

# Dynamic Tensor Shape Compilers

Yijia Diao

2022.11.15

# Outline

- Background
- Nimble
- DISC
- DietCode
- Summary

# Dynamic-Shape Workloads

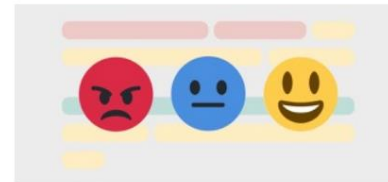
- Input sentences or audios have dynamic lengths.



Translation<sup>[1]</sup>



Speech Recognition<sup>[2]</sup>



Sentiment  
Analysis<sup>[3]</sup>

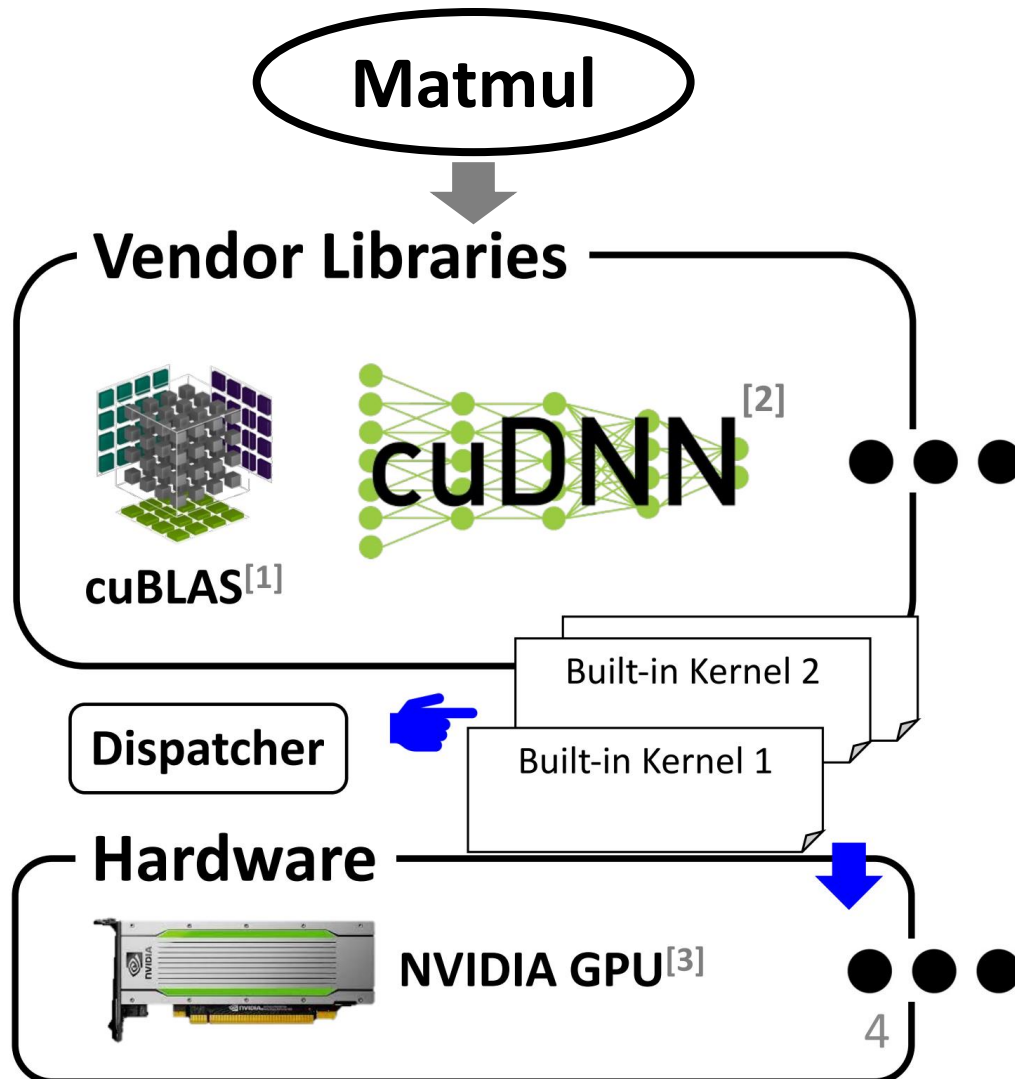
Text Auto-Complete<sup>[4]</sup>

[1] <https://translate.google.com/> [2] <https://github.com/NVIDIA/NeMo>

[3] J. Devlin et al. BERT. NAACL-HTL 2019 [4] A. Radford et al. GPT-2. 2019

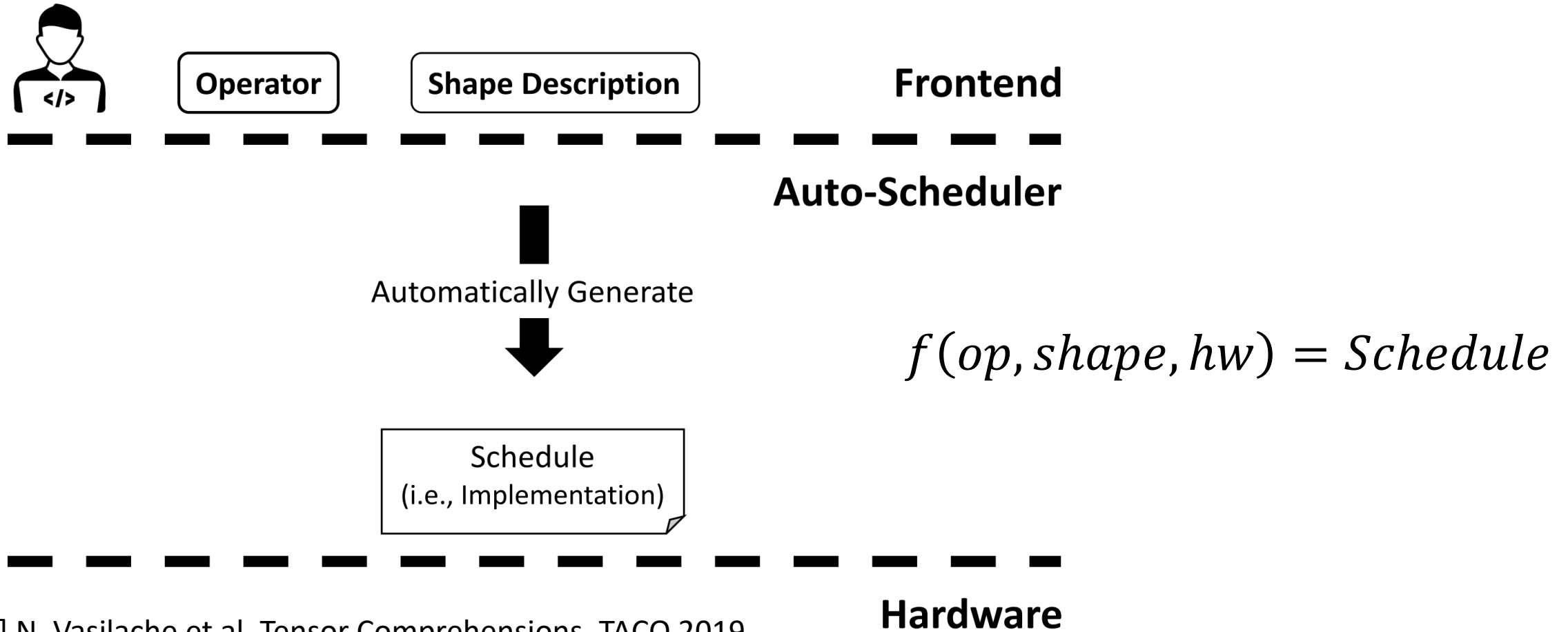
<https://mlsys.org/media/mlsys-2022/Slides/2175.pdf>

# Challenges: Vendor Libraries



- Huge engineering efforts
- Built-in kernels sometimes not optimal
  - If kernel for every shape optimal, the library would be huge 😬
- Unknown source code

# Challenges: Tensor Compilers



[1] N. Vasilache et al. Tensor Comprehensions. TACO 2019

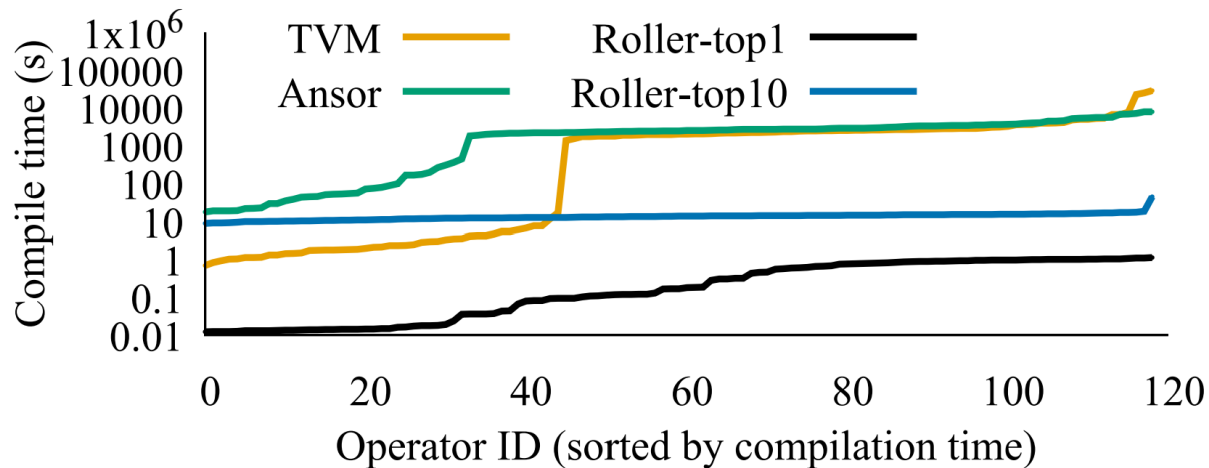
[2] L. Zheng et al. Ansor. OSDI 2020

[3] S. Zheng et al. FlexTensor. ASPLOS 2020

[4] H. Zhu et al. Roller. OSDI 2022

# Challenges: Tensor Compilers (Cont.)

- For tensor compilers like Ansor, it needs hour-level time to tune single shape for single op.
- The cost would be huge when input shape is dynamic.



	Tuning?	Complexity
Vendor Libraries	✗	-
Existing Auto-Schedulers	✓	$O( S )$

# Outline

- Background
- Nimble
- DISC
- DietCode
- Summary

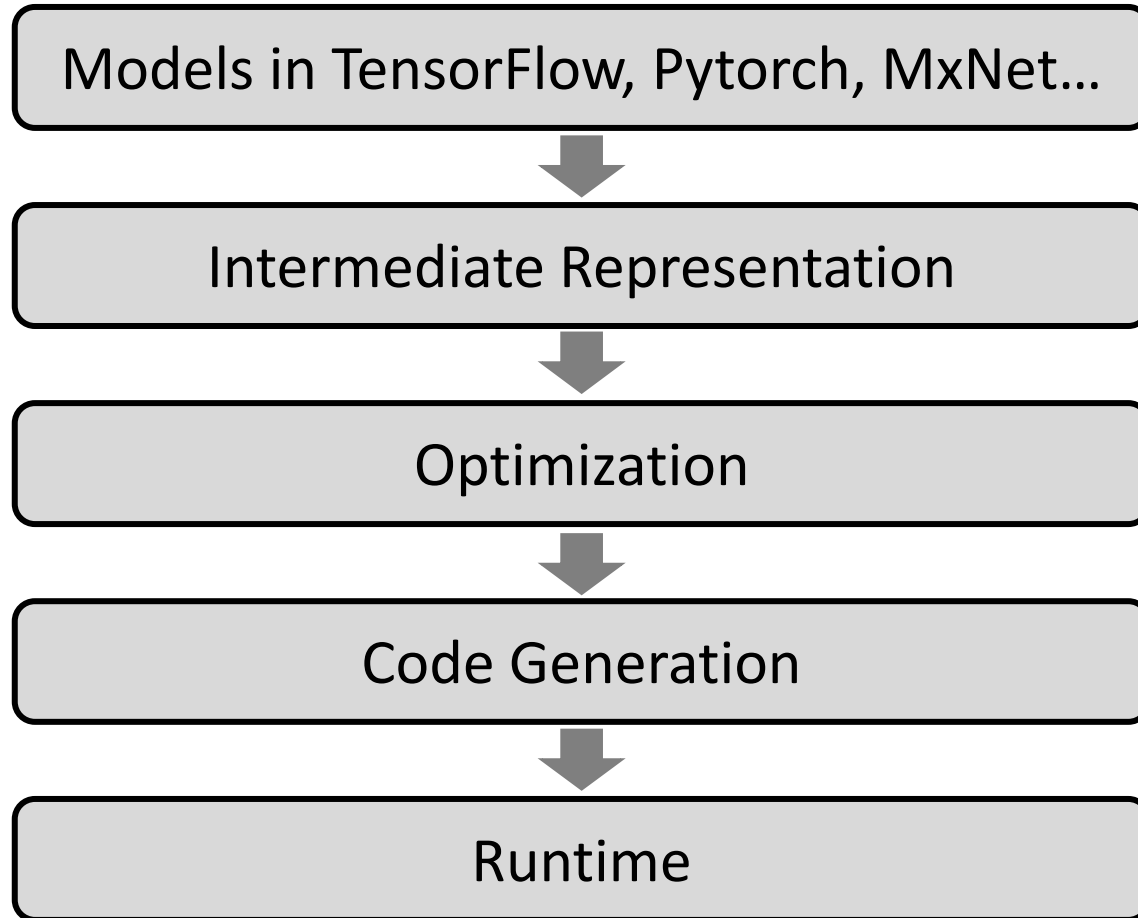
---

## **NIMBLE: EFFICIENTLY COMPILING DYNAMIC NEURAL NETWORKS FOR MODEL INFERENCE**

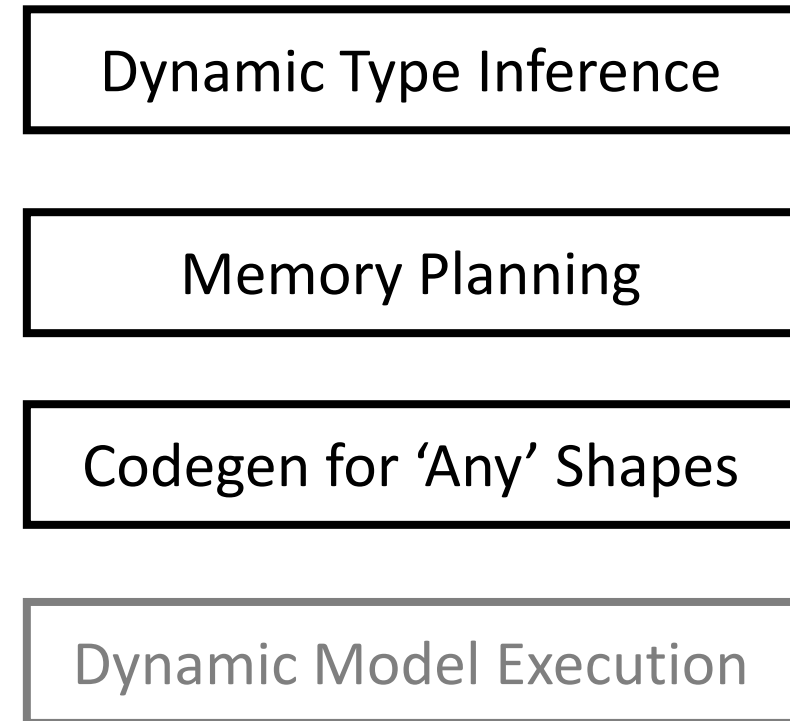
---

**Haichen Shen<sup>\*1</sup> Jared Roesch<sup>\*2</sup> Zhi Chen<sup>1</sup> Wei Chen<sup>1</sup> Yong Wu<sup>1</sup>  
Mu Li<sup>1</sup> Vin Sharma<sup>1</sup> Zachary Tatlock<sup>3</sup> Yida Wang<sup>1</sup>**

# Nimble: Overview



Static Tensor Compiler



Nimble



# Nimble: Dynamic Type Inference

- Symbolic shape:
  - Tensor[(1, 10, Any), float32]
- New propagation rules:
  - broadcast rel(Any, 1)  $\rightarrow$  Any
  - broadcast rel(Any, d)  $\rightarrow$  d
  - broadcast rel(Any, Any)  $\rightarrow$  Any

# Nimble: Memory Planning

- Static compilers: implicit memory planning
- Dynamic shape needs explicit memory planning

```
1 fn (x: Tensor<?, 2>, y: Tensor<1, 2>)
2     ->Tensor<?, 2> {
3     let in_sh0 = shape_of(x);
4     let in_sh1 = shape_of(y); // shape calculating function
5     let buf0 = alloc_storage(16, 64, ...); // allocate shape buffer & tensor
6     let out_sh0 = alloc_tensor(buf0, ...);
7     invoke_shape_func(concat, // invoke shape function and get actual shape
8         (in_sh0, in_sh1), (out_sh0, ), ...);
9 - - - - -
10    let buf1 = alloc_storage(...);
11    let out0 = alloc_tensor(
12        buf1, out_sh0, ...); // allocate data buffer & tensor
13    invoke_mut(concat, (x, y), (out0)); // invoke concat op
14    out_0
15 }
```

Host

Device

# Nimble: Codegen for 'Any' Shapes

- 3 Step to Codegen:
  - Any  $\rightarrow$  64, 128,... (Tensor[(1, 10, Any), float32])
  - Tune with static shapes, then pick TopK(100) results to test on other shapes
  - Pick the config: run best on other shapes
- Avoiding Boundary Check: duplicate kernels
  - Kernel tile = 8, then prepare 8 version kernels.
  - Automatically generate dispatch function.

# Comments for Nimble

- End-to-End Dynamic Workload Compiler
  - New IR and primitive is necessary
- Compilation Overhead?
- Dynamic kernel performance?

# Outline

- Background
- Nimble
- **DISC**
- DietCode
- Summary

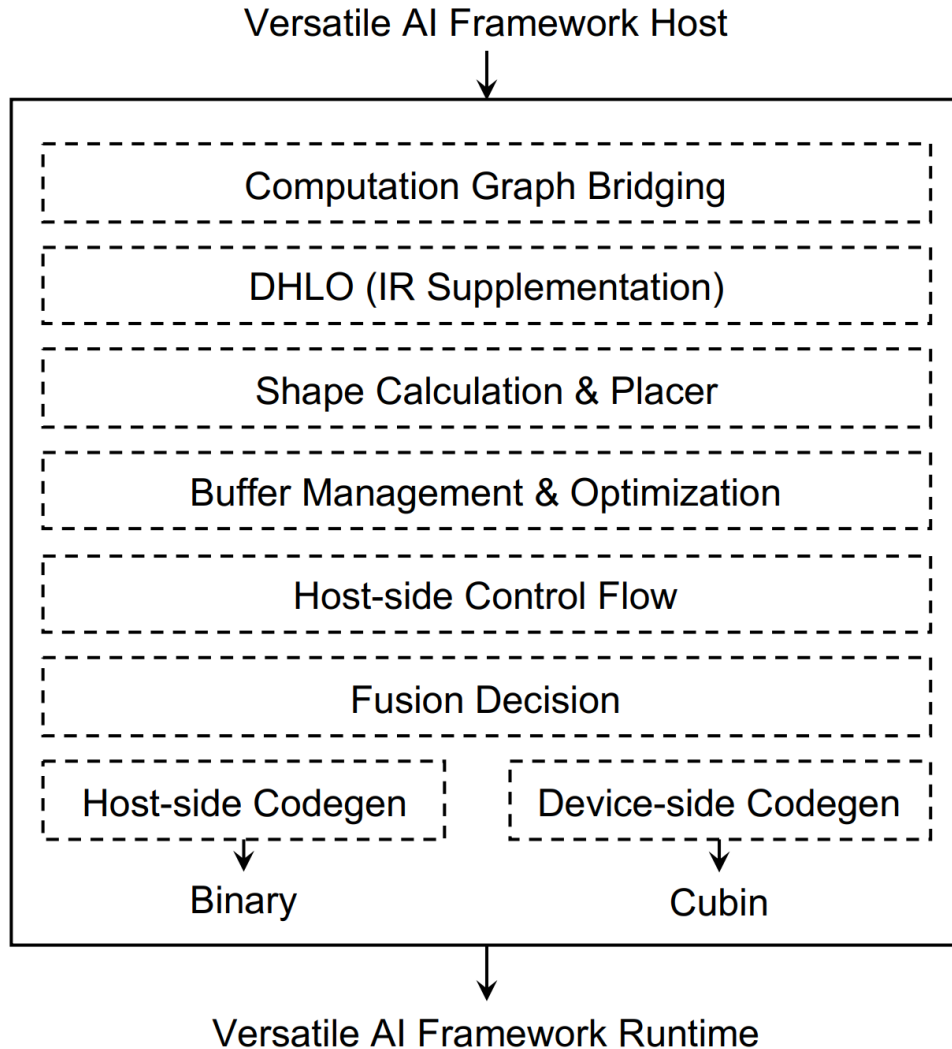
## ***DISC* : A Dynamic Shape Compiler for Machine Learning Workloads**

Kai Zhu, Wenyi Zhao, Zhen Zheng, Tianyou Guo, Pengzhan Zhao, Feiwen Zhu  
Junjie Bai, Jun Yang, Xiaoyong Liu, Lansong Diao, Wei Lin

Alibaba Group

{tashuang.zk, kevin.zwy, james.zz, tianyou.gty, pengzhan.zpz, feiwen.zfw,  
j.bai, muzhuo.yj, xiaoyong.liu, lansong.dls, weilin.lw}@alibaba-inc.com

# DISC Overview



- Based on MLIR
- Focus on Fusion of memory-bound OPs

# DISC: Fusion of Dynamic Kernels

- Fusion requires same shape
- Shape Propagation
  - ADD op: input tensor shape = output tensor shape
- Shape Constraints
  - Dimension equality: reduction dim
  - Tensor size equality: transpose op

# DISC: Which Kernel to Choose?

- Two Choices:
  - Choose 'Best' kernel from vendor libraries
  - Hand-tuned kernel for each shape
- Comments:
  - The offline workload is still huge for new ops & devices
  - Hand-tuned kernel is not a clear & general method



# Outline

- Background
- Nimble
- DISC
- DietCode
- Summary

---

## **DIETCODE: AUTOMATIC OPTIMIZATION FOR DYNAMIC TENSOR PROGRAMS**

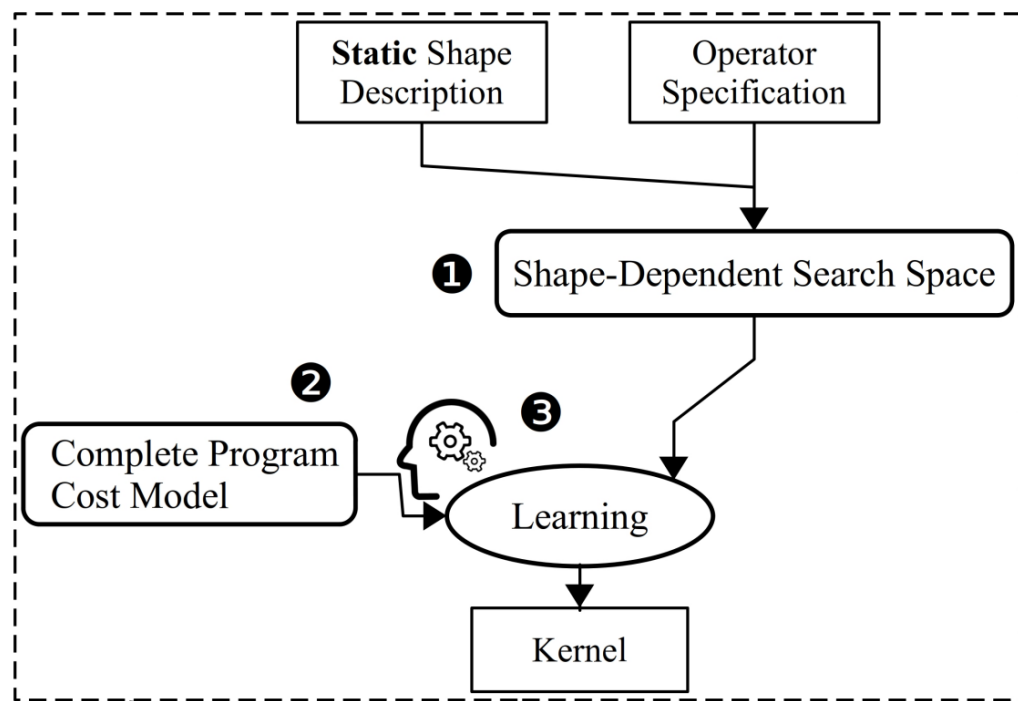
---

**Bojian Zheng**<sup>\* 1 2 3</sup> **Ziheng Jiang**<sup>\* 4</sup> **Cody Hao Yu**<sup>2</sup> **Haichen Shen**<sup>5</sup> **Josh Fromm**<sup>6</sup>  
**Yizhi Liu**<sup>2</sup> **Yida Wang**<sup>2</sup> **Luis Ceze**<sup>7 6</sup> **Tianqi Chen**<sup>8 6</sup> **Gennady Pekhimenko**<sup>1 2 3</sup>

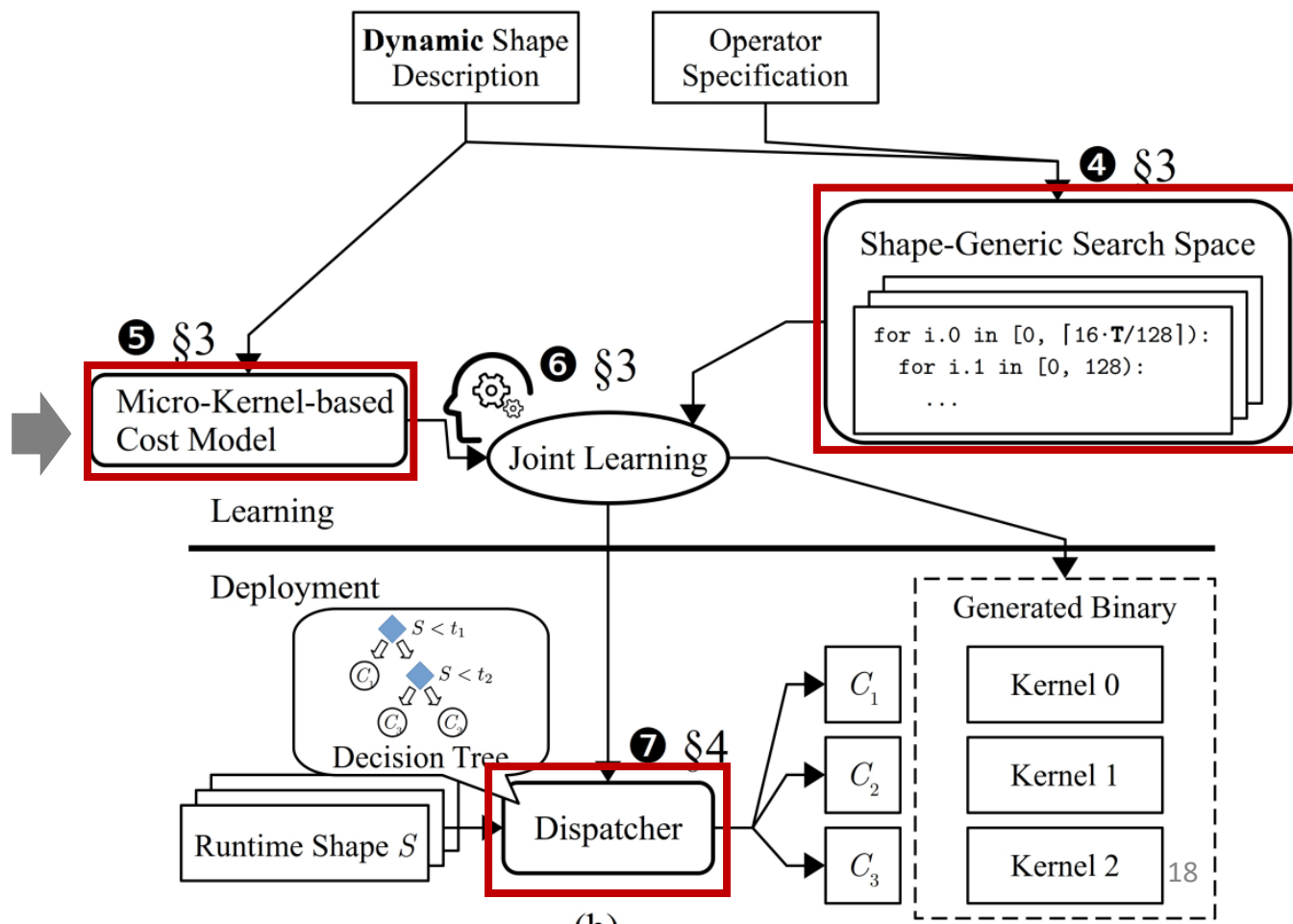
# DietCode: Overview

- Codegen for a batch of shapes.
- Codebase: Anso(TVM)

$$f(op, shape, hw) = Schedule$$



$$f(op, shapes, hw) = Schedules, Dispatcher$$



# Dynamic Shape Description


- Dynamic shape interface in TVM

```
51 def test_train_dynT():
52     B = 16
53     T = list(range(5, 128, 19))
54     T.append(128)
55     I = 768
56     H = 2304
57
58     wkl_insts = cross_product(T, (I, H))
59     logger.info("Workload Instances: {}".format(wkl_insts))
60
61     DynT, DynI, DynH = tir.DynShapeVar('T'), tir.DynShapeVar('I'), tir.DynShapeVar('H')
62
63     auto_scheduler.train(wkl_func=Dense,
64                          wkl_func_args=(B * DynT, DynI, DynH),
65                          shape_vars=[DynT, DynI, DynH], wkl_insts=wkl_insts,
66                          wkl_inst_weights=[1. for _ in wkl_insts],
67                          fcublas_fixture=cuBLASDenseFixture,
68                          sched_func_name_prefix='dense_{}xTx{}x{}'.format(B, I, H)
69                          )
```

# DietCode: Search Space for a Batch Shape

- Full Search space:  $t \in [2, +\infty)$
- Ansor Search space:  $t \in \{2, 5, 10, 25\}$
- DietCode: Any  $t < T$  is valid.
  - Bring branch overhead 😞

Loop Tiling Schedule:



```
for (int io = 0; io < [50/t]; ++io) {  
  for (int ii = 0; ii < t ++ii) {  
    if (io×t + ii < 50) A[io×t + ii] = ...  
  }  
}
```

# How to Reduce Branch Overhead?

Initial Code

```
for i.0 in [0, T):
  for i.1 in [0, t):
    if i.0*t+i.1 < T:
      X_local = X[...]
    if i.0*t+i.1 < T:
      Y_local = ...
    if i.0*t+i.1 < T:
      Y[...] = Y_local
```

Loop  
Partition

```
for i.0 in [0, T):
  for i.1 in [0, t):
    if i.0 <  $\lceil T/t \rceil$ :
      X_local = X[...]
      Y_local = ...
      Y[...] = Y_local
    else:
      if i.0*t+i.1 < T:
        X_local = X[...]
      if i.0*t+i.1 < T:
        Y_local = ...
      if i.0*t+i.1 < T:
        Y[...] = Y_local
```



```
X_pad = pad X to  $\lceil T/t \rceil \cdot t$ 
for i.0 in [0, T):
  for i.1 in [0, t):
    X_local = X_pad[...]
    Y_local = ...
    Y_pad[...] = Y_local
  slice Y_pad to [T]
```

Global  
Padding

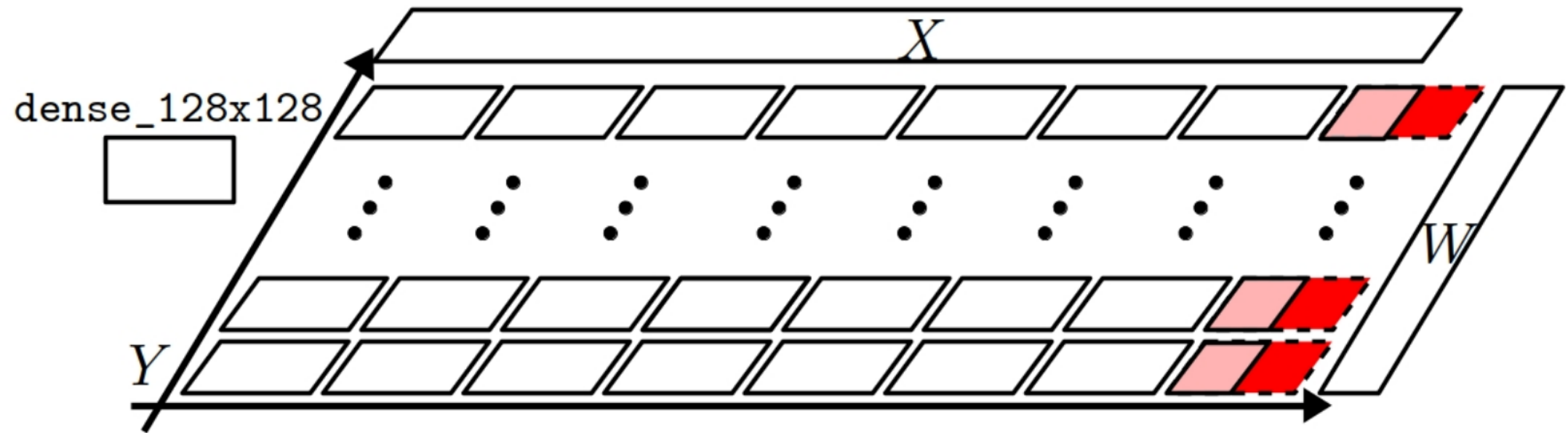
No branch

```
for i.0 in [0, T):
  for i.1 in [0, t):
    if i.0*t+i.1 < T:
      X_local = X[...]
    Y_local = ...
    if i.0*t+i.1 < T:
      Y[...] = Y_local
```

Local  
Padding

No branch Compute

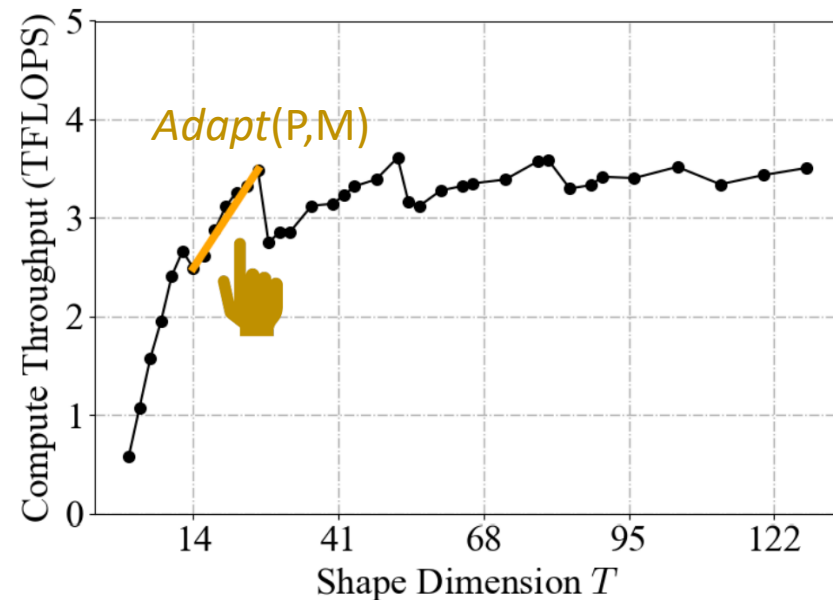
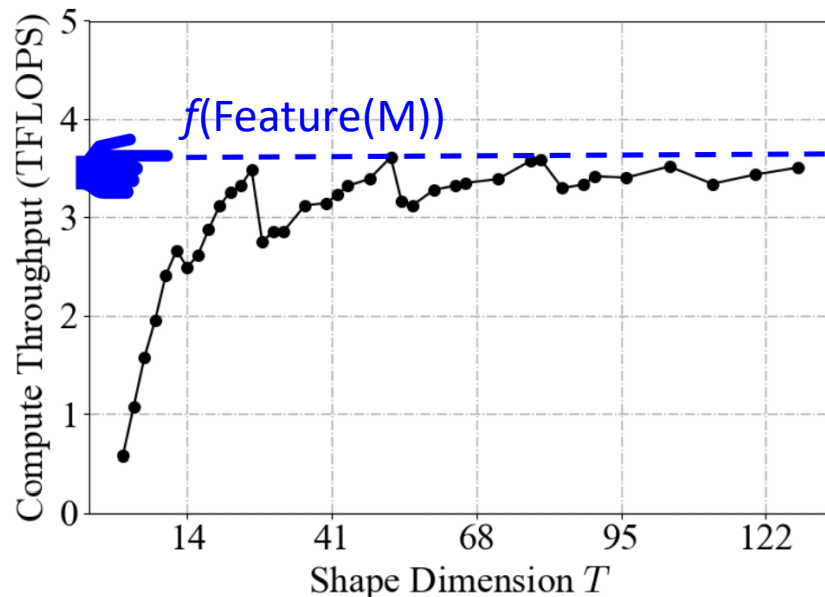
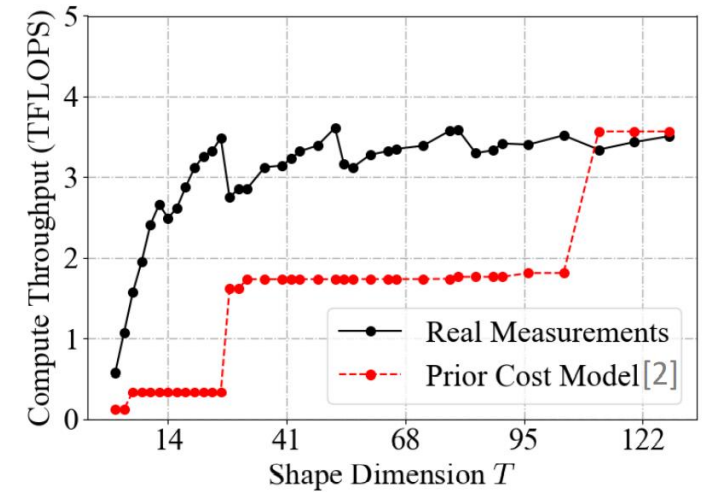
# Micro-Kernel in DietCode



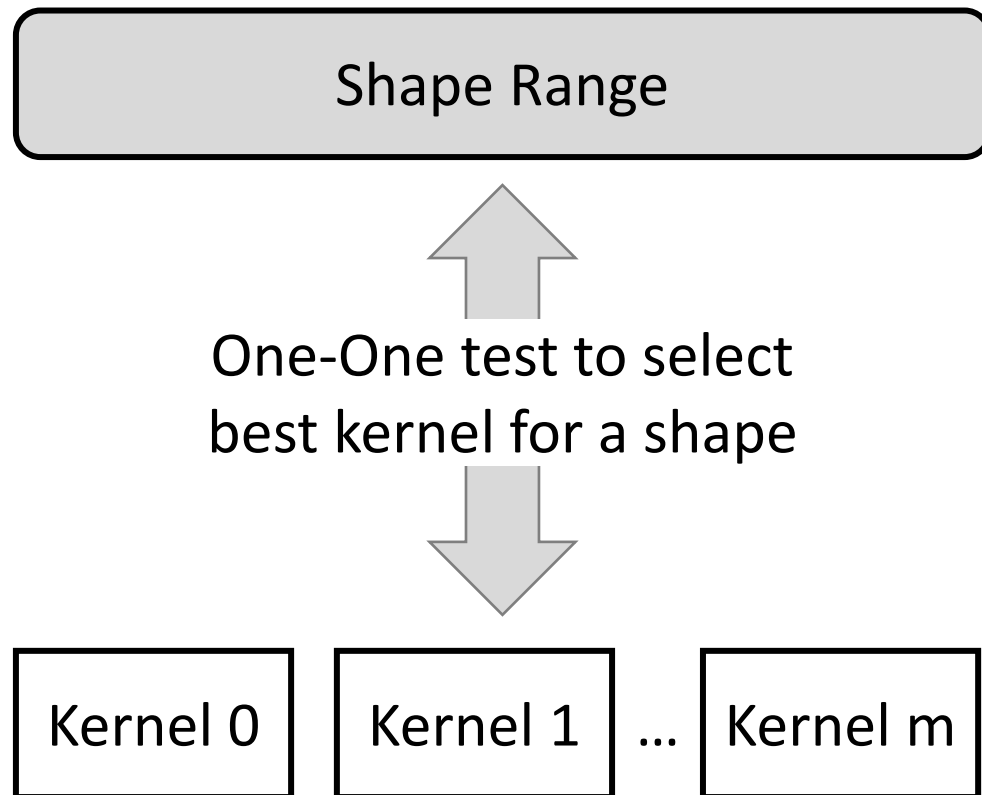
# DietCode Cost Model: Micro-Kernel Based

- $Y = XW^T, X: [16 \times T, 768], W = [2304, 768]$
- Ansor Cost Model:  $\text{Cost}(P) = f(\text{Feature}(P))$
- DietCode Cost Model:

$$\text{Cost}(P) = f(\text{Feature}(M)) \times \text{adapt}(P, M)$$



# Dispatcher of DietCode



- DietCode could generate a batch of micro-kernels.
- Dispatcher: Decision Tree(?)



# Summary

	<b>Nimble</b>	<b>DISC</b>	<b>DietCode</b>
End-To-End	Yes	Yes	No(TVM)
Backend	Virtual Machine	No VM	TVM
Dynamic Kernel	Off-line Tuning	Hand-tuned kernels + vendor	Search space + Cost Model
Dispatcher	Choose one kernel, duplicate, apply to all	?(maybe hardcode)	Decision Tree

Thanks!