# Learning to Optimize Tensor Programs
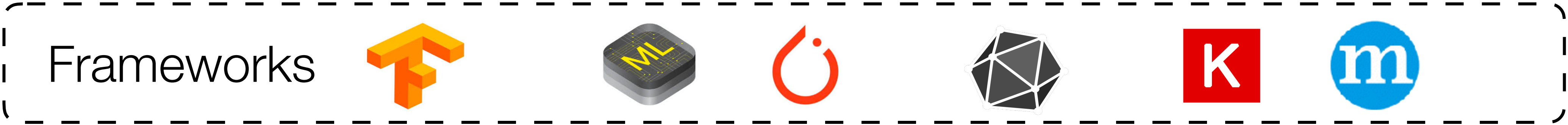
**Tianqi Chen,** Lianmin Zheng, Eddie Yan, Ziheng Jiang, Thierry Moreau,
Luis Ceze, Carlos Guestrin, Arvind Krishnamurthy

**W** PAUL G. ALLEN SCHOOL
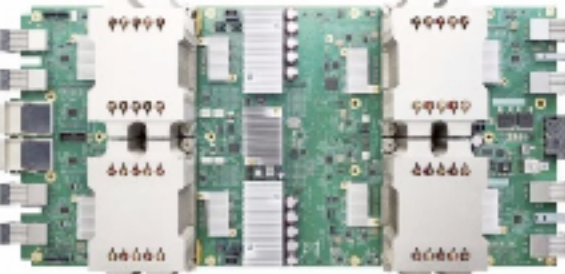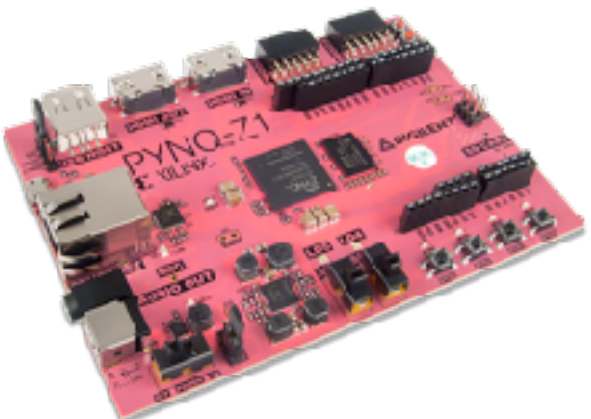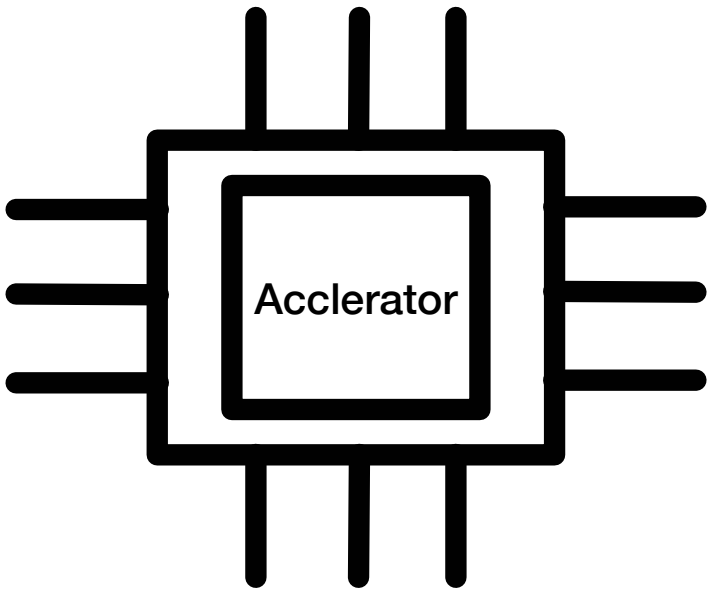OF COMPUTER SCIENCE & ENGINEERING

tvm.ai

sampl

# Goal: Deploy Deep Learning Everywhere
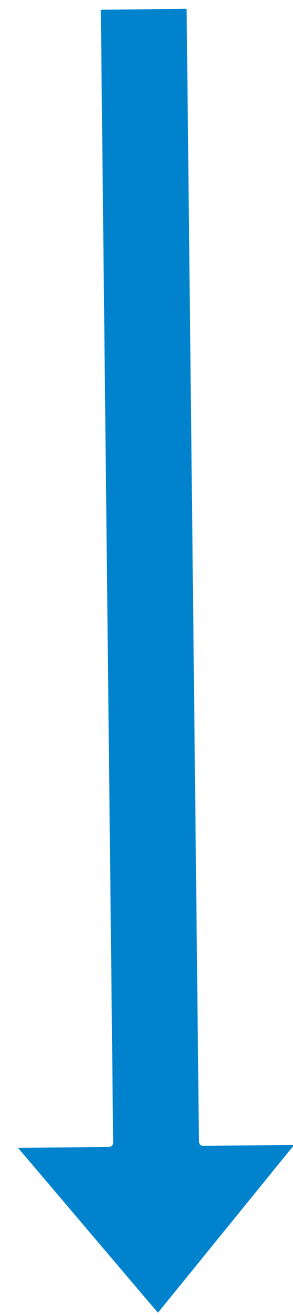
Frameworks

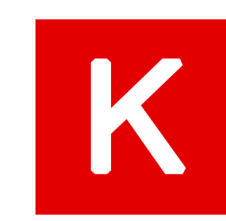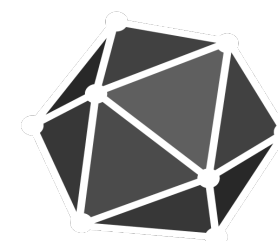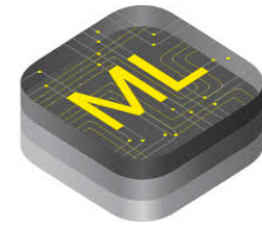**Explosion of models and frameworks**

**Huge gap between model/frameworks and hardware backends**

**Explosion of hardware backends**

Acclerator

# Existing Approach

Frameworks

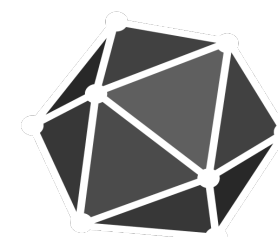Hardware

# Existing Approach

Frameworks

High-level data flow graph

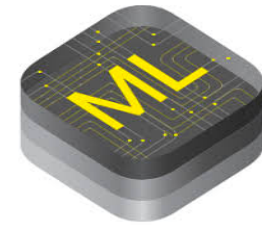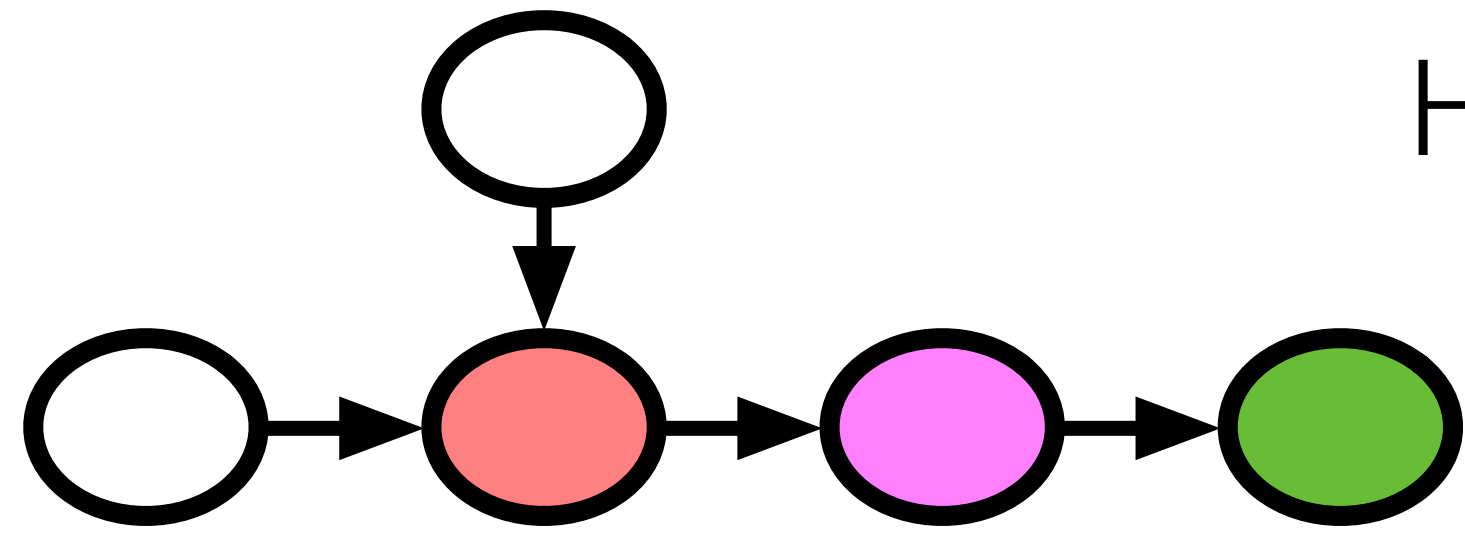Hardware

# Existing Approach

Frameworks

High-level data flow graph

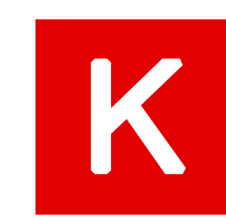Primitive Tensor operators such as Conv2D

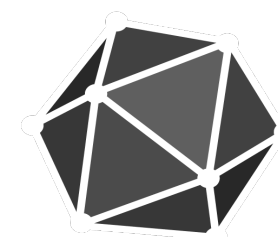Hardware

# Existing Approach

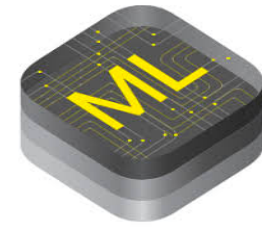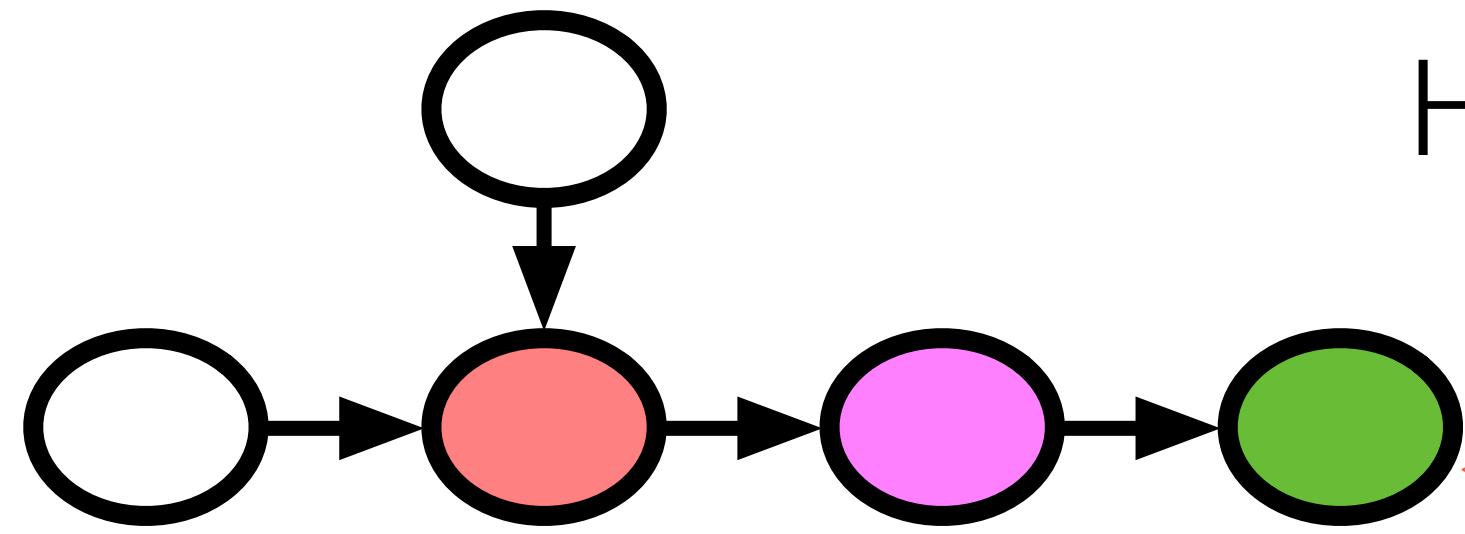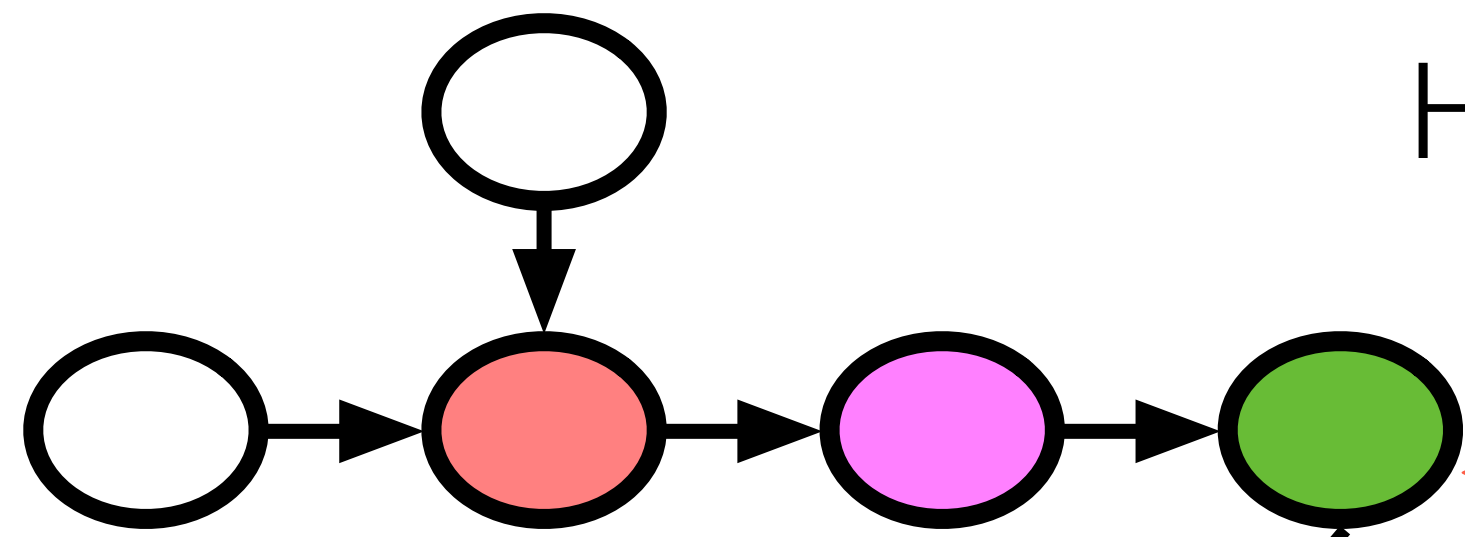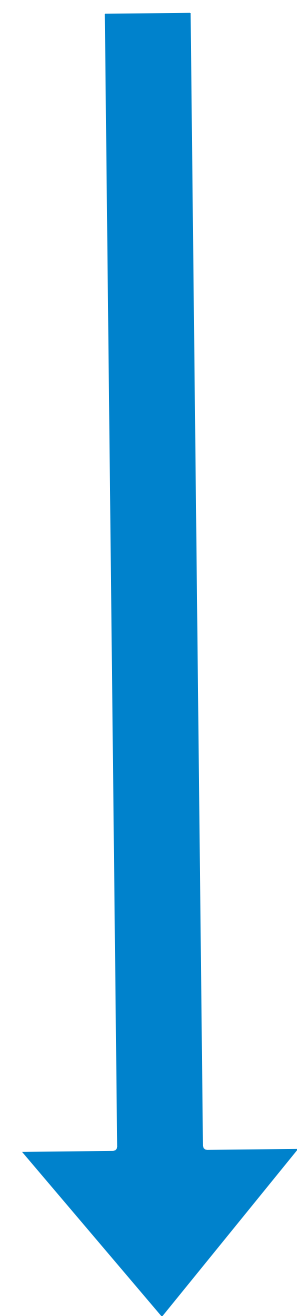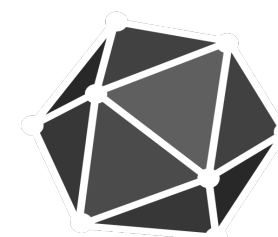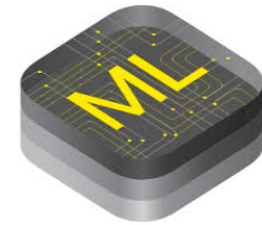Frameworks
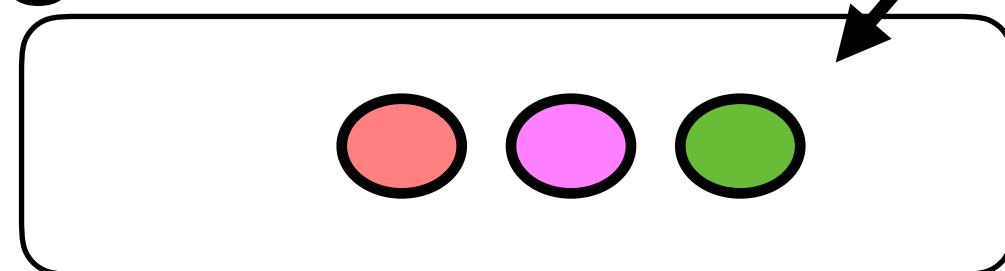
High-level data flow graph

Primitive Tensor operators such as Conv2D

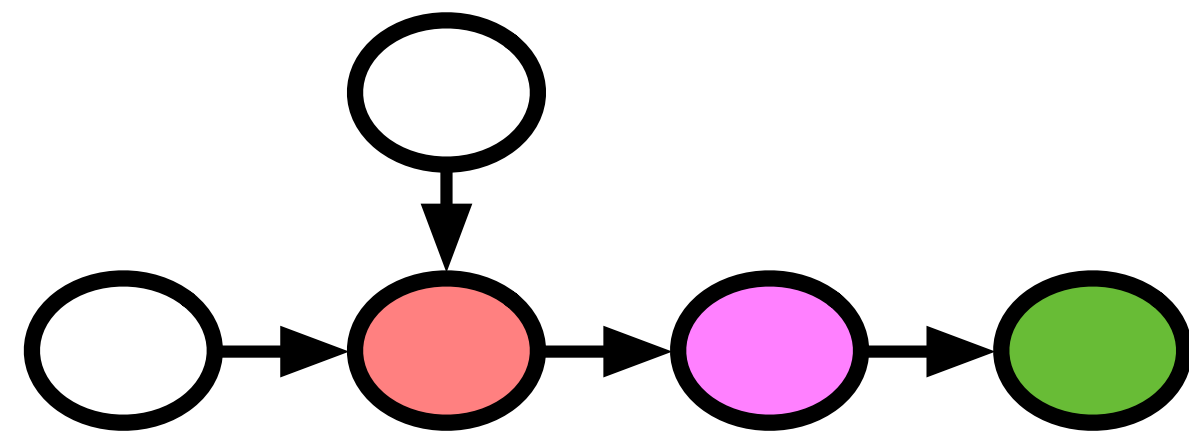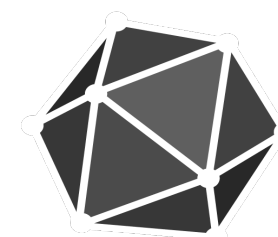**eg. cuDNN**

Offload to heavily optimized
DNN operator library

Hardware

# Limitations of Existing Approach

Frameworks



**cuDNN**

# Limitations of Existing Approach
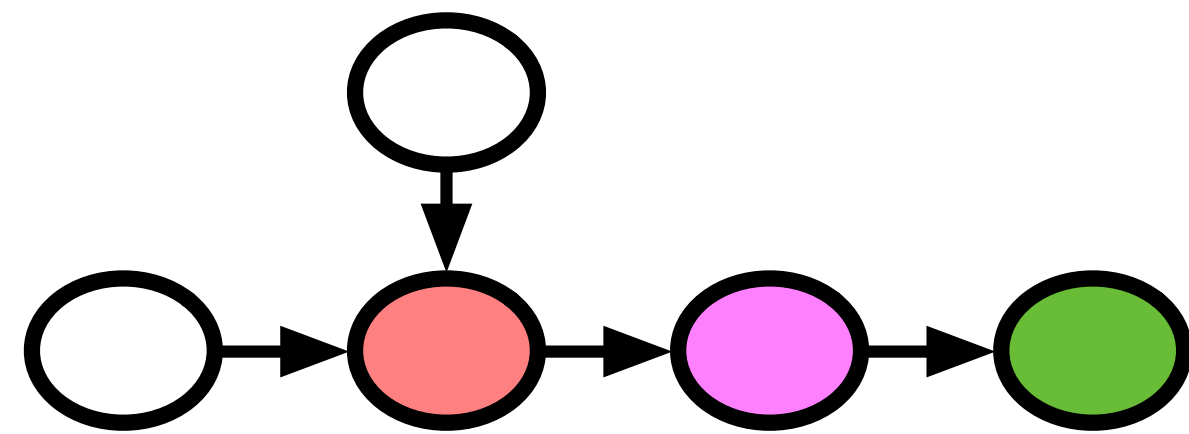
Frameworks

cuDNN

# Limitations of Existing Approach

Frameworks

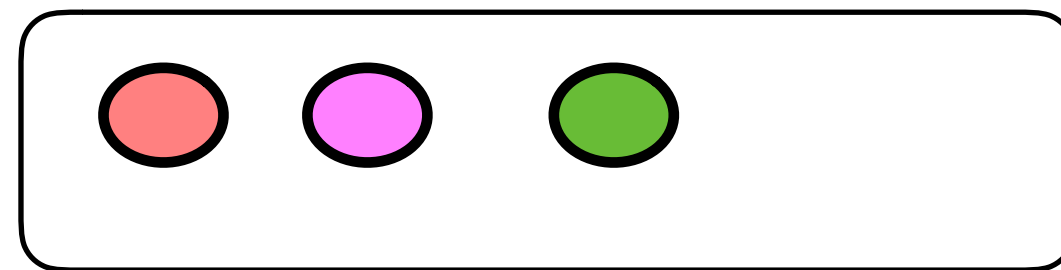**cuDNN**

# Limitations of Existing Approach



Frameworks

New operators

cuDNN

# Limitations of Existing Approach

Frameworks

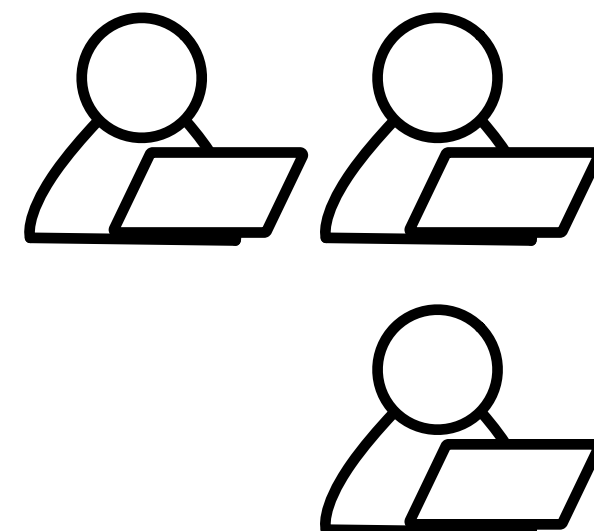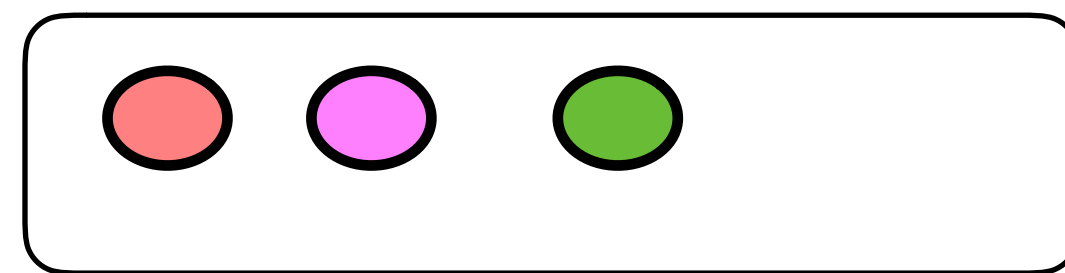**New operators**

**cuDNN**

# Limitations of Existing Approach

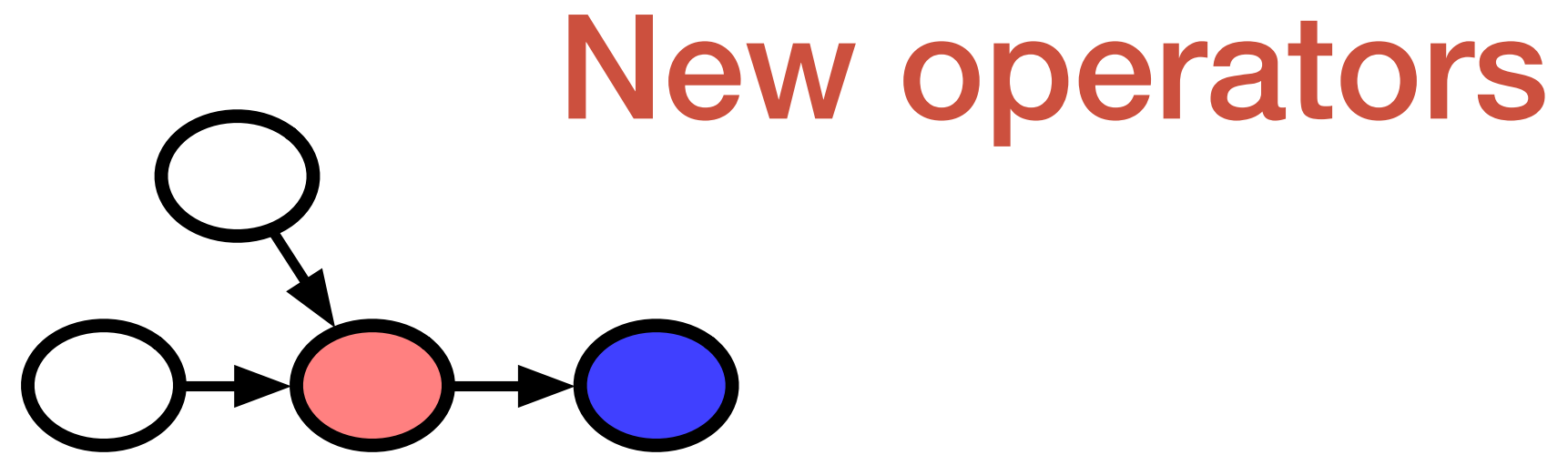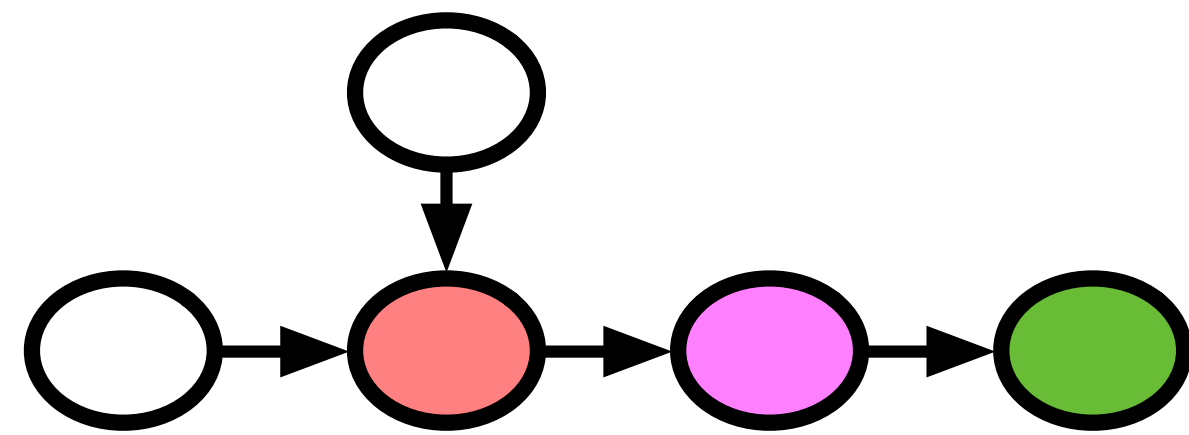Frameworks

**New operators**

**cuDNN**

nVIDIA

intel Xeon processor

# Limitations of Existing Approach

Frameworks

**New operators**

**Engineering intensive**

**cuDNN**

NVIDIA

# Learning to Optimize Tensor Programs

Frameworks

High-level data flow graph and optimizations

Hardware

# Learning to Optimize Tensor Programs

Frameworks

High-level data flow graph and optimizations

Hardware

# Learning to Optimize Tensor Programs

Frameworks

High-level data flow graph and optimizations

Machine Learning based Program Optimizer

Hardware

# Learning to Optimize Tensor Programs

# Search over Possible Program Transformations

**Compute Description**

```
C = tvm.compute((m, n),
    lambda y, x: tvm.sum(A[k, y] * B[k, x], axis=k))
```

| Loop Transformations | Thread Bindings | Cache Locality |
|---|---|---|
| Thread Cooperation | Tensorization | Latency Hiding |

Hardware

# Search over Possible Program Transformations



**Compute Description**

```
C = tvm.compute((m, n),
      lambda y, x: tvm.sum(A[k, y] * B[k, x], axis=k))
```

| Loop Transformations | Thread Bindings | Cache Locality |
|---|---|---|
| Thread Cooperation | Tensorization | Latency Hiding |

Hardware

# Search over Possible Program Transformations



**Compute Description**

```
C = tvm.compute((m, n),
    lambda y, x: tvm.sum(A[k, y] * B[k, x], axis=k))
```

**Billions of possible optimization choices**

| Loop Transformations | Thread Bindings | Cache Locality |
| Thread Cooperation | Tensorization | Latency Hiding |

Hardware

# Learning-based Program Optimizer



Program Optimizer → Code Generator → Program

# Learning-based Program Optimizer

**Program Optimizer** → Code Generator →$^{Program}$

**Training data** $\mathcal{D}$

# Learning-based Program Optimizer

# Learning-based Program Optimizer



**Unique Problem Characteristics**

- Relatively low experiment cost
- Domain-specific problem structure
- Large quantity of similar tasks

# Program-aware Cost Modeling

High-Level Configuration

# Program-aware Cost Modeling

**High-Level Configuration**

```
for y in range(8):
  for x in range(8):
    C[y][x]=0
    for k in range(8):
      C[y][x]+=A[k][y]*B[k][x]
```

Low-level Abstract Syntax Tree
(shared between tasks)

# Program-aware Cost Modeling

**High-Level Configuration**

touched memory

| | C | A | B |
|---|---|---|---|
| y | 64 | 64 | 64 |
| x | 8 | 8 | 64 |
| k | 1 | 8 | 8 |

outer loop length

| | |
|---|---|
| y | 1 |
| x | 8 |
| k | 64 |

statistical features

Boosted Tree Ensembles

```
for y in range(8):
  for x in range(8):
    C[y][x]=0
    for k in range(8):
      C[y][x]+=A[k][y]*B[k][x]
```

Low-level Abstract Syntax Tree
(shared between tasks)

# Program-aware Cost Modeling

**High-Level Configuration**

|   | touched memory | | | | outer loop length |
|---|---|---|---|---|---|
|   | C | A | B |   | |
| y | 64 | 64 | 64 | y | 1 |
| x | 8 | 8 | 64 | x | 8 |
| k | 1 | 8 | 8 | k | 64 |

→ Boosted Tree Ensembles

statistical features

```
for y in range(8):
  for x in range(8):
    C[y][x]=0
    for k in range(8):
      C[y][x]+=A[k][y]*B[k][x]
```

Low-level Abstract Syntax Tree
(shared between tasks)



TreeGRU

# Transfer Learning Among Different Workloads



Historical Optimization Tasks

Domain Invariant Program Representations

Transferable Models to speedup new tasks

Legend:
- GBT on Configuration $S$
- GBT on Flatten Loop Context $x$
- GBT on Context Relation $R$
- GBT No Transfer

(a) TITANX C7 in domain
(b) TITANX C1−C6 -> C7
(c) TITANX C1−C6 -> Matmul-1024
(d) Mali GPU C1−C6 -> A53 CPU C7

# State of Art Performance



**Nvidia GPU** — Legend: Tensorflow XLA, Tensorflow, MXNet, AutoTVM

**ARM CPU** — Legend: Tensorflow Lite, AutoTVM

**ARM GPU** — Legend: ARMComputeLib, AutoTVM

**tvm**.ai
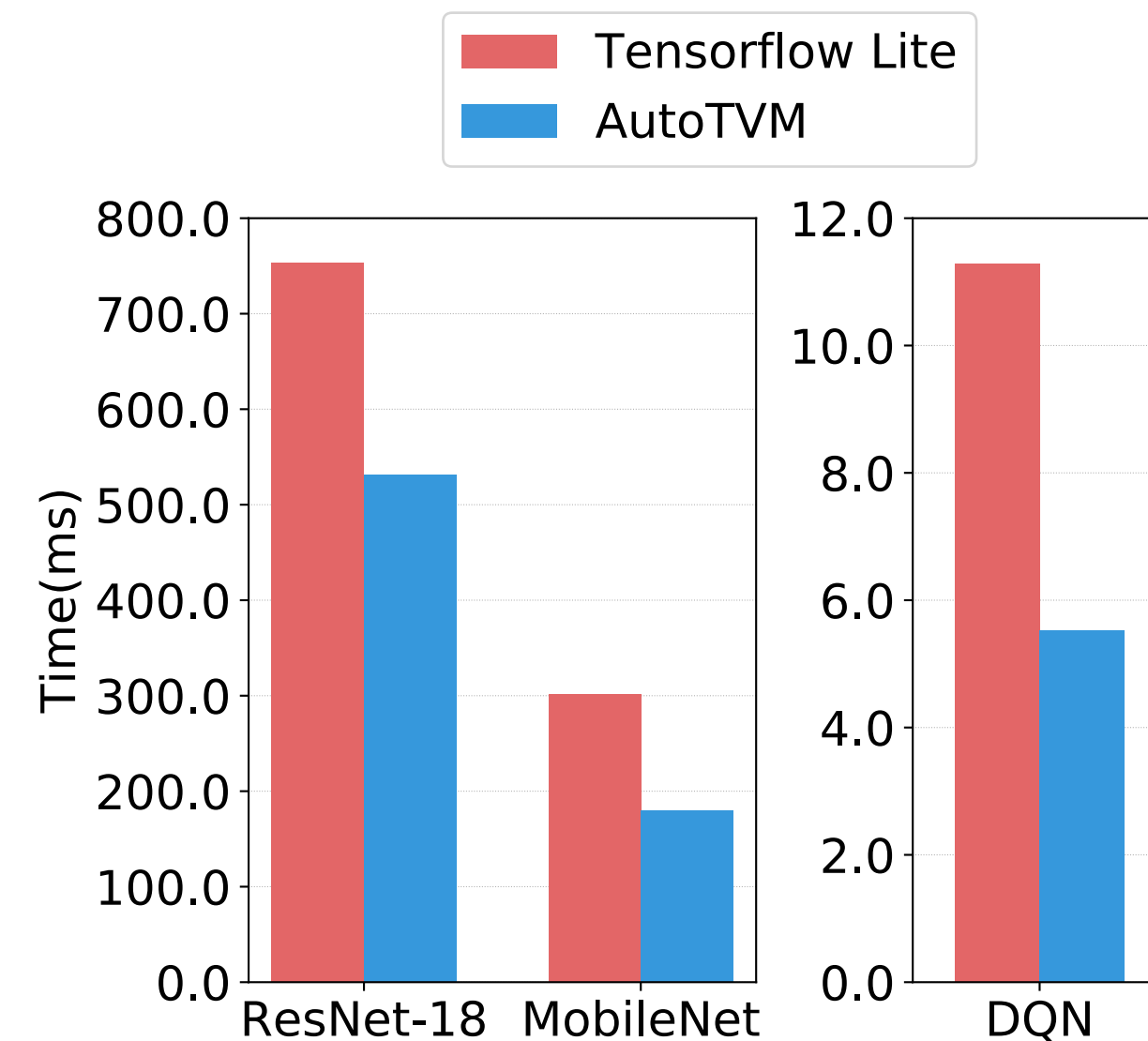
**Poster #104**

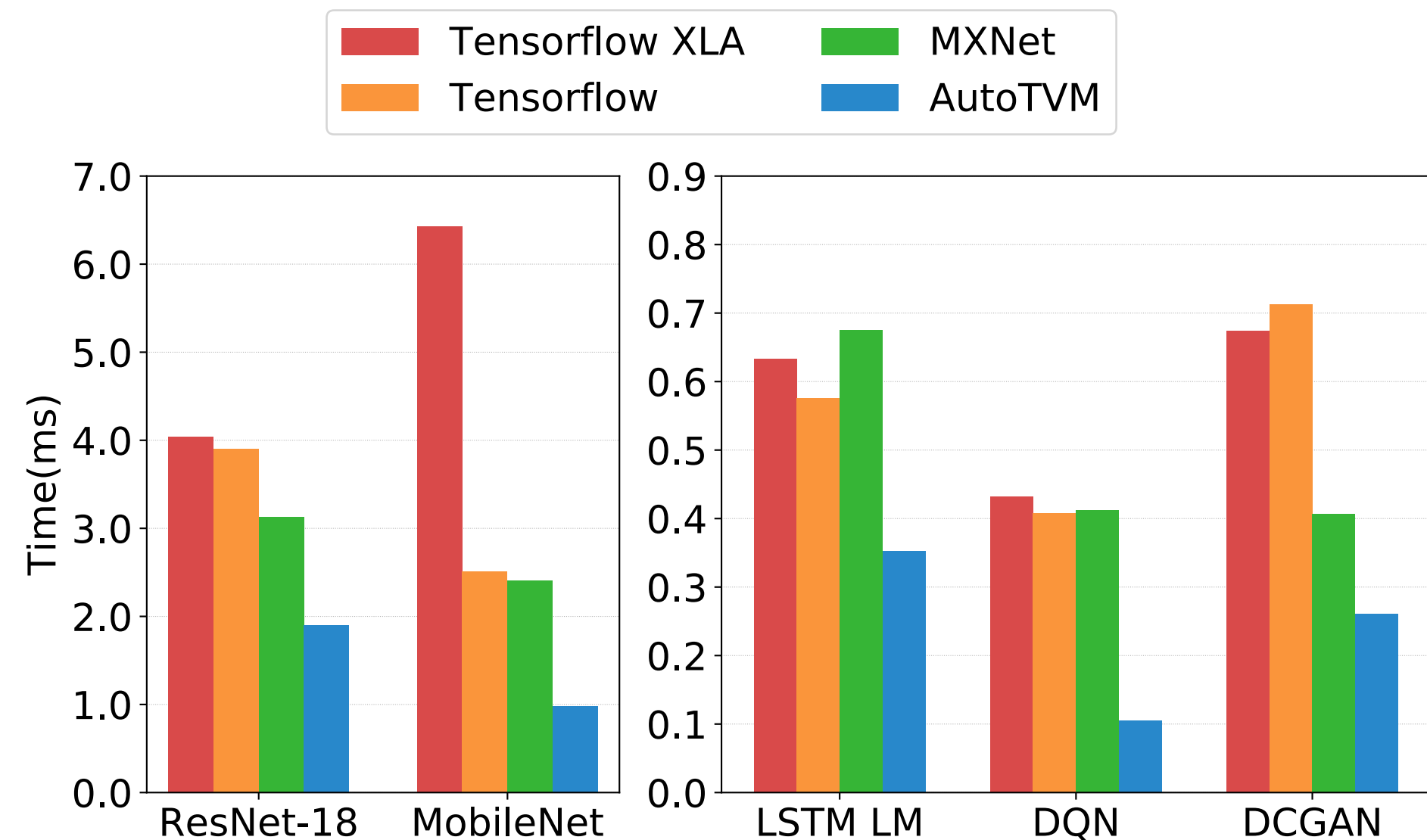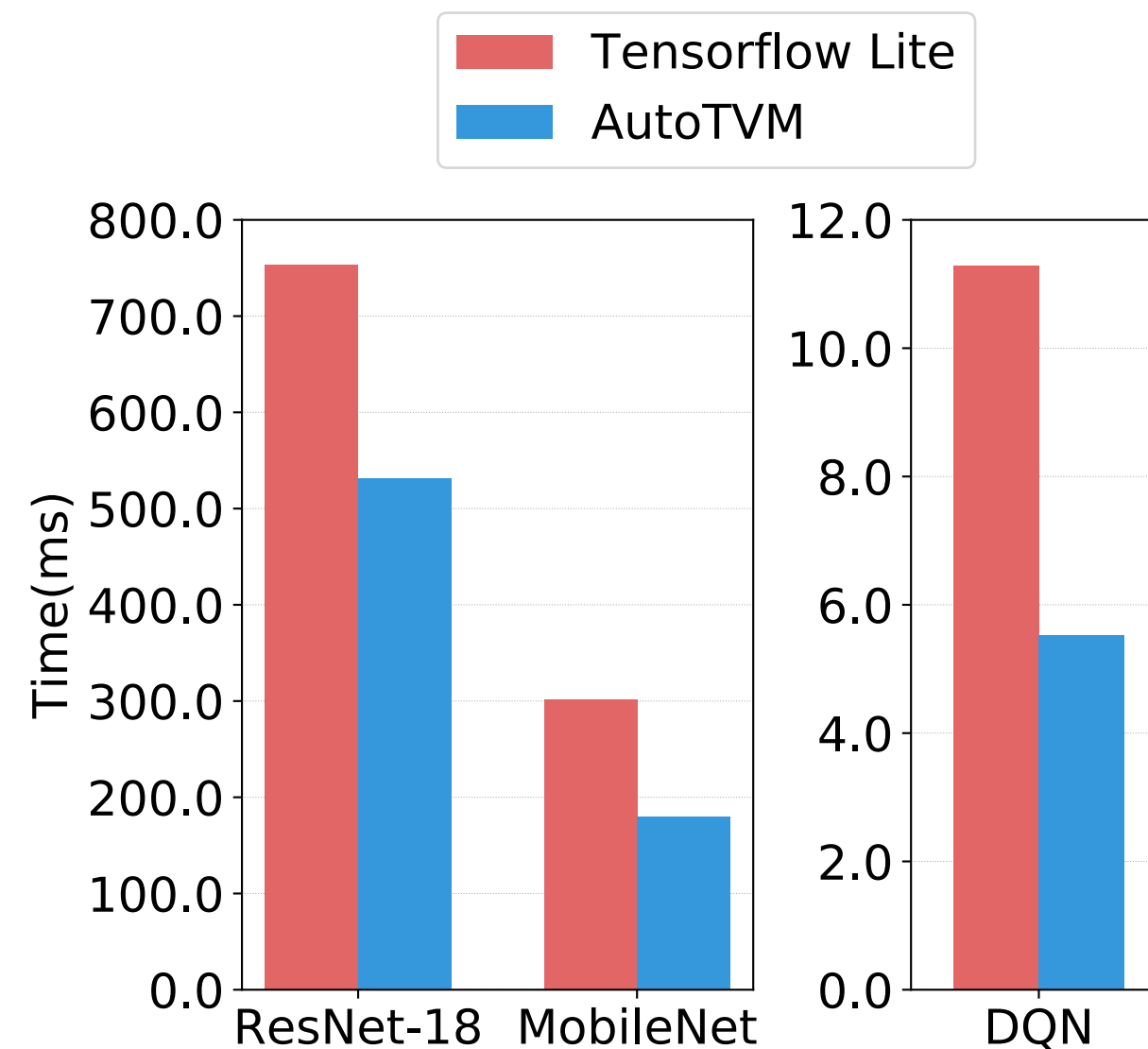# State of Art Performance
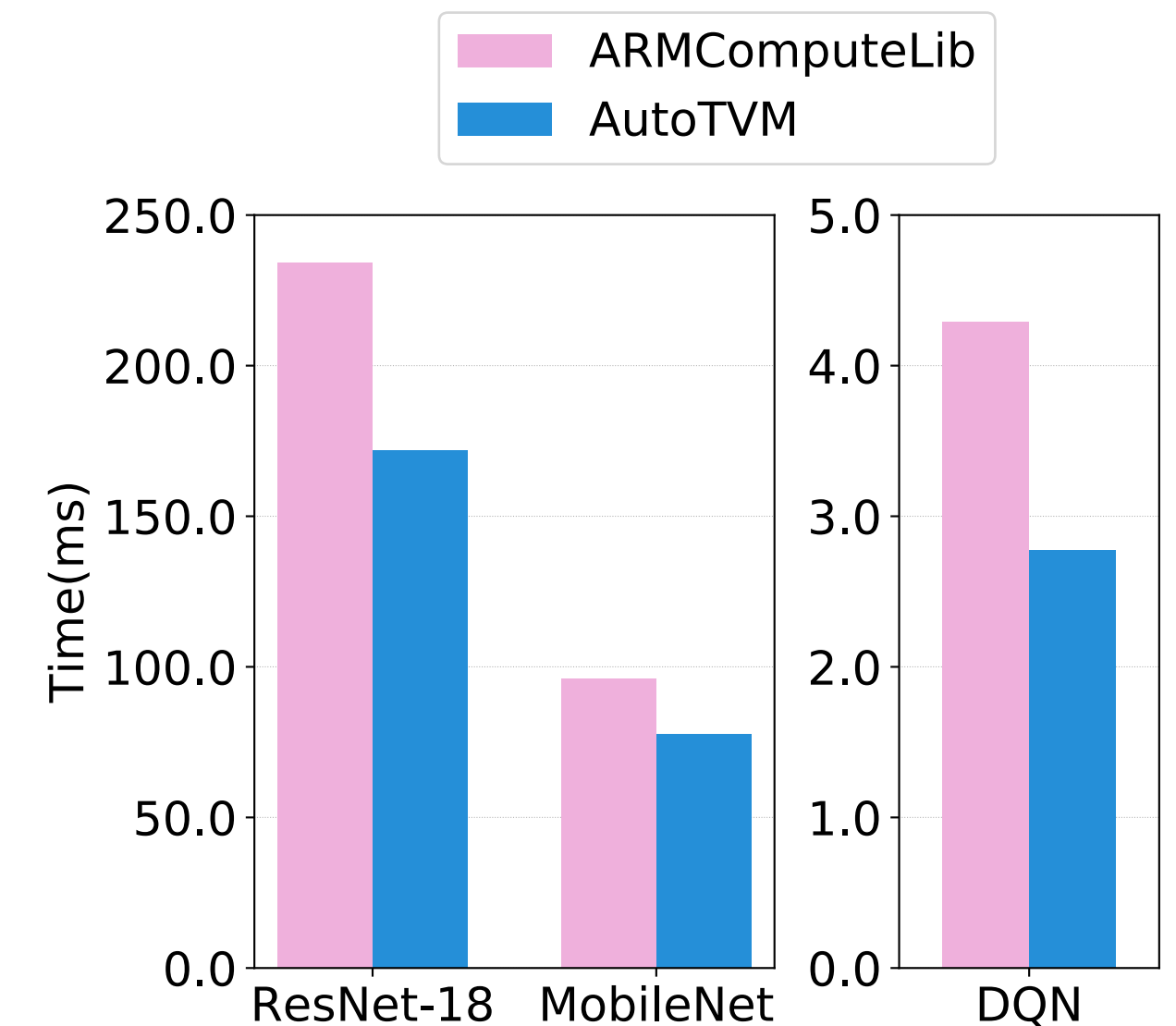
## Nvidia GPU



## ARM CPU



## ARM GPU



**In production use inside several major companies**

tvm.ai

**Poster #104**