# Hexagon-MLIR: An AI Compilation Stack For Qualcomm's Neural Processing Units (NPUs)

Mohammed Javed Absar[1], Muthu Baskaran[2], Abhikrant Sharma[2], Abhilash Bhandari[3], Ankit Aggarwal[2], Arun Rangasamy[3], Dibyendu Das[3], Fateme Hosseini[2], Franck Slama[2], Iulian Brumar[2], Jyotsna Verma[2], Krishnaprasad Bindumadhavan[3], Mitesh Kothari[2], Mohit Gupta[2], Ravishankar Kolachana[3], Richard Lethin[2], Samarth Narang[2], Sanjay Motilal Ladwa[3], Shalini Jain[3], Snigdha Suresh Dalvi[2], Tasmia Rahman[2], Venkat Rasagna Reddy Komatireddy[3], Vivek Vasudevbhai Pandya[3], Xiyue Shi[2], and Zachary Zipper[2]

[1]Qualcomm Technologies International, Ltd.
[2]Qualcomm Technologies, Inc.
[3]Qualcomm India Private Limited

## Abstract

In this paper, we present Hexagon-MLIR [17, 18], an open-source compilation stack that targets Qualcomm Hexagon Neural Processing Unit (NPU) [10] and provides unified support for lowering Triton kernels [24] and PyTorch models [15]. Built using the MLIR framework [11], our compiler applies a structured sequence of passes to exploit NPU architectural features to accelerate AI workloads. It enables faster deployment of new Triton kernels (handwritten or subgraphs from PyTorch 2.0 [2]), for our target by providing automated compilation from kernel to binary. By ingesting Triton kernels, we generate *mega-kernels* that maximize data locality in the NPU's Tightly Coupled Memory (TCM), reducing the bandwidth bottlenecks inherent in library-based approaches. This initiative complements our commercial toolchains by providing developers with an open-source MLIR-based compilation stack that gives them a path to advance AI compilation capabilities through a more flexible approach. Hexagon-MLIR is a work-in-progress, and we are continuing to add many more optimizations and capabilities in this effort.

## 1 Introduction

Hexagon-MLIR is a stack for compiling and executing Triton kernels [24] and PyTorch [15] graphs on the Qualcomm Hexagon Neural Processing Unit (NPU) [10]. Parts of the stack were open-sourced recently [17, 18]. Built using the MLIR framework [11], our compiler applies a structured sequence of conversion, optimization, and lowering passes to exploit NPU architectural features — including Hardware Vector eXtension (HVX); multi-threading that exploits hardware threads and multiple HVX contexts; memory hierarchy including Tightly Coupled Memory (TCM); interleaving of DMA transfers and computation; optimized and vectorized math libraries; and a high-throughput matrix multiplication engine — to accelerate AI workloads. Hexagon-MLIR adopts a Generative Approach: by treating fusion as a first-class compiler pass, we can generate specialized kernels for arbitrary long operator chains.

The traditional approach to ML compilation often relied on *operator libraries* which deliver peak performance but require time-consuming, careful manual implementation [4] [21]. With new operators and optimizations being written in Triton and other DSLs
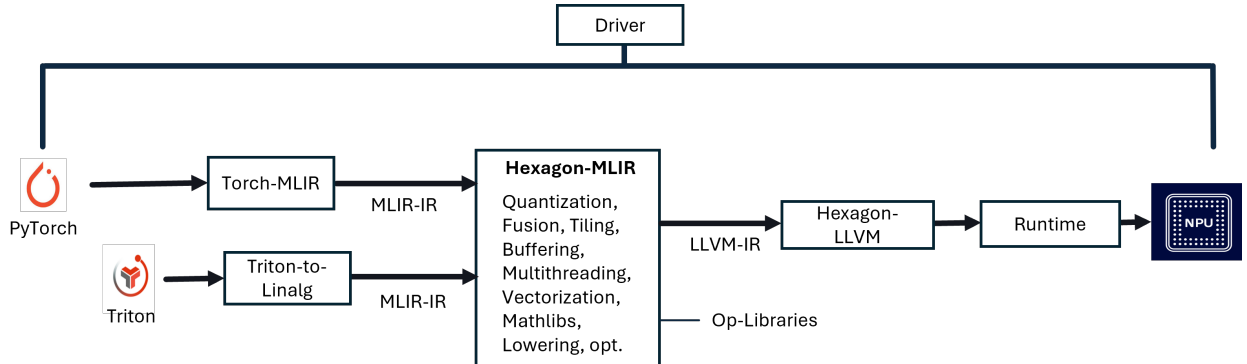
Figure 1: Hexagon-MLIR — AI Compilation Stack

[19, 25] continuously, a quick automated path from, e.g. PyTorch operators, hand-written or Inductor-generated Triton kernels, to reasonably optimized binaries provides a significant deployment speed advantage. Moreover, *operator libraries* act as opaque barriers to optimization. This fragmentation forces data to round-trip to DRAM between calls, creating a bandwidth bottleneck. Much like region compilation [7, 9], which offered locality and parallelism exploitation in the 1990s, Triton kernels via fusion [28] provide opportunities for locality and parallelism exploitation that operator libraries can miss. Hand-written assembly (peak-performance code) is faster for single operations, but automatic code generation is the only scalable solution for the rapid evolution of AI models. Moreover, the Triton ecosystem continues to evolve with higher-level abstractions such as Helion [16], which provides a PyTorch-centric DSL that compiles to Triton, further reducing the barrier to kernel development. Hexagon-MLIR provides the architectural foundation to ingest these large regions (e.g. via PyTorch Inductor [2, 15]) and compile them into binaries.

The challenge of efficiently compiling new large language models (LLMs) and recent Mixture-of-Experts architectures [6] remains an active area of research. Recent works have demonstrated that Triton can achieve state-of-the-art performance for critical LLM components such as attention mechanisms [5, 20, 26]. Additionally, novel activation functions such as poly-nomial composition activations (PolyCom) [31] have been proposed to enhance the expressivity and dynamics of transformers, further motivating the need for flexible compilation frameworks that can efficiently support emerging architectural innovations. With the PyTorch 2.0 Inductor path, Triton kernels can be generated automatically. Alternatively, new operators or algorithms can be quickly and efficiently written in the Triton language. The hand-optimized *operator library* path is therefore not scalable.

Our compiler is built using MLIR [11]. MLIR is not a specific compiler but a framework with which compilers can be constructed. One such AI compiler built using MLIR is IREE [27], which provides portability via a Hardware Abstraction Layer (HAL). Another MLIR-based compiler is TPP (Tensor Processing Primitives) [8], which targets Intel and AMD processors and, according to the authors, "achieves 90% of the performance of hand-optimized equivalent programs." Compiler-based approaches have also been applied to distributed training optimization, as demonstrated by DeepCompile [23], which uses profiling-guided optimization passes to improve communication-computation overlap in distributed settings. While wrapping hand-optimized libraries (e.g., Intel's approach with libxsmm) guarantees peak performance for standard operations, it fails to address the long tail of fusion patterns found in modern Transformers, and faces a combinatorial explosion of required kernels.

# 2 Organization

The remainder of this paper is organized as follows. Section 3 provides an overview of the Hexagon-MLIR compilation flow, formalizing key concepts including semantic preservation and the pass pipeline architecture. Section 4 presents detailed technical descriptions of our core compiler passes—including Triton-to-Linalg conversion, the `linalg.generic` abstraction, operator fusion, tiling for memory hierarchy, multi-threading, double buffering, math library integration, and additional optimizations. We use motivating kernel examples (softmax and GELU) throughout to illustrate the progressive transformation of code through the compilation pipeline. Section 5 presents experimental results and performance analysis, demonstrating the effectiveness of individual optimizations (vectorization, multi-threading, double buffering) and their interactions as kernels transition between compute-bound and memory-bound regimes. Section 6 discusses related work in deep learning compilers and optimization techniques. Finally, Section 7 concludes with a summary of our contributions and directions for future work.

# 3 Hexagon-MLIR: Overview

Figure 1 provides an overview of the compilation flow. PyTorch models convert to Linalg and other MLIR core ops using Torch-MLIR [1]. Similarly, Triton kernels convert to MLIR using triton-to-linalg [13].

**Domains.** Let $\mathcal{M}$ be the set of PyTorch models, $\mathcal{K}$ the set of Triton kernels, $\mathcal{L}$ the set of MLIR Linalg ops. Then we have:

$$f : \mathcal{M} \longrightarrow \mathcal{L} \quad \text{(Torch-MLIR conversion)}$$
$$g : \mathcal{K} \longrightarrow \mathcal{L} \quad \text{(Triton-to-Linalg conversion)}$$

**Semantics preservation.** For all $m \in \mathcal{M}$ and $k \in \mathcal{K}$, lowering paths to Linalg should preserve the denotational semantics.

$$[\![m]\!] = [\![f(m)]\!] \quad \text{and} \quad [\![k]\!] = [\![g(k)]\!]$$

While certain precision-related choices, e.g. quantization, can weaken strict semantic preservation, machine learning models are generally tolerant of such variations. Also, if custom named operations appear in the MLIR graph, we translate them to semantically equivalent MLIR core operations.

**Hexagon-MLIR Pass Pipeline.** The compilation process can be modeled as a sequence of morphisms between intermediate representations (IRs), where each pass transforms one IR into another:

$$\text{Obj}(\mathbf{IR}) = \{\text{IR}_\alpha \mid \alpha \text{ is an IR abstraction level}\}$$
$$\text{Hom}(\mathbf{IR}) = \{p : \text{IR}_\alpha \mapsto \text{IR}_\beta \mid p \in \text{pass-pipeline}\}$$

Passes could be classified into at least two types:

- **Transformation:** rewrites the IR to simplify (e.g. cse, canonicalization) or to improve performance. E.g. vectorization pass:

$$p : \text{IR}_{scalar} \mapsto \text{IR}_{vectorized}$$

- **Lowering:** lowers abstraction towards the target. E.g. structure-control-flow to unstructured control flow:

$$p : \text{IR}_{\text{scf}} \mapsto \text{IR}_{\text{cf}}$$

Hexagon-MLIR pass pipeline has many passes but in this paper we focus on some main ones

| Symbol | Description |
|--------|-------------|
| $C$ | Canonicalization, CSE, or CP |
| $F$ | Operator Fusion |
| $L$ | Layout transformation and propagation |
| $M$ | Multi-threading |
| $p_i$ | an optimization or lowering pass |
| $T$ | Tiling for memory hierarchy |
| $Q$ | Quantization |
| $V$ | Vectorization and related passes |

The Hexagon-MLIR pipeline is a composition of passes applied on the input IR:

$$\Pi = p_n \circ \cdots \circ T \circ F \circ \cdots \circ p_0 : \text{IR}_0 \longrightarrow \text{IR}_{\text{final}}.$$

The rest of the paper goes through in details on some of the important passes.

# 4 Hexagon-MLIR: Details

In this section we will deep-dive into some of the passes. We use a couple of examples to facilitate the explanation.

## 4.1 Triton to Linalg

**Softmax** function is a widely used function in machine learning. Mathematically, for a given $\mathbf{z} = (z_1, \ldots, z_K)$, softmax is:

$$\sigma(z_i) = \frac{\exp(z_i - \max_j z_j)}{\sum_{k=1}^{K} \exp(z_k - \max_j z_j)}.$$

Subtracting by $\max_j z_j$ improves numerical stability without changing the result.

```
@triton.jit
def softmax_kernel(output_ptr, input_ptr,
  ...
  row_minus_max = row - tl.max(row, axis=0)
  numerator = tl.exp(row_minus_max)
  denominator = tl.sum(numerator, axis=0)
  softmax_output = numerator / denominator
  ...
```

Listing 1: Extract from Triton Softmax Kernel

**Triton Kernel.** In Listing 1, we show an excerpt of a Triton kernel that implements a numerically stable row-wise softmax. The kernel subtracts the maximum value in each row from all elements to prevent overflow during exponentiation. Then it computes tl.exp and aggregates using a row-wise sum tl.sum to form the normalization factor. Finally, the softmax probabilities are obtained by dividing each exponential by this normalizing factor.

The Triton kernel in Listing 1 is lowered by the Triton compiler to Triton-IR and then to Linalg by the Triton-to-Linalg converter. Listing 2 shows the IR. It shows a sequence of linalg.generics. We take a small detour here to explain linalg.generic as it forms an important part of the compilation flow.

## 4.2 Linalg-Generic

The linalg.generic is a key operator in the Linalg (Linear Algebra) dialect of MLIR. As an op it has a number of constructs - indexing map, iterator type, ins and outs, and linalg body.

**Indexing maps:** Let $\mathcal{I} \subseteq \mathbb{Z}^d$ denote the *iteration domain*, an integer lattice subject to affine constraints. A linalg.generic operation specifies a finite collection of affine maps

$$\phi_\ell : \mathcal{I} \longrightarrow \mathbb{Z}^{n_\ell},$$

each of which defines a view into the coordinate space of an input or output tensor $X_\ell$. Collectively, these maps determine how points in the common iteration domain are projected onto the indexing sets of the respective operands.

**Iterator types:** Iteration domains labeled as parallel or reduction, encode interpretation of the computation along each axis in the body of the generic. The linalg body must be consistent with this declaration.

**Payload region:** The linalg.body which is a basic-block operating on scalars and yielding scalar results (outs) from scalar inputs. The body specifies the computation to be done at a single point in the iteration space of that linalg.generic. Each scalar is extracted from tensor operands using the indexing map and passed as block argument to the linalg.body.

The linalg.generic is a structured op that internalizes the geometry of iteration (index set, affine map, and iterator type) while delegating the computation to the linalg body. This design separates control over data traversal from computation semantics on the data, enabling principled transformations prior to lowering to explicit loops. For instance, every named op (e.g., linalg.add, linalg.conv_2d) uniquely lowers to linalg.generic. This is key to its usefulness which ensures uniform transformation, optimization and lowering machinery across the sea of ops in ML. Although you can have linalg.generic operating on buffers, the best practice is to use linalg-on-tensors,

```
1   %9 = linalg.generic
2          {indexing_maps = [affine_map<(d0) -> (d0)>, affine_map<(d0) -> (d0)>,
3           affine_map<(d0) -> (d0)>],
4           iterator_types = ["parallel"]}
5          ins(%3, %8 : tensor<16384xf32>, tensor<16384xf32>)
6          outs(%4 : tensor<16384xf32>) {
7        ^bb0(%in: f32, %in_5: f32, %out: f32):
8          %17 = arith.subf %in, %in_5 : f32
9          linalg.yield %17 : f32
10   } -> tensor<16384xf32>
11   %10 = tensor.empty() : tensor<16384xf32>
12   %11 = linalg.generic
13          {indexing_maps = [affine_map<(d0) -> (d0)>, affine_map<(d0) -> (d0)>],
14           iterator_types = ["parallel"]}
15          ins(%9 : tensor<16384xf32>) outs(%10 : tensor<16384xf32>) {
16        ^bb0(%in: f32, %out: f32):
17          %17 = math.exp %in : f32
18          linalg.yield %17 : f32
19   } -> tensor<16384xf32>
```

Listing 2: Linalg IR of the Softmax excerpt.

which we do, as it provides the benefits of SSA representation.

It must be mentioned on balance that many ops that have characteristics such as – shape-inconsistent outputs; data-dependent flow; non-affine maps etc – cannot be represented by a single `linalg.generic`. Important ML ops such as `softmax`, `prefix-sum`, `TopK` etc need other control flow ops along with `linalg.generic` to make them representable. Also, it must be mentioned that an intermediate abstraction layer was added to Linalg in the form of `linalg.contraction`, `linalg.elementwise`. These also lower uniquely to `linalg.generic`.

### 4.3 Operator Fusion

**Operator fusion** is a semantics-preserving transformation that combines multiple structured operations (e.g., `linalg.generic`) into a single operation, eliminating intermediate materialization and improving data locality.

Consider two operations:

$$X \xrightarrow{P} Y \xrightarrow{Q} Z,$$

where $P$ and $Q$ are `linalg.generic` operations with iteration domains $\mathcal{I}_P$ and $\mathcal{I}_Q$. The **producer** $P$ writes elements of $Y$ using:

$$\phi_{\text{out}}^P : \mathcal{I}_P \to \text{dom}(Y),$$

while the **consumer** $Q$ reads elements of $Y$ using:

$$\phi_{\text{in}}^Q : \mathcal{I}_Q \to \text{dom}(Y).$$

A case of fusion combines the two `linalg.generic` into one so that instead of writing back the entire tensor $Y$, in the body of the fused generic scalar value $x$ of $X$ is read to produce scalar value $y$ of $Y$ that is immediately consumed to produce $z$ of $Z$. Therefore, if

$$P = \texttt{linalg.generic}\big\{\mathcal{I}_P, \{\phi_\ell^P\}, \mathbf{y} := f(\mathbf{x})\big\},$$
$$Q = \texttt{linalg.generic}\big\{\mathcal{I}_Q, \{\phi_\ell^Q\}, \mathbf{z} := g(\mathbf{y})\big\},$$

then the fused operation is:

$$\texttt{linalg.generic}\big\{\mathcal{I}, \{\tilde{\phi}_\ell\}, \mathbf{z} := g\big(f(x)\big)\big\}.$$

Listing 3 shows the fused `subf` and `exp` generics from Listing 1. The example in Listing 3 is rather sim-

```
1    ...
2    %7 = linalg.generic
3            {indexing_maps = [affine_map<(d0) -> (d0)>,
4                              affine_map<(d0) -> (d0)>],
5             iterator_types = ["parallel"]}
6         ins(%3 : tensor<16384xf32>) outs(%4 : tensor<16384xf32>) {
7       ^bb0(%in: f32, %out: f32):
8           %11 = arith.subf %in, %extracted : f32
9           %12 = math.exp %11 : f32
10          linalg.yield %12 : f32
11      } -> tensor<16384xf32>
12   ...
```

Listing 3: IR of Softmax excerpt after Fusion

ple (identity-map) for illustration but in general fusion has to consider - broadcast, transposes, and re-computations.

## 4.4 Tiling for Memory Hierarchy

Tiling partitions large tensors that typically reside on the slower but large shared memory (DDR) into smaller `tiles` that are copied onto the smaller but faster local Tightly Coupled Memory (TCM). Computations such as vectorized HVX operations are then done on data that is on TCM. High reuse of data that is brought or generated on TCM is important to offset the latency of transfer latency between DDR and TCM. Operator fusion that we saw in the previous section helps improve vector register level reuse. Reuse across fused operators is facilitated be retaining tensor results in TCM.

Let us look at the tiling mechanism in Hexagon-MLIR in more details. Given a `linalg.generic` with iteration domain $\mathcal{I} \subseteq \mathbb{Z}^d$, a tiling vector $\mathbf{t} = (t_1, \ldots, t_d)$, and an optional interchange vector, tiling introduces two levels of iteration: an outer loop consisting of `scf::for` or `scf::forall`, and the body of the loop remains a `linalg.generic` that now operates on a tile on TCM.

To exploit memory hierarchy, the tiling pass in Hexagon-MLIR introduces two distinct memory spaces: DDR and TCM. Data transfers between DDR and TCM are performed via DMA operations. Tiling is applied at the `Linalg-on-Tensor`

level, where the pass inserts explicit data movement operations using `bufferization.alloc_tensor` with memory-space attributes. These attributes enable the compiler to later generate DMA transfers between DDR and TCM. Furthermore, the tiling pass incorporates loop interchange to position vectorizable dimensions innermost and annotates tiled loops to facilitate subsequent optimizations such as software pipelining. We will come back to softmax but for now let us turn to another interesting kernel for now - Gaussian Eror Linear Unit (GELU), to explore tiling, vectorization and double buffering.

$$\text{GELU}(x) \approx 0.5x \left( 1 + \tanh \left( \sqrt{\frac{2}{\pi}} (x + 0.044715x^3) \right) \right)$$

The GELU kernel after fusion and tiling is shown in Listing 4. Bufferization has not happended yet so the tiled generic still operates on tensors. Notice the `scf.for` loop, where in each iteration a slice of the larger input tensor is calculated using via `tensor.extract_slice`. That slice is then copied from the default memory space to TCM. The body of the `linalg.generic` operation computes on the data brought into TCM. The generic's body although presently looks like scalar code, after further transformations it will become multi-threaded HVX operations. Next, we will multi-thread this tiled generic.

## 4.5 Multi-threading

```
1   #map = affine_map<(d0) -> (d0)>
2   module attributes {llvm.target_triple = "hexagon"} {
3   func.func @gelu_kernel(%arg0: memref<*xf32> {tt.divisibility = 16 : i32}, ..) {
4     %cst_2 = arith.constant 7.978850e-01 : f32
5     ..
6     %4 = scf.for %arg8 = %c0 to %c1048576 step %c262144
7          iter_args(%arg9 = %0) -> (tensor<1048576xf32>) {
8        %extracted_slice = tensor.extract_slice %3[%arg8] [262144] [1]
9              : tensor<1048576xf32> to tensor<262144xf32>
10       %5 = bufferization.alloc_tensor()
11             copy(%extracted_slice) {memory_space = 1 : i64} : tensor<262144xf32>
12       ..
13       %7 = linalg.generic
14             {indexing_maps = [#map, #map], iterator_types = ["parallel"]}
15             ins(%5 : tensor<262144xf32>) outs(%6 : tensor<262144xf32>) {
16       ^bb0(%in: f32, %out: f32):
17         ..
18         %13 = arith.mulf %12, %cst : f32
19         %14 = math.exp %13 : f32
20         %15 = arith.addf %14, %cst_0 : f32
21         %16 = arith.subf %cst_0, %14 : f32
22         %17 = arith.divf %16, %15 : f32
23         ..
24         linalg.yield %20 : f32
25       } -> tensor<262144xf32>
26       %inserted_slice = tensor.insert_slice %7 into %arg9[%arg8] [262144] [1]
27             : tensor<262144xf32> into tensor<1048576xf32>
28       scf.yield %inserted_slice : tensor<1048576xf32>
29     } {all_parallel, tiled_generic}
30     ..
31   }
32 }
```

Listing 4: GELU after Fusion and Tiling

Qualcomm NPUs have hardware multi-threading to support parallel execution of Hexagon Vector eXtensions (HVX) instructions. A vector context is a set of resources that enables independent execution of HVX instructions, including a vector register file and a vector predicate file. Multiple hardware threads can run in parallel, each utilizing a different vector context. In this section we will show multi-threading of the GELU kernel we tiled in previous section. We have two multi-threading schemes. One uses scf::for starting point and the other uses linalg.generic. We will show the latter one here. Our multi-threading uses MLIR's Async dialect which provides a structured fork–join IR:

- async.execute creates an async region, returning a token (and optionally values) with async.yield.

- async.await waits on a token/value by suspending the current co-routine.

- async.create_group, async.add_to_group, and async.await_all implement barriers over multiple tokens.

The asynchronous representation lowers to

```
1   ...
2   %12 = async.create_group %11 : !async.group
3   scf.for %arg9 = %c0_10 to %c262144_11 step %c8_12 {
4     %token = async.execute {
5       %subview_14 = memref.subview %2[%arg9] [8] [1]
6               : memref<262144xf32, 1> to memref<8xf32, strided<[1], offset: ?>, 1>
7       ..
8       scf.for %arg10 = %c0 to %c8 step %c1 {
9         %14 = memref.load %subview_14[%arg10] : memref<8xf32, strided<[1], offset: ?>, 1>
10        %15 = arith.mulf %14, %cst_1 fastmath<fast> : f32
11        ...
12        %21 = math.exp %20 fastmath<fast> : f32
13        %22 = arith.addf %21, %cst_0 fastmath<fast> : f32
14        %23 = arith.subf %cst_0, %21 fastmath<fast> : f32
15        %24 = arith.divf %23, %22 fastmath<fast> : f32
16        ...
17        memref.store %27, %subview_15[%arg10] : memref<8xf32, strided<[1], offset: ?>, 1>
18      }
19      async.yield
20    }
21    %13 = async.add_to_group %token, %12 : !async.token
22  }
23  async.await_all %12
24  ...
```

Listing 5: GELU after async-threads creation

LLVM coroutines (`llvm.coro.id`, `llvm.coro.begin`, `llvm.coro.suspend`, `llvm.coro.end`) together with the MLIR async runtime. Hexagon-MLIR exploits HVX multi-threading with a two-stage pipeline:

1. *Form-Virtual-Threads*: partitions work and encodes parallel structure using `scf.forall`.

2. *Form-Async-Threads*: lowers that structure to MLIR's Async dialect. This approach keeps parallel semantics declarative and analyzable.

**Stage-1:** Analyzes linalg.generic to decide if parallelization is profitable using a polytope size heuristic. It performs a tiling over parallel iterators to produce scf.forall (virtual threads) using block distribution when ranges divide evenly; using block-cyclic for thread size balancing.

**Stage-2:** Each `scf.forall` is rewritten into an explicit asynchronous fork–join pattern. First, `async.create_group(N)` establishes a barrier for `N` threads. A surrounding `scf.for` then iterates over the tiles, and each iteration spawns:

```
async.execute { /* body of the forall tile
*/ } → token
```

The resulting token is added to the group using:

```
async.add_to_group(group, token)
```

Finally, `async.await_all(group)` enforces completion before any dependent work proceeds. This approach preserves the original parallel semantics while preparing for coroutine-based code generation.

Listing 5 shows the IR of theGELU kernel after multi-threading. Its input IR was the tiled GELU kernel of Listing 4.

## 4.6 Double Buffering

Double buffering in software pipelining is a technique of overlapping computation with data movement, thereby mitigating memory latency in hierarchical memory systems. This approach leverages two buffers (ping and pong) — while one is actively engaged in computation the other is concurrently being filled (and then the roles swap) - enabling data transfers and computations to happen in parallel.

In AI/ML workloads, where tensor operations dominate and data sizes often exceed the capacity of fast local memory, double buffering becomes critical. The interleaving of communication and computation transforms memory-bound kernels towards latency-tolerant pipelines, exploiting parallelism across memory and compute units. Similar overlapping techniques have been explored in distributed AI systems, where compiler-based approaches enable native overlapping optimizations for distributed workloads [29].

The implementation in hexaogn-mlir introduces double buffering in two layered passes that separate structural software pipelining from asynchronous data movement.

**Stage 1: Structural Transformation** Stage 1 finds canonical single-buffered, tiled loops `scf.for` that are in "normal form". The loops have attributes indicating they originate from tiled `linalg.generic`. The tiling leaves per-tile schedule triplets

$$\texttt{memref.subview} \rightarrow \texttt{memref.alloc} \rightarrow \texttt{memref.copy},$$

along with the compute region and the write-back sequence. Double buffering pass synthesizes a ping-pong flow: a guarded prologue that prefetches the first tile into the ping buffers, and a rebuilt kernel with two mutually exclusive sub-kernels (annotated `db_ping_kernel`/`db_pong_kernel`). Each sub-kernel issues a prefetch for the next iteration, clones the compute with `alloc` remapped to the current ping/pong buffers, and reconstructs the stores by re-materializing subviews at the current induction variable. A memory-resident boolean `toggle` drives ping↔pong selection and is updated each iteration, while attributes like `db_generic` (unique ID) and `db_prologue`/`db_prefetch` provide IR anchors for

downstream passes. This keeps the transformation deterministic and hazard-free, preserving the original tiling semantics and establishing the overlap between communication (prefetch) and computation without yet committing to a specific asynchronous transport.

**Stage 2: Asynchronous DMA Integration** Stage 2 materializes true asynchronous transfers and synchronization. It parses the Stage 1 IR using the injected attributes to recover the prologue, the kernel, and the ping/pong sub-schedules. . The pass then rewrites all preload and prefetch `memref.copy` operations into `memref.dma_start`—allocating distinct ping/pong tags so that each phase can be waited on independently—and inserts `memref.dma_wait` immediately after the corresponding prefetch blocks to ensure data is resident before compute. For the store-back path, it similarly replaces each `memref.copy` with `dma_start` followed by `dma_wait` (using separately allocated store tags), and emits balanced deallocations for all tags after the kernel. This design decouples legality/scheduling (Stage 1) from transport semantics (Stage 2), enabling latency hiding through explicit DMA start/wait without perturbing the compute IR.

Listing 6 shows the IR after double-buffering. The IR is annotated in C++ style comments and some variables given custom names for easier comprehension.

We have been illustrating our flow with GELU IRs so lets recap here. Listing 3 showed at an earlier stage, right after fusion. Fusion happens on `linalg-on-tensor`.Then we did tiling to bring slices from DDR to TCM in Listing 4. The tiled code forms bridge between tensors and memrefs (buffers). Some other passes such as bufferization run subsequently. Then we showed multi-threading for HVX in the IR in Listing 5. After that comes double buffering which we showed in this section. There is still more lowering that happens after this leading finally to LLVM-IR with runtime calls.

## 4.7 Math Library

In AI/ML workloads, replacing transcendental functions like `exp` with low-degree polynomial approximations is a common way to trade controlled numerical

error for gains in throughput. We invoke Qualcomm Hexagon Library (QHL) – a set of math libraries pre-optimized for vectorized computations.

MLIR also provides `polynomial approximation` expansions of math operations for fast approximations that do not rely on any custom library functions. The general approach e.g. to approximate $e^x$, firstly express it as

$$e^x = e^{\epsilon + b} = e^{\epsilon + n \log 2} = e^\epsilon \cdot 2^n.$$

Then, as $\epsilon$ is chosen to be small, Taylor series up to some terms can be used to compute $e^\epsilon$. Horner's formula is often used to reduce the number of multiplications and additions to compute the polynomial.

In addition to above, there are other optimizations and algorithms used to produce fast implementations of math operations. One example is **fast inverse square root**.

## 4.8   Optimizations and Lowering

The flow discussed in previous section leads to successful compilation of PyTorch and Triton kernels by our tool. We did not cover all the passes and optimizations that are in place or in development. We also did not discuss MLIR core passes that we use e.g. bufferization, deallocation, and lowering to llvm as they are quite well known. We identify convolutions and matrix multiplications earlier on in our flow and generally they go to a dedicated engine for matrmix multiplication. The data needs to be laid out in a particular way and we use `linalg.pack/unpack` to do the data-space transformations. To reduce the number of packs/unpacks we use layout-propagation.

## 5   Results and Analysis

The results in this section are organized to map the compiler optimizations described in previous sections to the performance improvements shown in the accompanying tables and bar charts. We enable individually and in combination set of passes in the Hexagon-MLIR pipeline: fusion (F) removes intermediate materialization and improves locality; tiling (T) moves working sets into TCM; HVX vectorization

(V) exploits wide SIMD execution; multi-threading (M) distributes tiled work across hardware HVX contexts; and double buffering (DB) overlaps DMA transfers with computation. This structured presentation allows each optimization described in the technical sections to be correlated with its observed runtime effect, and highlights how these passes interact—sometimes constructively, sometimes competitively—as kernels transition from compute-bound to memory-bound regimes.

## 5.1   Vectorization Speedups

Table 1 summarizes speedups from HVX vectorization across some kernels. GELU achieves $63.9\times$ on `float16` (16,384 elements) and $16.1\times$ on `float32`, indicating that HVX width utilization and data type packing are primary drivers of throughput.

Gains from vectorization of RMS-norm reaches $46.5\times$ on $127{\times}513$ shape, confirming that stencil linalg.elementwise dominated computations, even with some reductions, gain from vectorization. Just to recap, RMS-Norm is defined as:

$$\mathrm{RMSNorm}(x) = \frac{x}{\sqrt{\frac{1}{d}\sum_{i=1}^{d} x_i^2 + \epsilon}} \cdot g$$

The computation part of its triton kernel is shown in Listing 7.

```
sq = tl.sum(x * x, axis=0) / NUM_COLS
rms = tl.sqrt(sq + EPSILON)
y = (x / rms) * g
```

Listing 7: Extract from Triton RMS Norm Kernel

Table 1 also includes results for the **Sigmoid Linear Unit (SiLU)** activation. SiLU is defined as:

$$\mathrm{SiLU}(x) = x \cdot \sigma(x) = x \cdot \frac{1}{1 + e^{-x}},$$

We get good $5 - 8x$ gains. The computation is dominated by the term $e^x$. For `float16`, performance decreases when operations are internally promoted to `float32`.

Table 1: Speedups by vectorization.

| Test | Speedup | Notes |
|---|---|---|
| Gaussian Error Linear Unit (GELU) | 63.9x | float16, elements = 16384 |
| Gaussian Error Linear Unit (GELU) | 16.1x | float32, elements = 16384 |
| Double Buffered GELU | 15.4x | float32, element= 16 x 16384, double buffered |
| Flash Attention | 4.7x | float32, N_CTX=1024, DIM_HEAD=64, BLOCK_N=64 |
| Root Mean Square Norm (RMS-norm) | 46.5x | float16, elements = 127 x 513, |
| Sigmoid Linear Unit (SiLU) | 4.8 | float16, elements = 16384 |
| Sigmoid Linear Unit (SiLU) | 7.1 | float32, elements = 16384 |
| Vec-Add-2D | 40.6x | float32, elements = 64 x 16384 |

## 5.2 Single vs. Multi-threading

Figure 2 compares single-threaded execution against HVX multi-threading over HVX contexts across increasing problem sizes (8K $\rightarrow$ 1M elements). At small sizes ($\leq$ 16K), multi-threaded runtimes are higher due to thread start-up, token management, and barrier overheads. Beyond 32K elements, multi-threading consistently outperforms single-threaded execution with speedups ranging from 2.28$\times$ (32K) to 3.95$\times$ (512K), and 3.40$\times$ at 1M elements (see the speedup bar chart). These gains emerge when the parallel payload amortizes the fixed costs of `async.execute` and `async.await_all`. At 1M elements, speedup dips slightly compared to 512K, reflecting increasing pressure on TCM and DMA queues. In such case re-tuning could be beneficial. Figure 3 presents the speedup, defined as the ratio of single-threaded execution time to multi-threaded execution time, in a more interpretable format.

## 5.3 Double Buffering

Figure 5 shows improvements with vectorization, multi-threading (MT), and double-buffering (DB). Vectorization gives significant gains, but if we take it as a given (base case), then MT and DB improve performance further. Figure 7 shows the gains from double buffering and other optimizations for GELU. Similarly, Figure 6 shows improvements for the compute-intensive exponent series kernel.

## 5.4 Multi-pass Interaction

We now turn to multi-pass interactions. Techniques like vectorization and multi-threading mainly reduce computation time, which makes memory transfer overhead more pronounced. Double buffering helps by overlapping data movement with computation, boosting overall performance. However, the net effect can be that an optimized kernel becomes memory-bound. Figure 4 illustrates an idealized case: when memory transfers completely dominate (100%), double buffering offers no benefit; likewise, if computation dominates entirely, the technique is again ineffective. In practice, most workloads fall between these extremes and so double-buffering provides gain, and secondary effects—such as overheads and scheduling — play an important role. Applying optimization passes selectively provides additional insight into kernel characteristics and these secondary interactions.

## 6 Related Work

Li et al. [12] provide a comprehensive survey of deep learning (DL) compilers, analyzing their design architectures and optimization strategies. Although published five years ago, the concepts remain relevant today. The paper explains how DL compilers transform models from various frameworks into optimized code for heterogeneous hardware, focusing on multi-level intermediate representations (IR) and both frontend and backend optimizations. It discusses key techniques such as operator fusion, memory planning, layout transformation, hardware mapping, and auto-

tuning.

Yifan et al [28], present Neptune, a tensor compiler that supports loop fusion algorithms for reduction operators, that satisfies data dependency by recomputation and automatically derives required algebraic repairs. Two instances, Rolling Update Fusion and Split-K Fusion, are particularly suited for optimizing attention-like operators. They implemented Neptune on top of the Apache TVM schedule tensor compiler and the Triton tile tensor compiler. A more recent paper [30] by the same author on *Strategies for Graph Optimization in Deep Learning Compilers* delves into front-end optimization techniques, which according to the author are instrumental in refining the structure of computational graphs, thereby significantly bolstering the efficiency of neural network training and inference phases.

Chin et al, in the paper [14] *HloEnv: A Graph Rewrite Environment for Deep Learning Compiler Optimization Research*, analyze the impact on the performance of an optimization pass/pipeline on the HLO dataset from two perspectives: the proportion of the dataset affected by that pass, and the average change in performance as a result of that pass (runtime ratio w/ and w/o the pass).

Optimizing both computation speed and memory size and transfer requirements is important for an AI compilation package. An XLA compiler extension that adjusts the computational data-flow representation of an algorithm according to a user-specified memory limit is discussed in [3].Benoit et al. [22] present MODeL, an ILP approach that optimizes the lifetime and memory location of the tensors used to train neural networks to reduce the memory usage of models.

## 7  Conclusion

This work presented Hexagon-MLIR, an MLIR-based compilation stack designed for Qualcomm's Hexagon NPU for AI workloads. By combining MLIR's extensible infrastructure and hardware-aware passes, Hexagon-MLIR delivers a systematic lowering pipeline for PyTorch models and Triton kernels. Key optimizations—such as operator fu-

sion, tiling, HVX vectorization, multi-threaded execution, and double buffering—enable efficient utilization of Tightly Coupled Memory (TCM), asynchronous DMA transfers, and specialized math libraries. Experimental results demonstrate performance gains across representative kernels, validating the effectiveness of our approach. Beyond performance, this work illustrates how MLIR's design principles—structured IRs, declarative transformations, and staged lowering—provide a strong foundation for building compilers targeting AI workloads. While the current implementation achieves promising results, ongoing efforts focus on expanding coverage for emerging models, refining scheduling heuristics, and improve feature support and performance. Ultimately, Hexagon-MLIR aims to deliver a scalable, high-performance solution for deploying next-generation AI workloads on Hexagon NPUs.

## References

[1] Torch-mlir: Bridging pytorch and mlir ecosystems. GitHub repository, 2021–2025. `https://github.com/llvm/torch-mlir`.

[2] J. Ansel et al. Pytorch 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*, 2024.

[3] Artem Artemev, Yuze An, Tilman Roeder, and Mark Van der Wilk. Memory safe computations with xla compiler, 2022.

[4] Tianqi Chen, Thierry Moreau, Ziheng Jiang, Lianmin Zheng, et al. Tvm: An automated end-to-end optimizing compiler for deep learning, 2018.

[5] Tri Dao, Daniel Haziza, Francisco Massa, and Grigory Sizov. Flash-decoding for long-context inference. `https://crfm.stanford.edu/2023/10/12/flashdecoding.html`, 2023.

[6] DeepSeek-AI. Deepseek-v3 technical report, 2024. Model checkpoints: `https://github.com/deepseek-ai/DeepSeek-V3`.

[7] J. A. Fisher. Trace scheduling: A technique for global microcode compaction. *IEEE Transactions on Computers*, 1981.

[8] Renato Golin, Lorenzo Chelini, Adam Siemieniuk, Kavitha Madhu, Niranjan Hasabnis, Hans Pabst, Evangelos Georganas, and Alexander Heinecke. Towards a high-performance ai compiler with upstream mlir, 2024.

[9] W. M. Hwu et al. The superblock: An effective technique for vliw and superscalar compilation. *The Journal of Supercomputing*, 1993.

[10] Qualcomm Technologies Inc. Unlocking on-device generative ai with an npu and heterogeneous computing. `https://www.qualcomm.com/processors/ai-engine`, 2024.

[11] Chris Lattner, Mehdi Amini, Uday Bondhugula, Albert Cohen, Andy Davis, Jacques Pienaar, River Riddle, Tatiana Shpeisman, Nicolas Vasilache, and Oleksandr Zinenko. Mlir: A compiler infrastructure for the end of moore's law, 2020.

[12] Mingzhen Li, Yi Liu, Xiaoyan Liu, Qingxiao Sun, Xin You, Hailong Yang, Zhongzhi Luan, Lin Gan, Guangwen Yang, and Depei Qian. The deep learning compiler: A comprehensive survey. *IEEE Transactions on Parallel and Distributed Systems*, 32(3), 2021. Available at: `https://arxiv.org/abs/2002.03794`.

[13] Microsoft. triton-shared: Shared middle-layer for triton compilation. GitHub repository, 2019–2025. `https://github.com/microsoft/triton-shared`.

[14] Chin Yang Oh, Kunhao Zheng, Bingyi Kang, Xinyi Wan, Zhongwen Xu, Shuicheng Yan, Min Lin, and Yangzihao Wan. Hloenv: A graph rewrite environment for deep learning compiler optimization research. In *36th Conference on Neural Information Processing Systems (NeurIPS 2022) ML for Systems Workshop*, 2022.

[15] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, et al. Pytorch: An imperative style, high-performance deep learning library, 2019.

[16] PyTorch Team. Helion: A high-level dsl for triton kernel development. `https://pytorch.org/blog/helion/`, 2025.

[17] Qualcomm Technologies, Inc. Compile triton & pytorch for hexagon npu with open source hexagon-mlir. `https://www.qualcomm.com/developer/blog/2026/02/`, February 2026.

[18] Qualcomm Technologies, Inc. Hexagon-mlir: Open-source mlir-based compilation stack for hexagon npu. `https://github.com/qualcomm/hexagon-mlir`, 2026.

[19] J. Ragan-Kelley, C. Barnes, A. Adams, S. Paris, F. Durand, and S. Amarasinghe. Halide: a language and compiler for optimizing parallelism, locality, and recomputation in image processing pipelines. *ACM SIGPLAN Notices*, 48(6), 2013. Available at: `https://arxiv.org/pdf/1802.04730`.

[20] Burkhard Ringlein, Jan van Lunteren, Radu Stoica, and Thomas Parnell. The anatomy of a triton attention kernel, 2025.

[21] Nadav Rotem et al. Glow: A graph lowering compiler techniques for neural networks, 2018.

[22] Benoit Steiner, Mostafa Elhoushi, Jacob Kahn, and James Hegarty. Model: memory optimizations for deep learning. In *Proceedings of the 40th International Conference on Machine Learning (ICML'23)*, 2023.

[23] Masahiro Tanaka, Du Li, Umesh Chand, Ali Zafar, Haiying Shen, and Olatunji Ruwase. Deepcompile: A compiler-driven approach to optimizing distributed deep learning training, 2025.

[24] P. Tillet, H.-T. Kung, and D. D. Cox. 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 (MAPL@PLDI)*, pages 10–19, Phoenix, AZ, USA, 2019. ACM.

[25] Nicolas Vasilache et al. Tensor comprehensions: Framework-agnostic high-performance machine learning abstractions, 2018.

[26] Ashish Vaswani et al. Attention is all you need. In *Advances in Neural Information Processing Systems 30 (NIPS 2017)*, Long Beach, CA, USA, 2017.

[27] Hanhan Wang. Data-tiling in iree: Achieving high performance through compiler design. LLVM Developers Meeting, June 2025. Available at: `https://llvm.org/devmtg/2025-06/slides/technical-talk/wang-data-tilling.pdf`.

[28] Yifan Zhao, Egan Johnson, Prasanth Chatarasi, Vikram Adve, and Sasa Misailovic. Neptune: Advanced ml operator fusion for locality and parallelism on gpus, 2024.

[29] Size Zheng, Wenlei Bao, Xin Liu, et al. Triton-distributed: Compiler-based overlapping optimizations for distributed ai systems, 2025. Project page: `https://github.com/ByteDance-Seed/Triton-distributed`.

[30] Yijian Zheng. Strategies for graph optimization in deep learning compilers. In *2024 International Conference on Interactive Intelligent Systems and Techniques (IIST)*, 2024.

[31] Zhijian Zhuo, Ya Wang, Yutao Zeng, Xiaoqing Li, Xun Zhou, and Jinwen Ma. Polynomial composition activations: Unleashing the dynamics of large language models, 2024. School of Mathematical Sciences, Peking University.

```
1  // -----// IR Dump After hexagon-double-buffer-generic-s2   // ------------- //
2  func.func @gelu_kernel(%arg0: memref<*xf32> {tt.divisibility = 16 : i32}, ..) {
3    %false = arith.constant false
4    // allocate ping-pong buffers on TCM
5    %ping_buffer = memref.alloc() {alignment = 2048 : i64} : memref<262144xf32, 1>
6    %pong_buffer = memref.alloc() {alignment = 2048 : i64} : memref<262144xf32, 1>
7    %dma_ping_tag = memref.alloc() : memref<1xi32>
8    %dma_pong_tag = memref.alloc() : memref<1xi32>
9    ...
10   /// Double Buffering - Prologue
11   scf.if %loop_executes_at_least_once {
12     %subview = memref.subview %Input[0] [262144] [1]
13                : memref<1048576xf32, ..>> to memref<262144xf32,..>>
14     memref.dma_start %subview[%c0], %ping_buffer[%c0],%c262144,%dma_ping_tag[%c0]
15           : memref<262144xf32,..>>, memref<262144xf32,1>, memref<1xi32>
16     ...
17   } {db_generic = 0 : i64, db_prologue}
18   /// Double Buffering - Kernel (main loop)
19   scf.for %arg8 = %c0 to %c1048576 step %c262144 {
20     %is_ping_stage = memref.load %ping_pong_toggle[] : memref<i1>
21     ...
22     scf.if %is_ping_stage {
23       scf.if %is_not_last_iteration {
24         // Prefetch to pong buffers for next iteration.
25         %subview_17 = memref.subview %Input[%4] [262144] [1]
26                        : memref<1048576xf32, ..>> to memref<262144xf32, ..>>
27         memref.dma_start %subview_17[%c0], %pong_buffer[%c0],
28                  %c262144, %dma_pong_tag[%c0]
29               : memref<262144xf32, ..>>, memref<262144xf32, 1>, memref<1xi32>
30         ...
31       } {db_prefetch}
32       memref.dma_wait %dma_ping_tag[%c0], %c262144 : memref<1xi32>
33       ...
34       // Do multi-threaded HVX execution on ping buffers on TCM
35       scf.forall (%arg9) = (0) to (262144) step (65536) {
36         %subview_17 = memref.subview %ping_buffer[%arg9] [65536] [1]
37                        : memref<262144xf32, 1> to memref<65536xf32, ...>, 1>
38         ...
39         scf.for %arg10 = %c0 to %c65536 step %c32 {
40           %subview_19 = memref.subview %subview_17[%arg10] [32] [1]
41               : memref<65536xf32,..>, 1> to memref<32xf32, ..>, 1>
42           ...
43           %15 = math.exp %14 fastmath<fast> : vector<32xf32>
44           ...
45           vector.transfer_write ...   : vector<32xf32>, memref<32xf32, ..>, 1>
46         }
47       }
48       memref.dma_start ... // write-back ping buffer results.
49       memref.store %false, %ping_pong_toggle[] : memref<i1> // ping-pong state
50     } {db_ping_kernel}
51     scf.if %is_pong_stage {
52       ... // pong-stage
53     } {db_pong_kernel}
54     ... // deallocations etc.
55 }
```

15
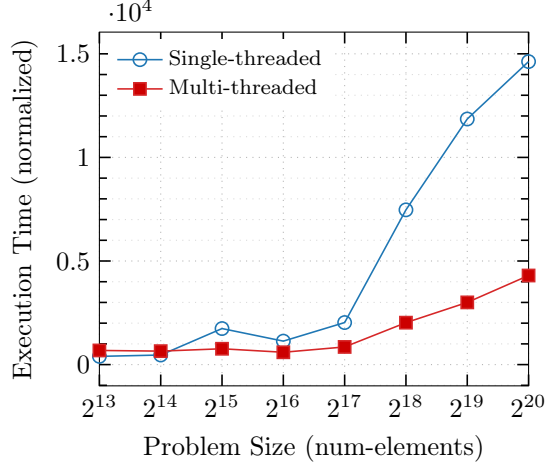
Listing 6: GELU after double-buffering

Figure 2: Execution times of single-threaded vs multi-threaded GELU kernel implementations.
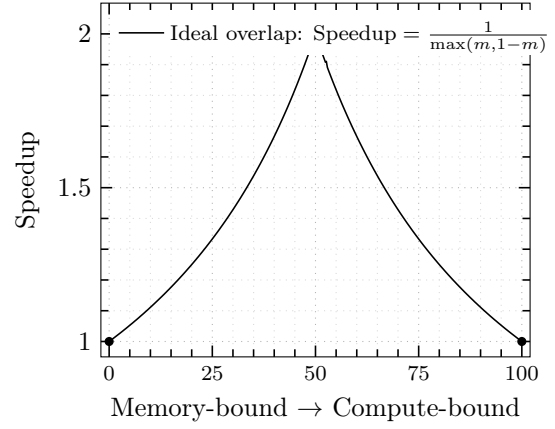


Figure 4: Speedup for double buffering with perfect overlap as memory fraction varies.
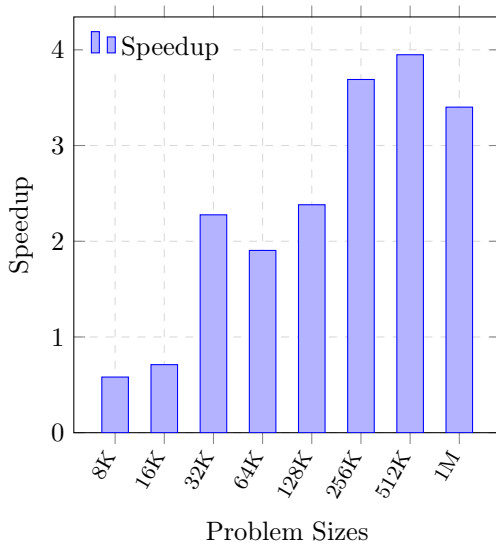


Figure 3: Speedup of multi-threaded GELU over increasing problem-sizes.
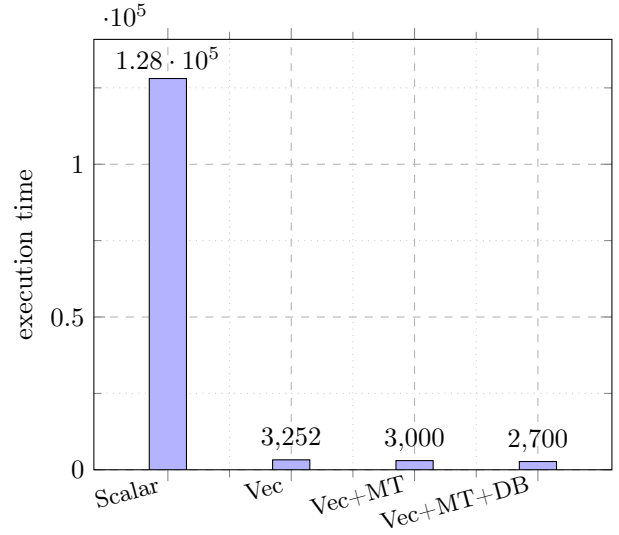


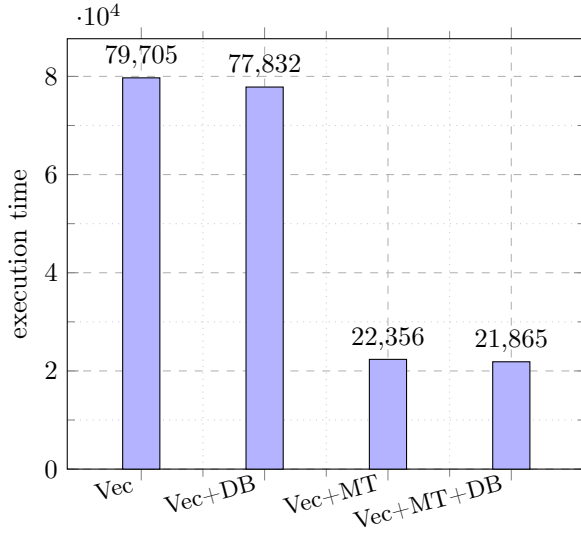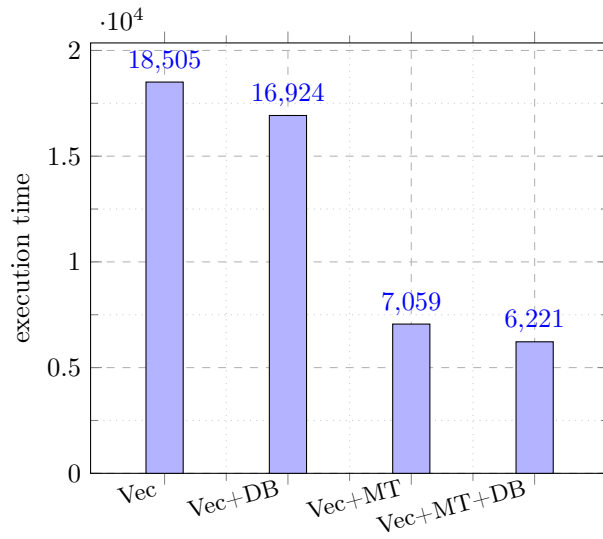Figure 5: Vector-Add 2D performance across optimization passes.

Figure 6: Multi-pass analysis of Exponent Series.



Figure 7: Multi-pass analysis of GELU.