

## ArcSysU Weekly Paper Sharing

# Hidet: Task-Mapping Programming Paradigm for Deep Learning Tensor Programs

Yaoyao Ding<sup>\*†</sup>  
University of Toronto  
Toronto, Canada  
yaoyao@cs.toronto.edu

Yizhi Liu  
Amazon Web Services  
Santa Clara, USA  
yizhiliu@amazon.com

Cody Hao Yu  
Amazon Web Services  
Santa Clara, USA  
hyuz@amazon.com

Yida Wang  
Amazon Web Services  
Santa Clara, USA  
wangyida@amazon.com

Bojian Zheng<sup>†</sup>  
University of Toronto  
Toronto, Canada  
bojian@cs.toronto.edu

Gennady Pekhimenko<sup>†</sup>  
University of Toronto  
Toronto, Canada  
pekhimenko@cs.toronto.edu

## ASPLOS'23

## Contents

1 Background

2 Motivation

## 3 Design

## 4 Evaluation

5 Conclusion

中山大學  
SUN YAT-SEN UNIVERSITY

## About the Authors

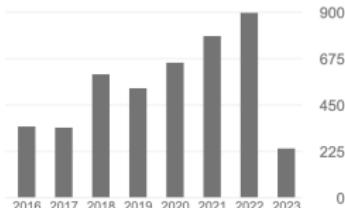
Yaoyao Ding

- Homepage: <https://www.cs.toronto.edu/~yaoyao/>
  - 1st Year PhD at University of Toronto
  - Bachelor from Shanghai Jiao Tong University (at ACM Honors Class)
  - MLSys'21: IOS: Inter-Operator Scheduler for CNN Acceleration

Gennady Pekhimenko

- Homepage: <http://www.cs.toronto.edu/~pekhimenko/>

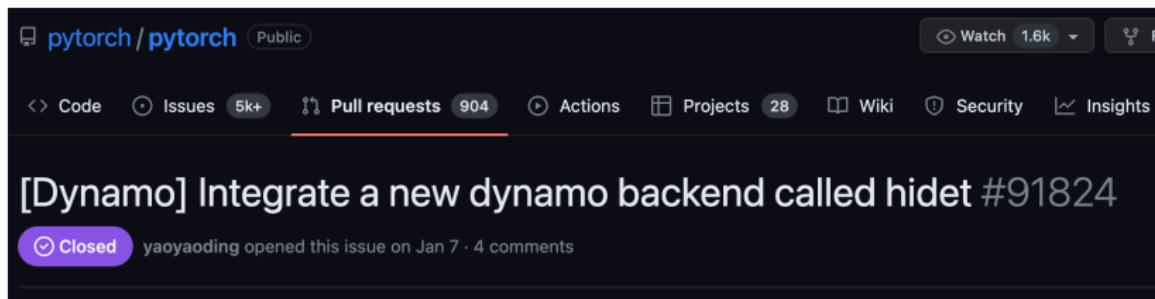
	All	Since 2018
Citations	4712	3716
h-index	32	32
i10-index	50	48



中山大學  
SUN YAT-SEN UNIVERSITY

# About the Project

- Repository: <https://github.com/hidet-org/hidet>.
- Plug-and-play with PyTorch<sup>1</sup> ( 20k LOC in Python and C++):  
`torch.compile(model, backend='hidet')`.
- Proposed a PR to PyTorch, but rejected.



A screenshot of a GitHub repository page for "pytorch/pytorch". The repository is public, has 1.6k stars, and 904 pull requests. The "Pull requests" tab is selected. A prominent pull request is shown with the title "[Dynamo] Integrate a new dynamo backend called hidet #91824". The status is "Closed" by yaoyaoding on Jan 7, with 4 comments. The repository has 5k+ issues, 28 projects, and 1 Wiki page.



PyTorch 2.0 experimental preview released on 2022/12/02, the Hidet version on artifact evaluation used PyTorch 1.11. By now it has integrated with PyTorch 2.0's compile feature

# TVM's Loop-Oriented Scheduling

## Computation

### ① Generate Default Tensor Program

## Default Program

```
C = compute([M, N], lambda i, j: reduce([K], A[i, k] * B[k, j]))  
  
for i in range(M):  
    for j in range(N):  
        for k in range(K):  
            C[i, j] += A[i, k] * B[k, j]
```

### ② Apply Declarative Scheduling Primitives

Schedule Primitives	Original Program	Scheduled Program
fuse(i, j)	for i in range(128): for j in range(4): body(i, j)	for ij in range(512): body(ij / 4, ij % 4)
split(i, 128)	for i in range(512): body(i)	for oi in range(4): for ii in range(128): body(oi * 128 + ii)
reorder(i, j)	for i in range(128): for j in range(4): body(i, j)	for j in range(4): for i in range(128): body(i, j)
bind(i, threadIdx.x)	for i in range(128): body(i)	body(threadIdx.x)

# TVM's Loop-Oriented Scheduling (cont.)

```
for i in range(M):
    for j in range(N):
        for k in range(K):
            C[i, j] += A[i, k] * B[k, j]
```

Default Program



```
for io in range(M/64):
    for ii in range(64):
        for oj in range(N/64):
            for ij in range(64):
                for k in range(K):
                    i, j = io * 64 + ii, jo * 64 + ij
                    C[i, j] += A[i, k] * B[k, j]
```

reorder(io, oj, ii, ij)

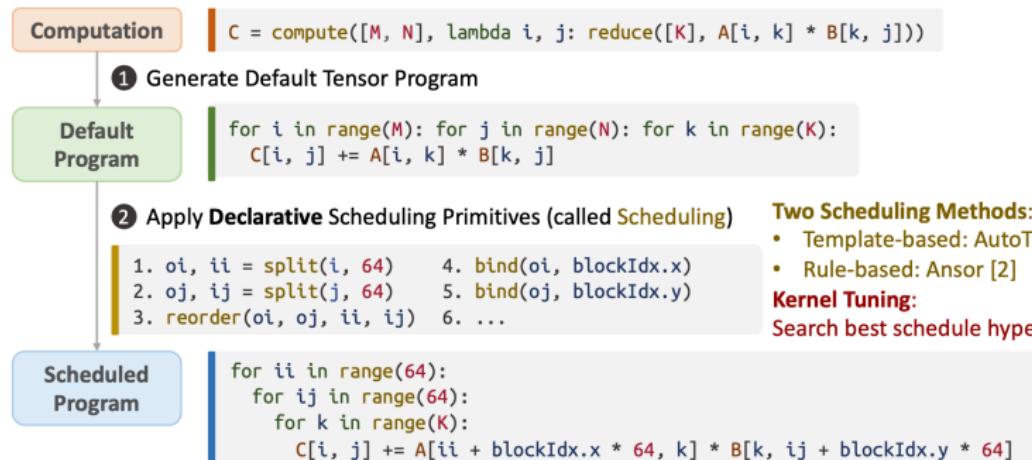
```
for ii in range(64):
    for ij in range(64):
        for k in range(K):
            i = blockIdx.x * 64 + ii
            j = blockIdx.y * 64 + ij
            C[i, j] += A[i, k] * B[k, j]
```

Scheduled Program



```
bind(oi, blockIdx.x)
bind(oj, blockIdx.y)
```

# TVM's Loop-Oriented Scheduling (cont.)



**Two Scheduling Methods:**

- Template-based: AutoTVM [1]
- Rule-based: Ansor [2]

**Kernel Tuning:**  
Search best schedule hyper-parameter

## Declarative Loop-Oriented Scheduling

Background  
oooooooo

Motivation  
●oooooooo

Design  
oooooooooooo

Evaluation  
oooooooo

Conclusion  
oooo

# Contents

## 1 Background

## 2 Motivation

## 3 Design

## 4 Evaluation

## 5 Conclusion

# Non Loop-Oriented Optimization

An example of tiled matrix multiply:

```
1 def matmul(A: fp32[M, K], B: fp32[K, N], C: fp32[M, N]):  
2     SmemA, SmemB = shared fp32[64, 8], fp32[8, 64]  
3     RegsC = local fp32[...]  
4     ... # Step ①: Calculate sub-task offset  
5     for k0 in range(128): # Iterate each K tile  
6         # Step ②: Load A and B frag. to SmemA and SmemB  
7         SmemA, SmemB = cooperative_load(A, B, k0)  
8         sync_threads()  
9         # Step ③: Block MMA (RegsC = SmemA * SmemB + RegsC)  
10        RegsC = block_mma(SmemA, SmemB, RegsC)  
11        sync_threads()  
12        ... # Step ④: Write back C (RegsC => C)
```

# Non Loop-Oriented Optimization (cont.)

Use double buffer to overlap computation and data transfer:

```
1 RegsA, RegsB = register fp32[...], fp32[...]
2 SmemA, SmemB = shared fp32[2, 64, 8], fp32[2, 8, 64]
3 ... Two Buffers for A & B
4 SmemA[0], SmemB[0] = cooperative_load(A, B, 0)
5 sync_threads()
6 for k0 in range(127):
7     p, q = k0 % 2, (k0 + 1) % 2
8     RegsA, RegsB = cooperative_load(A, B, k0 + 1)
9     RegsC = block_mma(SmemA[p], SmemB[p], RegsC)
10    SmemA[q], SmemB[q] = RegsA, RegsB
11    sync_threads() Store Next Tile of A/B into Shared Memory
12    RegsC = block_mma(SmemA[0], SmemB[0], RegsC)
13    ...
```

The diagram illustrates the double-buffering optimization. It shows two parallel processes: "Preloading Next Tile of A/B into Regs" (orange box) and "Computation of Current Tile" (blue box). Arrows indicate data flow from memory to registers and registers to memory, showing how computation and data transfer overlap.

# Non Loop-Oriented Optimization (cont.)

More on non *loop-oriented* optimization:

- Thread block swizzle<sup>2</sup>.
- Tensor Core MMA primitive.
- Multi-stage asynchronous prefetching<sup>3</sup>.

## Challenge 1

TVM's existing loop-oriented scheduling primitives cannot manipulate the loop body at **fine granularity**.

---

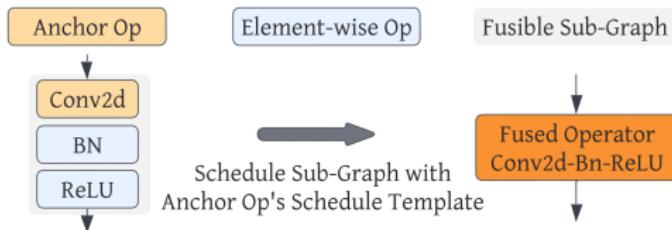
<sup>2</sup>Louis Bavoil. 2020. Optimizing Compute Shaders for L2 Locality using Thread Group ID Swizzling.

<https://developer.nvidia.com/blog/optimizing-compute-shaders-for-l2-locality-using-thread-group-id-swizzling>



<sup>3</sup>Andrew Kerr, Duane Merrill, Julien Demouth, et. al. 2018. CUTLASS: CUDA Template Library for Dense Linear Algebra at All Levels and Scales.

# Dedicated Schedule Template for Fusion



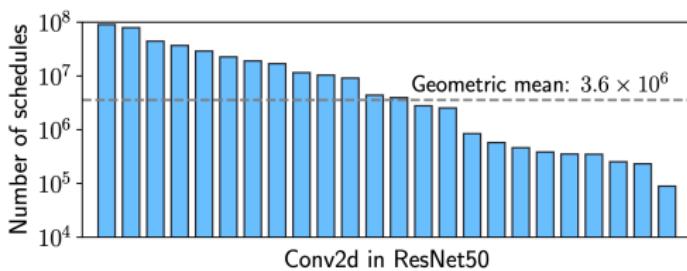
## Workflow of TVM's subgraph fusion:

- Identify anchor operator, which is usually the most compute-intensive one.
- Use anchor operator's schedule template to schedule the whole subgraph.

### Challenge 2

Fusion rules needs to be general enough to support all existing operators (and even new operators).

## Long Tuning Time



Tiling factor is the majority of search space:

- Constrained to perfect tile size, e.g.,  $N = 2039$  will fail.
  - Input size grows → Schedule space grows exponentially.
  - Long tuning time → Fail to co-optimize with graph-level opt and neural architecture search.

## Challenge 3

Large search space, long tuning time, input size sensitive.



# Contribution

- **C1: Non loop-oriented optimization**  
Introduce fine-grained scheduling primitive<sup>4</sup> for tensor programs.
- **C2: Dedicated schedule template for fusion**  
Based on that primitive, post-scheduling fusion is straight-forward yet general.
- **C3: Long tuning time**  
Based on that primitive, search space could be narrowed down to hardware-centric.
- **Programming Effort**  
Implemented with PyTorch 2.0, conduct solid experiments.



<sup>4</sup>The authors named it as "task-mapping oriented programming paradigm".

Background  
oooooo

Motivation  
oooooooo

Design  
●oooooooooooo

Evaluation  
oooooooo

Conclusion  
oooo

# Contents

1 Background

2 Motivation

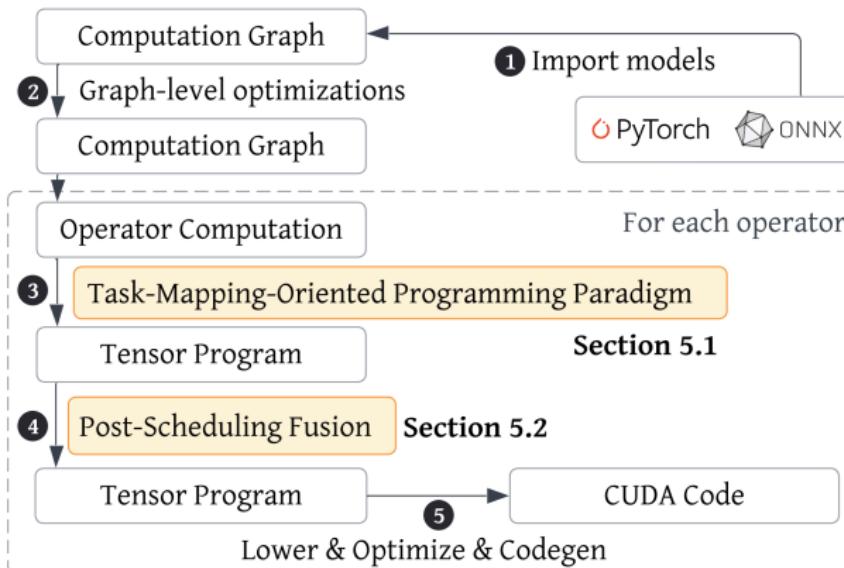
3 Design

4 Evaluation

5 Conclusion

# Design Overview

- Task mapping: CUDA-like thread block abstraction.
- Tuning: hardware-centric search space.
- Fusion: pre-defined rules.
- Model → Graph-level IR → Tensor-level IR → CUDA C++.

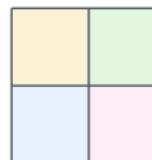


# Key Idea: Task Mapping

- Similar to CUDA thread-block abstraction.
- Task map: worker ID  $\rightarrow$  loop/array index.

0	1
2	3

Task execution order



Task Assignment

Legend

0	1
2	3

repeat(2, 2)

$$f(w) = [(0, 0), (0, 1), (1, 0), (1, 1)]$$

Assign all tasks to a single worker

(a)

0	0
0	0

spatial(2, 2)

$$f(w) = [(w / 2, w \% 2)]$$

Assign each worker a task (4 workers in this case)

(b)

# Key Idea: Task Mapping (cont.)

- Can be combined and nested.
- Represents CUDA's `blockDim`, `threadDim` or grid-strided loop.

---

$$\begin{array}{|c|c|c|} \hline 0 & 1 & 2 \\ \hline \end{array} \times \begin{array}{|c|c|} \hline 0 & 0 \\ \hline 0 & 0 \\ \hline \end{array} = \begin{array}{|c|c|c|c|c|c|} \hline 0 & 0 & 1 & 1 & 2 & 2 \\ \hline 0 & 0 & 1 & 1 & 2 & 2 \\ \hline \end{array}$$

(a)  $\text{repeat}(1, 3) \times \text{spatial}(2, 2)$ 

---

$$\begin{array}{|c|c|} \hline 0 & 0 \\ \hline 0 & 0 \\ \hline \end{array} \times \begin{array}{|c|c|c|} \hline 0 & 1 & 2 \\ \hline \end{array} = \begin{array}{|c|c|c|c|c|c|} \hline 0 & 1 & 2 & 0 & 1 & 2 \\ \hline 0 & 1 & 2 & 0 & 1 & 2 \\ \hline \end{array}$$

(b)  $\text{spatial}(2, 2) \times \text{repeat}(1, 3)$ 

---

$$\begin{array}{|c|c|} \hline 0 & 0 \\ \hline \end{array} \times \begin{array}{|c|c|} \hline 0 & 1 \\ \hline \end{array} = \begin{array}{|c|c|} \hline 0 & 1 \\ \hline \end{array}$$
  
$$\begin{array}{|c|c|c|c|} \hline 0 & 1 & 0 & 1 \\ \hline \end{array} \times \begin{array}{|c|c|} \hline 0 & 0 \\ \hline \end{array} = \begin{array}{|c|c|c|c|} \hline 0 & 0 & 1 & 1 \\ \hline 0 & 0 & 1 & 1 \\ \hline \end{array}$$
  
$$\begin{array}{|c|c|c|c|c|c|c|} \hline 0 & 0 & 1 & 1 & 0 & 0 & 1 & 1 \\ \hline \end{array}$$
  
$$\begin{array}{|c|c|} \hline 0 & 2 \\ \hline 1 & 3 \\ \hline \end{array}$$

(c)  $\text{spatial}(2) \times \text{repeat}(2) \times \text{spatial}(2)$ (d)  $\text{repeat}(1, 2) \times \text{repeat}(2, 1)$

# Key Idea: Task Mapping (cont.)

- Can be combined and nested.
- Represent CUDA's blockDim, threadDim or grid-strided loop.

```
1 # spatial(4, 2) * repeat(2, 2) * spatial(4, 8) * repeat(4, 4)
2
3 __global__ void kernel(){
4     for bi in range(2):                                # repeat(2, 2)
5         for bj in range(2):                            # repeat(2, 2)
6             for ti in range(4):                          # repeat(4, 4)
7                 for tj in range(4):                      # repeat(4, 4)
8                     body(bi, bj, ti, tj);    # task body
9 }
10
11 int main(){
12     dim3 blockDim(4, 2);
13     dim3 threadDim(4, 8);
14     kernel<<<blockDim, threadDim>>>();
15 }
```

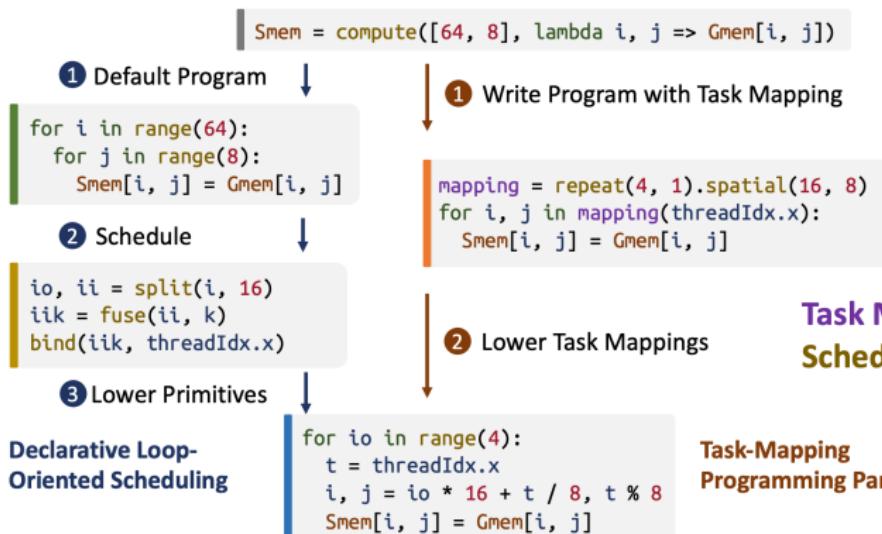
# From Loop-Oriented Scheduling to Task Mapping

TVM's loop scheduling primitives

Schedule Primitives	Original Program	Scheduled Program
<code>fuse(i, j)</code>	<code>for i in range(128):     for j in range(4):         body(i, j)</code>	<code>for ij in range(512):     body(ij / 4, ij % 4)</code>
<code>split(i, 128)</code>	<code>for i in range(512):     body(i)</code>	<code>for oi in range(4):     for ii in range(128):         body(oi * 128 + ii)</code>
<code>reorder(i, j)</code>	<code>for i in range(128):     for j in range(4):         body(i, j)</code>	<code>for j in range(4):     for i in range(128):         body(i, j)</code>
<code>bind(i, threadIdx.x)</code>	<code>for i in range(128):     body(i)</code>	<code>body(threadIdx.x)</code>

# From Loop-Oriented Schd. to Task Mapping (cont.)

Task map: worker ID → loop/array index



# From Loop-Oriented Schd. to Task Mapping (cont.)

## Putting it together

```

SmemA = compute([64, 8], lambda i, k: A[i, k])

1 Generate default tensor program
for i in range(64):
    for k in range(8):
        SmemA[i, k] = A[i, k]

2 Apply declarative schedule primitives
io, ii = split(i, 16)
iik = fuse(ii, k)
bind(iik, threadIdx.x)

3 Lower schedule primitives
def cooperative_load_A(A: fp32[64, 8]):
    SmemA = shared fp32[64, 8]
    t = threadIdx.x
    for io in range(4):
        i, k = io * 16 + t / 8, t % 8
        SmemA[i, k] = A[i, k]

```

Declarative Loop-Oriented Scheduling

Task-Mapping-Oriented Programming Paradigm

① Scheduling directly in the tensor program with **task mapping**.

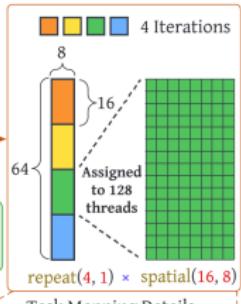
```

def cooperative_load_A(A: fp32[64, 8]):
    SmemA = shared fp32[64, 8]
    task_map = repeat(4, 1) × spatial(16, 8)
    for i, k in task_map(threadIdx.x):
        SmemA[i, k] = A[i, k]

```

② Lower task mapping

③ Implement task (i, k)



Task Mappings	Task Shape	Workers	Mapping (worker id w => assigned tasks)
repeat(4, 1)	(4, 1)	1	[(0, 0), (1, 0), (2, 0), (3, 0)]
spatial(16, 8)	(16, 8)	128	[(w / 8, w % 8)]
repeat(4, 1) × spatial(16, 8) (× task mapping composition)	(64, 8)	128	[(w / 8, w % 8), (w / 8 + 16, w % 8), (w / 8 + 32, w % 8), (w / 8 + 48, w % 8)]

# Complex Example: MMA (Matrix Multiply Accumulate)

- Task:  $16 \times 16$  grid
- Warp: worker
- Each warp compute 4 WMMA

---

```
1 def block_mma(SmemA: fp32[64, 8], SmemB: fp32[8, 64],  
2             RegsC: fp32[4, 4, 4]):  
3     RegsA, RegsB = register fp32[4], fp32[4]  
4     task_map = spatial(2, 2) * repeat(2, 2)  
5     worker_id = threadIdx.x / 32      # warp index  
6     for i, j in task_map(worker_id):  
7         wmma_load_a(&SmemA[i * 16, 0], RegsA)  
8         wmma_load_b(&SmemB[0, j * 16], RegsB)  
9         wmma_mma(RegsA, RegsB, RegsC[i, j])
```

---

# Scheduling Mechanism

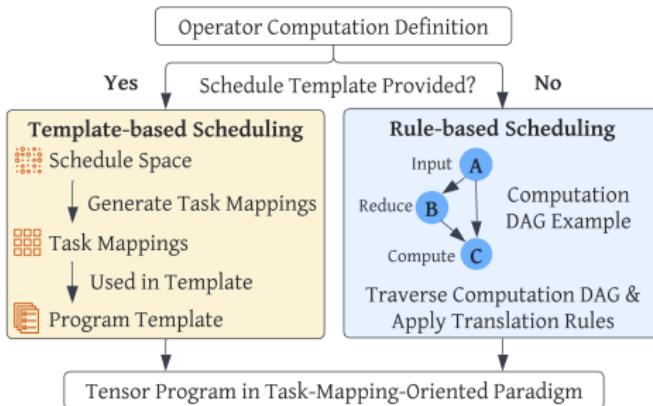
Hidet use similar scheduling mechanism with AutoTVM(template-based) and Ansor(rule-based)<sup>5</sup>

	AutoTVM Workflow	Auto-scheduler Workflow
<b>Step 1:</b> Write a compute definition  (relatively easy part)	# Matrix multiply  C = te.compute((M, N), lambda x, y: te.sum(A[x, k] * B[k, y], axis=k))  <b># 20-100 lines of tricky DSL code</b>	# The same
<b>Step 2:</b> Write a schedule template  (difficult part)	# Define search space cfg.define_split("tile_x", batch, num_outputs=4) cfg.define_split("tile_y", out_dim, num_outputs=4) ...  # Apply config into the template bx, txz, tx, xi = cfg["tile_x"].apply(s, C, C.op.axis[0]) by, tyz, ty, yi = cfg["tile_y"].apply(s, C, C.op.axis[1]) s[C].reorder(by, bx, tyz, txz, ty, tx, yi, xi) s[CC].compute_at(s[C], tx) ...	# Not required
<b>Step 3:</b> Run auto-tuning (automatic search)	tuner.tune(...)	task.tune(...)



# Scheduling Mechanism (cont.)

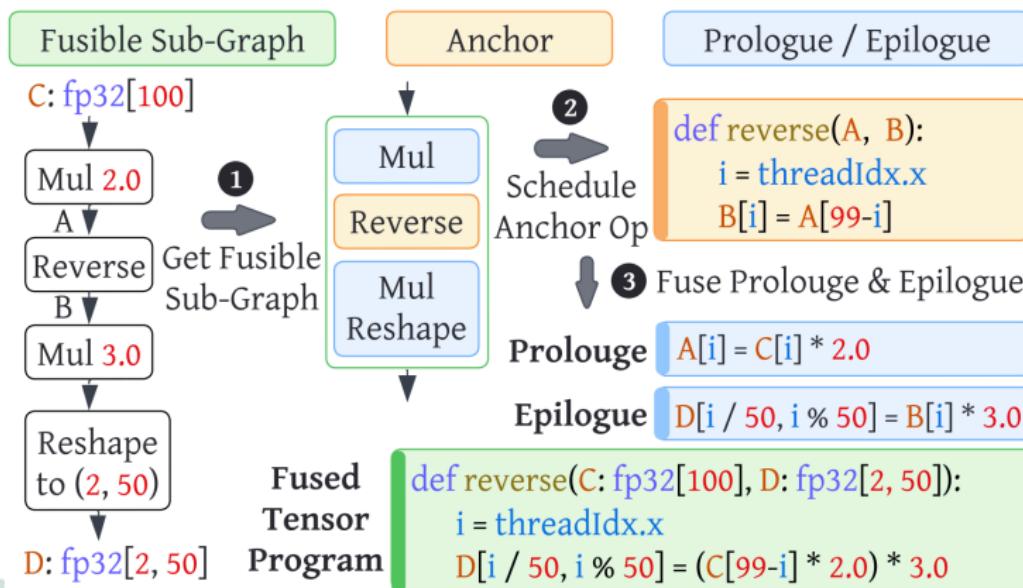
- **Rule-based:** Ansor-like, automatically translate DAG node by pre-defined rules, suitable for most of the operators (element-wise).
- **Template-based:** AutoTVM-like, requires manual effort, suitable for extreme-performance operators (GEMM, conv).
- Hidet only implemented GEMM and Reduce template, and they are able to cover all the operators in the evaluation.



# Post-Scheduling Fusion

Rules to decide if an operator is fusible:

- Prologue: no reduction.
- Epilogue: no reduction and is element-wise.



Background  
oooooo

Motivation  
oooooooo

Design  
oooooooooooo

Evaluation  
●oooooooo

Conclusion  
oooo

# Contents

1 Background

2 Motivation

3 Design

4 Evaluation

5 Conclusion

# Setup

## Platform

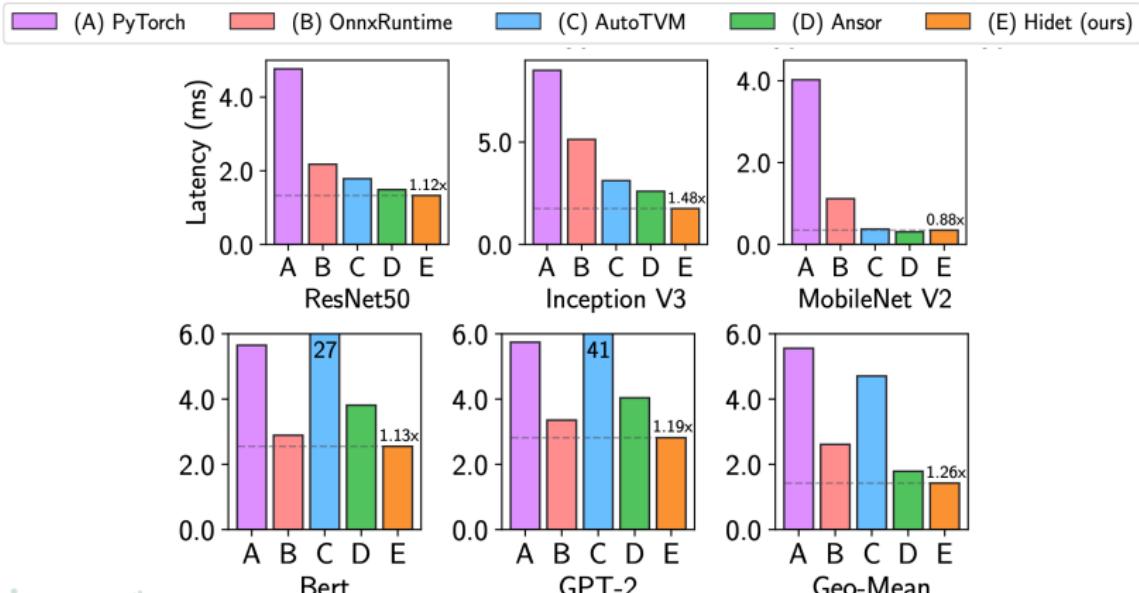
- Intel i9-12900K CPU (16-core, 24-thread, 64GB RAM)
- 1 × NVIDIA RTX 3090 GPU
- Ubuntu LTS 20.04, NVIDIA driver 510.73.08, CUDA 11.6

## Workload

- Models from TorchVision, exported to ONNX
- AutoTVM/Ansor tuning trial: 1000/800 (official's suggestion)
- Baseline: Ansor
- Metric: Inference latency, tuning cost

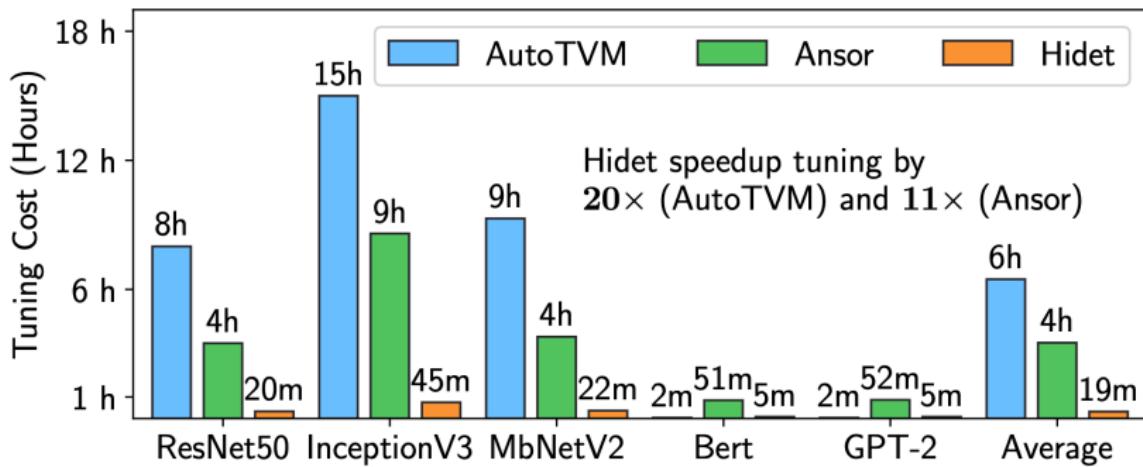
# End to End Performance

- MobileNet: Ansor has better schedule on *deepwise convolution*.
- Bert and GPT: AutoTVM is lack of schedule template (< 20).



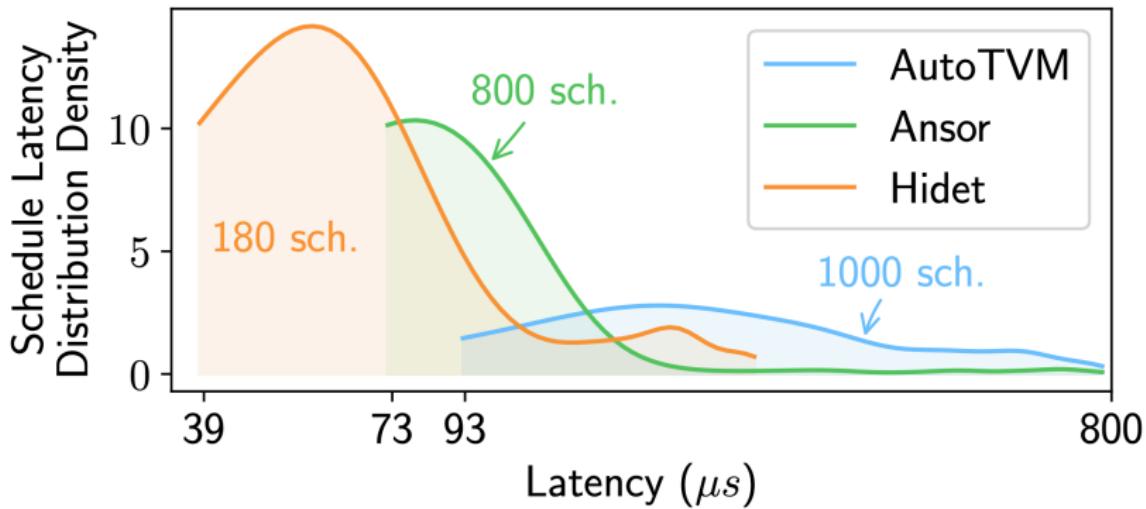
# Tuning Cost

- Schedule space: AutoTVM  $10^8$ , Ansor  $10^5$
- Hidet has small schedule space (180 for matrix multiplication), so can exhaust the schedule space while guarantee the performance.



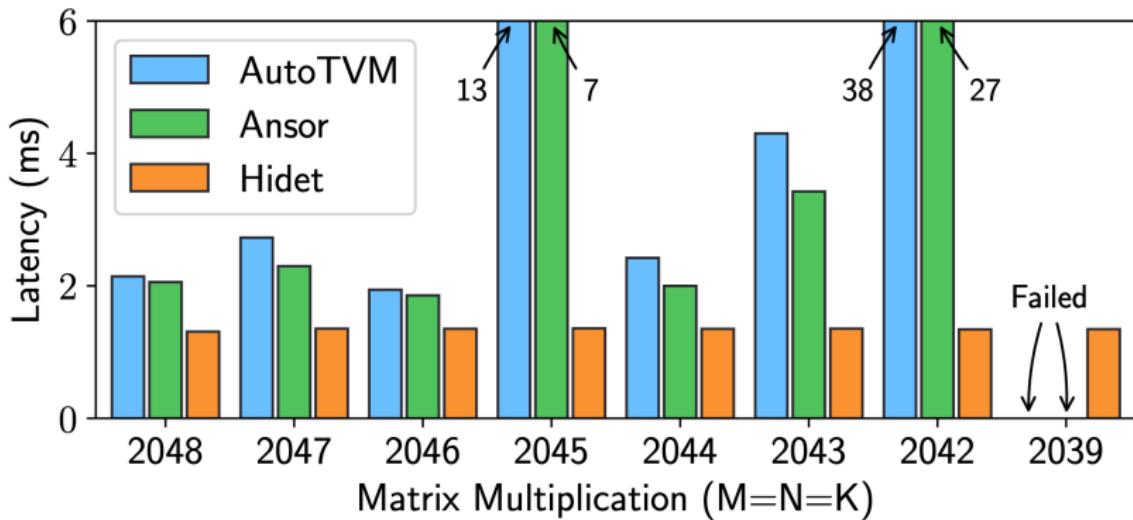
# Quality of Schedule Space

- Workload: a conv in ResNet50, batch 1, input  $28 \times 28$ , channels 256, kernel  $3 \times 3$ , padding 1, stride 2.
- Figure below: distribution of the performance in the schedule space.



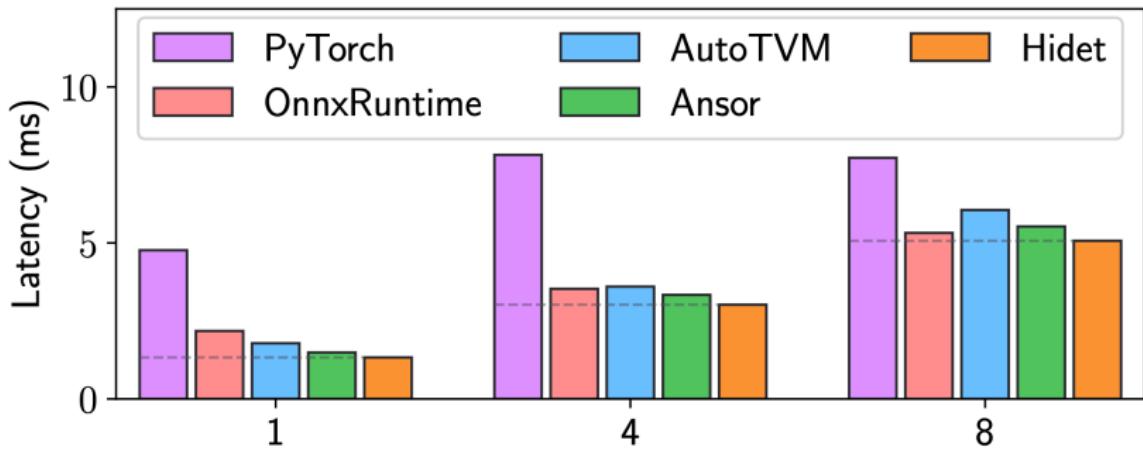
# Performance Sensitivity Over Input Size

- AutoTVM and Ansor's scheduling is *input-centric*, while Hidet is *hardware-centric*.
- AutoTVM and Ansor fails for prime number 2039.



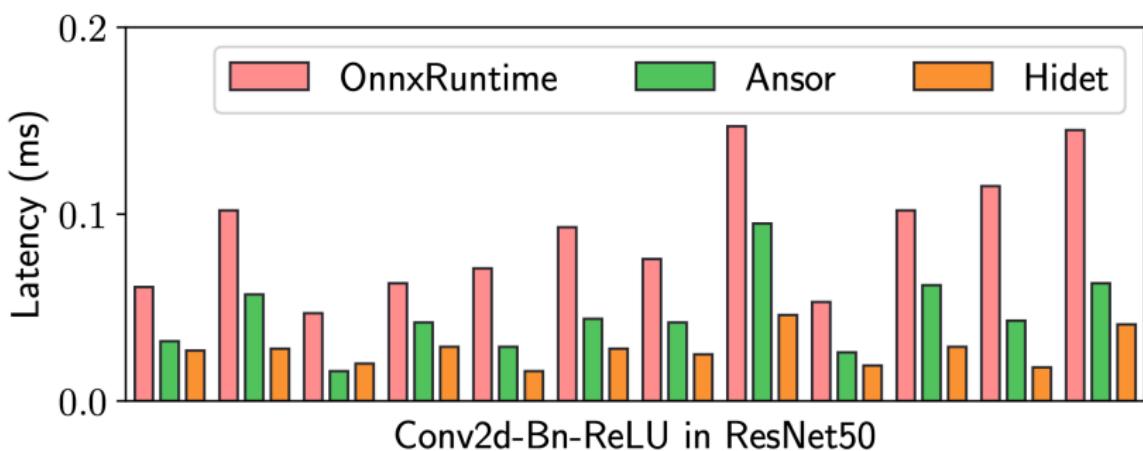
# Performance Sensitivity Over Batch Size

- AutoTVM and Ansor: lack of non loop-oriented optimization.
- Why Hidet outperform Onnx?



# Operator Fusion Performance

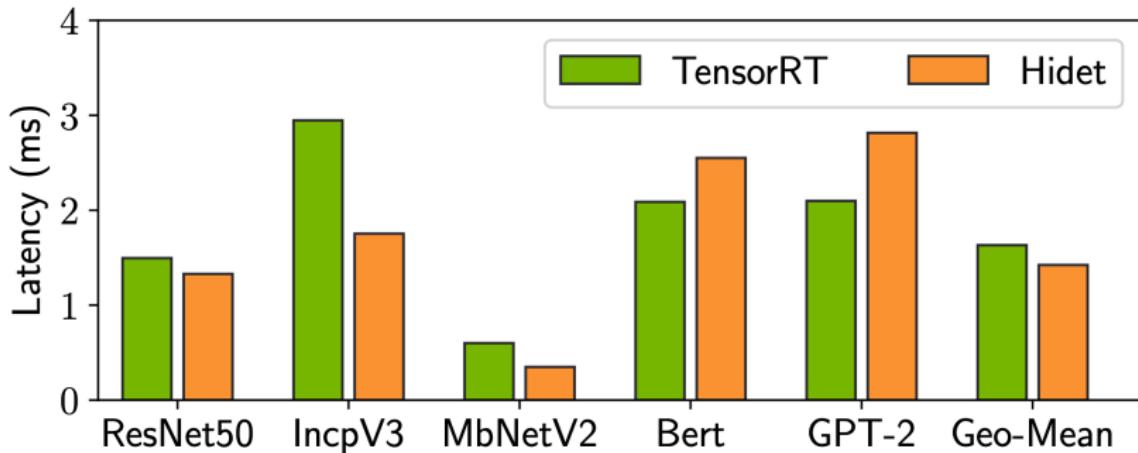
- Hidet implement conv by img2col<sup>6</sup> algorithm, using GEMM and reduce.
- Such implicit implementation could be co-optimized higher level optimization, and saturate hardware resource.



<sup>6</sup>Kumar Chellapilla, Sidd Puri, and Patrice Simard. 2006. High performance convolutional neural networks for document processing. In Tenth international workshop on frontiers in handwriting recognition. Suvisoft.

# Performance vs. TensorRT

- Hidet is able to fuse operator and tune for given input size.
- TensorRT may have dedicated optimization for Transformer?



Background  
oooooooo

Motivation  
oooooooo

Design  
oooooooooooo

Evaluation  
oooooooo

Conclusion  
●ooo

# Contents

1 Background

2 Motivation

3 Design

4 Evaluation

5 Conclusion

## Summary

- **Observation:** limited expressiveness of existing schedule method.
  - **Proposed method:** task-mapping programming paradigm.
    - Non loop-oriented optimizations.
    - Reduce tuning time.
  - **Implementation:** From scratch and integrate with PyTorch 2.0.
  - **Results:** 1.4x speedup and 20x tuning time reduction.



## Reference

- Paper: <https://dl.acm.org/doi/10.1145/3575693.3575702>
  - Repository: <https://github.com/hidet-org/hidet>
  - Talk on TVMCon'23: <https://www.youtube.com/watch?v=-vMSGZS9-DA>
  - Part of official slide: [https://centml.github.io/asplos23-tutorial/hidet\\_tutorial.pdf](https://centml.github.io/asplos23-tutorial/hidet_tutorial.pdf)
  - AutoTVM vs Ansor: <https://tvm.apache.org/2021/03/03/intro-auto-scheduler>

