] 2024. Nvidia CUTLASS- a collection of CUDA C++ template abstractions. Retrieved Aug 24, 2024 from https://github.com/NVIDIA/cutlass
- [6] Jason Ansel, Edward Yang, Horace He, Natalia Gimelshein, Animesh Jain, Michael Voznesensky, Bin Bao, Peter Bell, David Berard, Evgeni Burovski, Geeta Chauhan, Anjali Chourdia, Will Constable, Alban Desmaison, Zachary DeVito, Elias Ellison, Will Feng, Jiong Gong, Michael Gschwind, Brian Hirsh, Sherlock Huang, Kshiteej Kalambarkar, Laurent Kirsch, Michael Lazos, Mario Lezcano, Yanbo Liang, Jason Liang, Yinghai Lu, C. K. Luk, Bert Maher, Yunjie Pan, Christian Puhrsch, Matthias Reso, Mark Saroufim, Marcos Yukio Siraichi, Helen Suk, Shunting Zhang, Michael Suo, Phil Tillet, Xu Zhao, Eikan Wang, Keren Zhou, Richard Zou, Xiaodong Wang, Ajit Mathews, William Wen, Gregory Chanan, Peng Wu, and Soumith Chintala. 2024. Py-Torch 2: Faster Machine Learning Through Dynamic Python Bytecode Transformation and Graph Compilation. In Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2 (La Jolla, CA, USA) (ASPLOS '24). Association for Computing Machinery, New York, NY, USA, 929-947. https://doi.org/10.1145/3620665.3640366
- [7] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient Memory Management for Large Language Model Serving with PagedAttention. In Proceedings of the 29th Symposium on Operating Systems Principles (Koblenz, Germany) (SOSP '23). Association for Computing Machinery, New York, NY, USA, 611–626. https://doi.org/10.1145/3600006.3613165
- [8] Albert Gu and Tri Dao. 2023. Mamba: Linear-time sequence modeling with selective state spaces. arXiv preprint arXiv:2312.00752 (2023).
- [9] Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. 2020. DeepSpeed: System Optimizations Enable Training Deep Learning Models with Over 100 Billion Parameters. In Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (Virtual Event, CA, USA) (KDD '20). Association for Computing Machinery, New York, NY, USA, 3505–3506. https://doi.org/10.1145/3394486.3406703
- [10] Chris Lattner, Mehdi Amini, Uday Bondhugula, Albert Cohen, Andy Davis, Jacques Pienaar, River Riddle, Tatiana Shpeisman, Nicolas Vasilache, and Oleksandr Zinenko. 2020. MLIR: A Compiler Infrastructure for the End of Moore's Law. arXiv:2002.11054 [cs.PL] https://arxiv.org/abs/2002.11054
- [11] 2024. Triton Dialect Definition. Retrieved Aug 24, 2024 from https://github.com/triton-lang/triton/blob/main/include/triton/ Dialect/Triton/IR/TritonOps.td
- [12] Aart Bik, Penporn Koanantakool, Tatiana Shpeisman, Nicolas Vasilache, Bixia Zheng, and Fredrik Kjolstad. 2022. Compiler Support for

- Sparse Tensor Computations in MLIR. ACM Trans. Archit. Code Optim. 19, 4, Article 50 (sep 2022), 25 pages. https://doi.org/10.1145/3544559
- [13] 2024. Triton Layout Encoding. Retrieved Aug 24, 2024 from https://github.com/triton-lang/triton/blob/main/include/triton/ Dialect/TritonGPU/IR/TritonGPUAttrDefs.td
- [14] 2021. Introducing Triton: Open-source GPU programming for neural networks. Retrieved Aug 24, 2021 from https://openai.com/index/triton
- [15] 2023. Triton matmul kernel using block-pointer. Retrieved Aug 24, 2024 from https://github.com/triton-lang/triton/blob/v2.1.0/python/ tutorials/08-experimental-block-pointer.py
- [16] 2020. Intel VC Intrisic Definition. Retrieved Aug 24, 2024 from https://github.com/intel/vc-intrinsics/blob/v0.19.0/GenXIntrinsics/ include/llvm/GenXIntrinsics/Intrinsic_definitions.py#L2255
- [17] 2018. Intel GenISA Intrisic Definition. Retrieved Aug 24, 2024 from https://github.com/intel/intel-graphics-compiler/blob/releases/ igc-1.0.17384/IGC/GenISAIntrinsics/Intrinsic_definitions.py#L2701
- [18] 2022. Triton flashAttention kernel. Retrieved Aug 24, 2024 from https://github.com/triton-lang/triton/blob/v2.1.0/python/ tutorials/06-fused-attention.py
- [19] Maxim Milakov and Natalia Gimelshein. 2018. Online normalizer calculation for softmax. CoRR abs/1805.02867 (2018). arXiv:1805.02867 http://arxiv.org/abs/1805.02867
- [20] Jay Shah, Ganesh Bikshandi, Ying Zhang, Vijay Thakkar, Pradeep Ramani, and Tri Dao. 2024. FlashAttention-3: Fast and Accurate Attention with Asynchrony and Low-precision. arXiv:2407.08608 [cs.LG] https://arxiv.org/abs/2407.08608
- [21] 2023. Flash-decoding for long-context inference. Retrieved Aug 24, 2024 from https://crfm.stanford.edu/2023/10/12/flashdecoding.html
- [22] 2024. SYCL 2020 Specification. Retrieved Aug 24, 2024 from https://registry.khronos.org/SYCL/specs/sycl-2020/html/sycl-2020.html
- [23] 2023. XeTLA- Intel Xe Templates for Linear Algebra. Retrieved Aug 24, 2024 from https://github.com/intel/xetla/releases/tag/v0.3.6