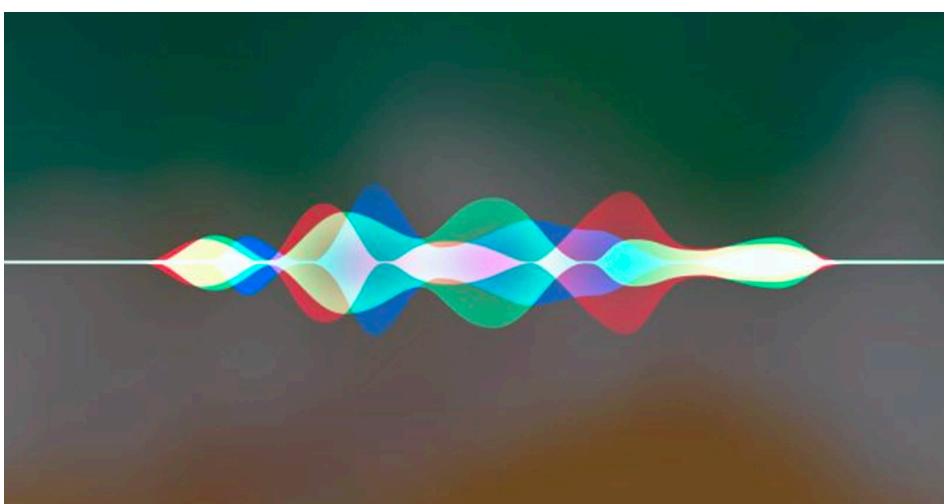


# TVM: An End to End Automated Deep Learning Compiler

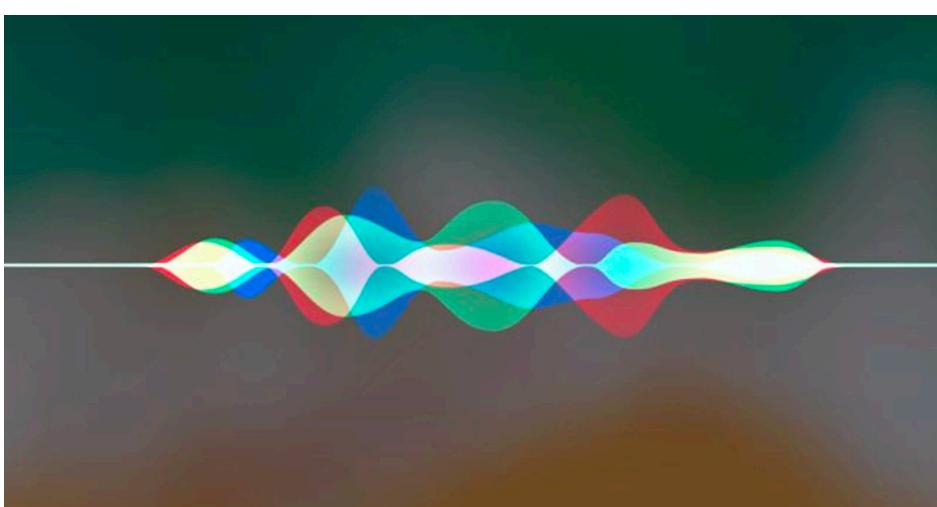
Tianqi Chen

# Machine Learning is Everywhere

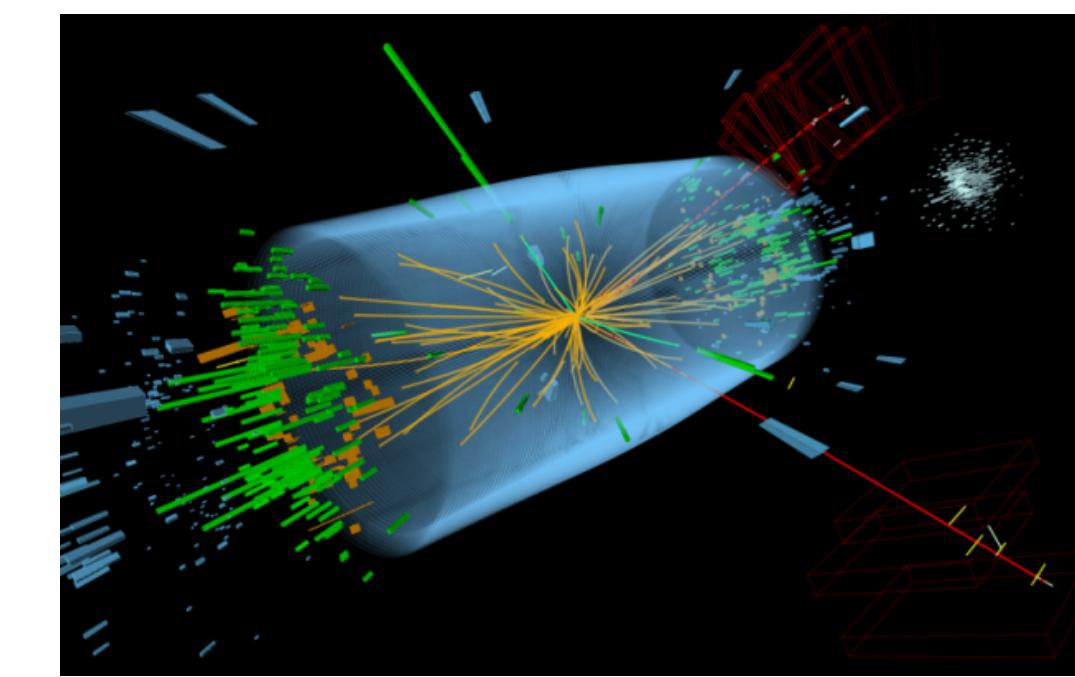
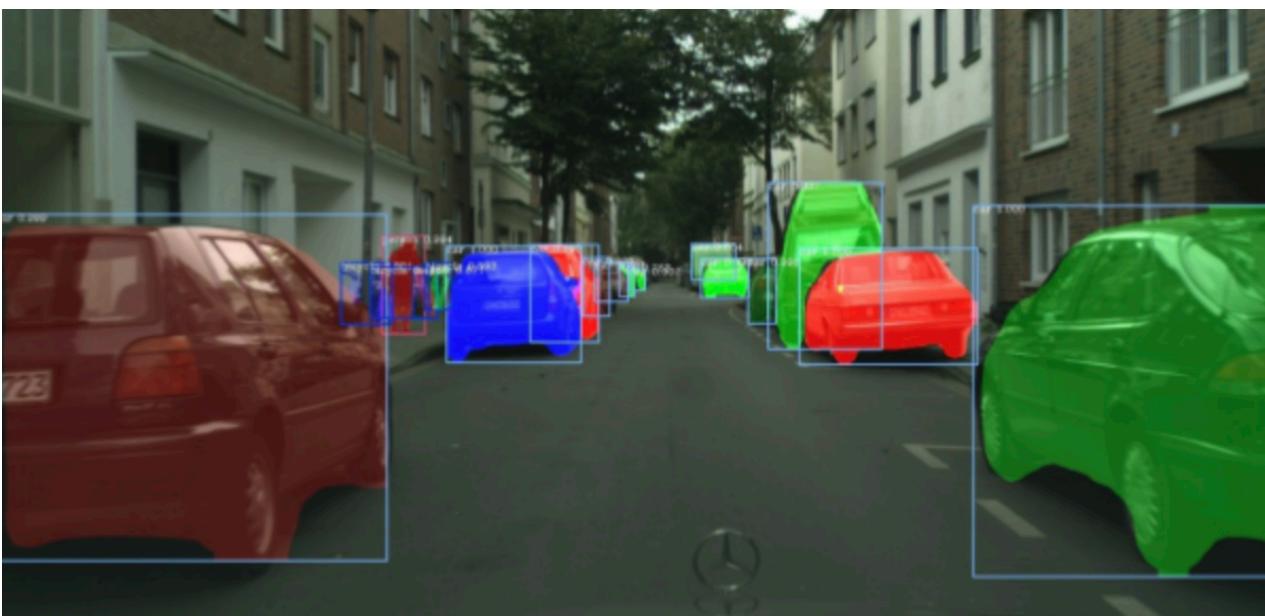
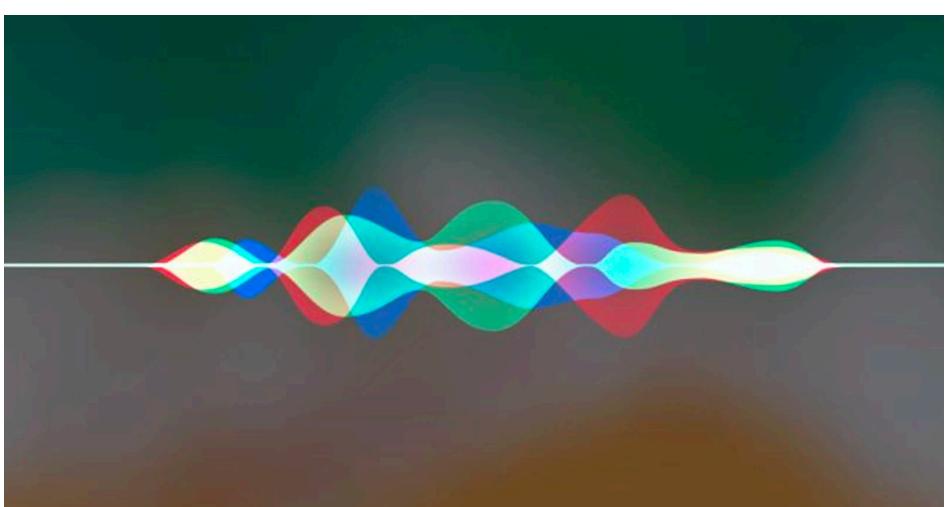
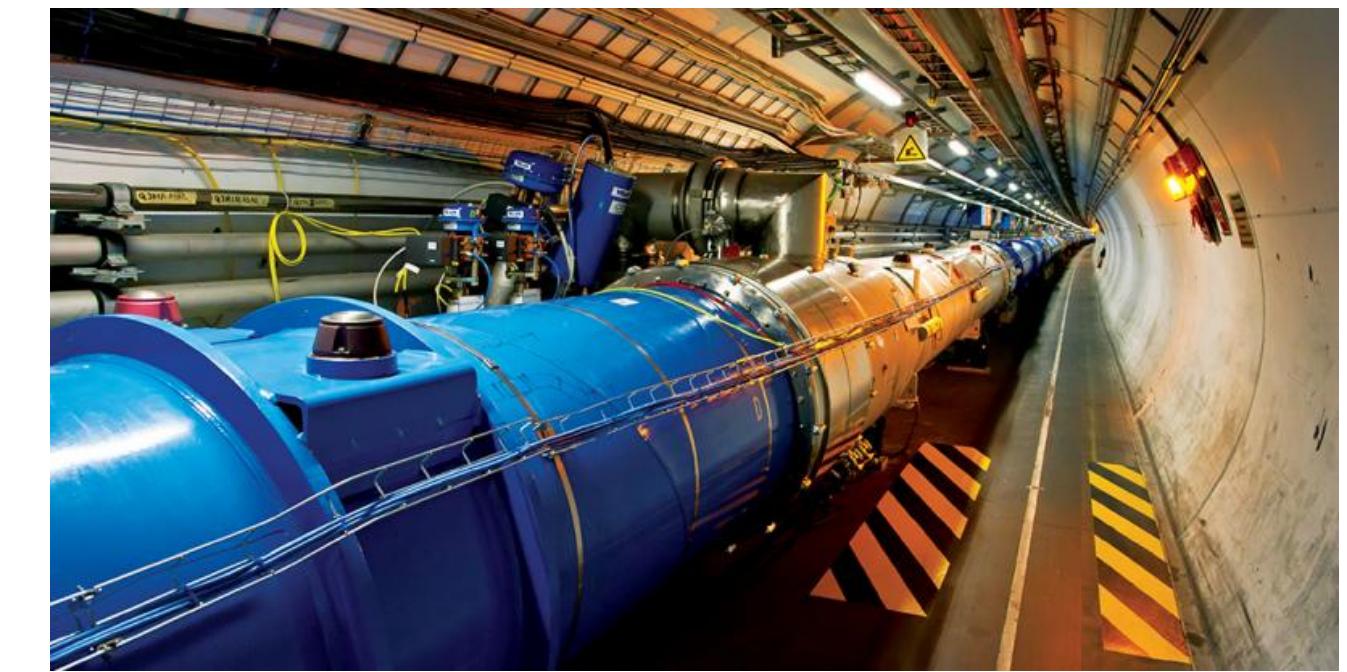
# Machine Learning is Everywhere



# Machine Learning is Everywhere



# Machine Learning is Everywhere



Higgs challenge

# Abbreviated History of Machine Learning

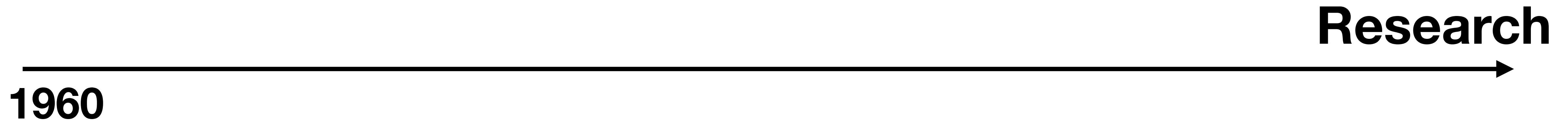
Based on personal view

# Abbreviated History of Machine Learning

**1960**

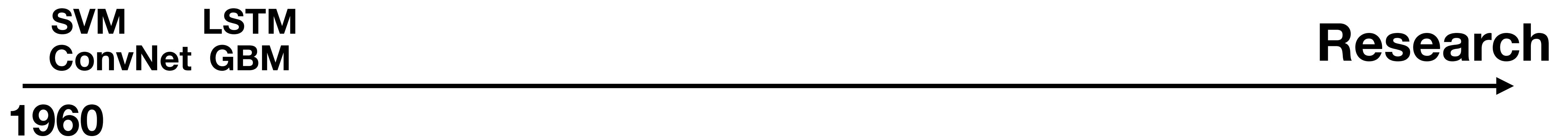
Based on personal view

# Abbreviated History of Machine Learning



Based on personal view

# Abbreviated History of Machine Learning



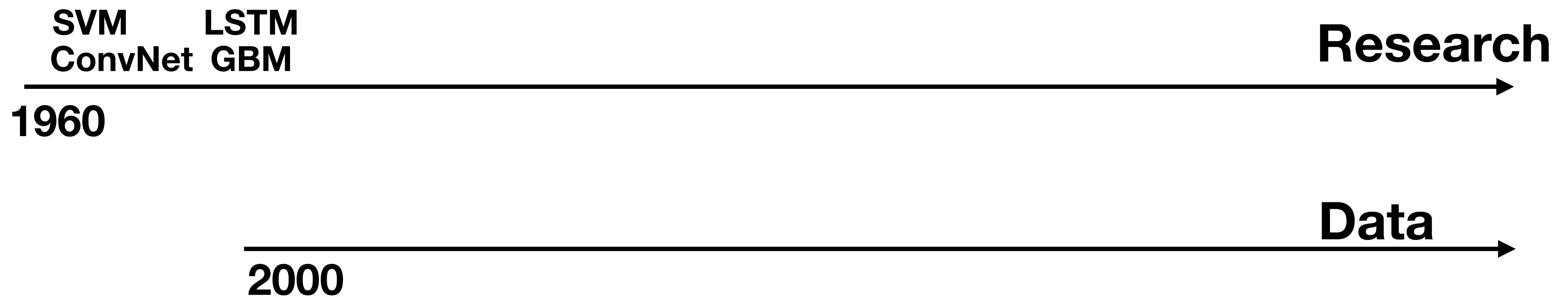
Based on personal view

# Abbreviated History of Machine Learning



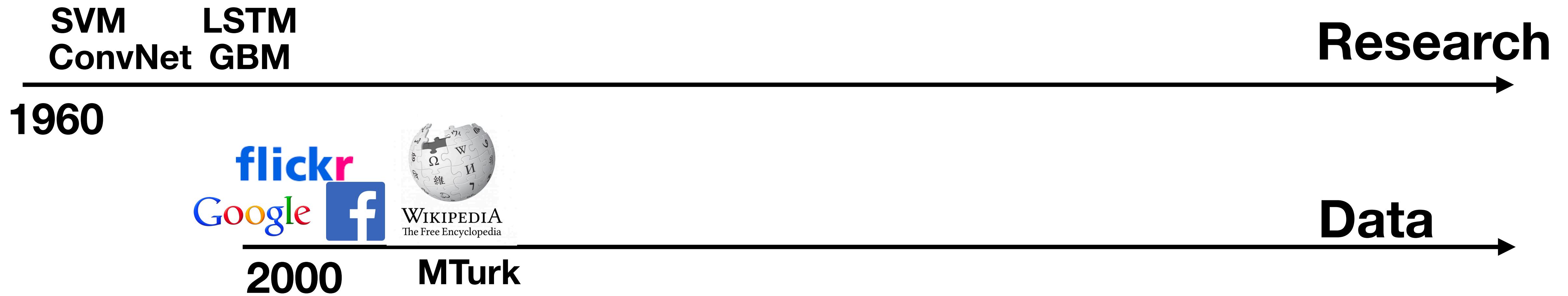
Based on personal view

# Abbreviated History of Machine Learning



Based on personal view

# Abbreviated History of Machine Learning



Based on personal view

# Abbreviated History of Machine Learning



Based on personal view

# Abbreviated History of Machine Learning



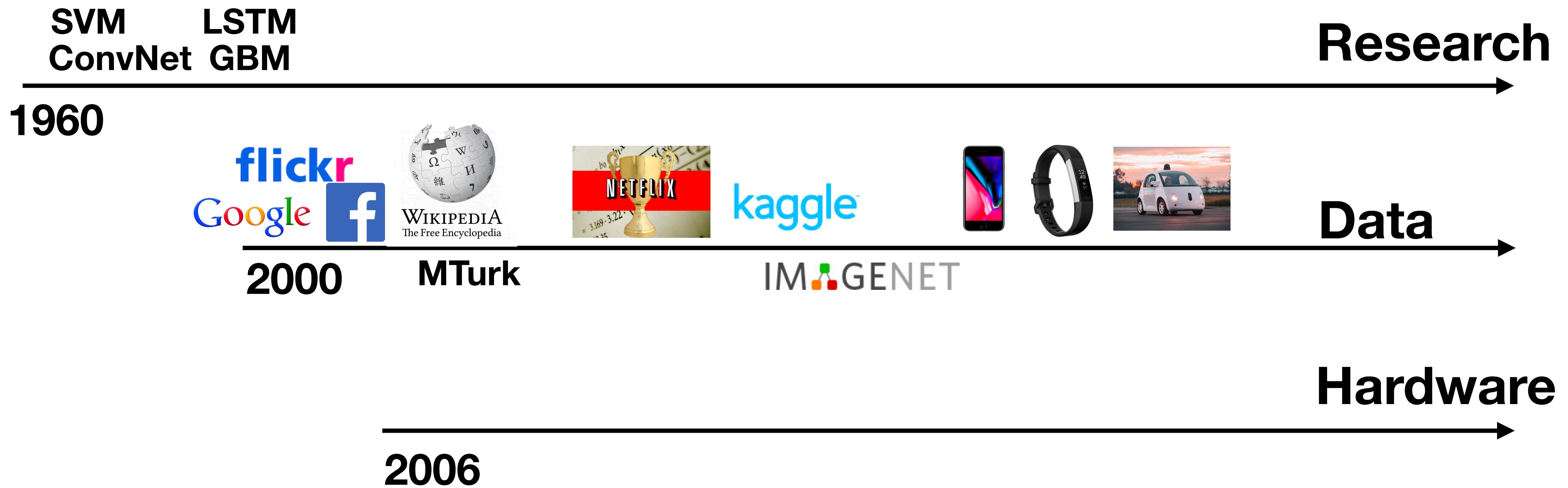
Based on personal view

# Abbreviated History of Machine Learning



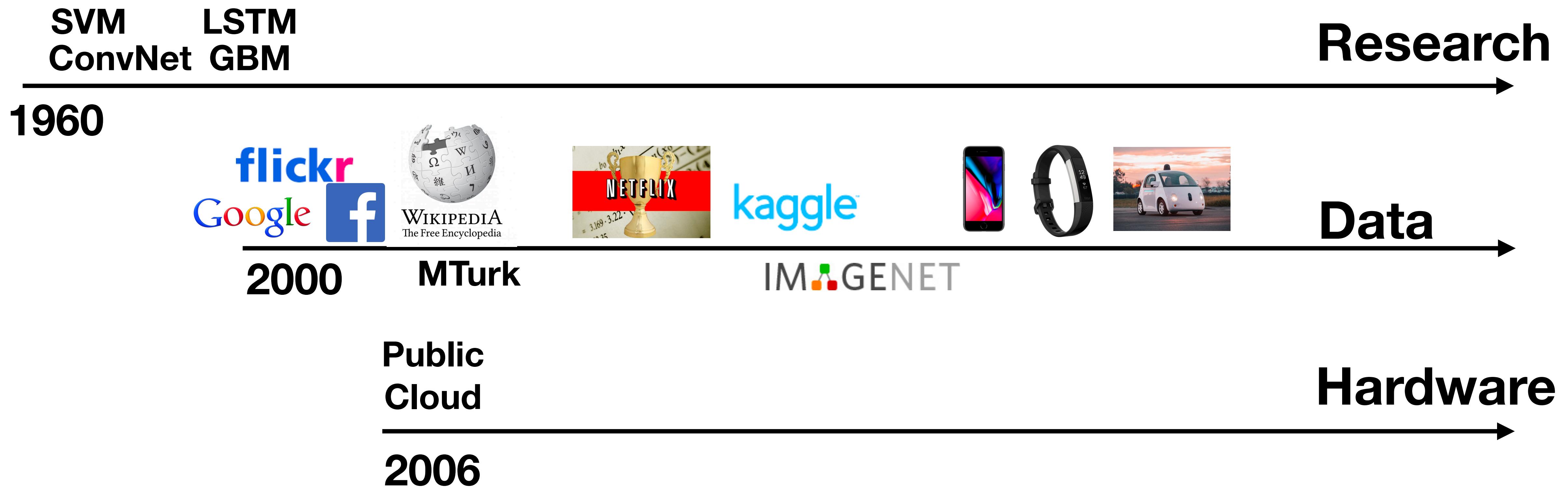
Based on personal view

# Abbreviated History of Machine Learning



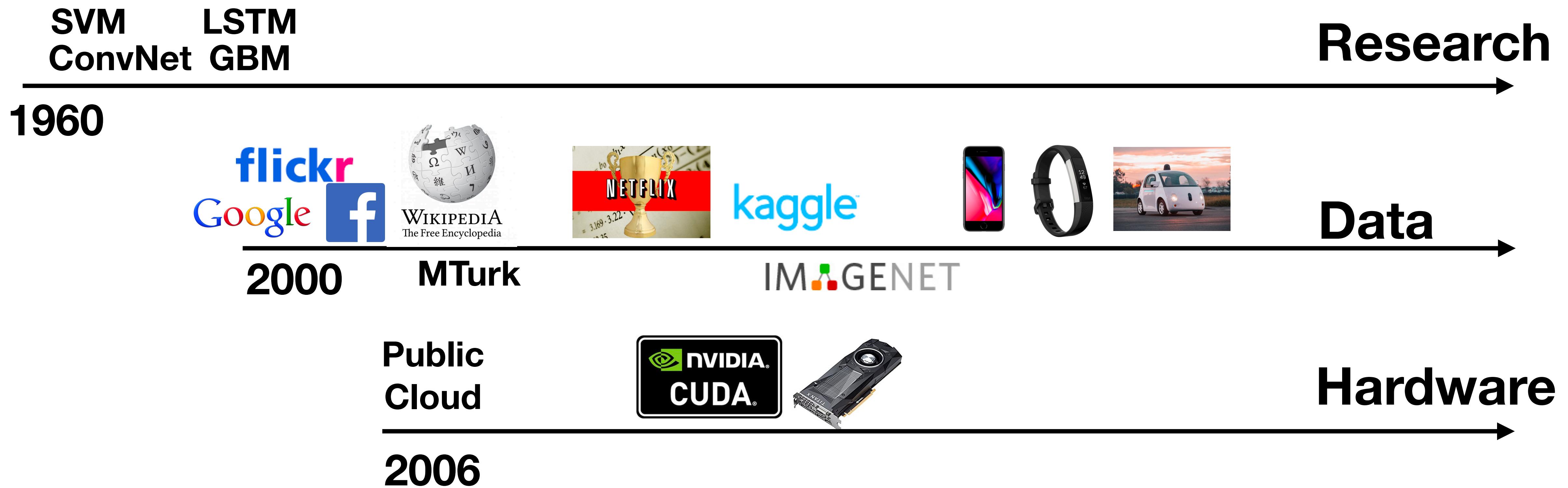
Based on personal view

# Abbreviated History of Machine Learning



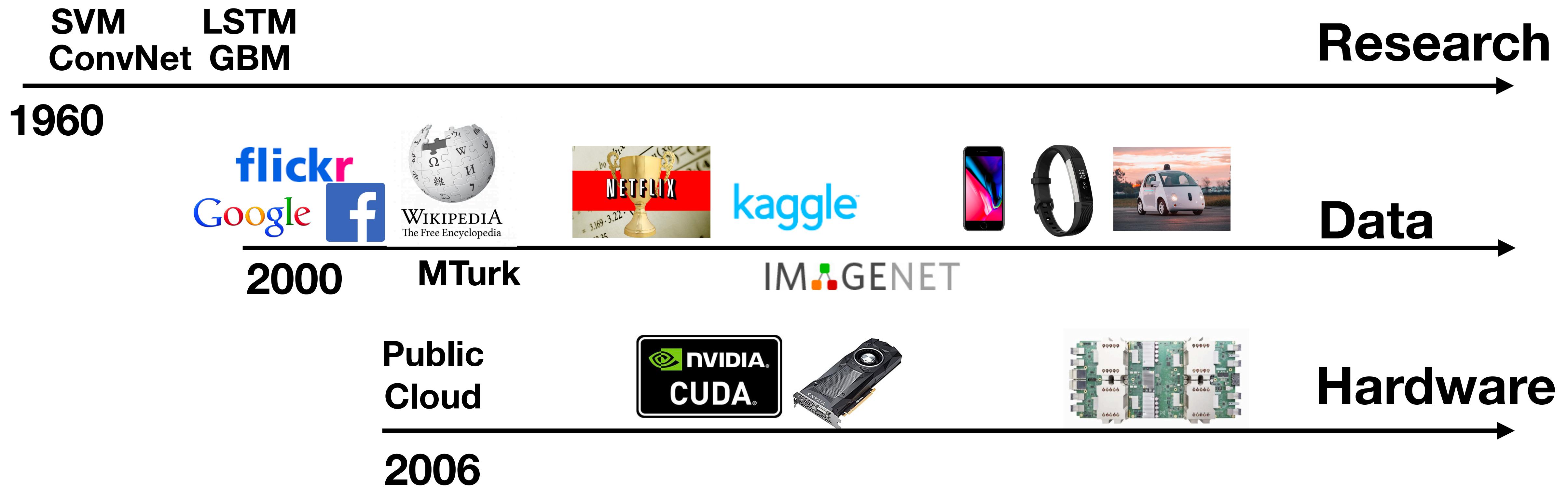
Based on personal view

# Abbreviated History of Machine Learning



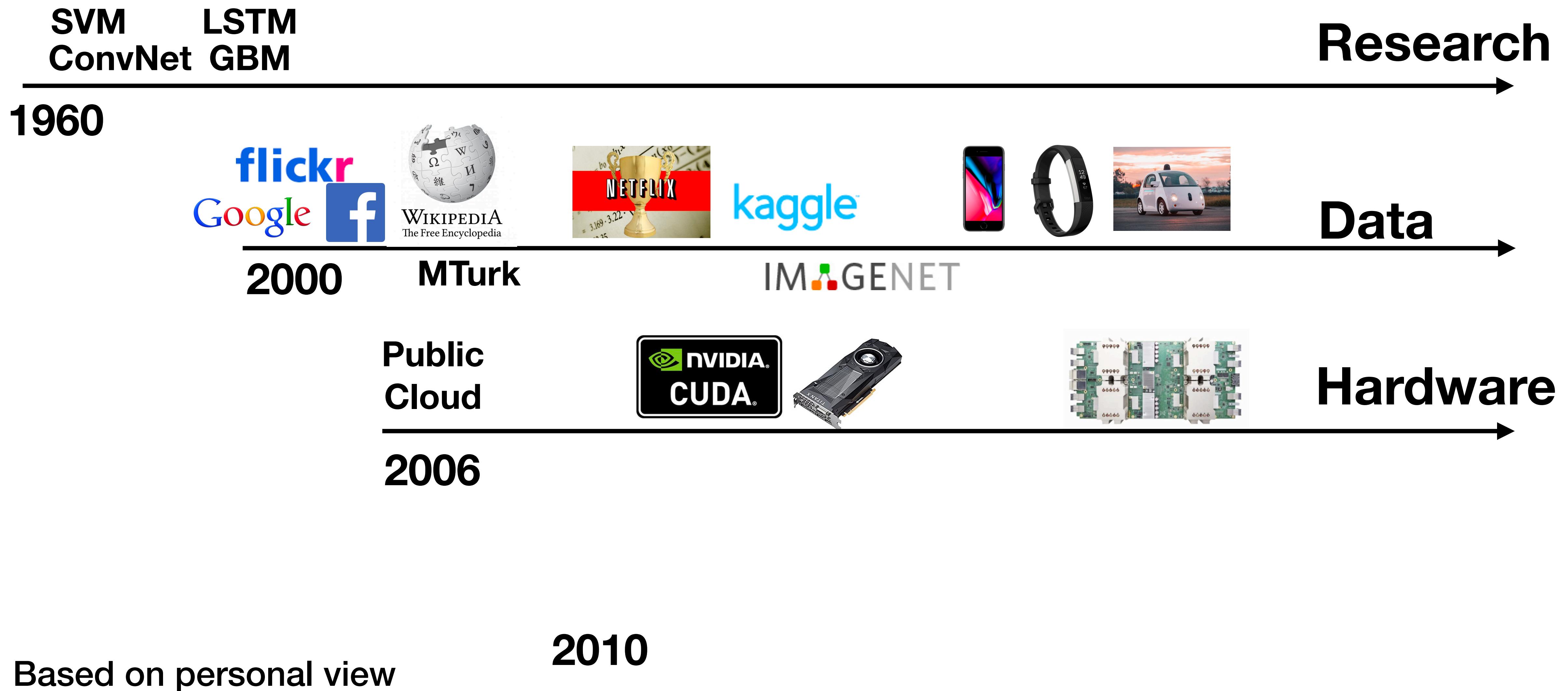
Based on personal view

# Abbreviated History of Machine Learning

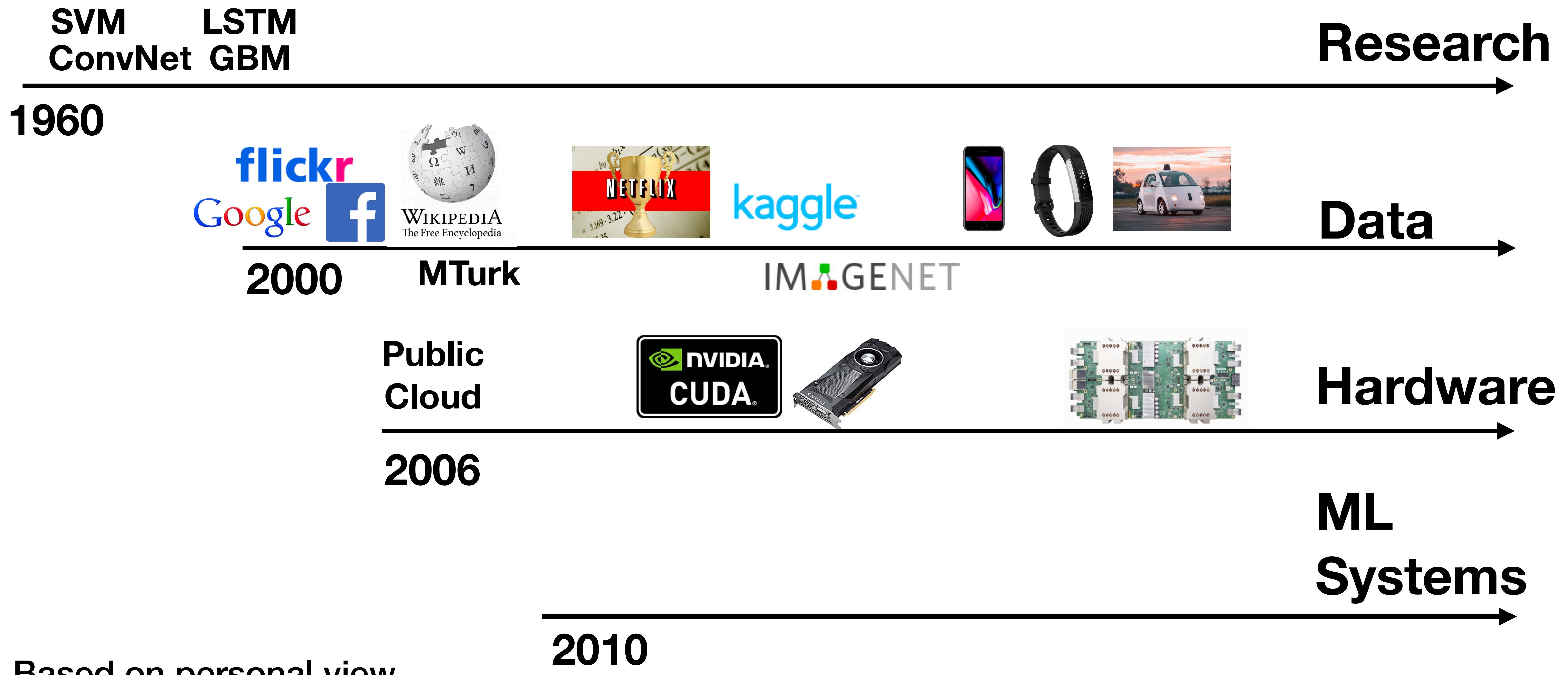


Based on personal view

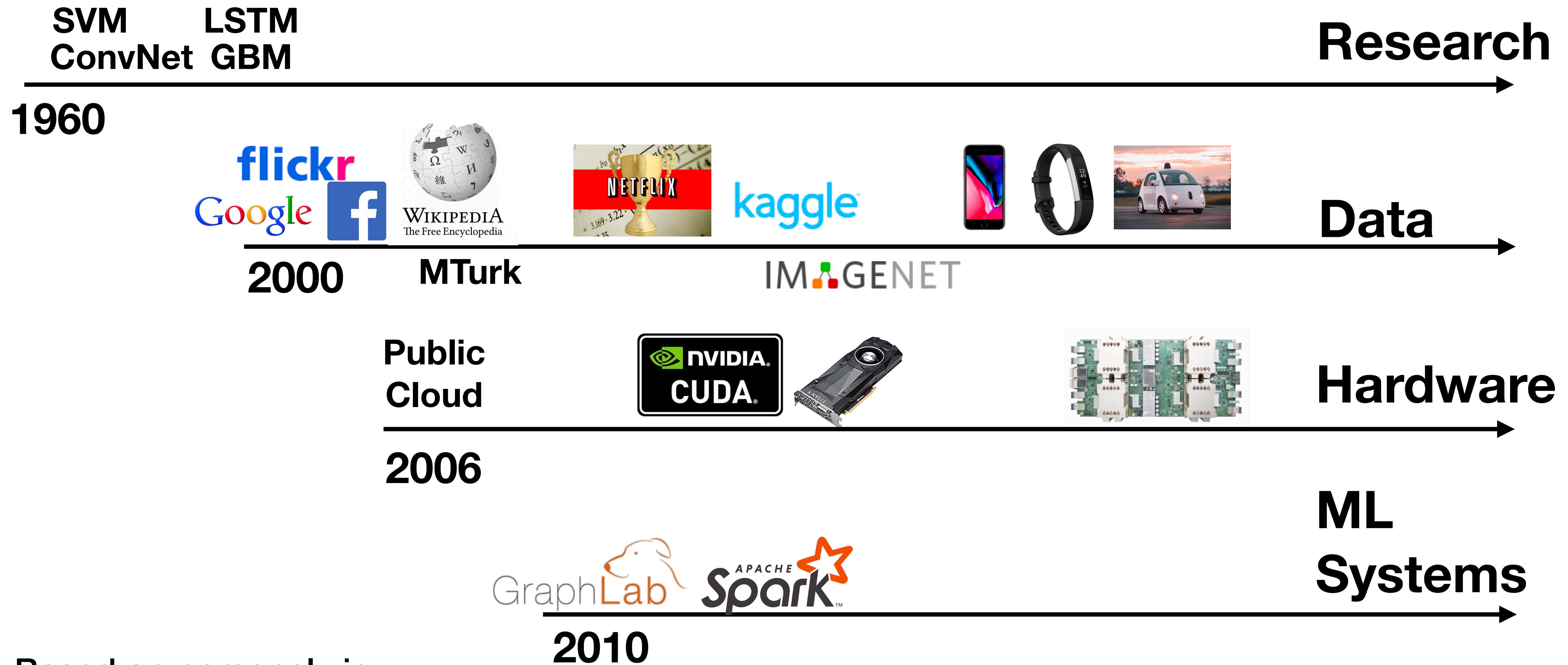
# Abbreviated History of Machine Learning



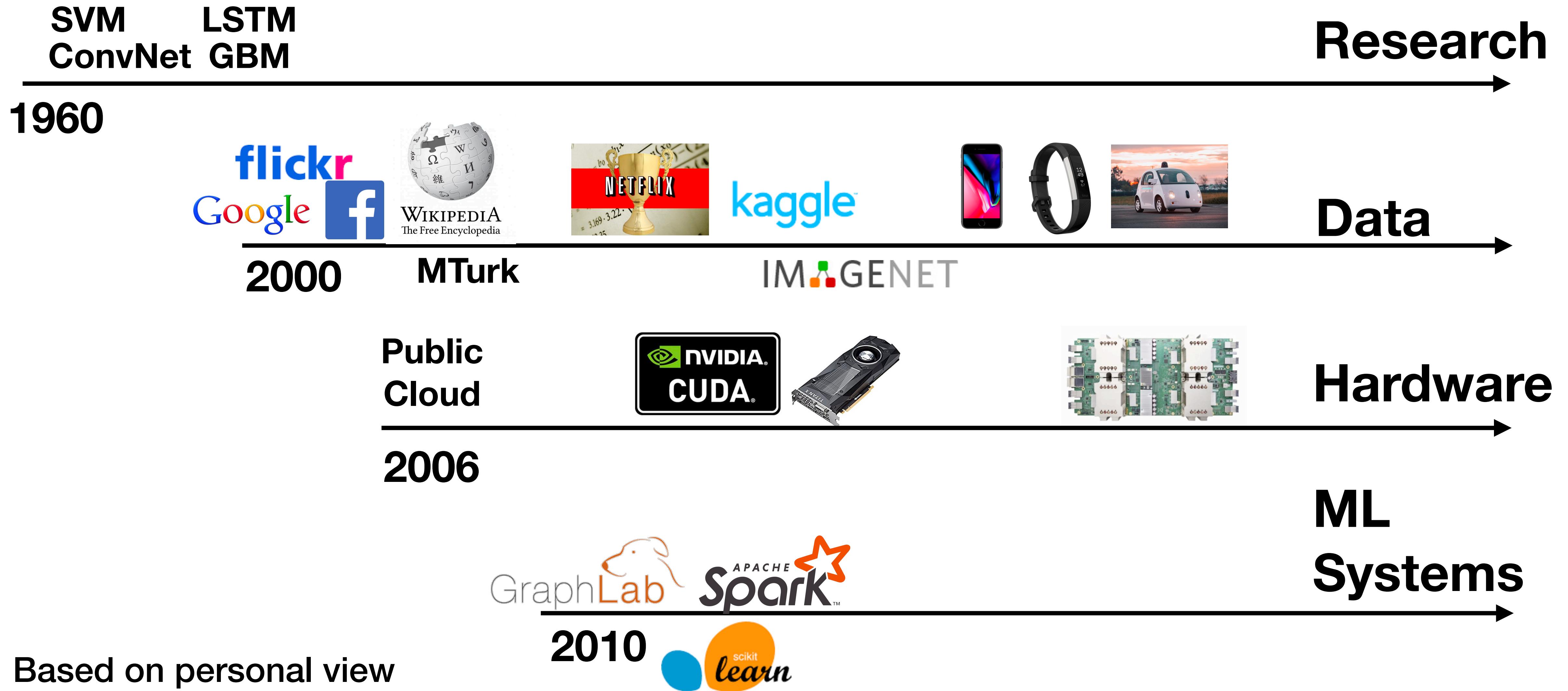
# Abbreviated History of Machine Learning



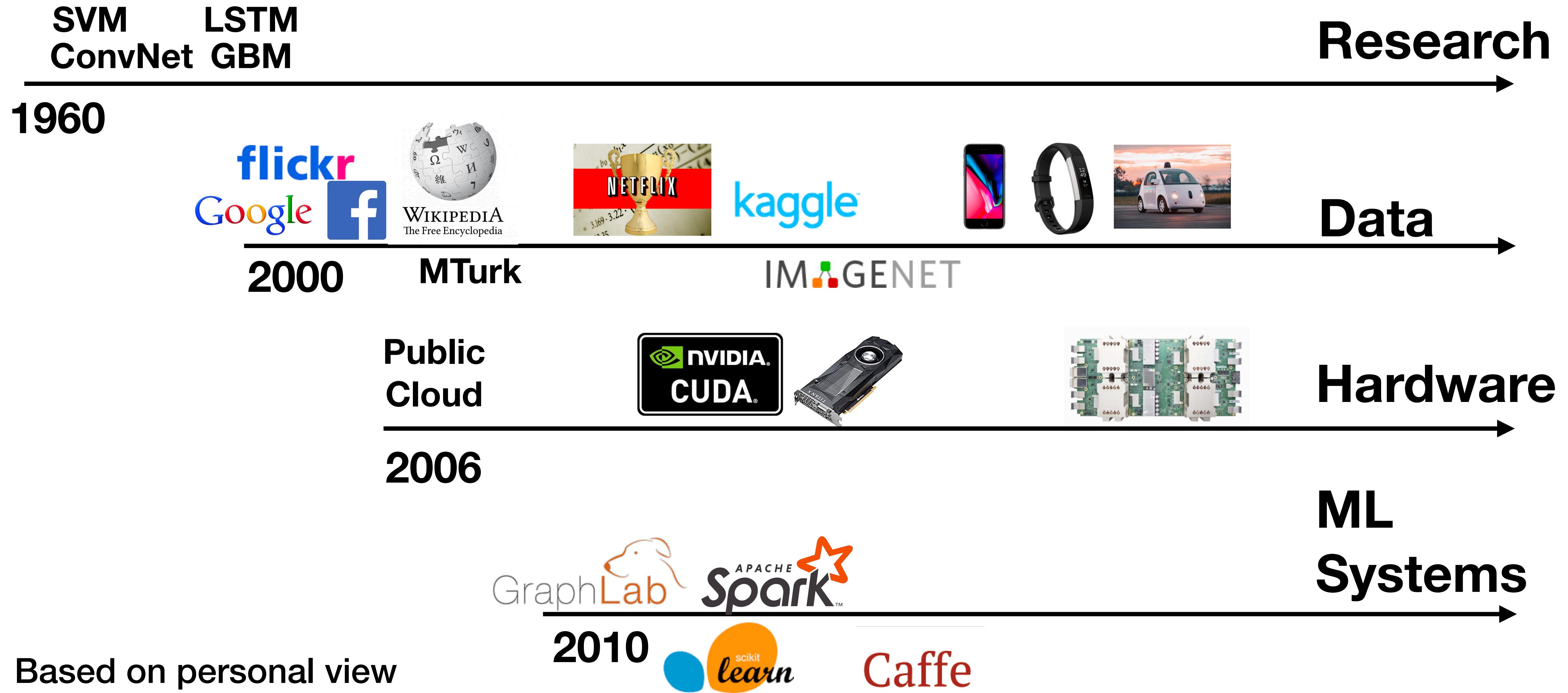
# Abbreviated History of Machine Learning



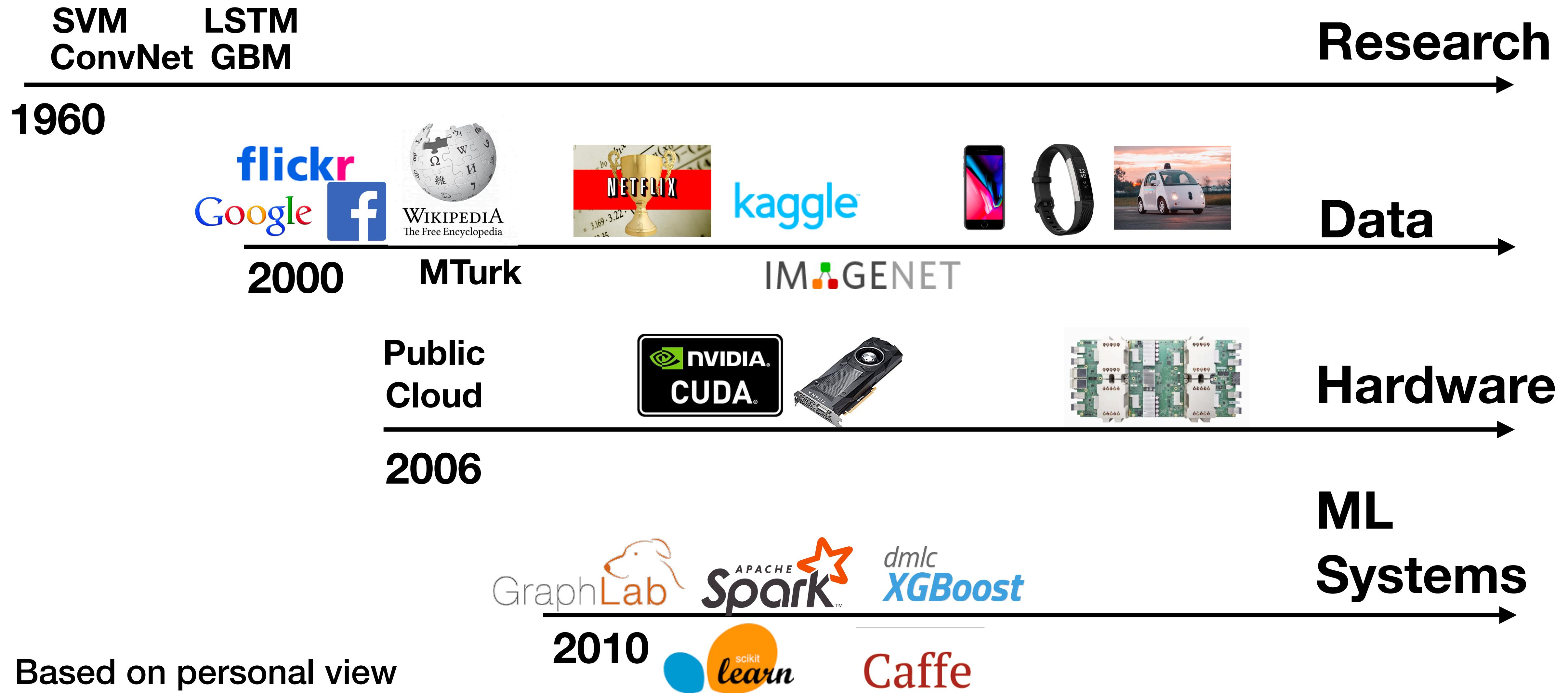
# Abbreviated History of Machine Learning



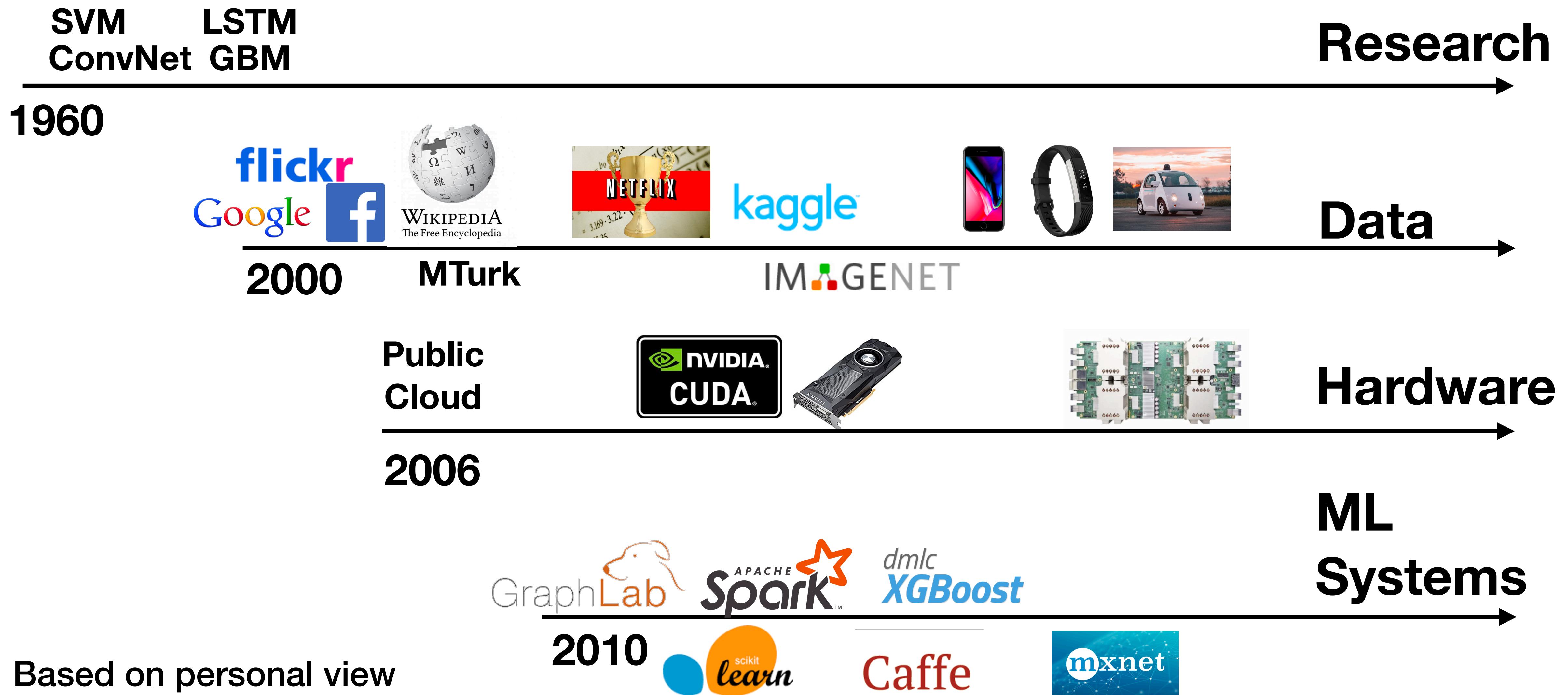
# Abbreviated History of Machine Learning



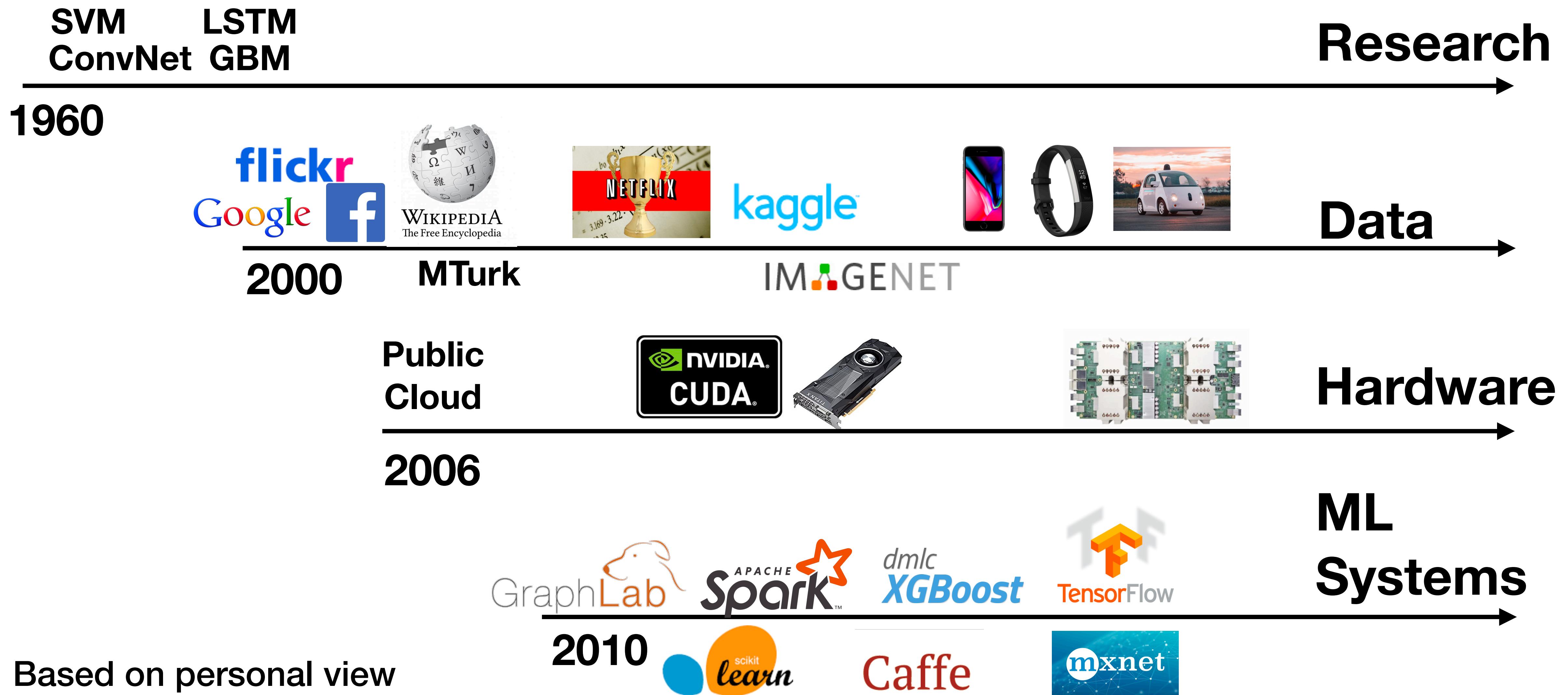
# Abbreviated History of Machine Learning



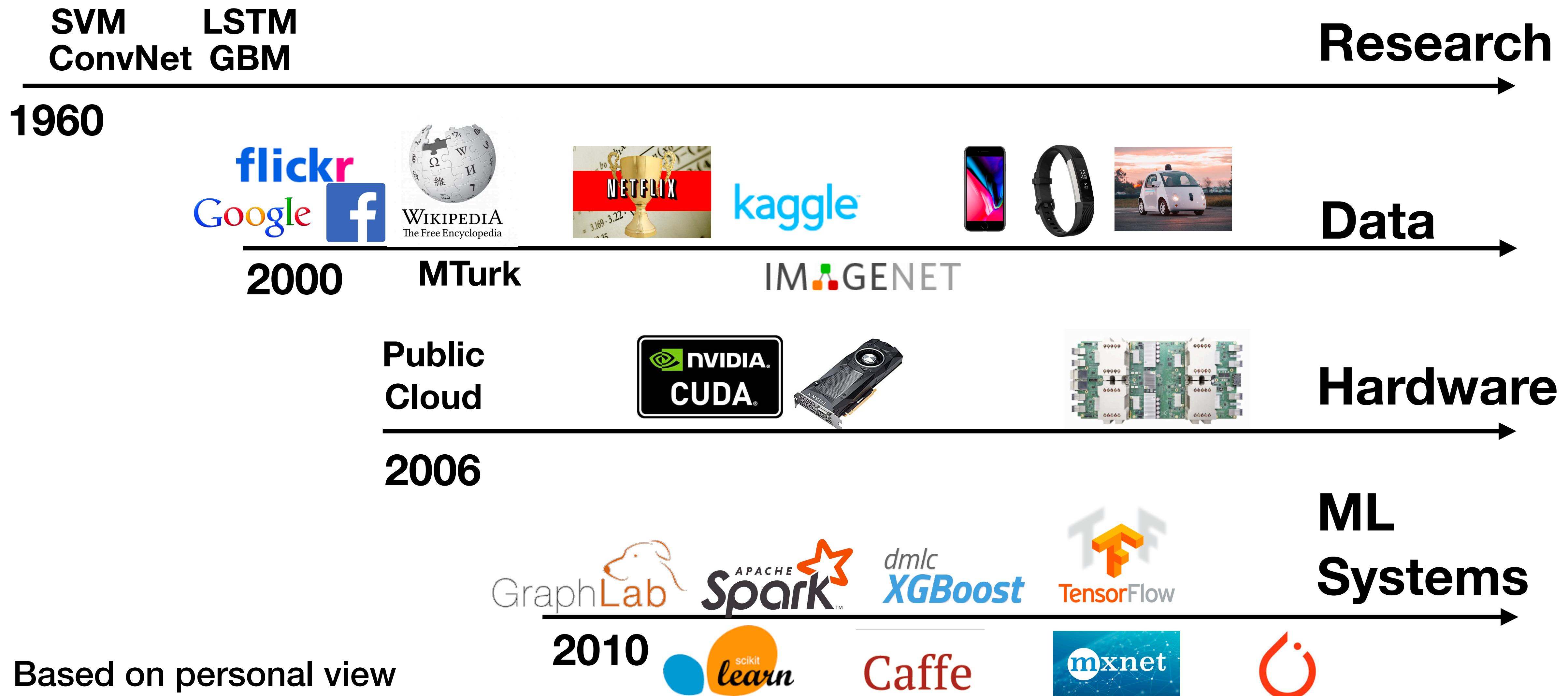
# Abbreviated History of Machine Learning



# Abbreviated History of Machine Learning

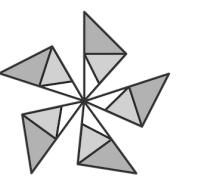
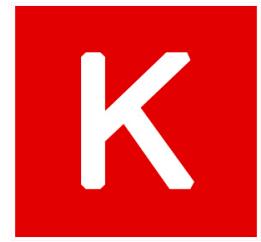
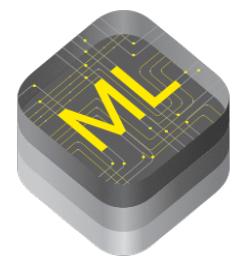


# Abbreviated History of Machine Learning



# Current Deep Learning Systems ecosystem

Frameworks and  
Inference engines



ONNX  
RUNTIME



DL Compilers



Kernel Libraries

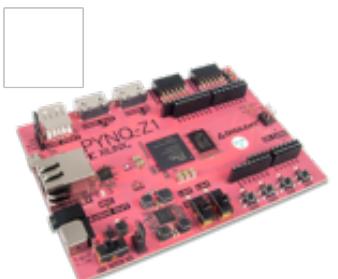
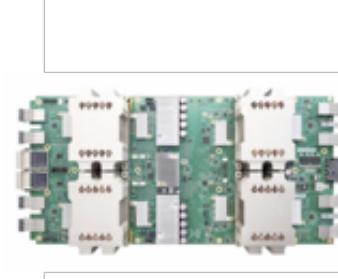
CuDNN

NNPack

MKL-DNN

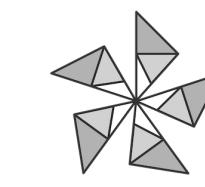
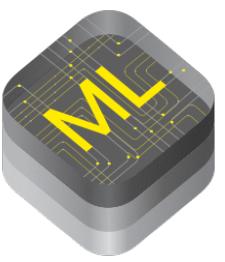
Hand optimized

Hardware

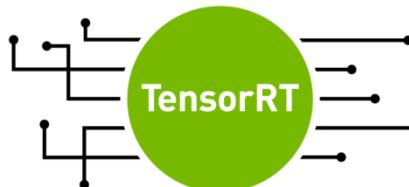


# Current Deep Learning Systems ecosystem

Frameworks and  
Inference engines



ONNX  
RUNTIME



DL Compilers



Kernel Libraries

CuDNN

NNPack

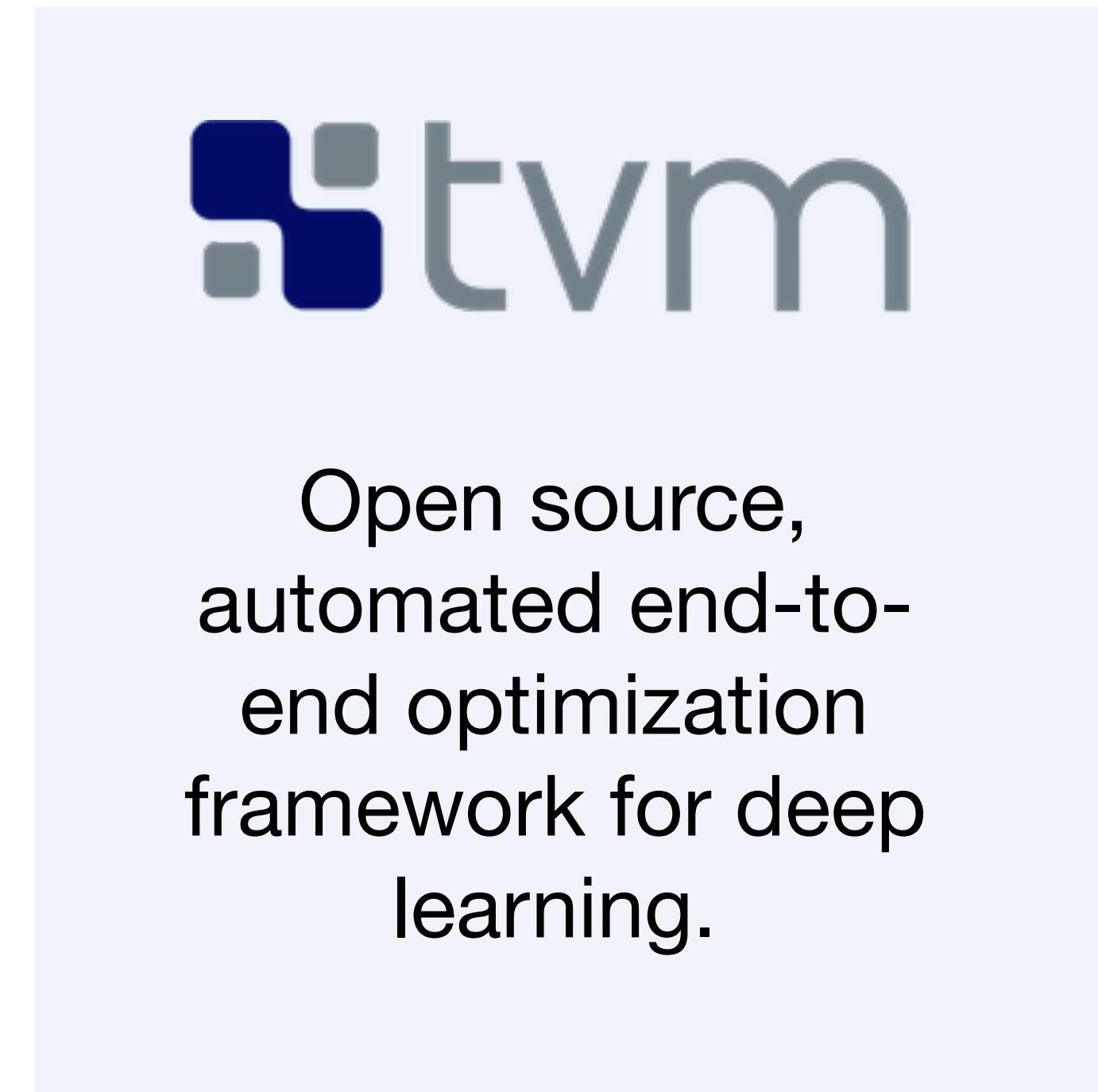
MKL-DNN

Hand optimized

Hardware

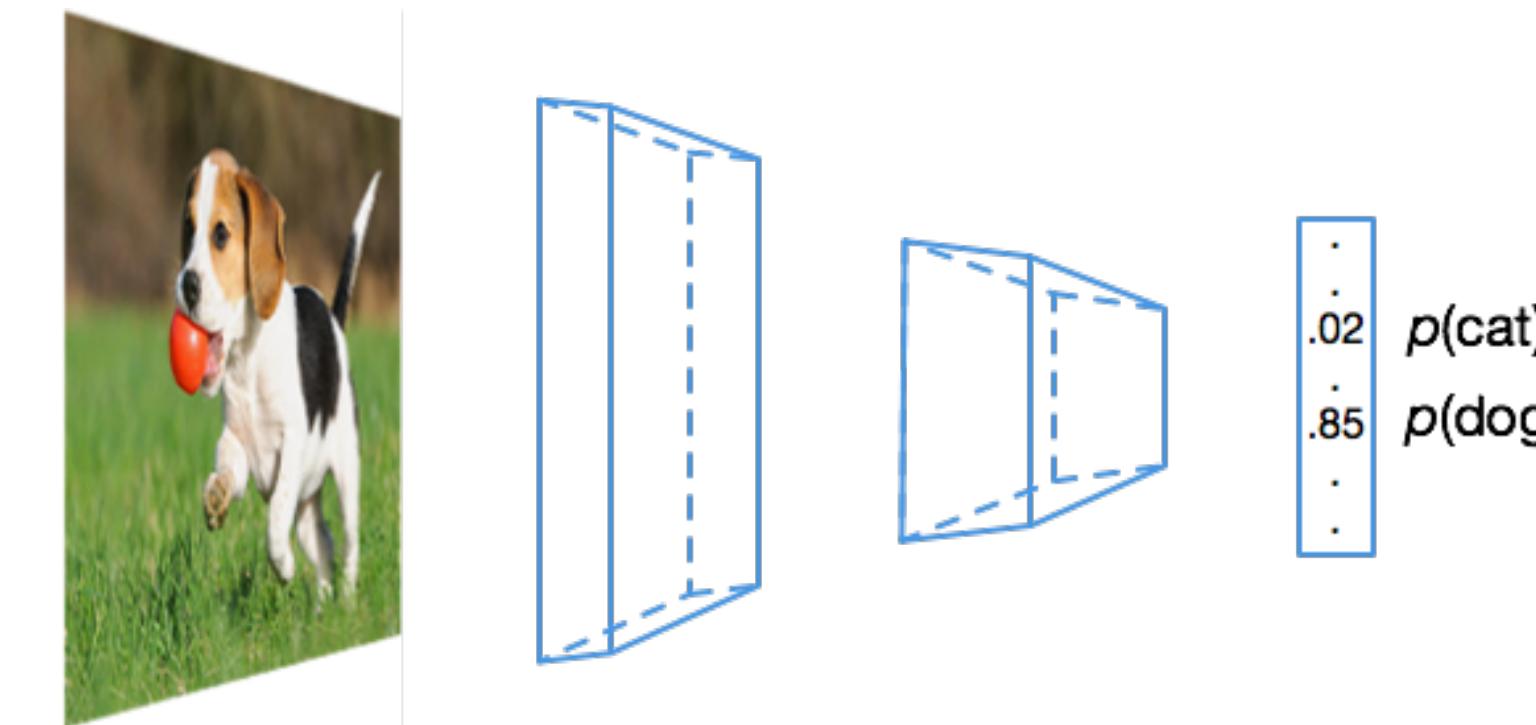


Open source,  
automated end-to-  
end optimization  
framework for deep  
learning.



# Problem: Deep Learning Deployment

Model

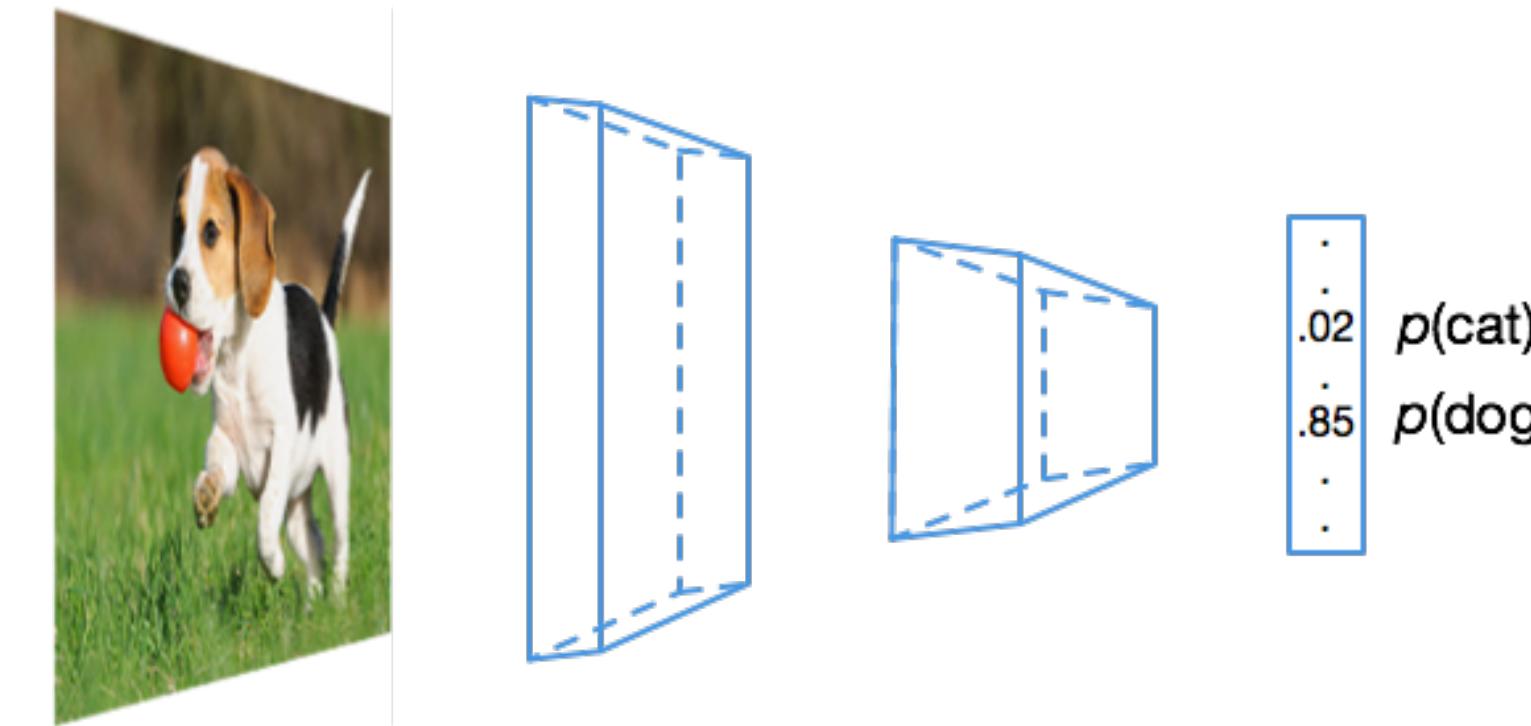


Hardware  
Backends

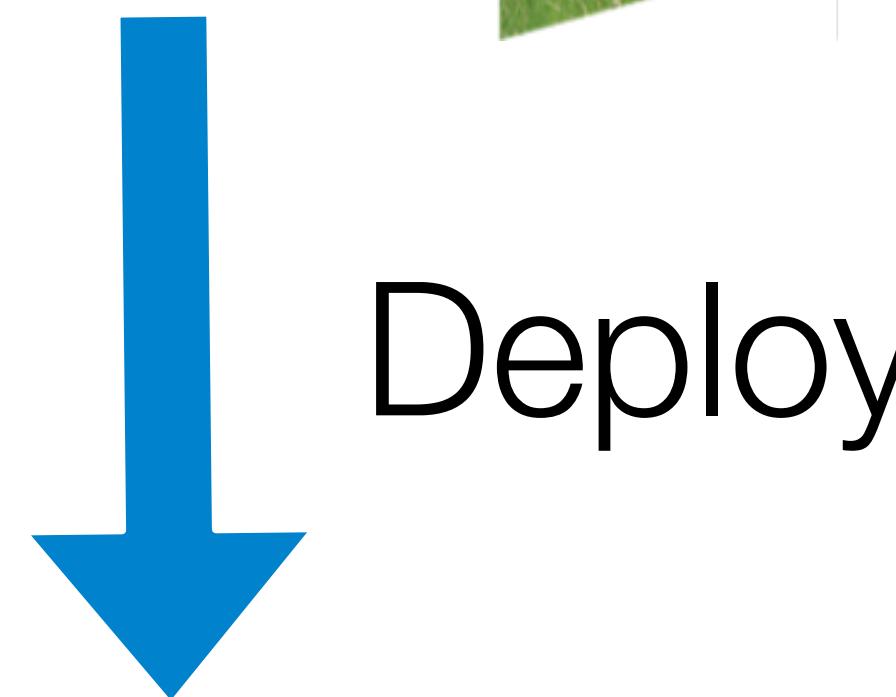


# Problem: Deep Learning Deployment

Model



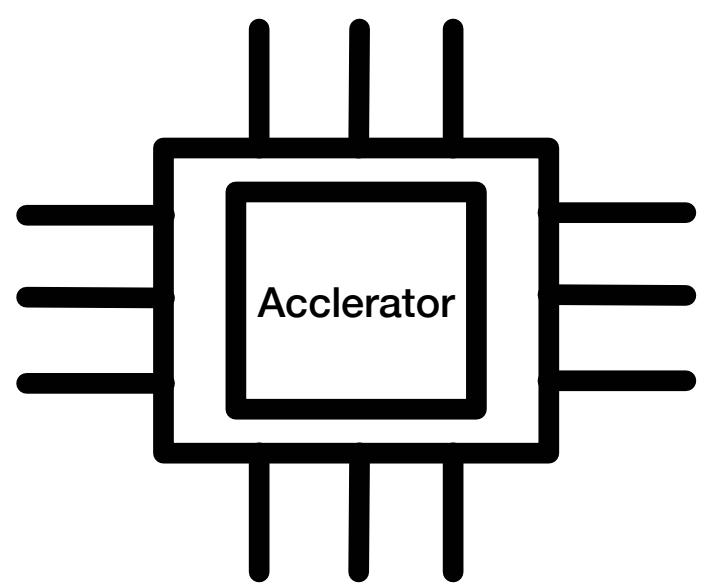
Hardware  
Backends



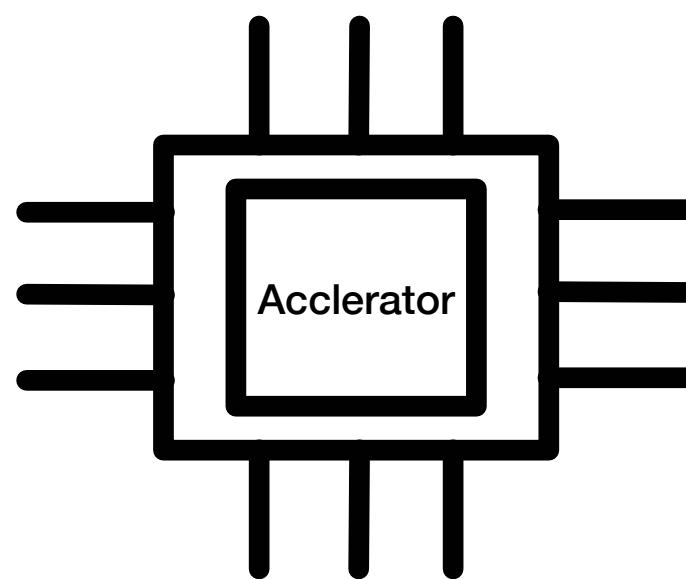
Deploy

# Original Motivating Problem

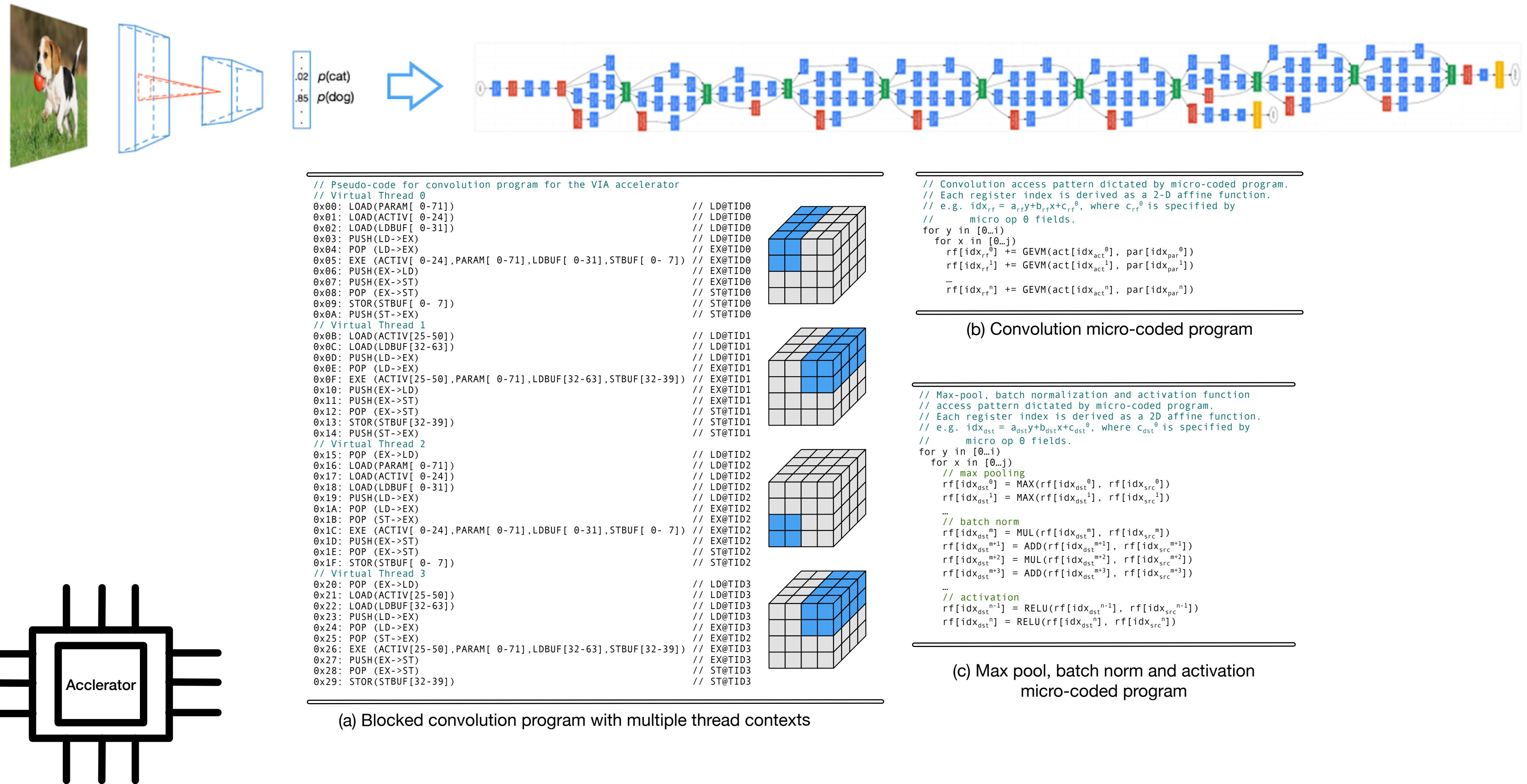
# Original Motivating Problem



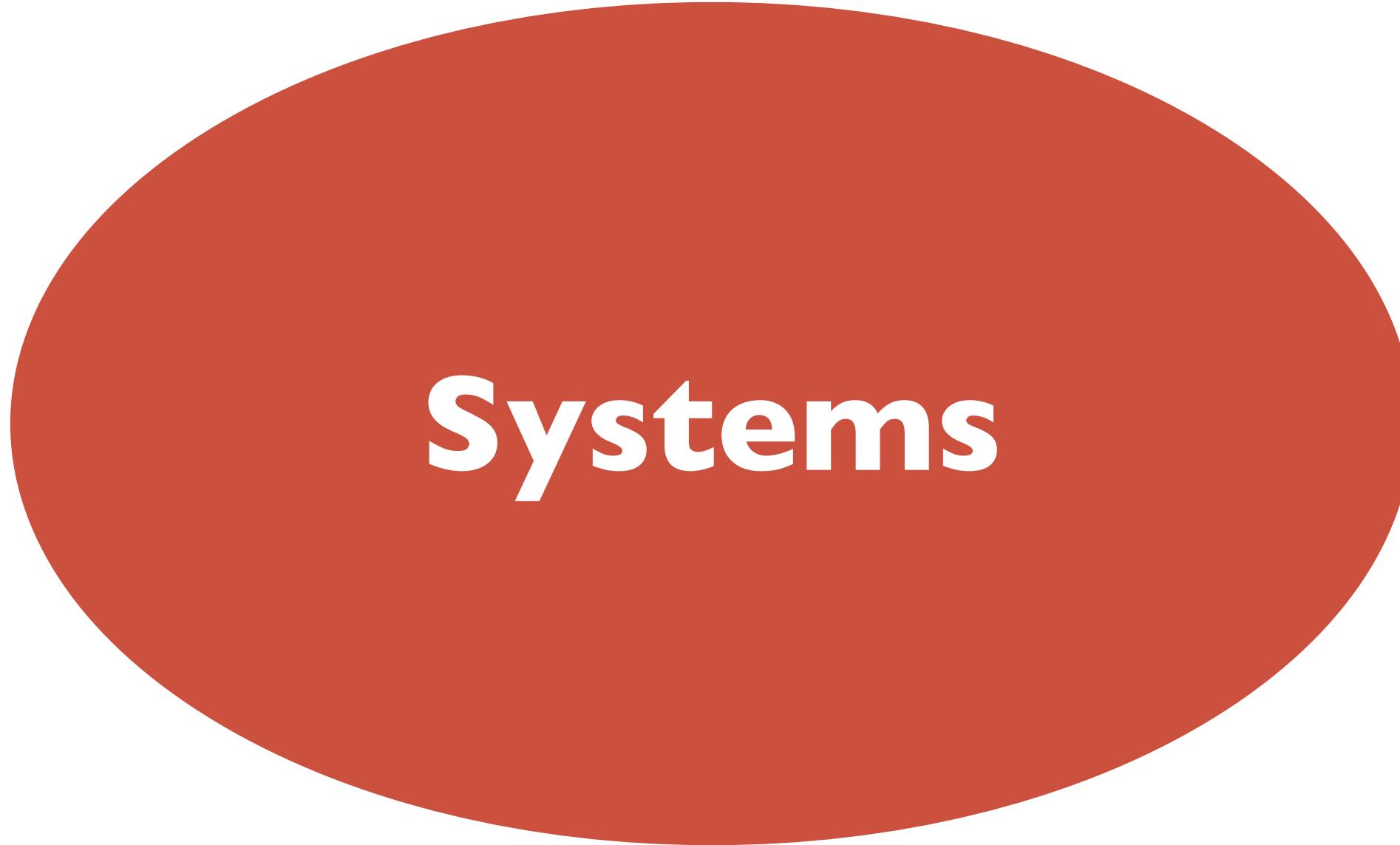
# Original Motivating Problem



# Original Motivating Problem

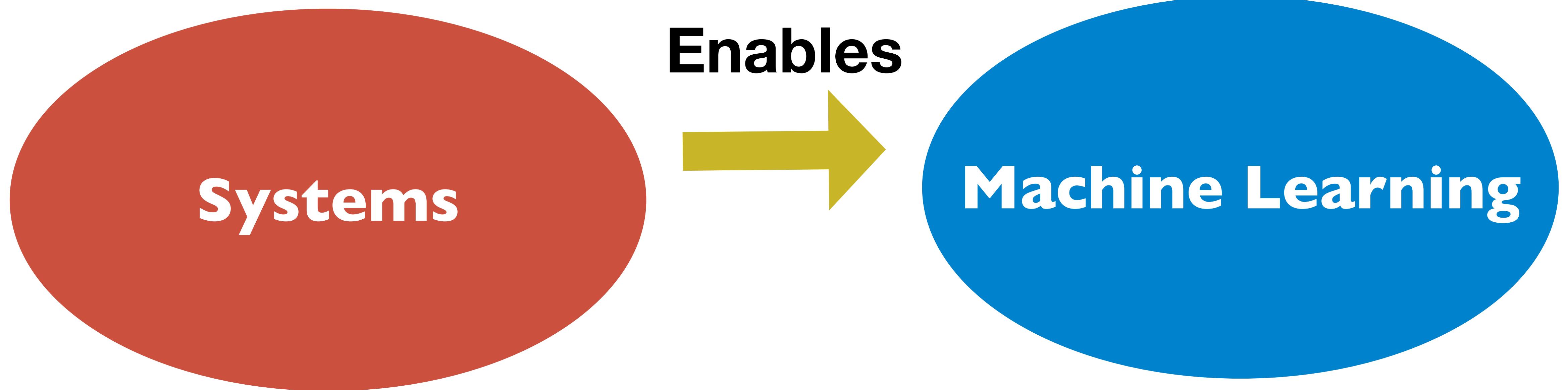


# Current Learning Systems

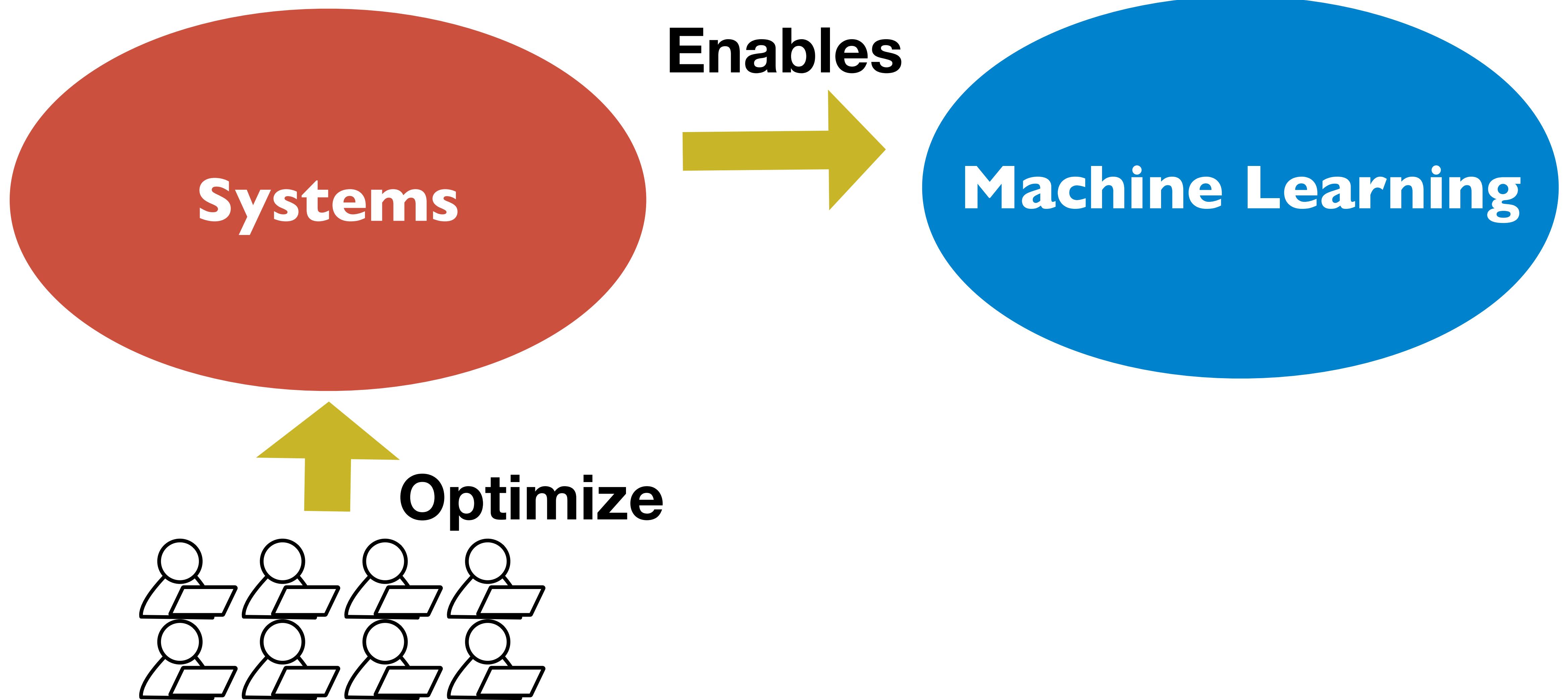


**Systems**

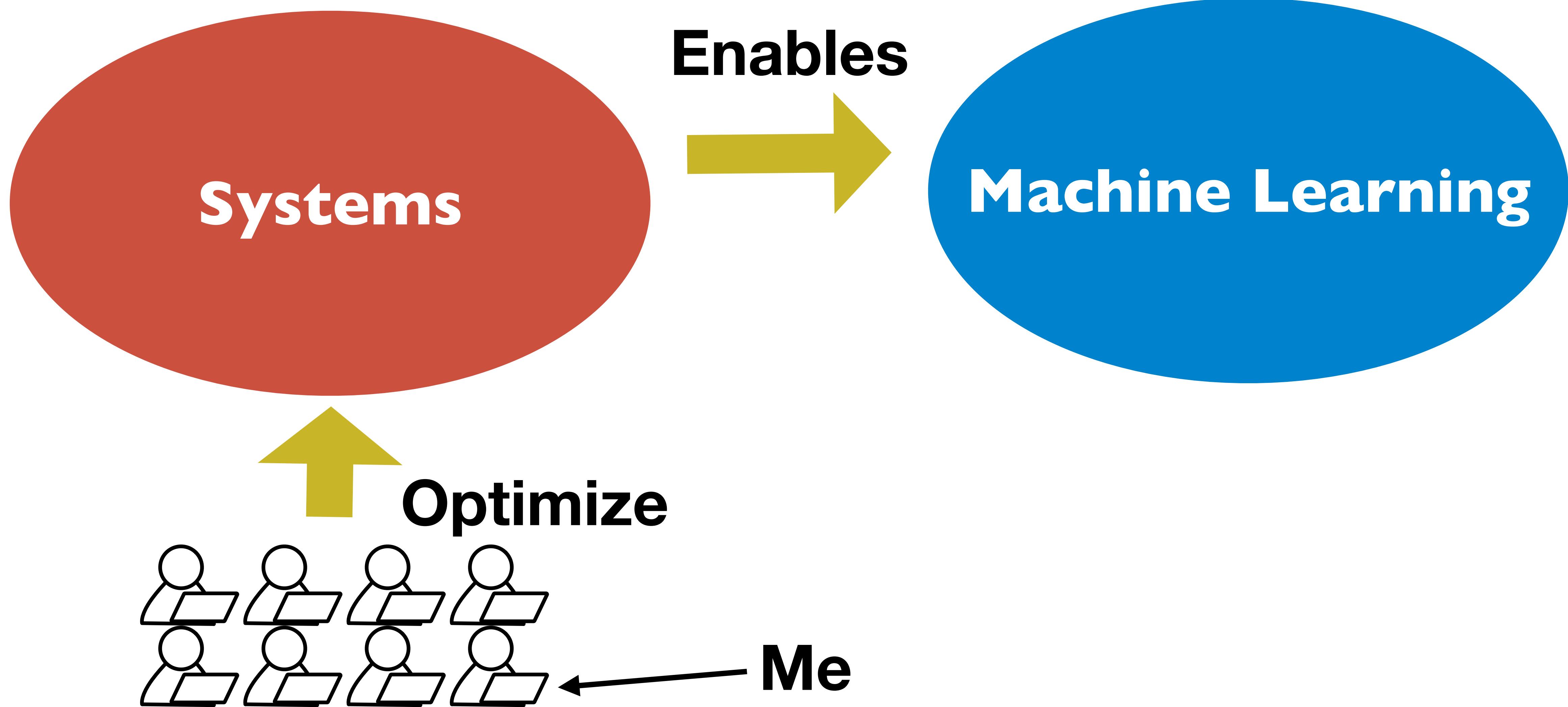
# Current Learning Systems



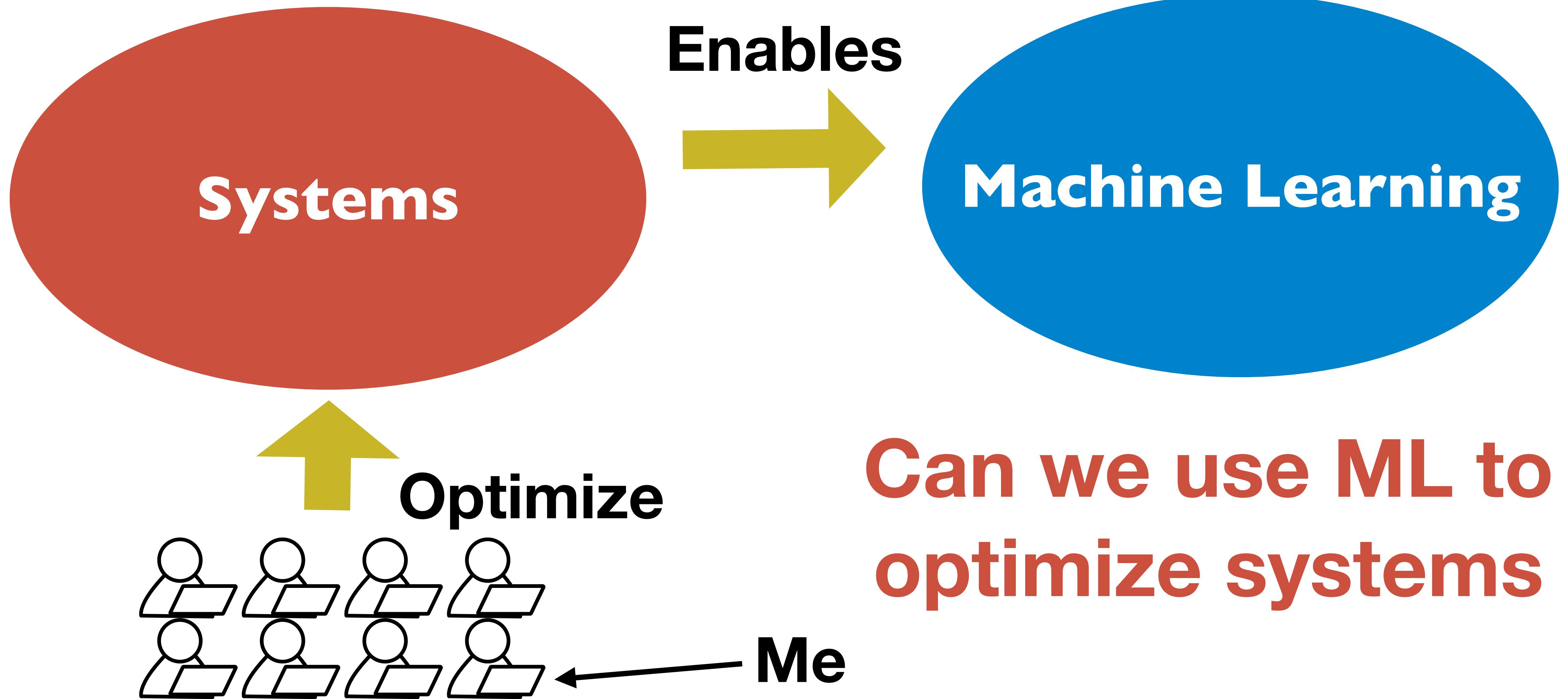
# Current Learning Systems



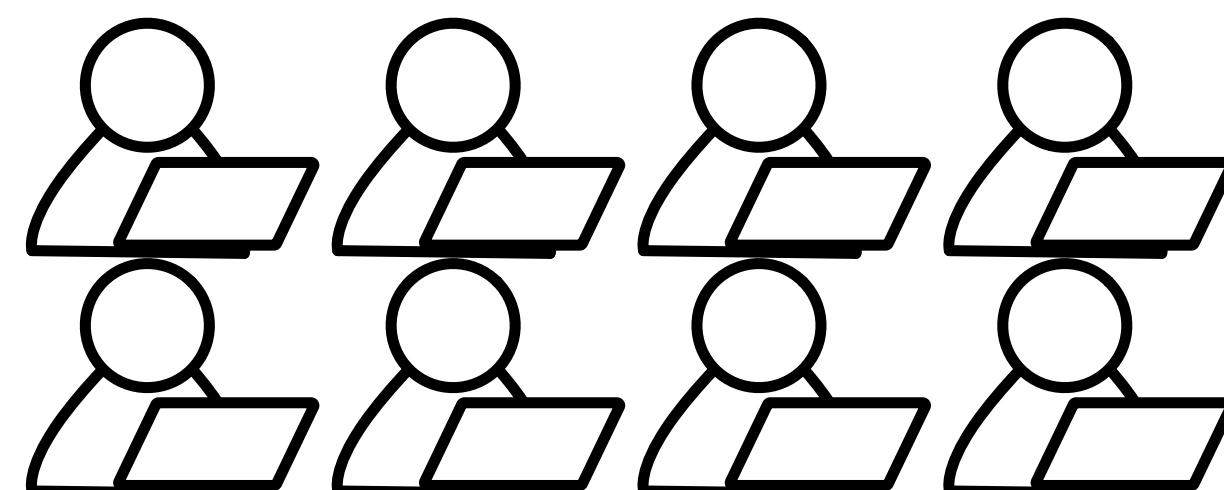
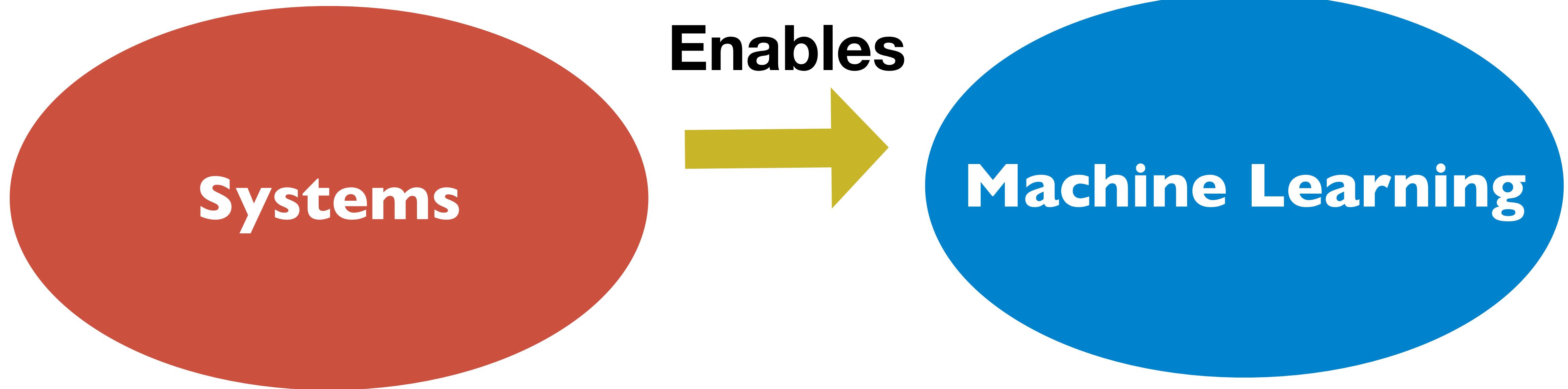
# Current Learning Systems



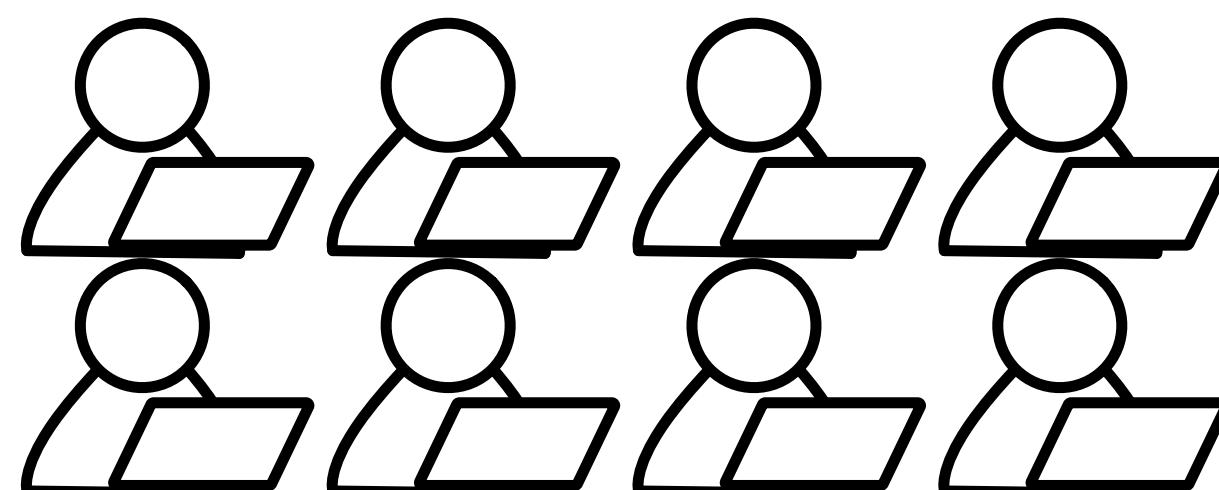
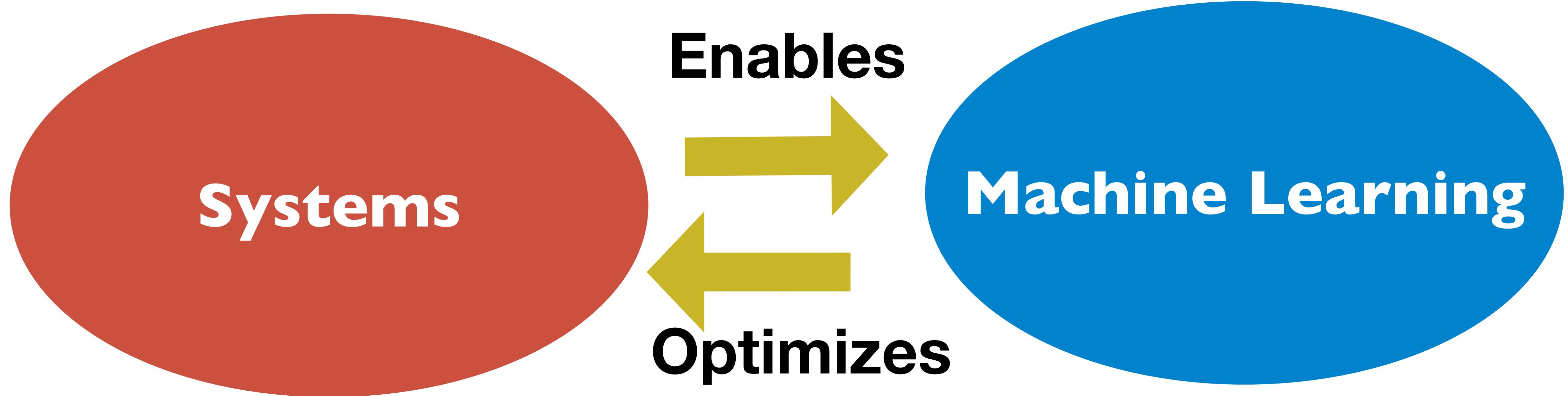
# Current Learning Systems



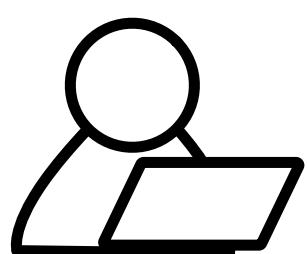
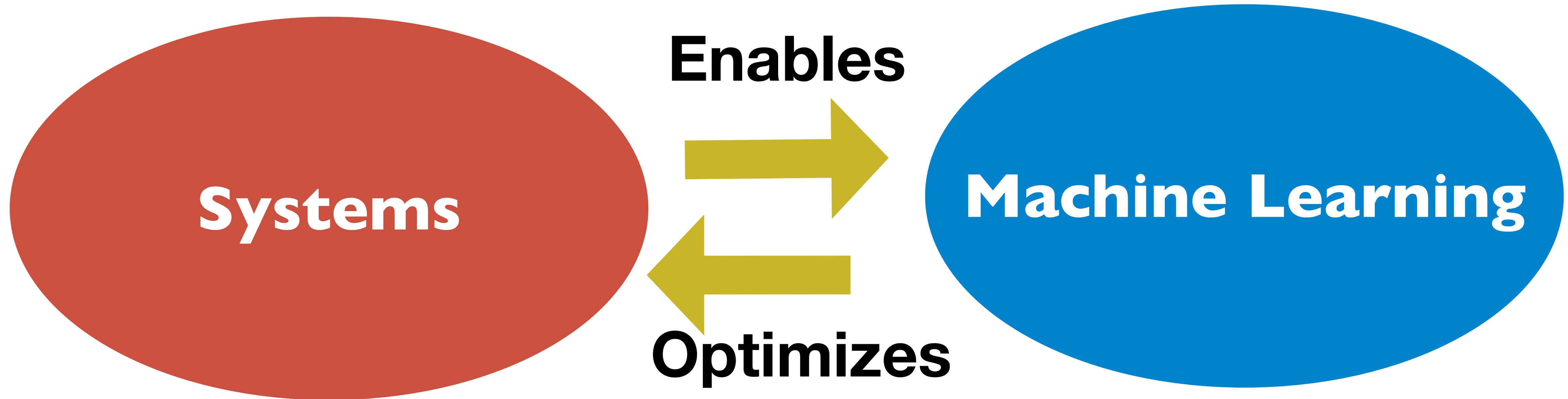
# Learning-based Learning Systems



# Learning-based Learning Systems



# Learning-based Learning Systems



# TVM: Learning-based Learning System

Why do we need machine learning for systems

How to build intelligent systems with learning

End to end learning-based learning system stack

# TVM: Learning-based Learning System

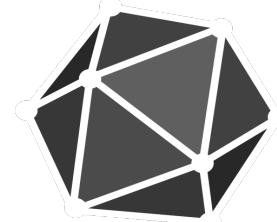
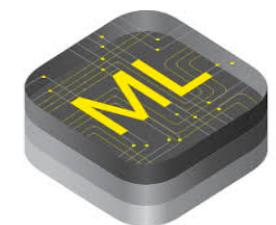
**Why do we need machine learning for systems**

How to build intelligent systems with learning

End to end learning-based learning system stack

# Deploy Deep Learning Everywhere

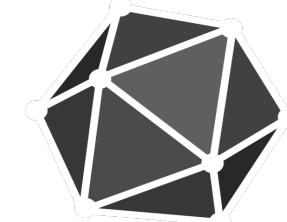
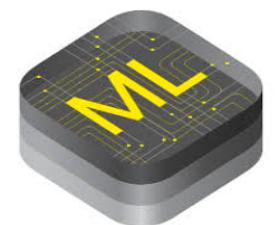
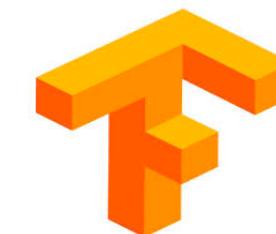
Frameworks



Hardware

# Deploy Deep Learning Everywhere

Frameworks

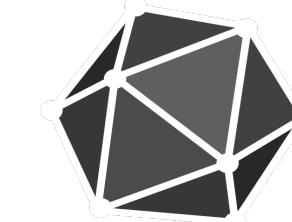
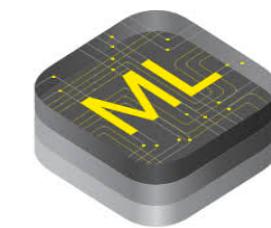


Hardware



# Deploy Deep Learning Everywhere

Frameworks

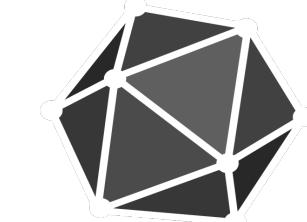
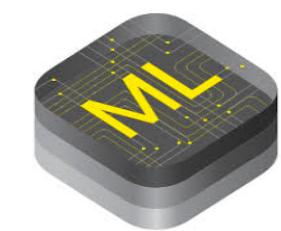


Hardware



# Deploy Deep Learning Everywhere

Frameworks

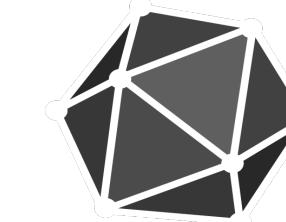
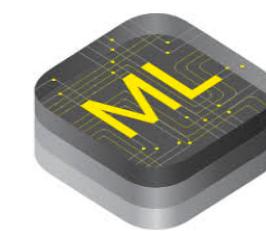


Hardware



# Deploy Deep Learning Everywhere

Frameworks

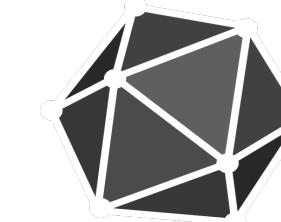
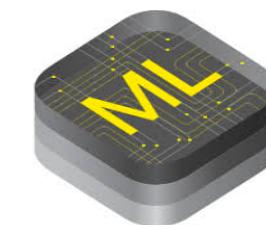


Hardware

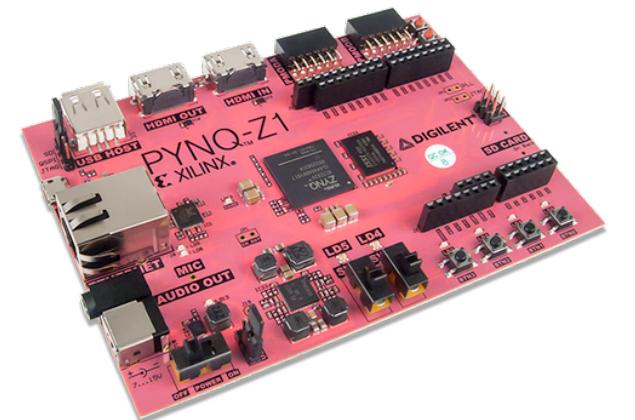


# Deploy Deep Learning Everywhere

Frameworks

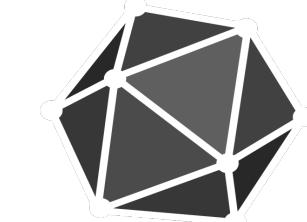
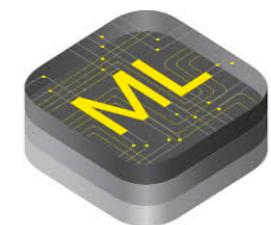


Hardware

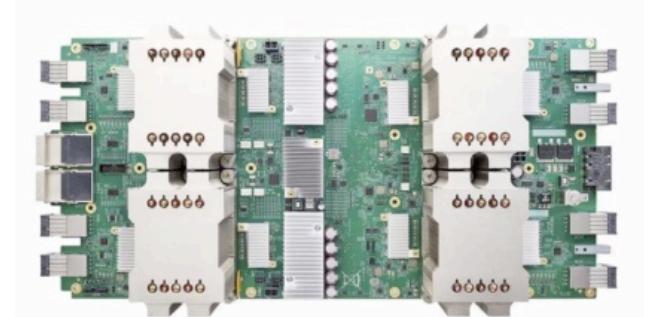
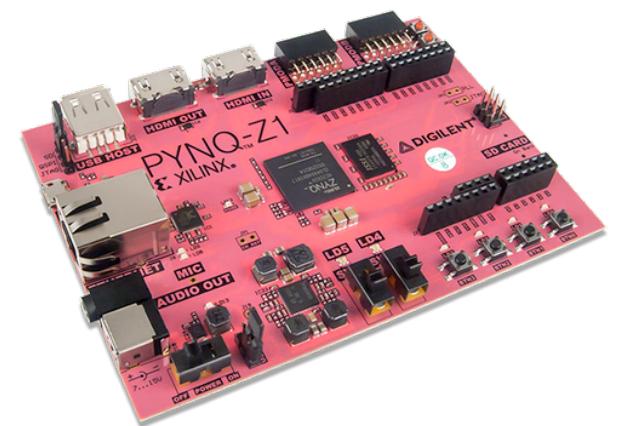


# Deploy Deep Learning Everywhere

Frameworks

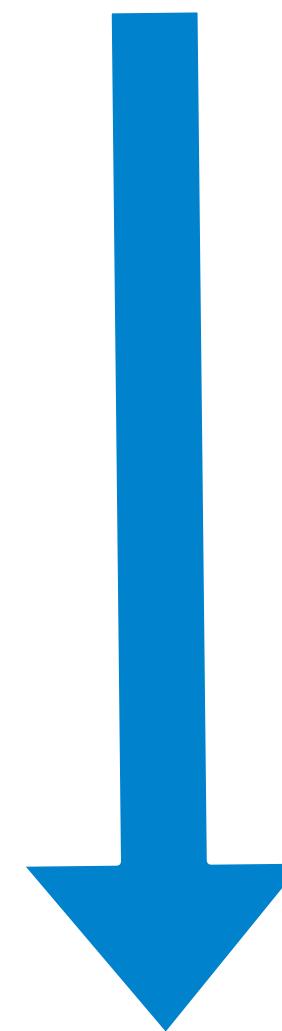
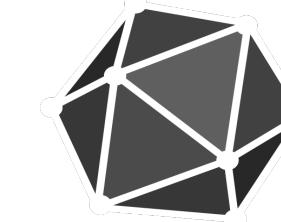
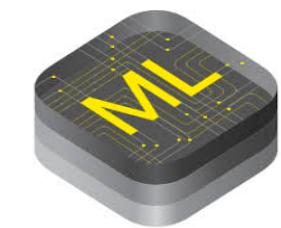


Hardware



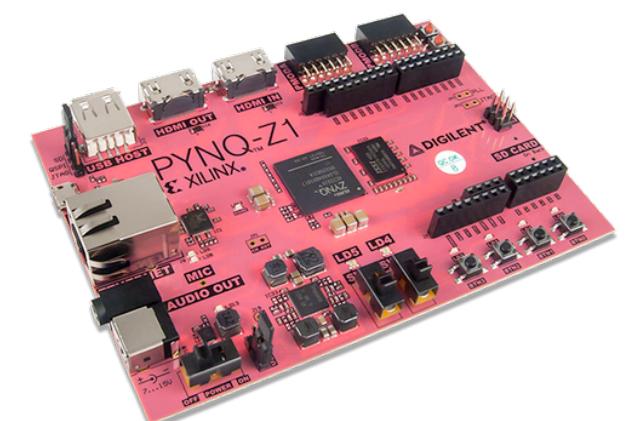
# Deploy Deep Learning Everywhere

Frameworks

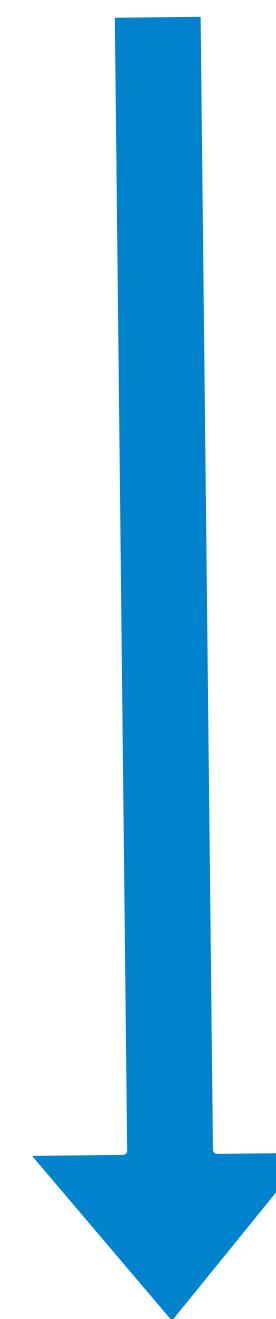
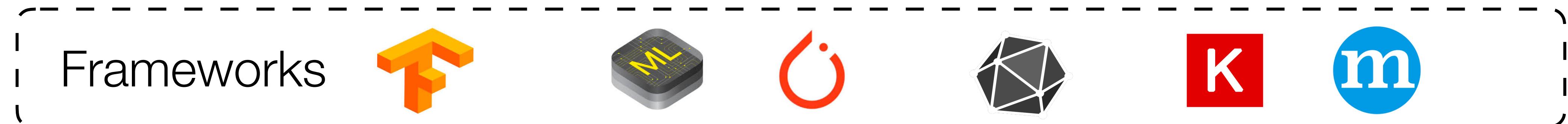


Huge gap between model/frameworks and hardware backends

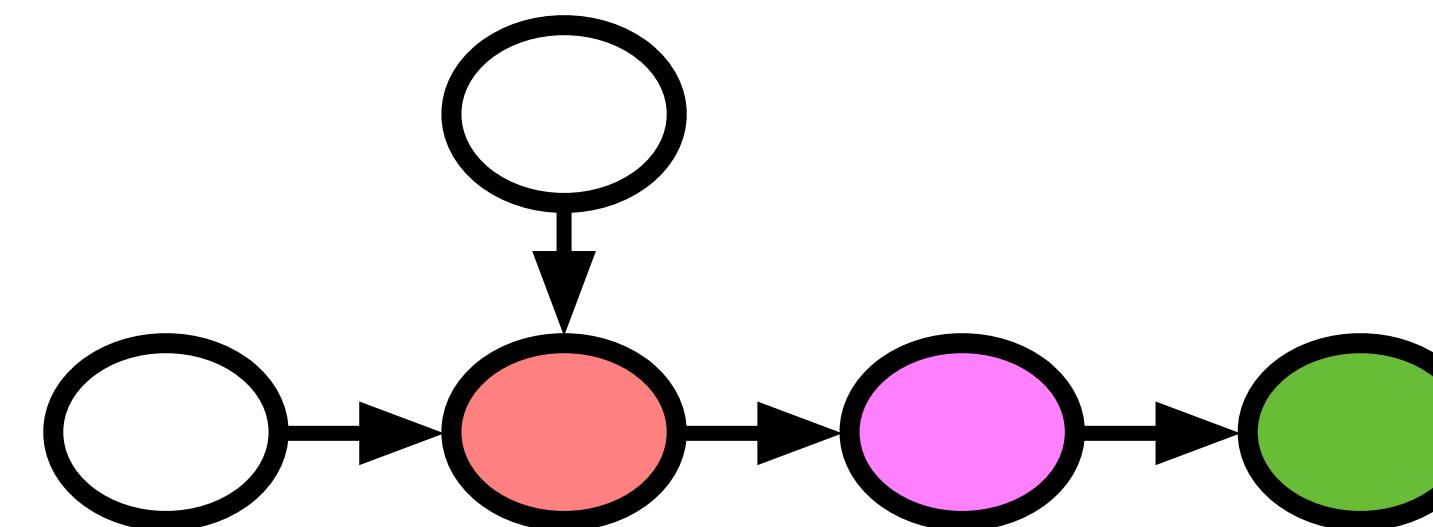
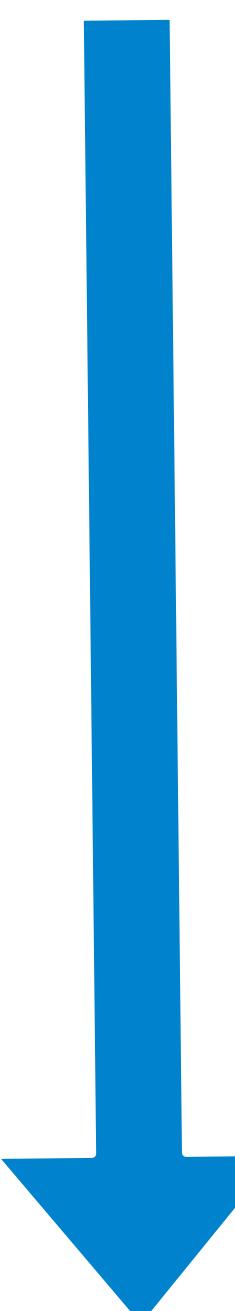
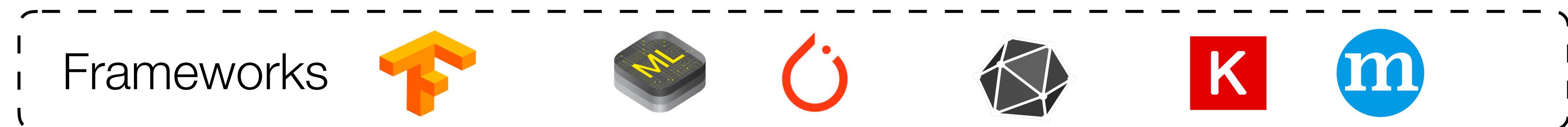
Hardware



# Existing Deep Learning Frameworks



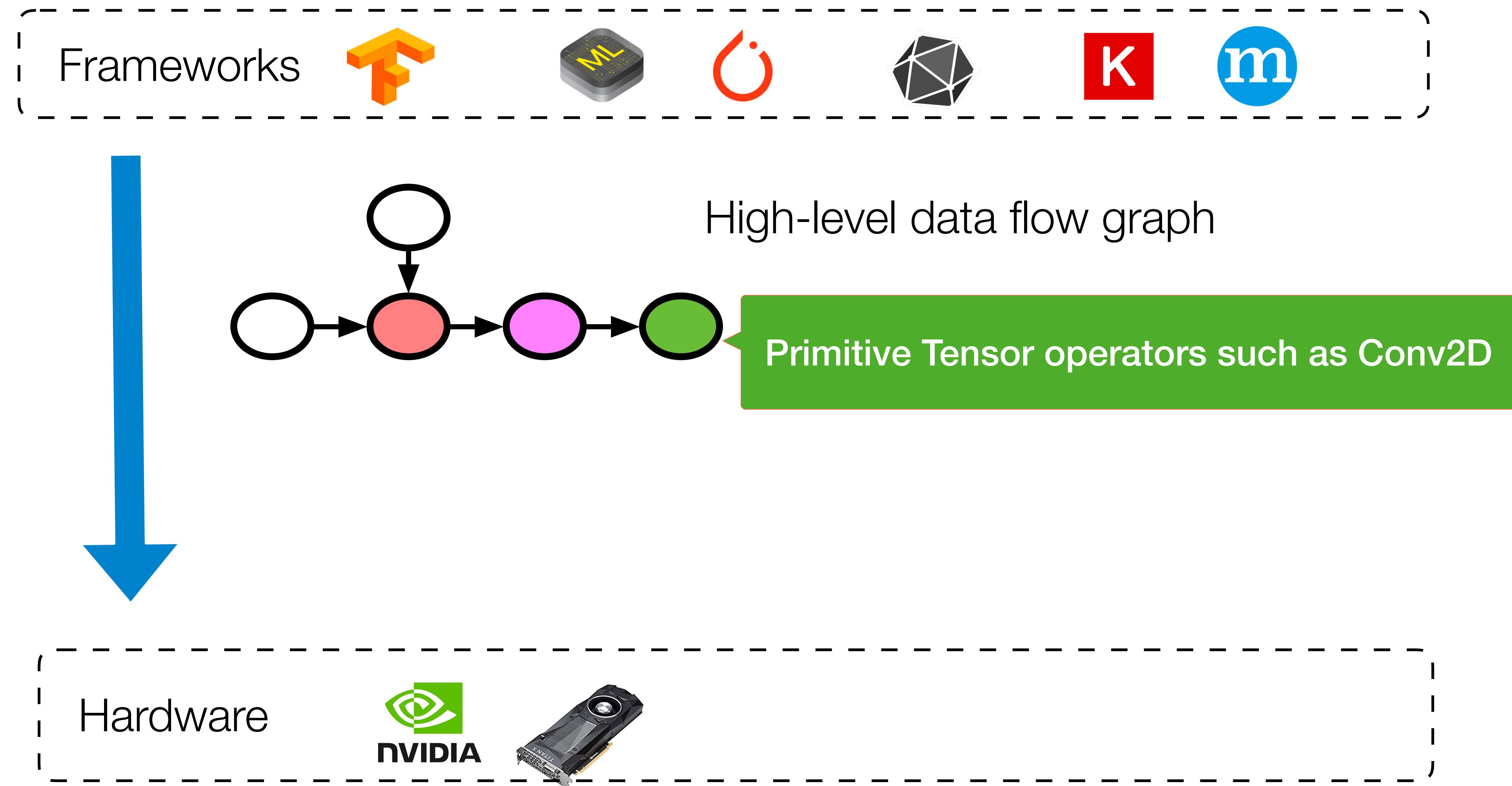
# Existing Deep Learning Frameworks



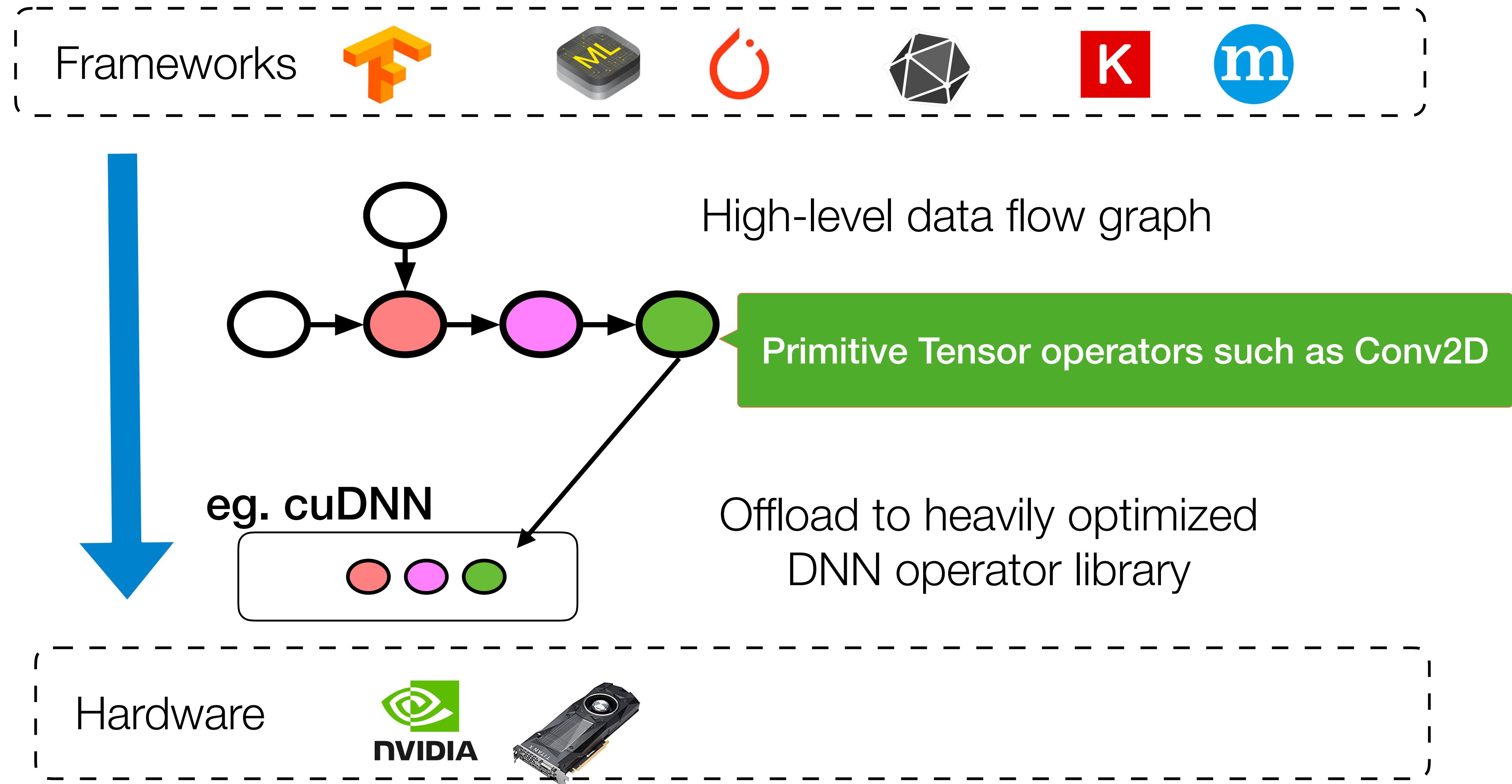
High-level data flow graph



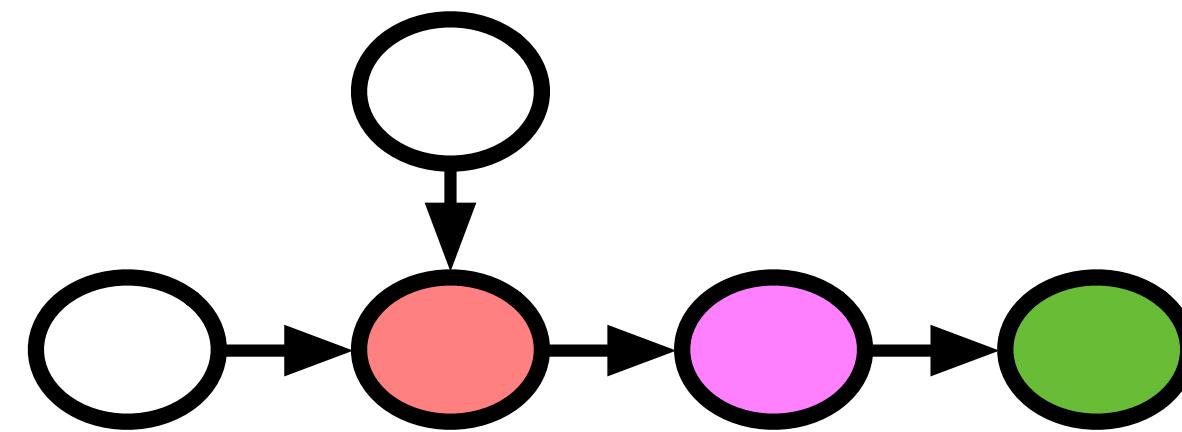
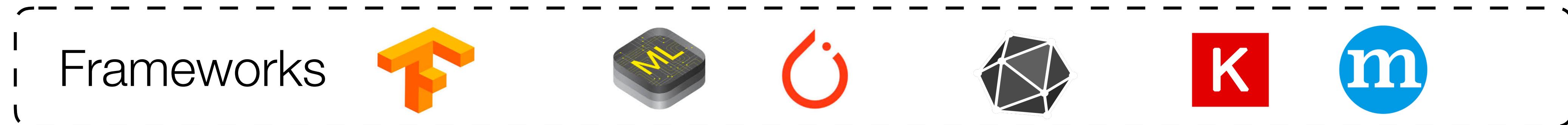
# Existing Deep Learning Frameworks



# Existing Deep Learning Frameworks



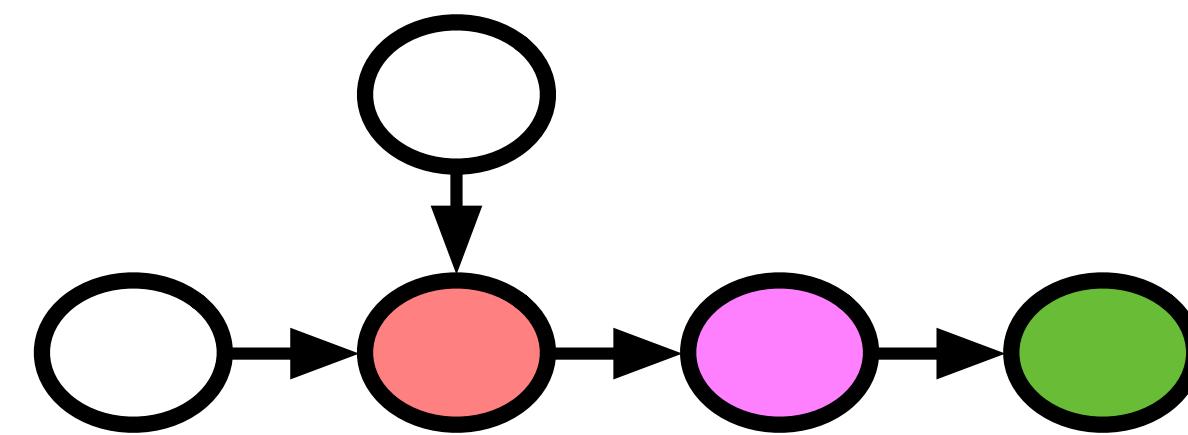
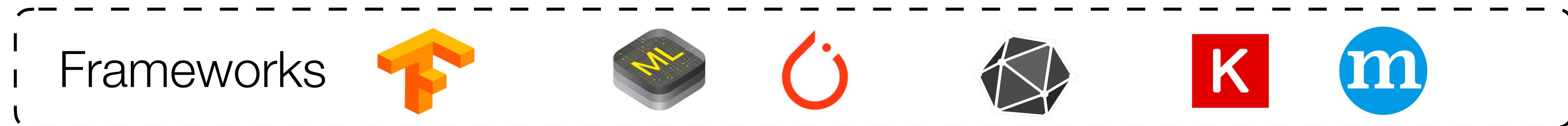
# Limitations of Existing Approach



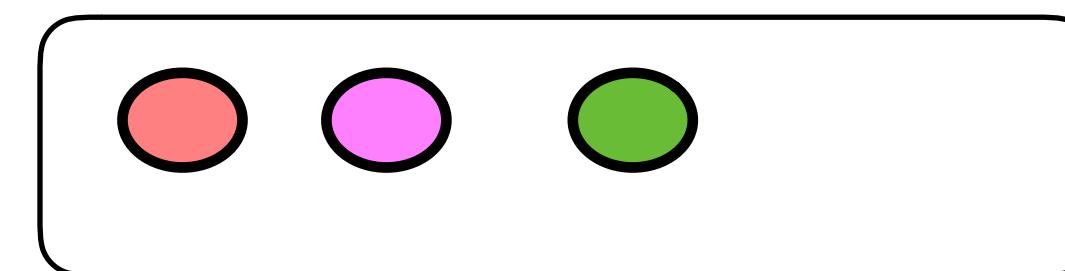
**cuDNN**



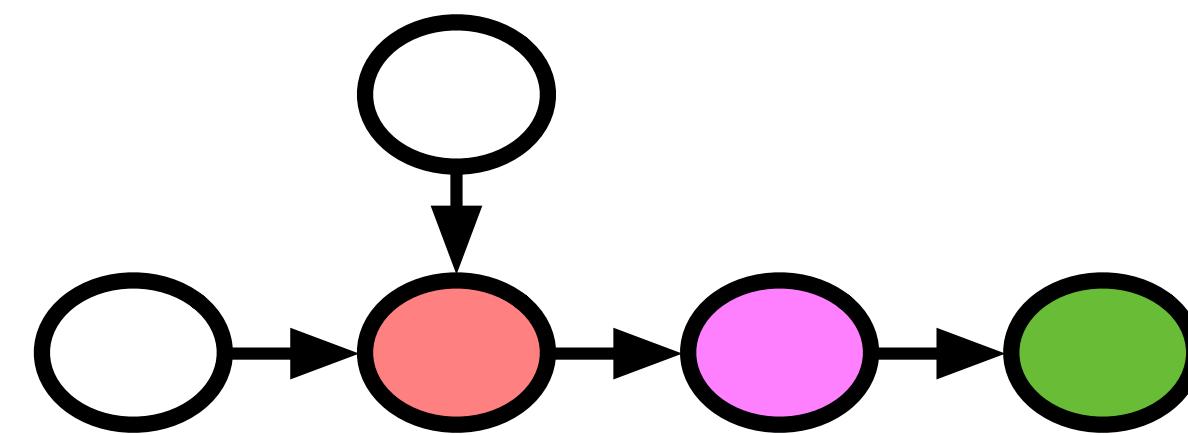
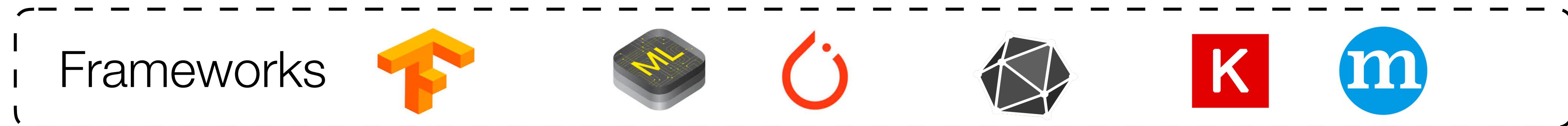
# Limitations of Existing Approach



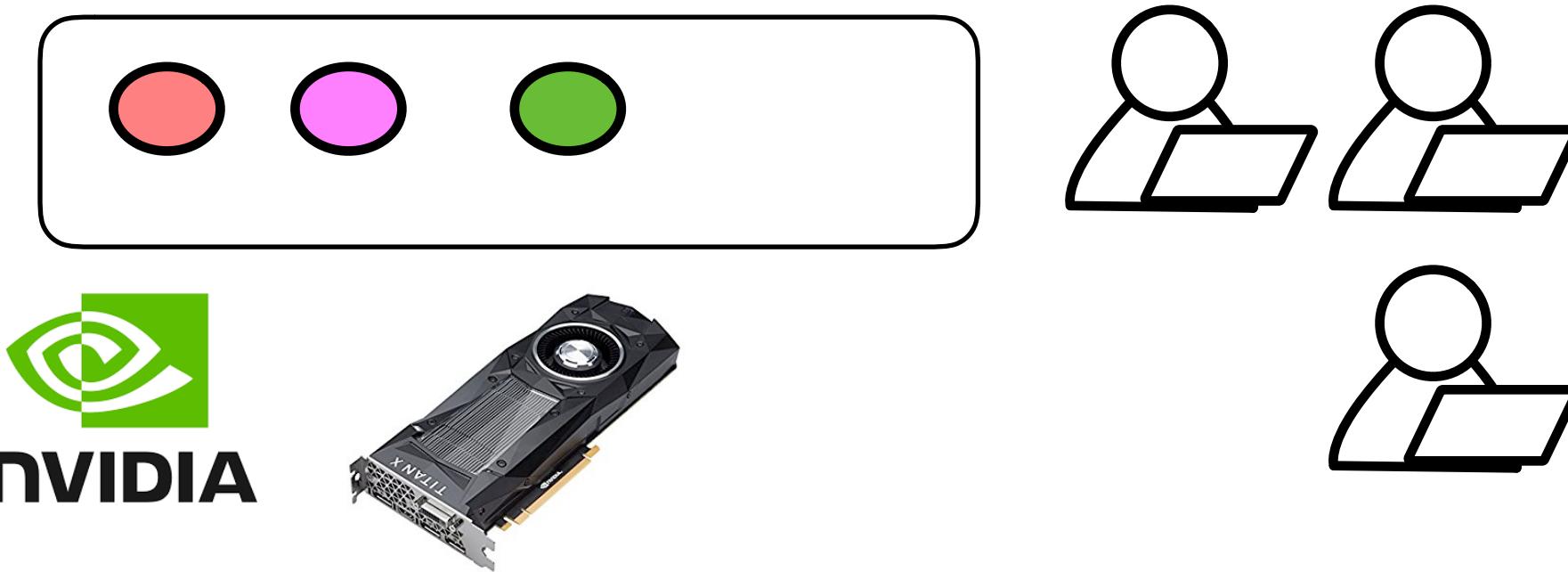
**cuDNN**



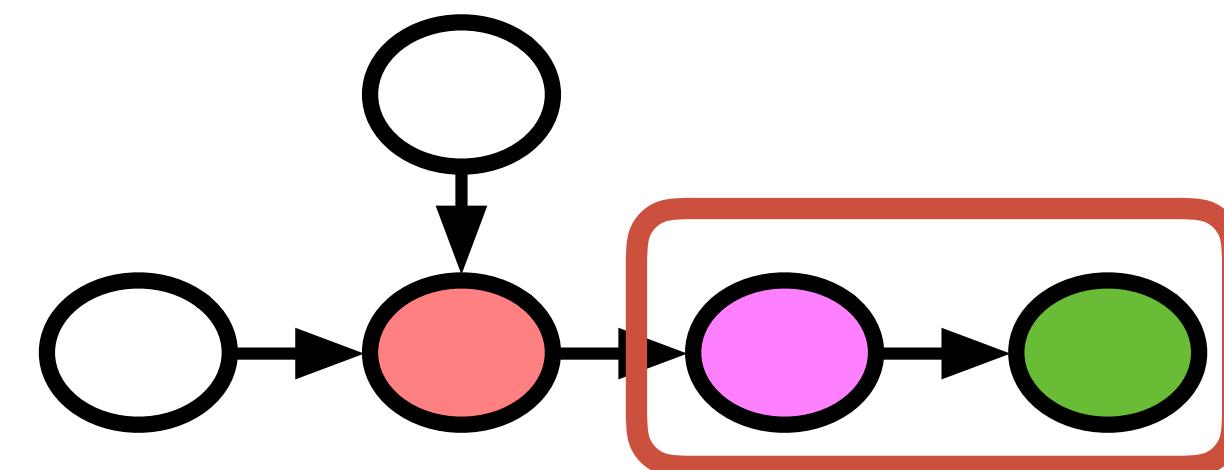
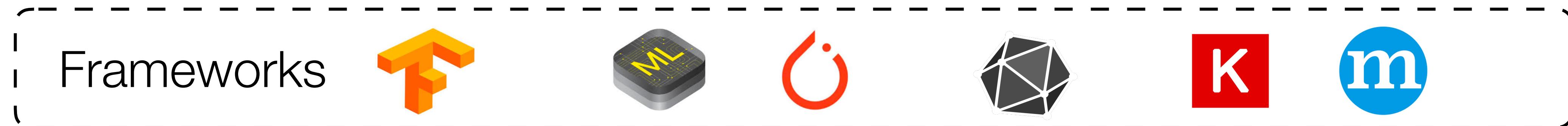
# Limitations of Existing Approach



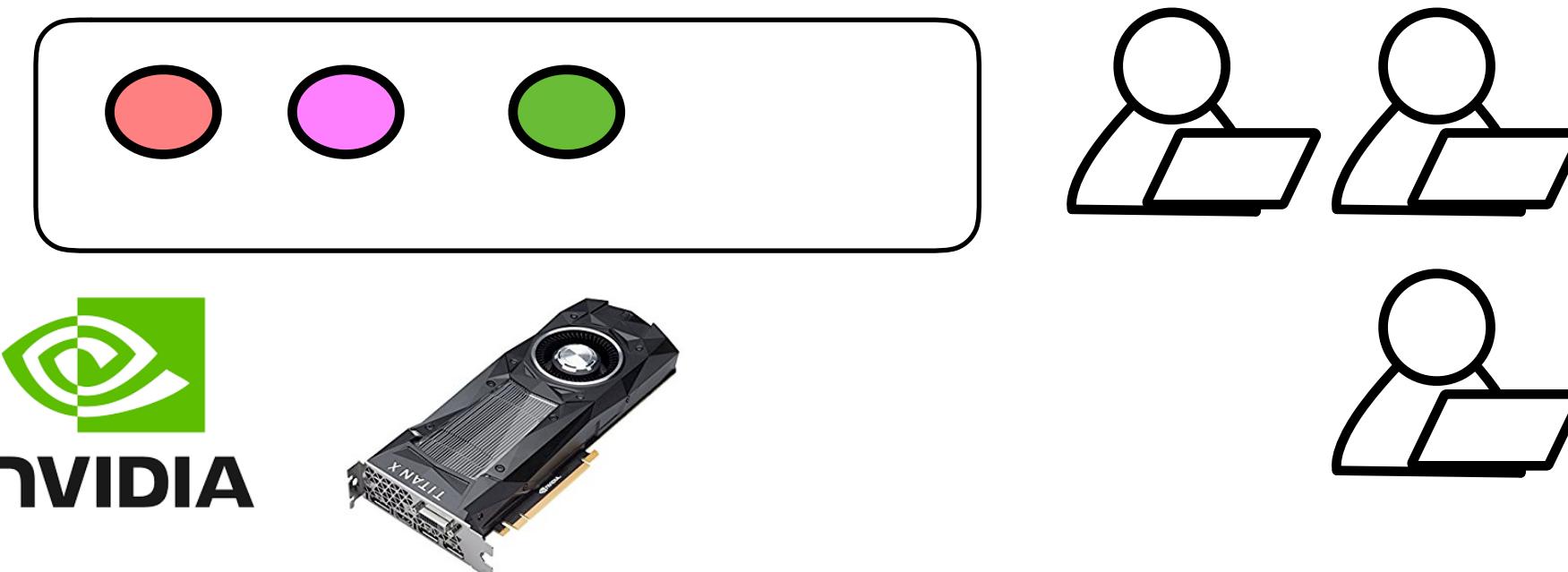
**cuDNN**



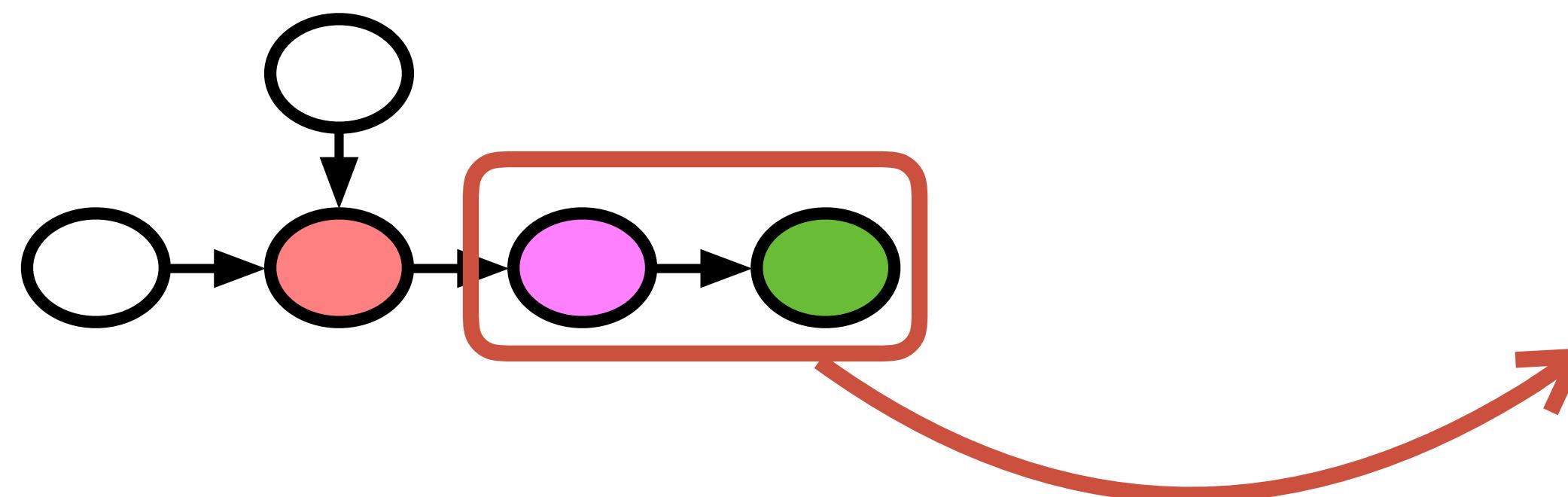
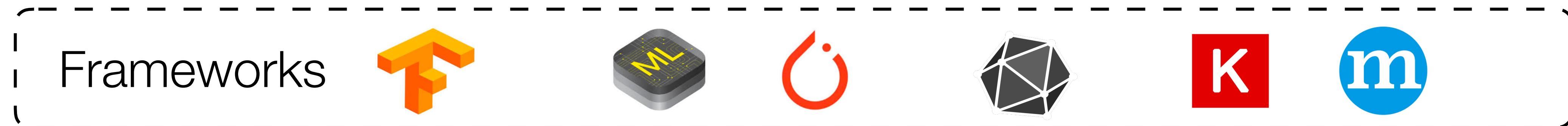
# Limitations of Existing Approach



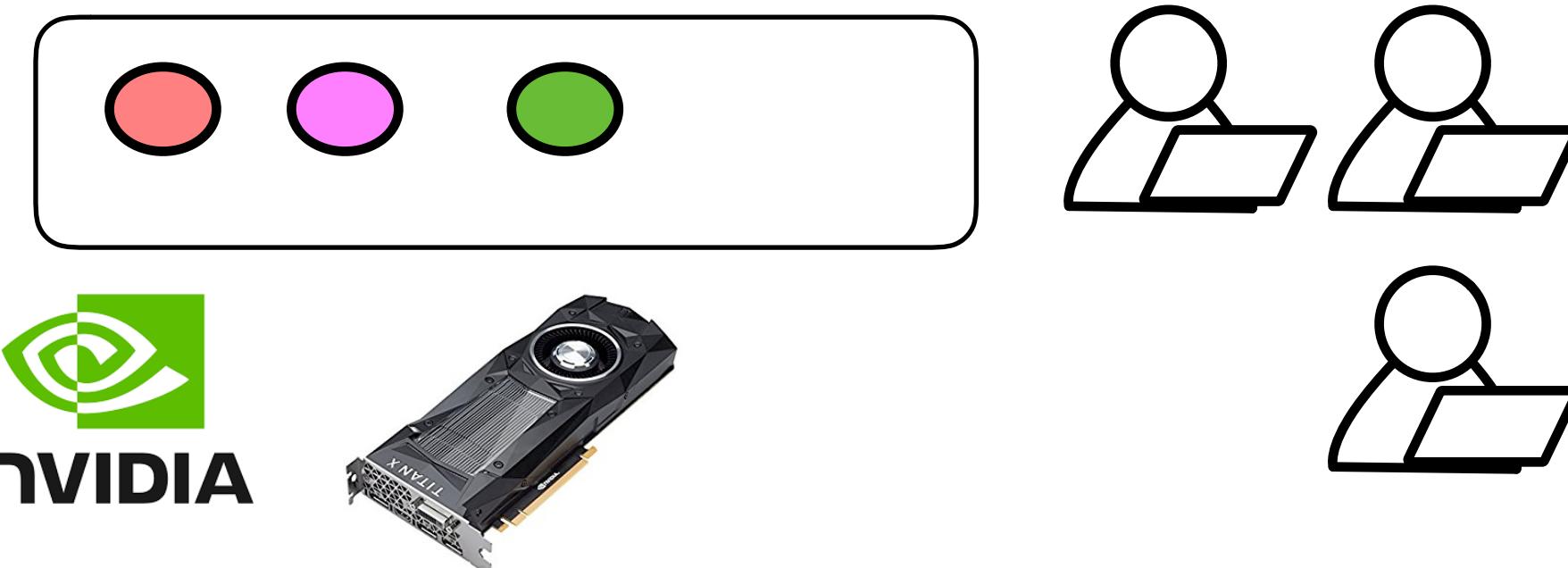
cuDNN



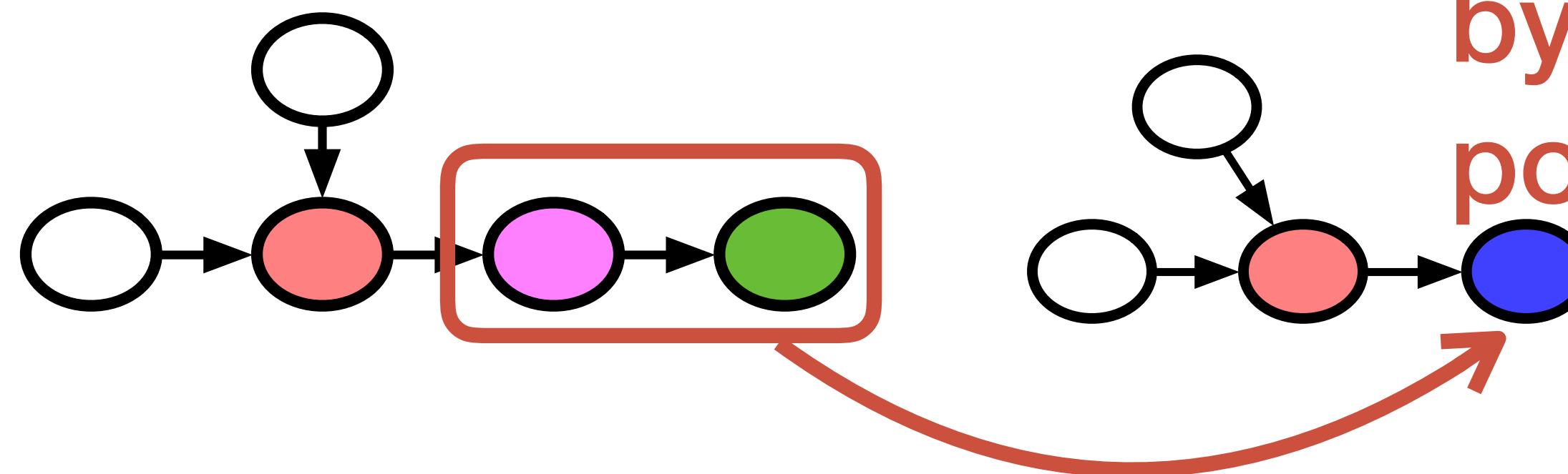
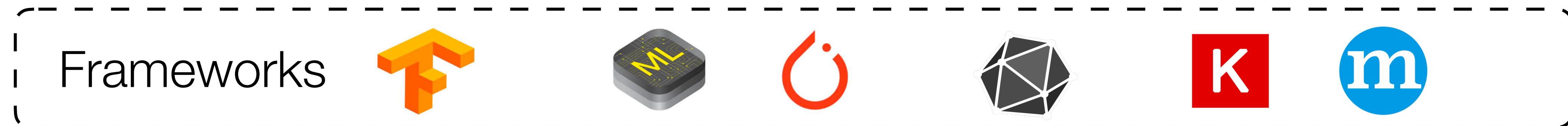
# Limitations of Existing Approach



**cuDNN**

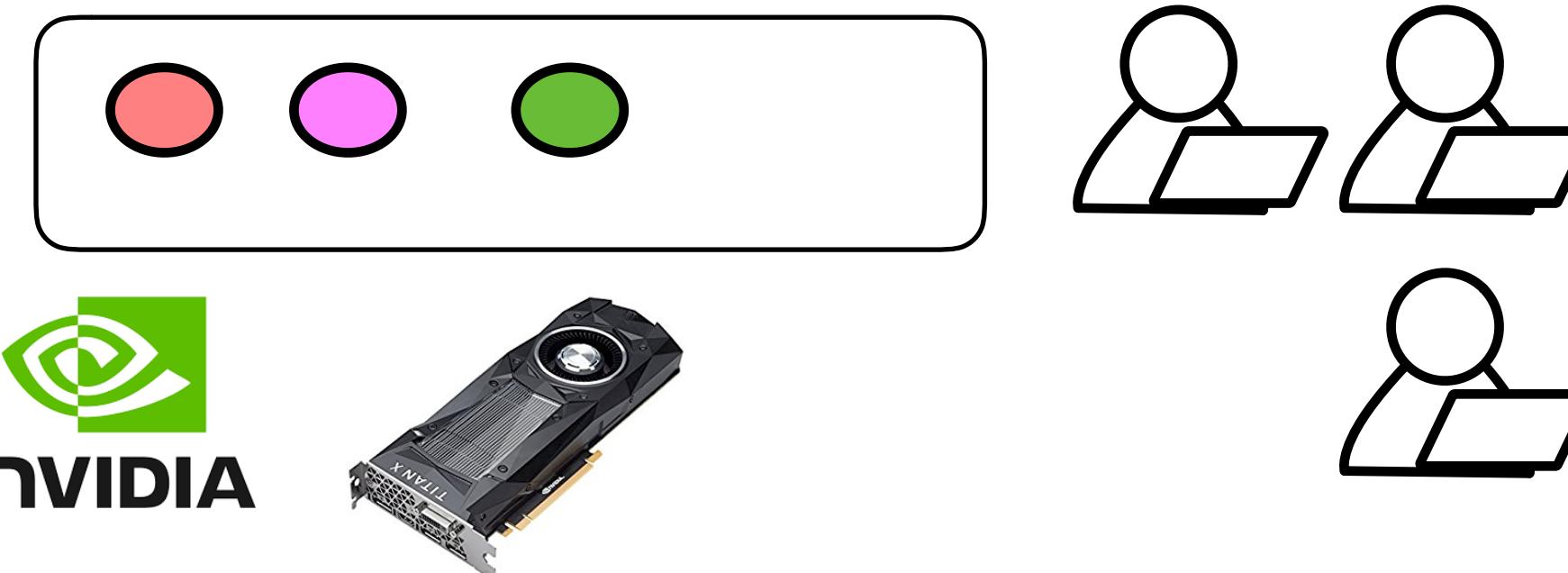


# Limitations of Existing Approach

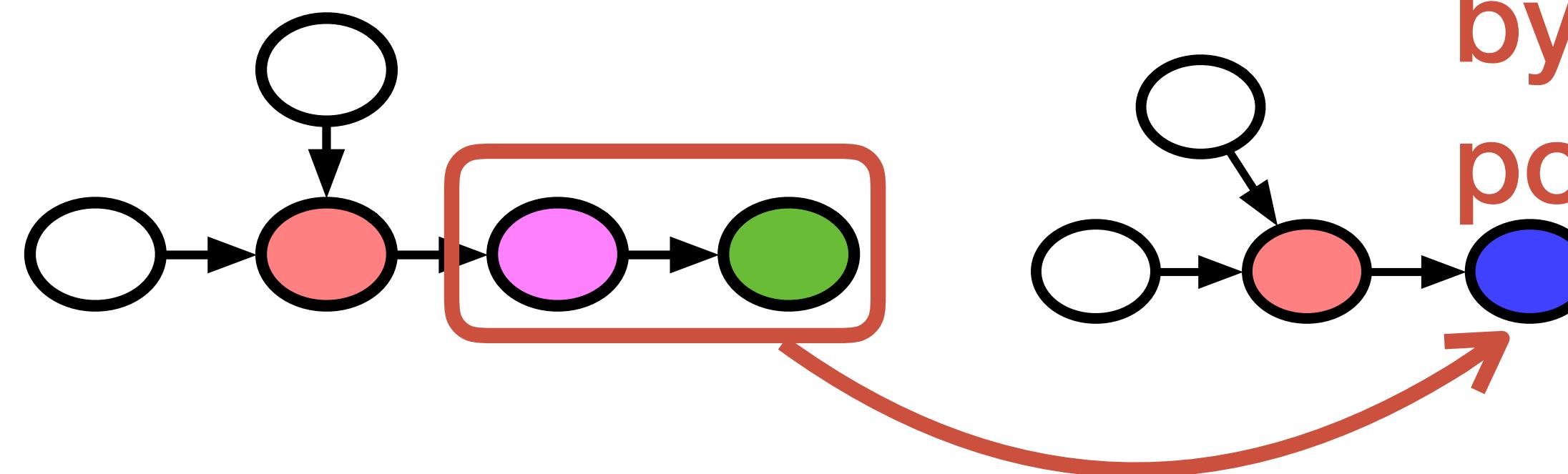
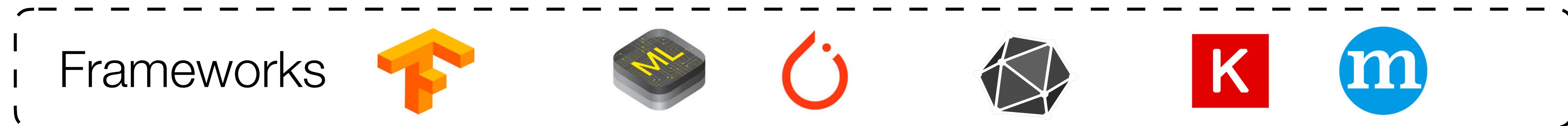


New operator introduced  
by operator fusion optimization  
potential benefit: 1.5x speedup

cuDNN

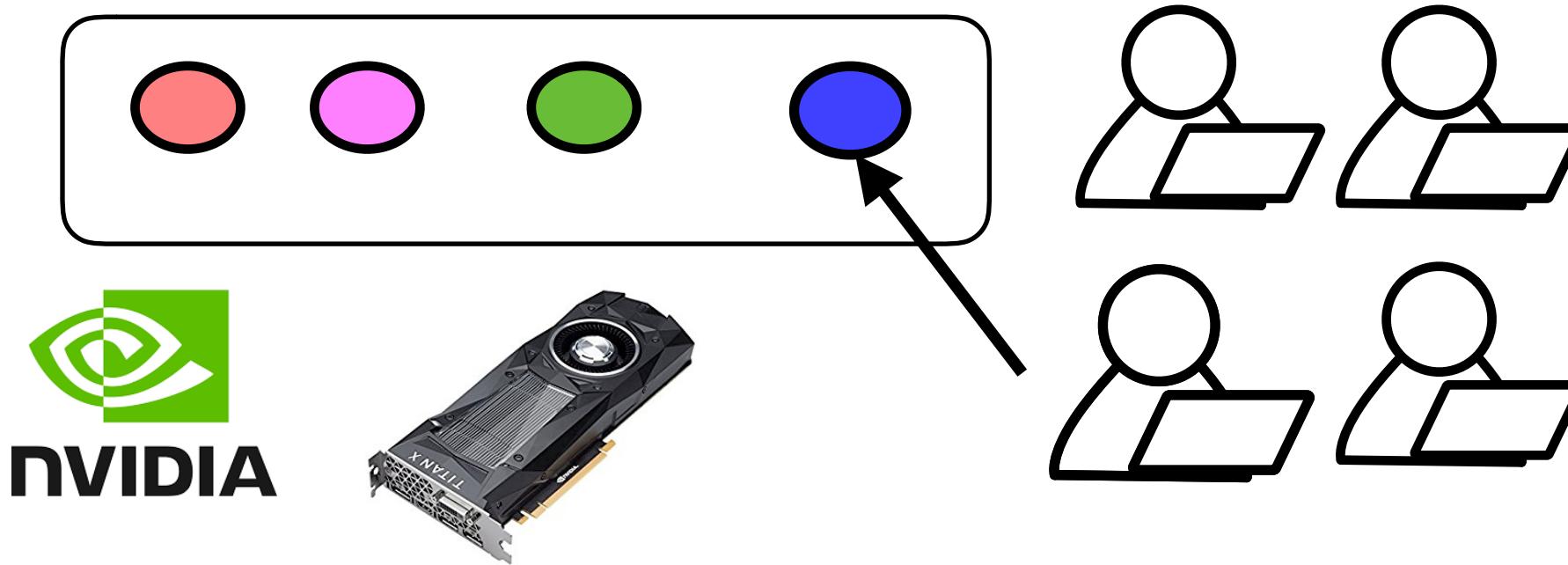


# Limitations of Existing Approach

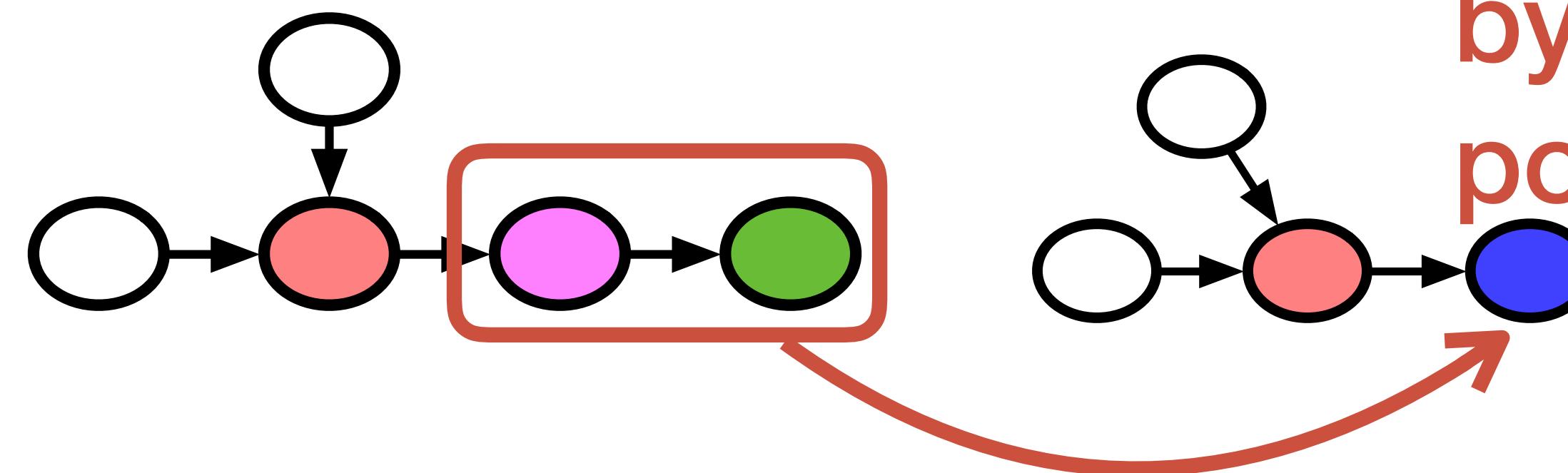
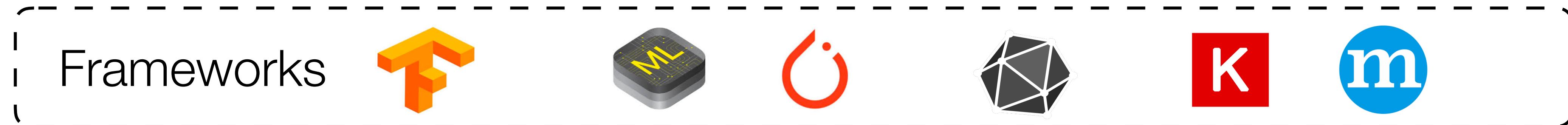


New operator introduced  
by operator fusion optimization  
potential benefit: 1.5x speedup

cuDNN

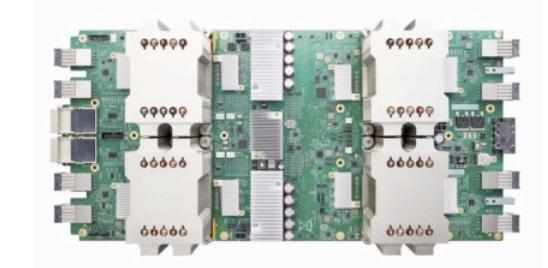
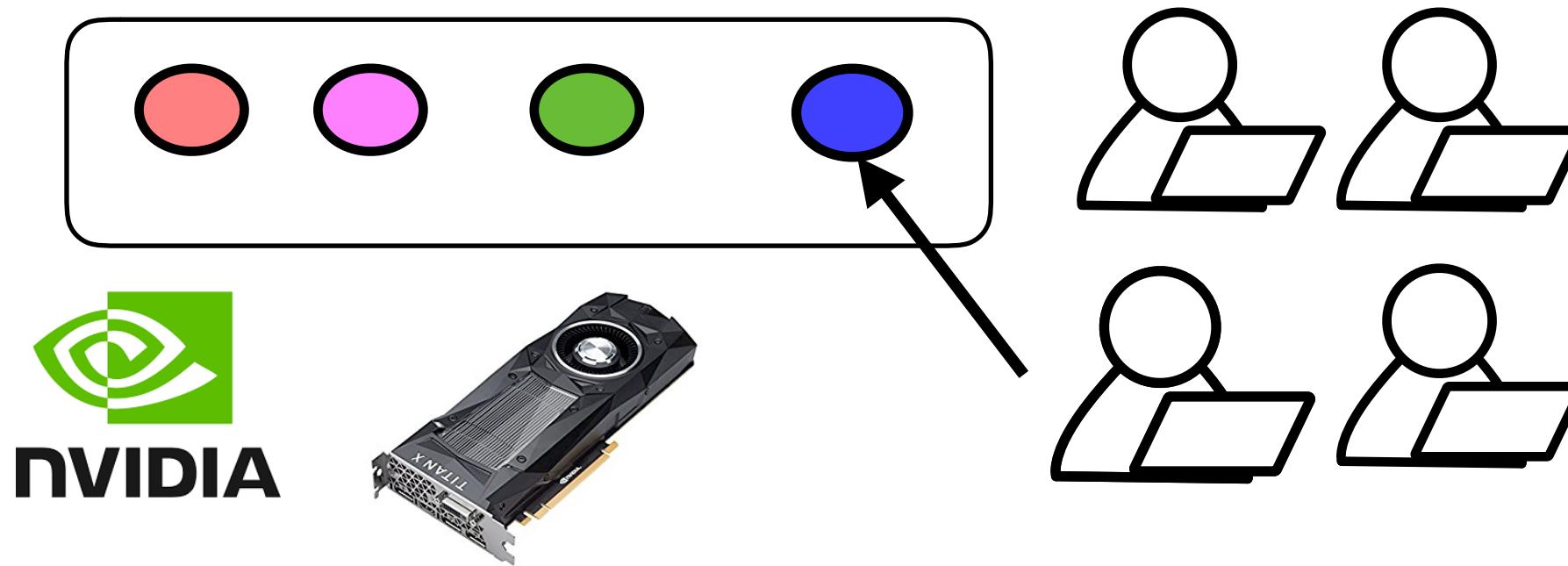


# Limitations of Existing Approach

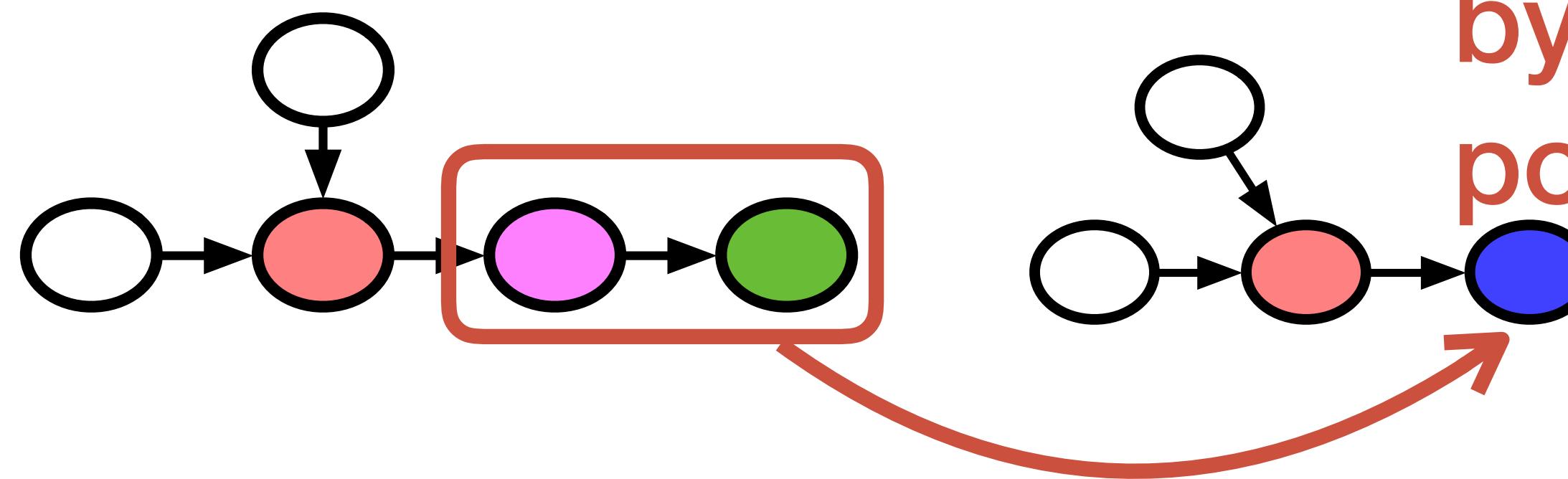
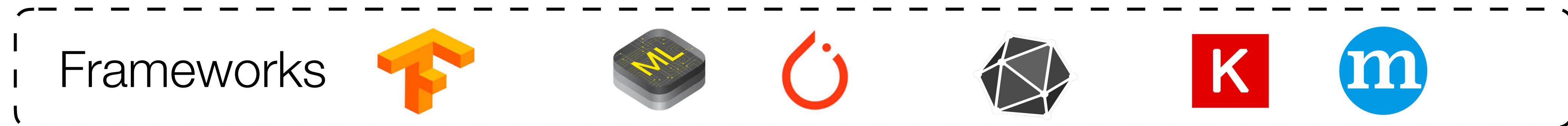


New operator introduced  
by operator fusion optimization  
potential benefit: 1.5x speedup

cuDNN



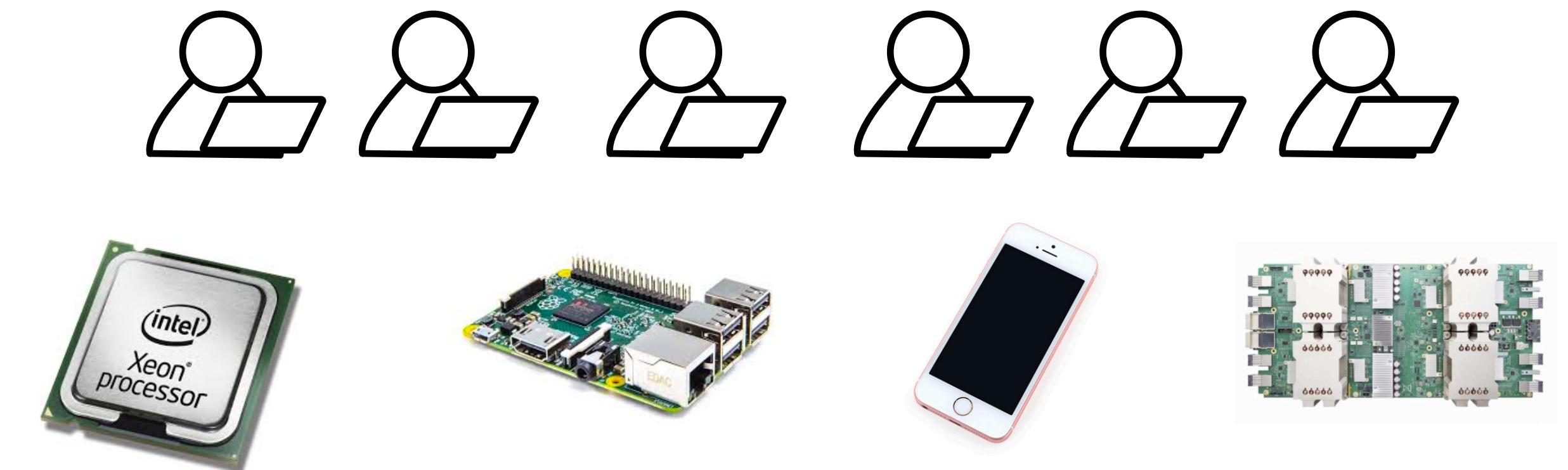
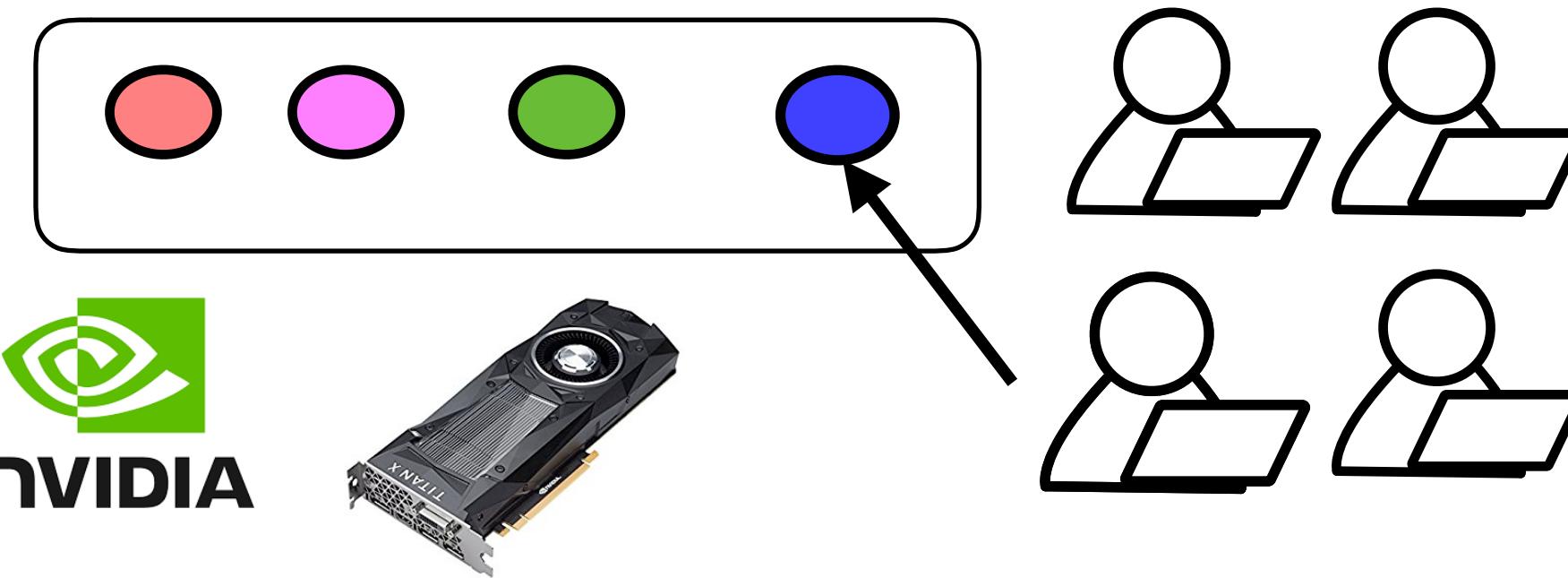
# Limitations of Existing Approach



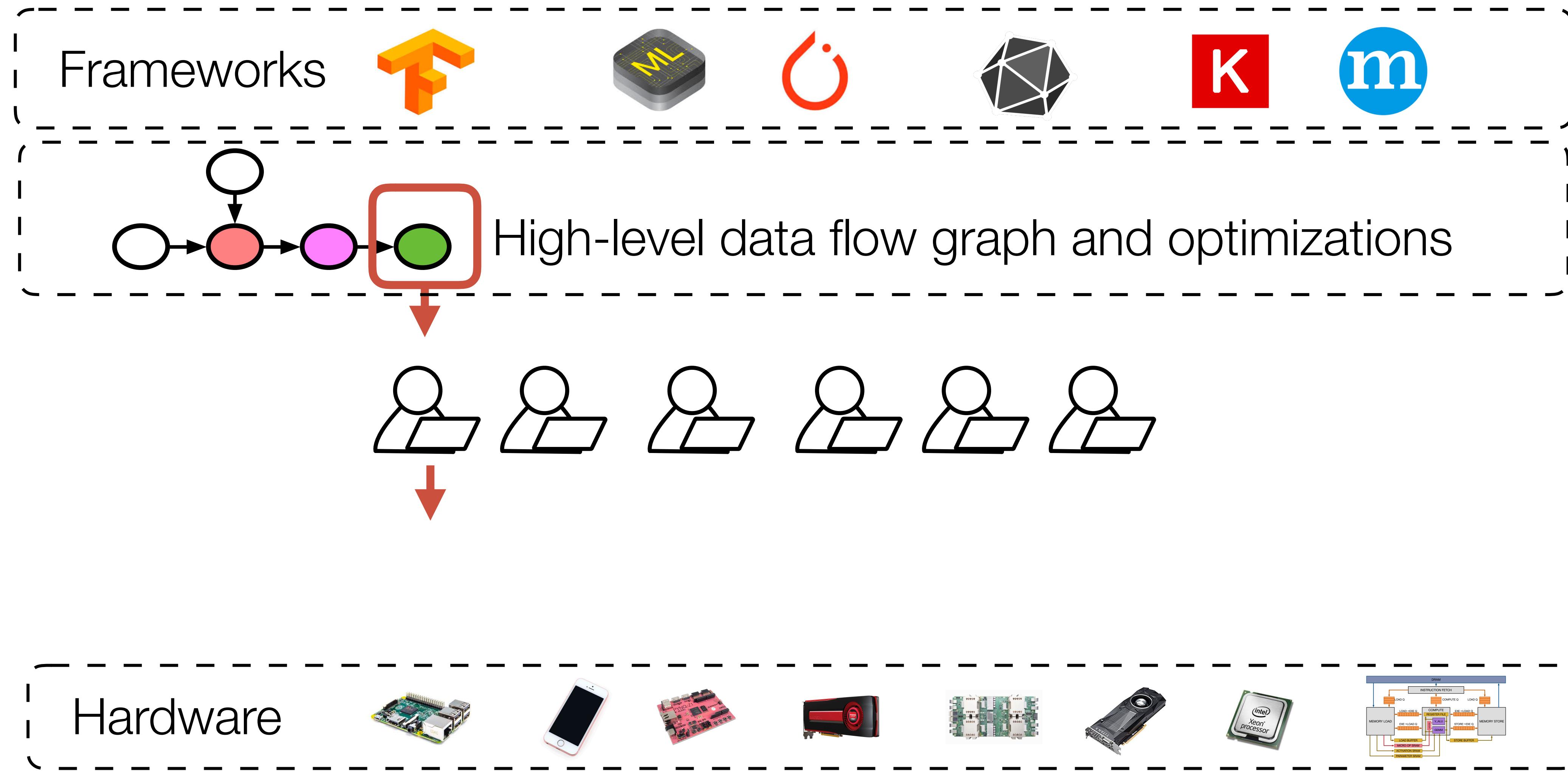
New operator introduced  
by operator fusion optimization  
potential benefit: 1.5x speedup

Engineering intensive

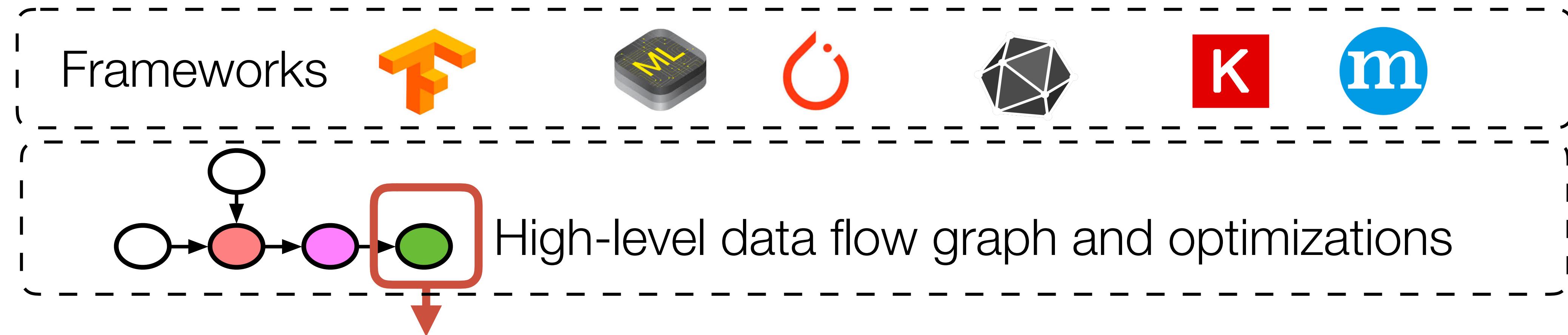
cuDNN



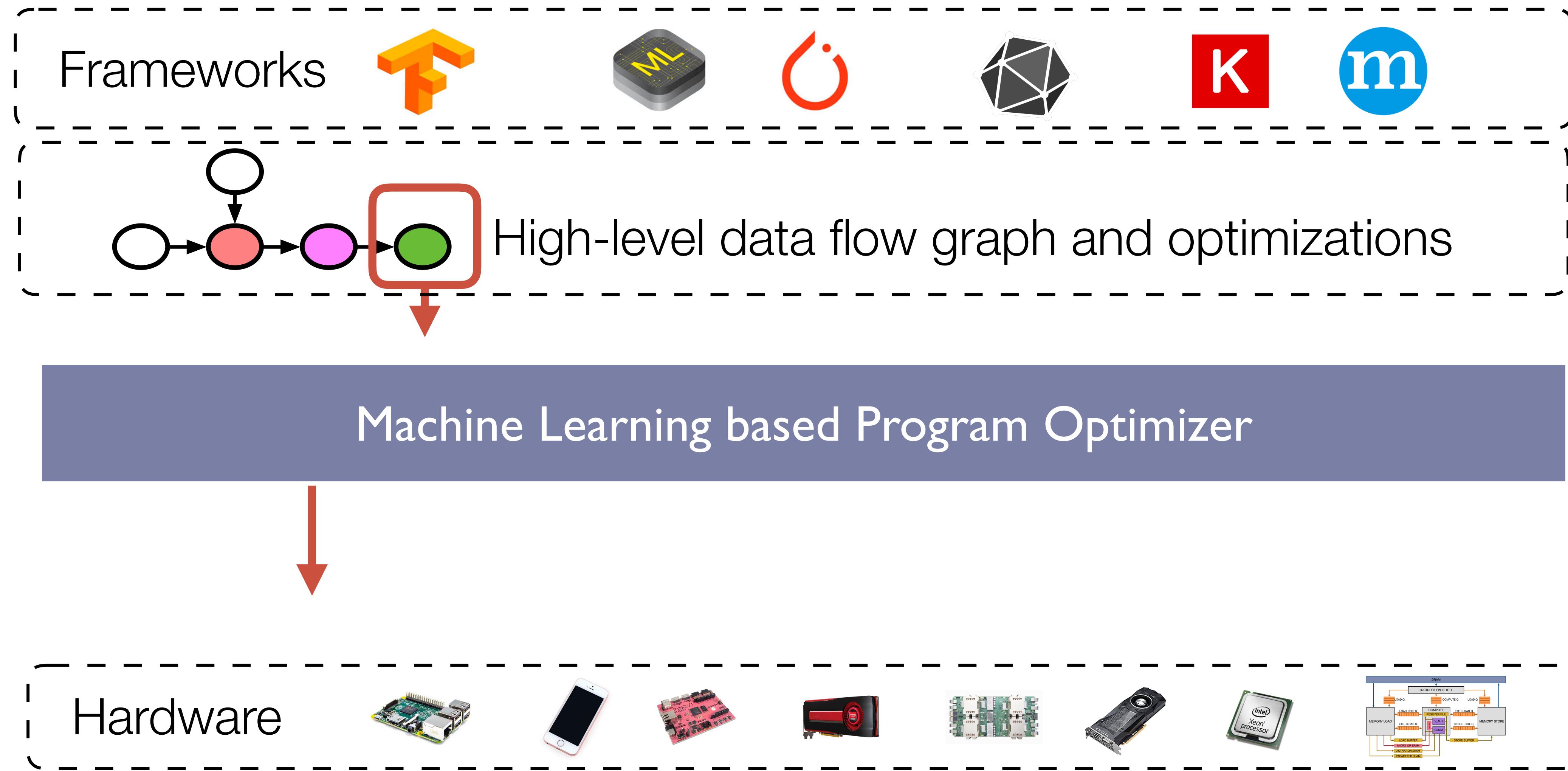
# Learning-based Learning System



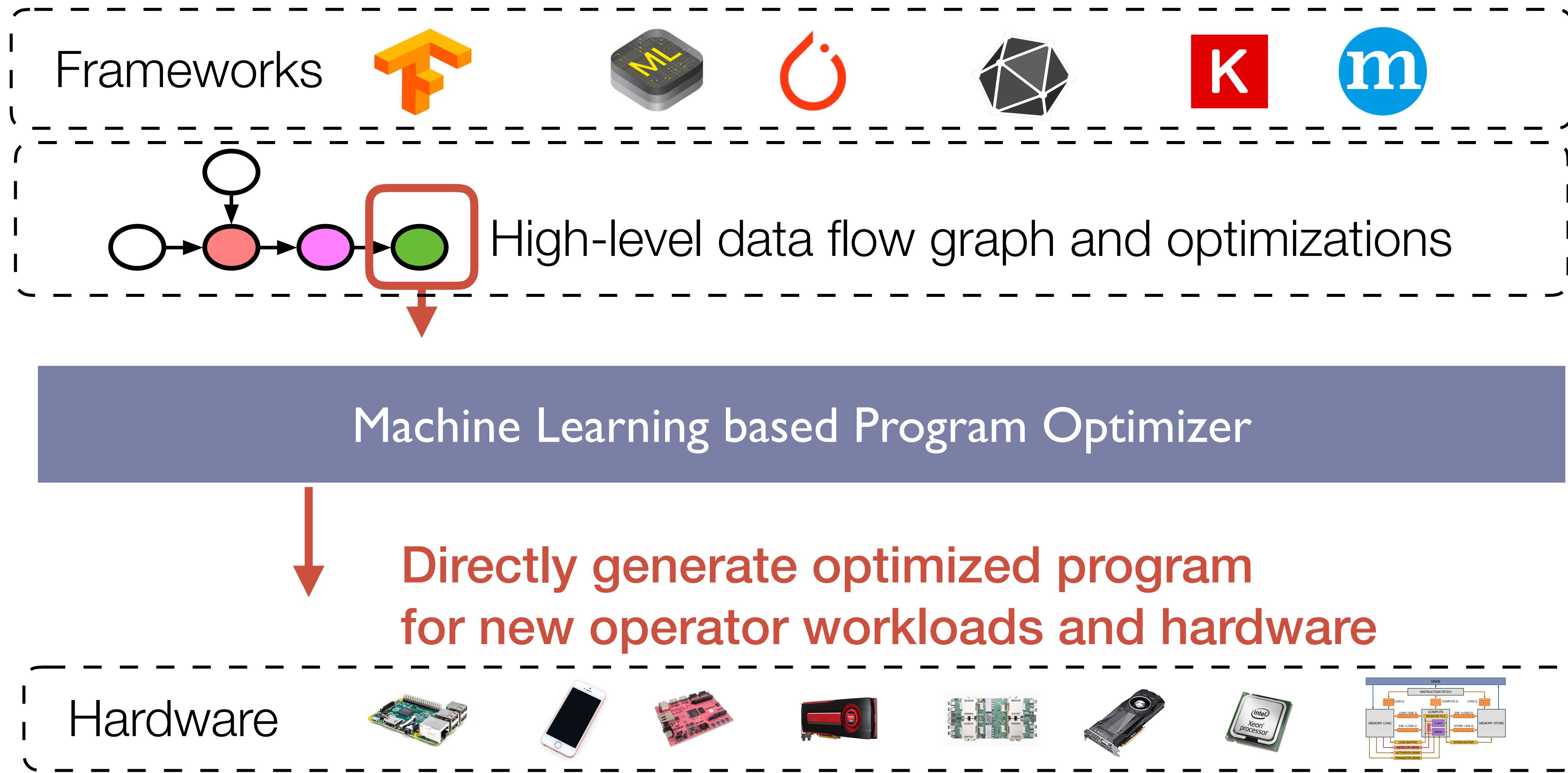
# Learning-based Learning System



# Learning-based Learning System



# Learning-based Learning System



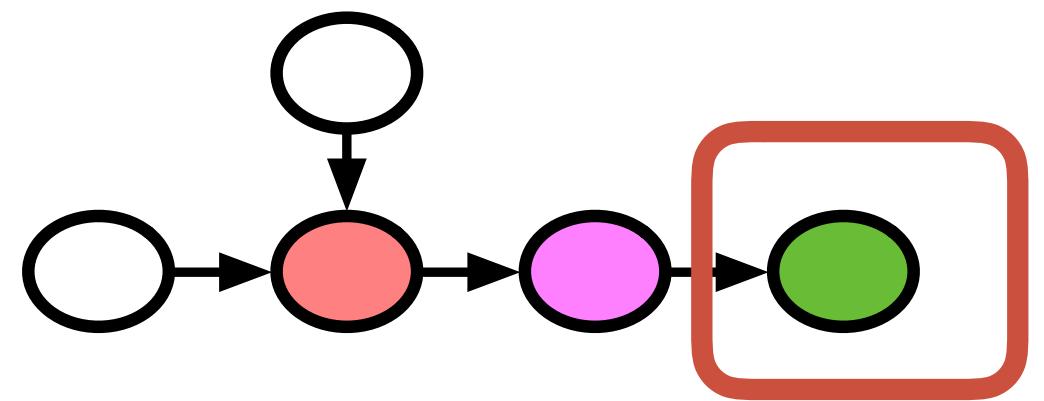
# TVM: Learning-based Learning System

Why do we need machine learning for systems

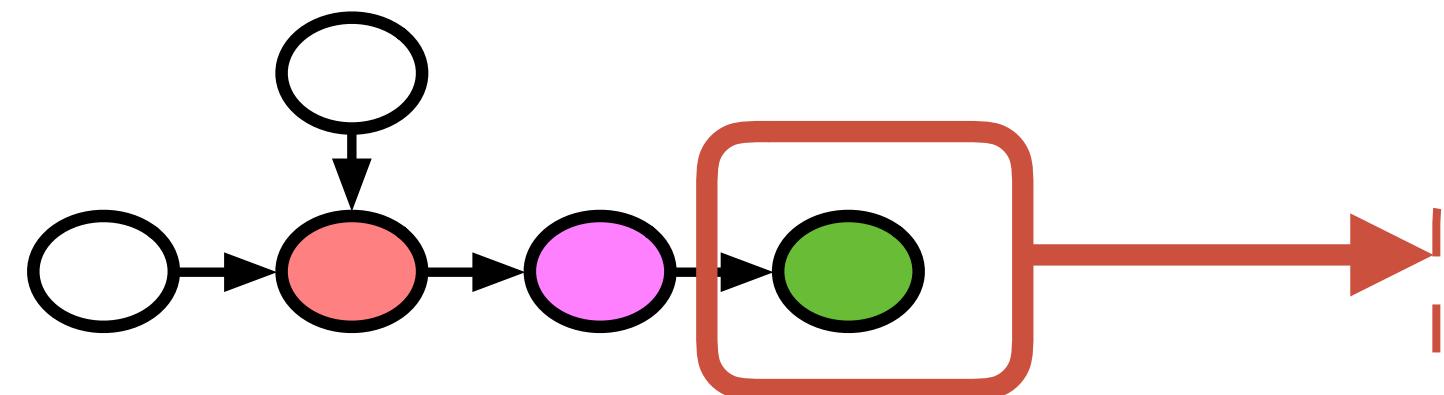
**How to build intelligent systems with learning**

End to end learning-based learning system stack

# Problem Setting



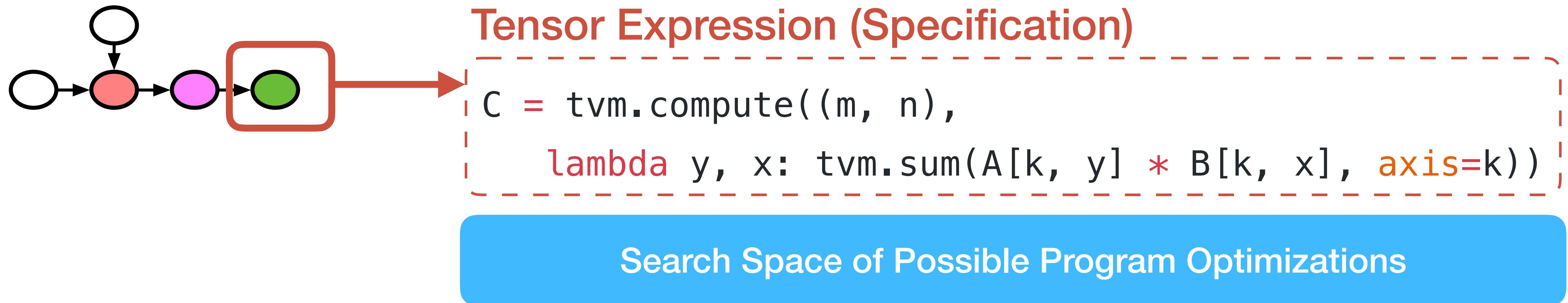
# Problem Setting



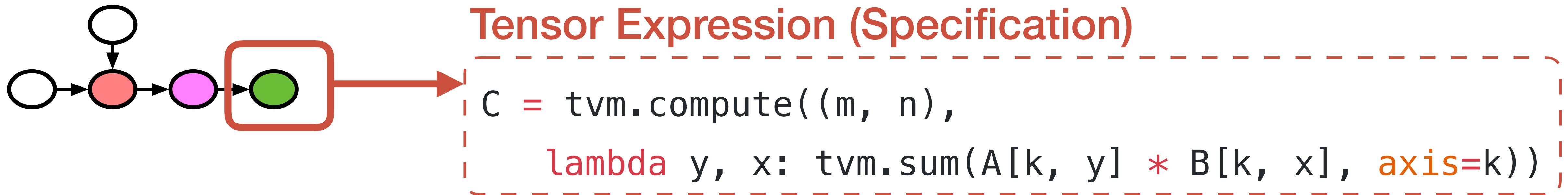
## Tensor Expression (Specification)

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

# Problem Setting

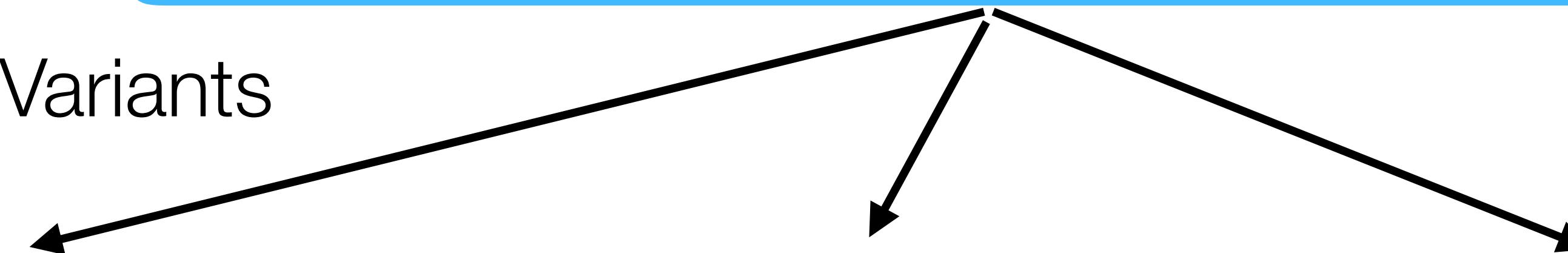


# Problem Setting



Search Space of Possible Program Optimizations

Low-level Program Variants

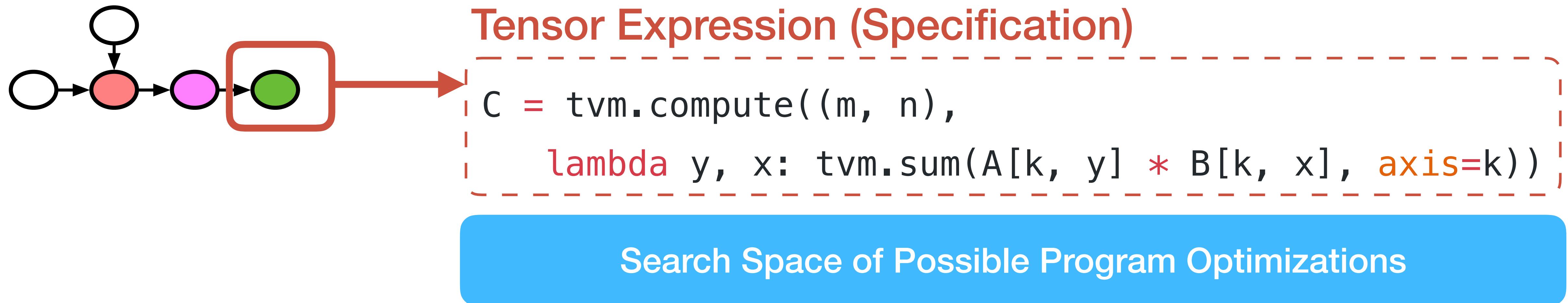


```
inp_buffer AL[8][8], BL[8][8]  
acc_buffer CL[8][8]  
for yo in range(128):  
    for xo in range(128):  
        vdla.fill_zero(CL)  
        for ko in range(128):  
            vdla.dma_copy2d(AL, A[ko*8:ko*8+8][yo*8:yo*8+8])  
            vdla.dma_copy2d(BL, B[ko*8:ko*8+8][xo*8:xo*8+8])  
            vdla.fused_gemm8x8_add(CL, AL, BL)  
            vdla.dma_copy2d(C[yo*8:yo*8+8,xo*8:xo*8+8], CL)
```

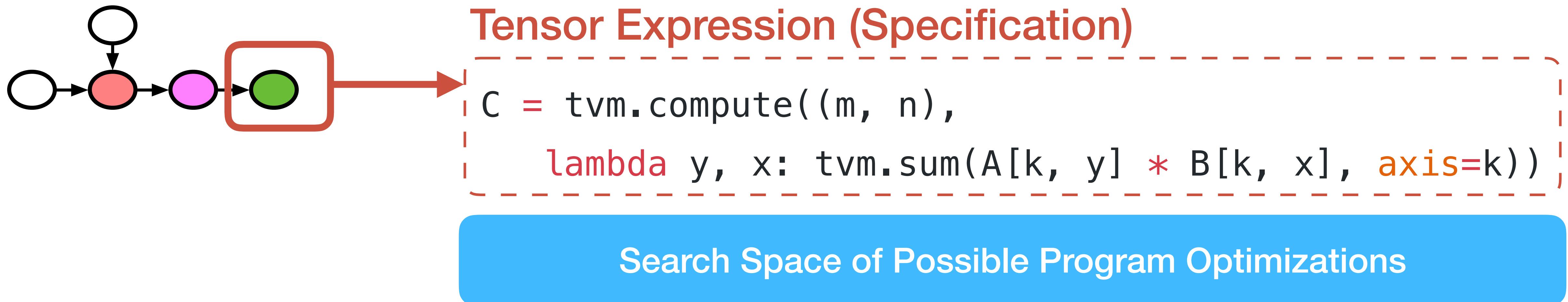
```
for yo in range(128):  
    for xo in range(128):  
        C[yo*8:yo*8+8][xo*8:xo*8+8] = 0  
        for ko in range(128):  
            for yi in range(8):  
                for xi in range(8):  
                    for ki in range(8):  
                        C[yo*8+yi][xo*8+xi] +=  
                            A[ko*8+ki][yo*8+yi] * B[ko*8+ki][xo*8+xi]
```

```
for y in range(1024):  
    for x in range(1024):  
        C[y][x] = 0  
        for k in range(1024):  
            C[y][x] += A[k][y] * B[k][x]
```

# Example Instance in a Search Space



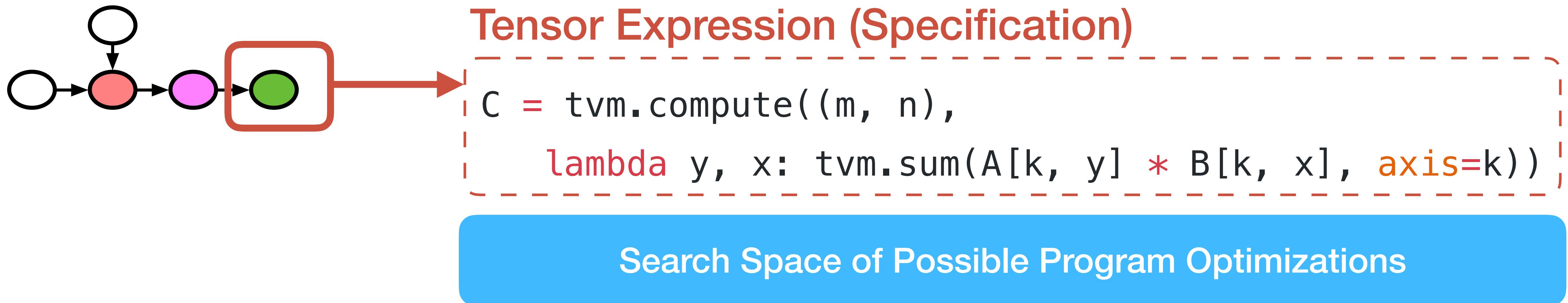
# Example Instance in a Search Space



Vanilla Code

```
for y in range(1024):  
    for x in range(1024):  
        C[y][x] = 0  
        for k in range(1024):  
            C[y][x] += A[k][y] * B[k][x]
```

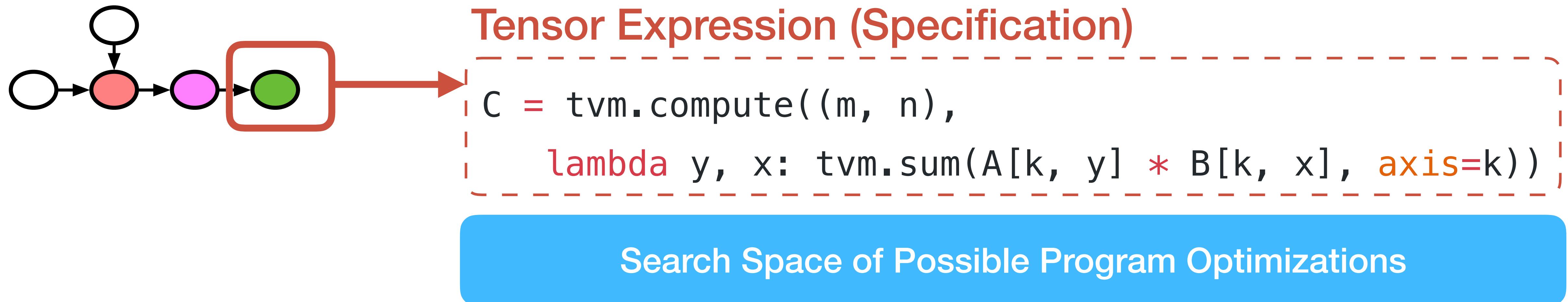
# Example Instance in a Search Space



## Loop Tiling for Locality

```
for yo in range(128):  
    for xo in range(128):  
        C[yo*8:yo*8+8][xo*8:xo*8+8] = 0  
        for ko in range(128):  
            for yi in range(8):  
                for xi in range(8):  
                    for ki in range(8):  
                        C[yo*8+yi][xo*8+xi] +=  
                            A[ko*8+ki][yo*8+yi] * B[ko*8+ki][xo*8+xi]
```

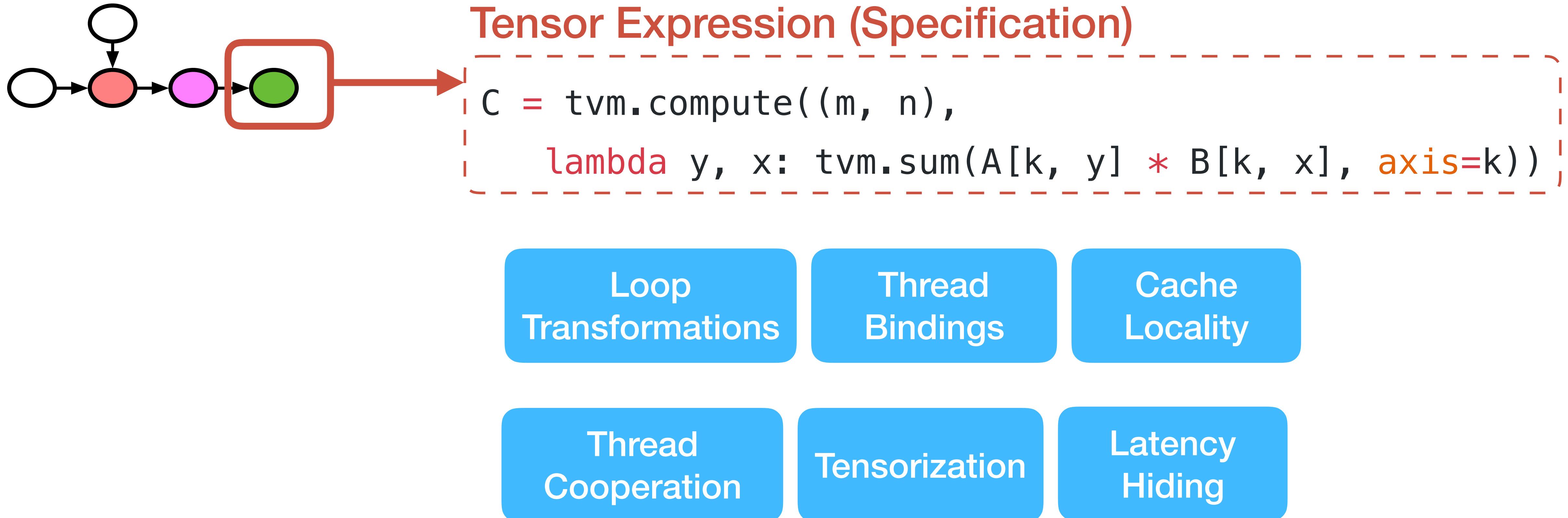
# Example Instance in a Search Space



## Map to Accelerators

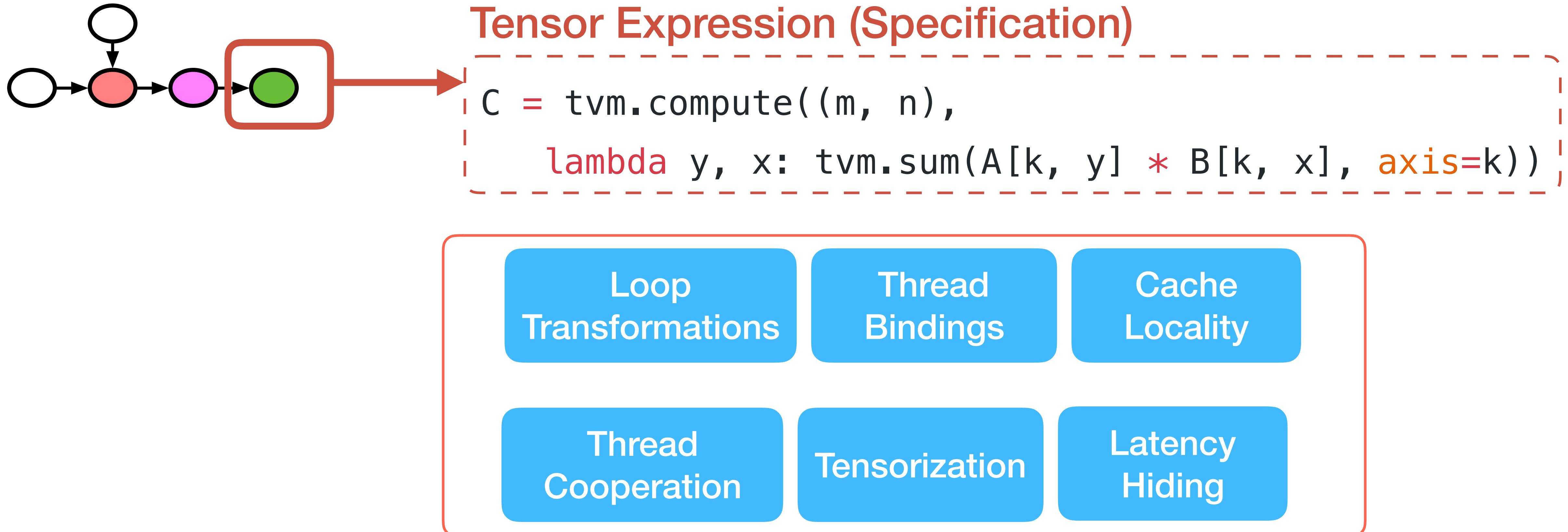
```
inp_buffer AL[8][8], BL[8][8]  
acc_buffer CL[8][8]  
for yo in range(128):  
    for xo in range(128):  
        vdla.fill_zero(CL)  
        for ko in range(128):  
            vdla.dma_copy2d(AL, A[ko*8:ko*8+8][yo*8:yo*8+8])  
            vdla.dma_copy2d(BL, B[ko*8:ko*8+8][xo*8:xo*8+8])  
            vdla.fused_gemm8x8_add(CL, AL, BL)  
        vdla.dma_copy2d(C[yo*8:yo*8+8,xo*8:xo*8+8], CL)
```

# Optimization Choices in a Search Space



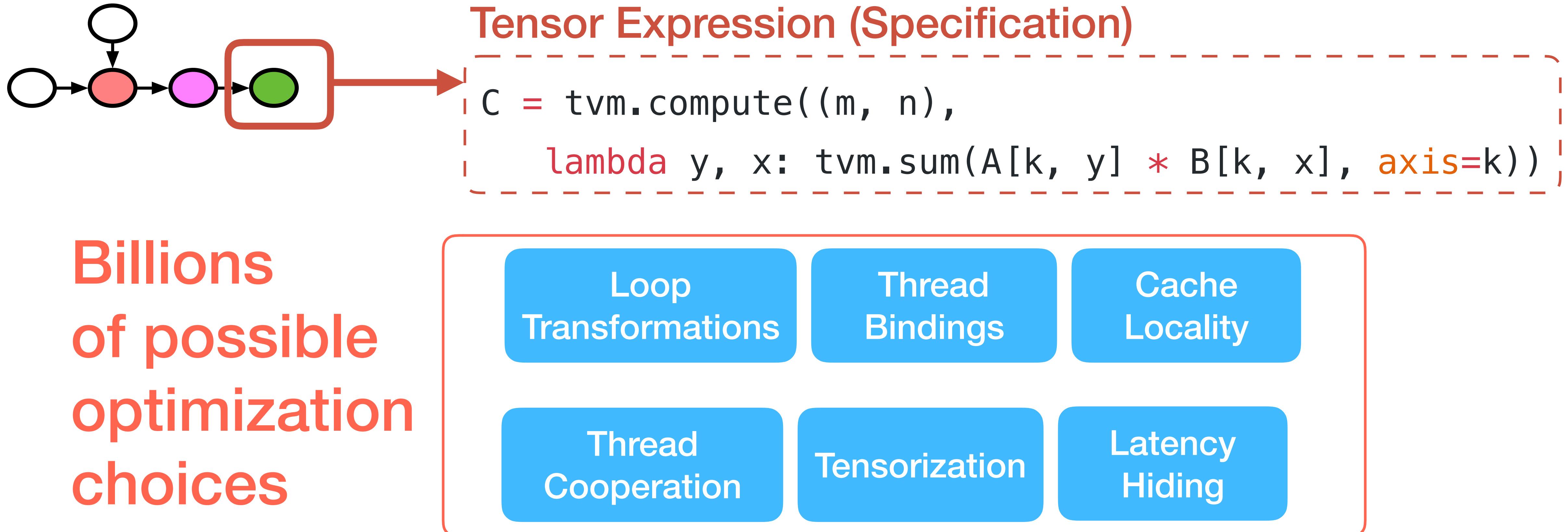
More details about the search space  
in the second half of the talk

# Optimization Choices in a Search Space



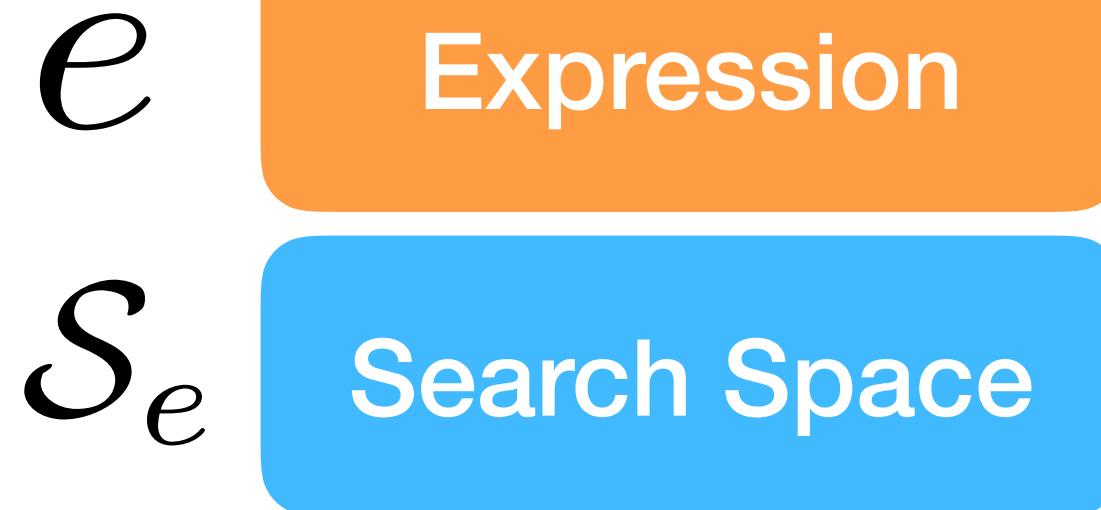
More details about the search space  
in the second half of the talk

# Optimization Choices in a Search Space

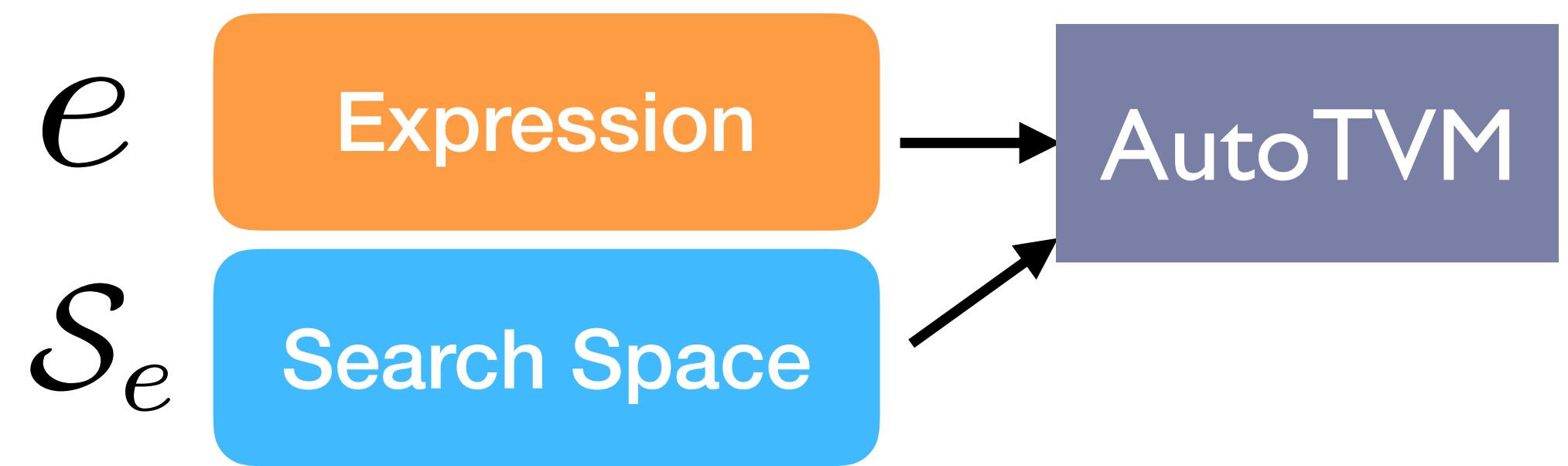


More details about the search space  
in the second half of the talk

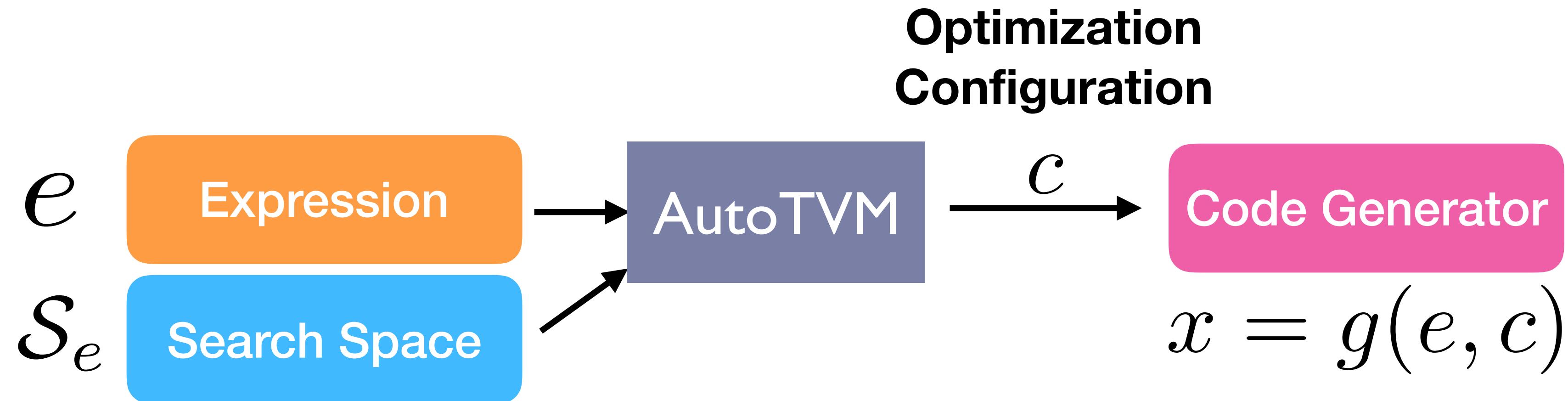
# Problem Formalization



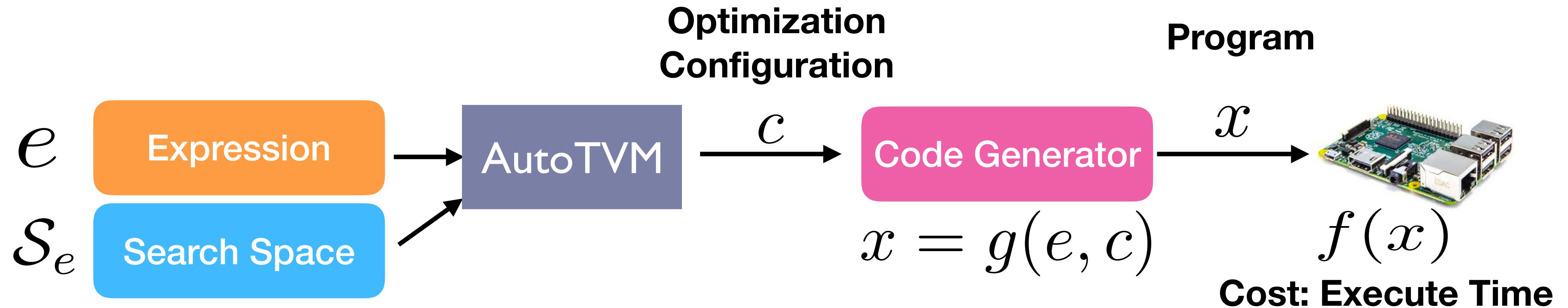
# Problem Formalization



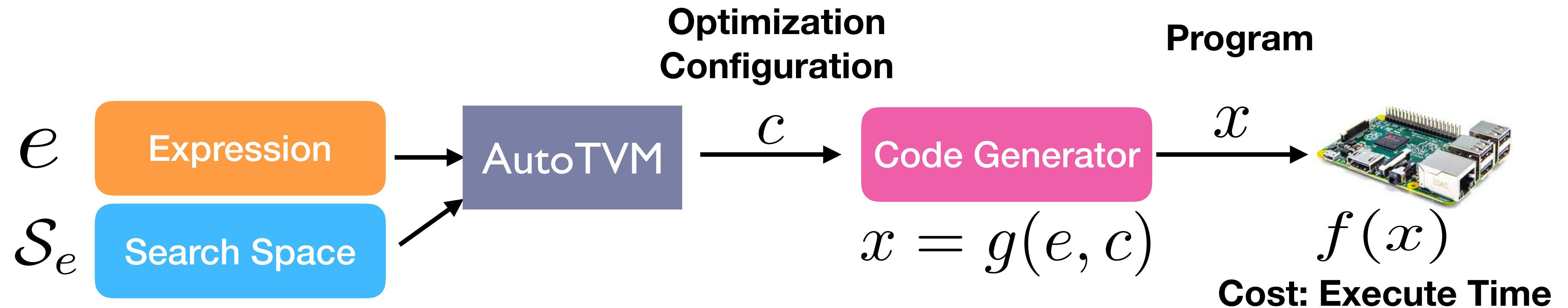
# Problem Formalization



# Problem Formalization



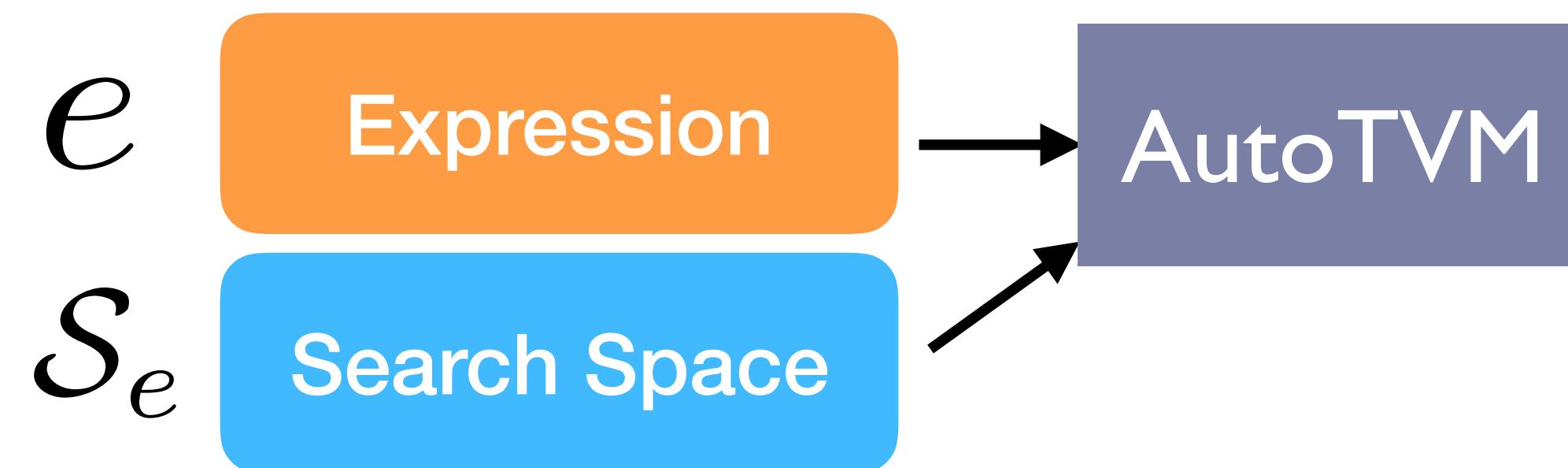
# Problem Formalization



**Objective**  $\operatorname{argmin}_{c \in S_e} f(g(e, c))$

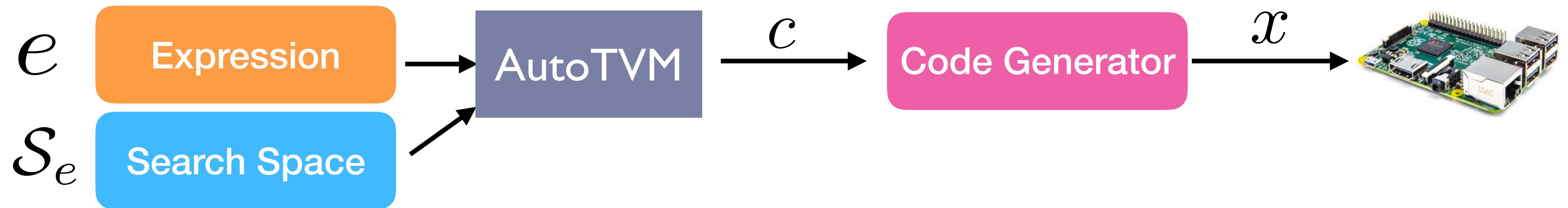
# Black-box Optimization

Try each configuration  $c$  until we find a good one



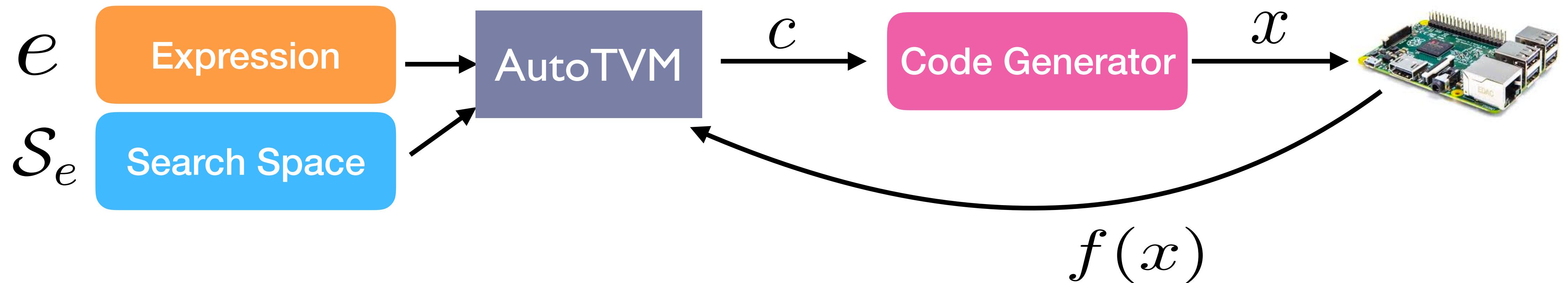
# Black-box Optimization

Try each configuration  $c$  until we find a good one



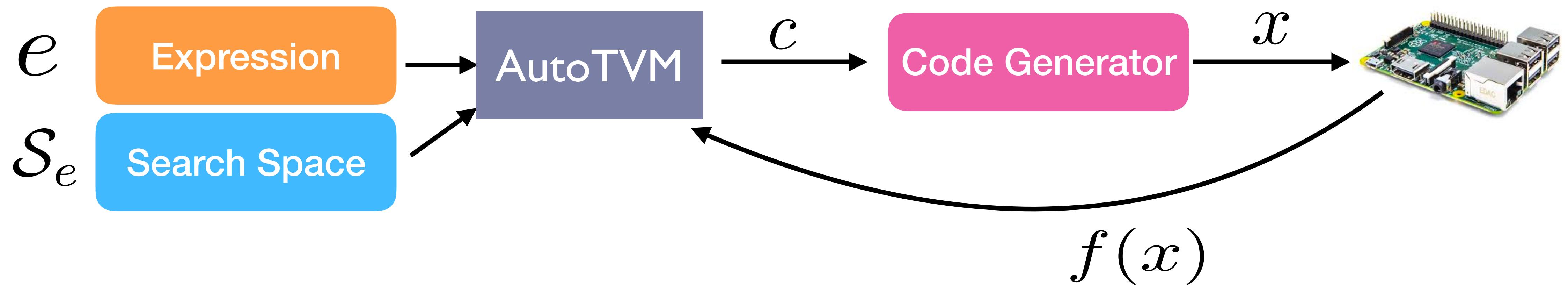
# Black-box Optimization

Try each configuration  $c$  until we find a good one



# Black-box Optimization

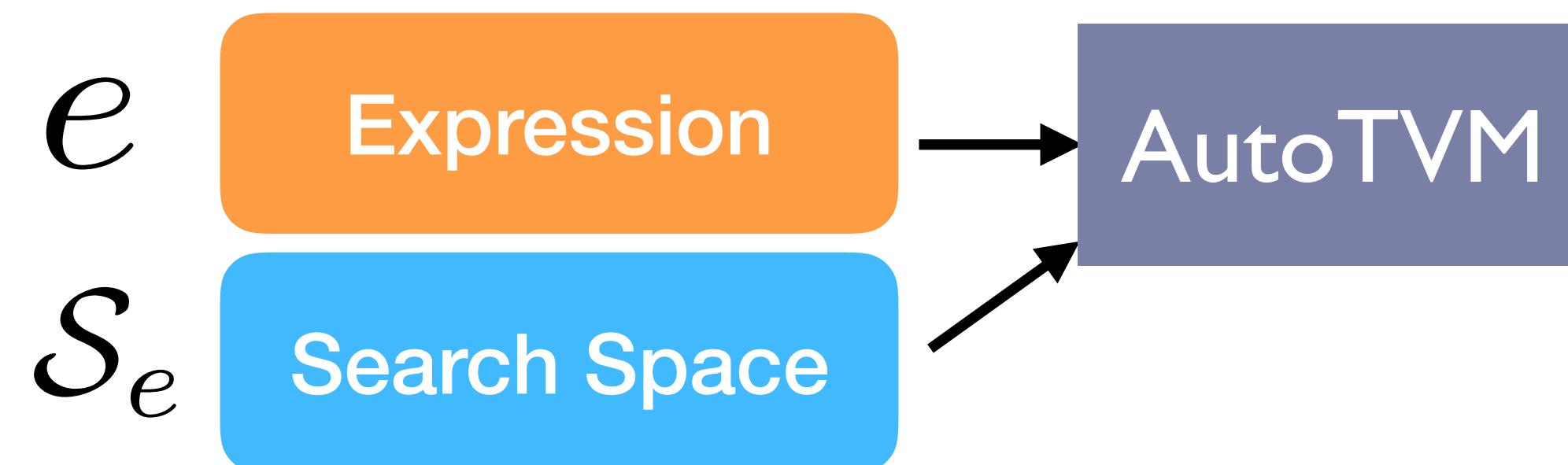
Try each configuration  $c$  until we find a good one



**Challenge:** lots of experimental trials, each trial costs ~1 second

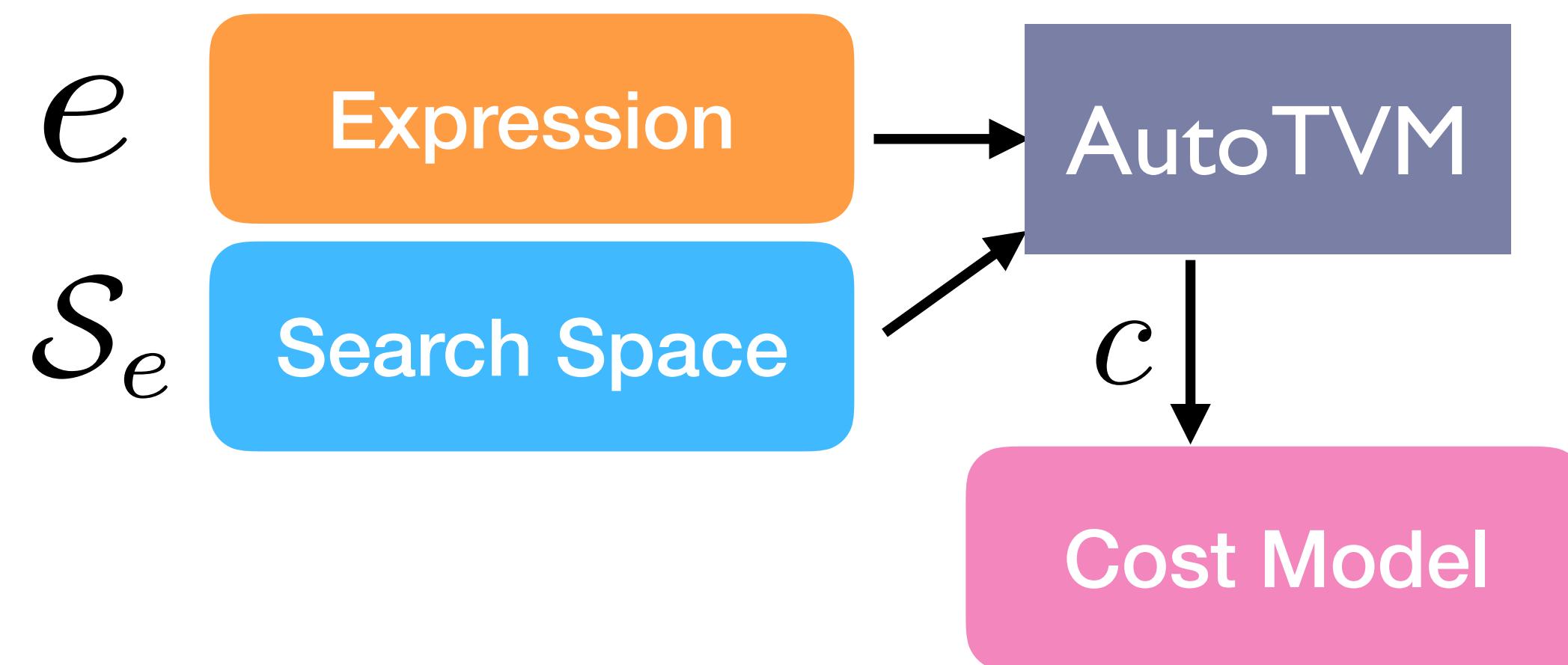
# Cost-model Driven Approach

Use cost model to pick configuration



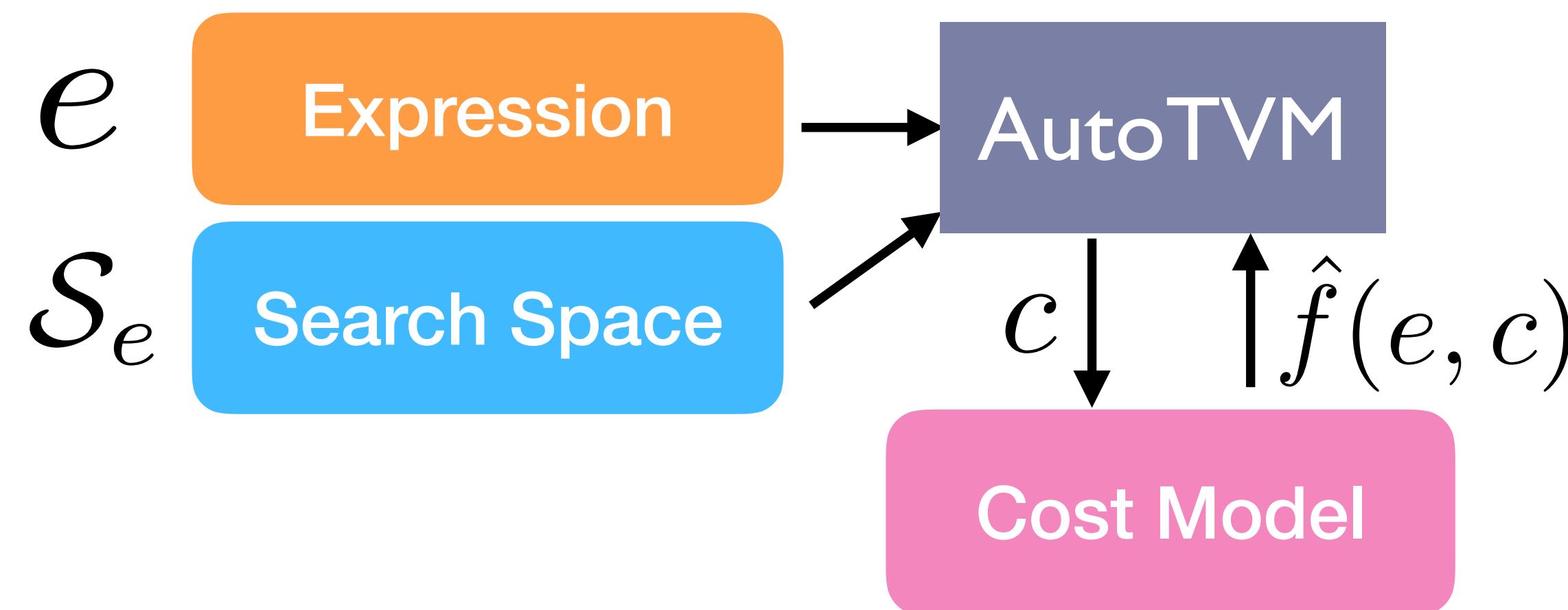
# Cost-model Driven Approach

Use cost model to pick configuration



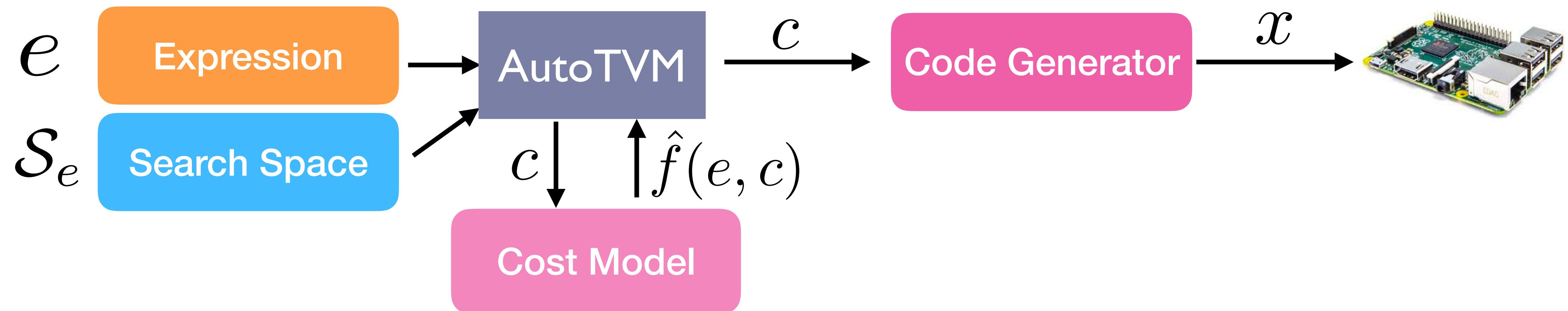
# Cost-model Driven Approach

Use cost model to pick configuration



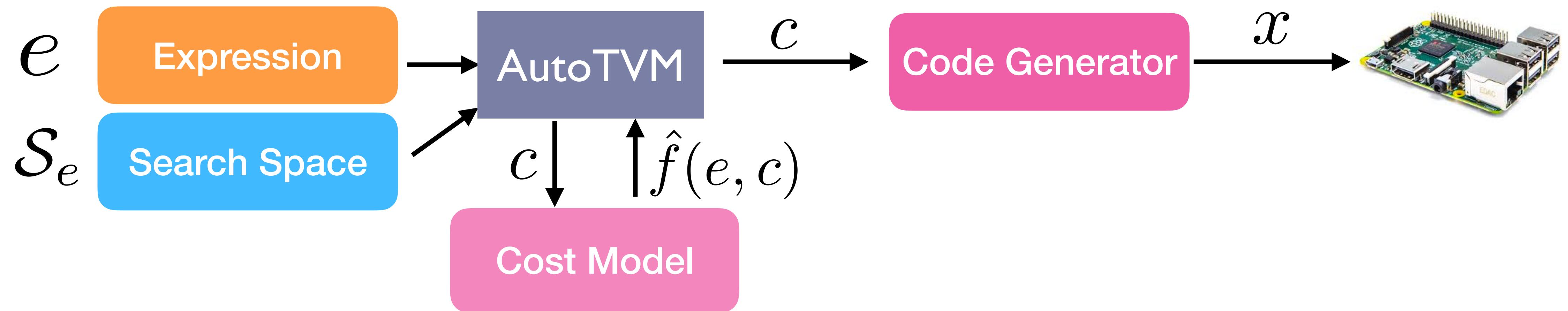
# Cost-model Driven Approach

Use cost model to pick configuration



# Cost-model Driven Approach

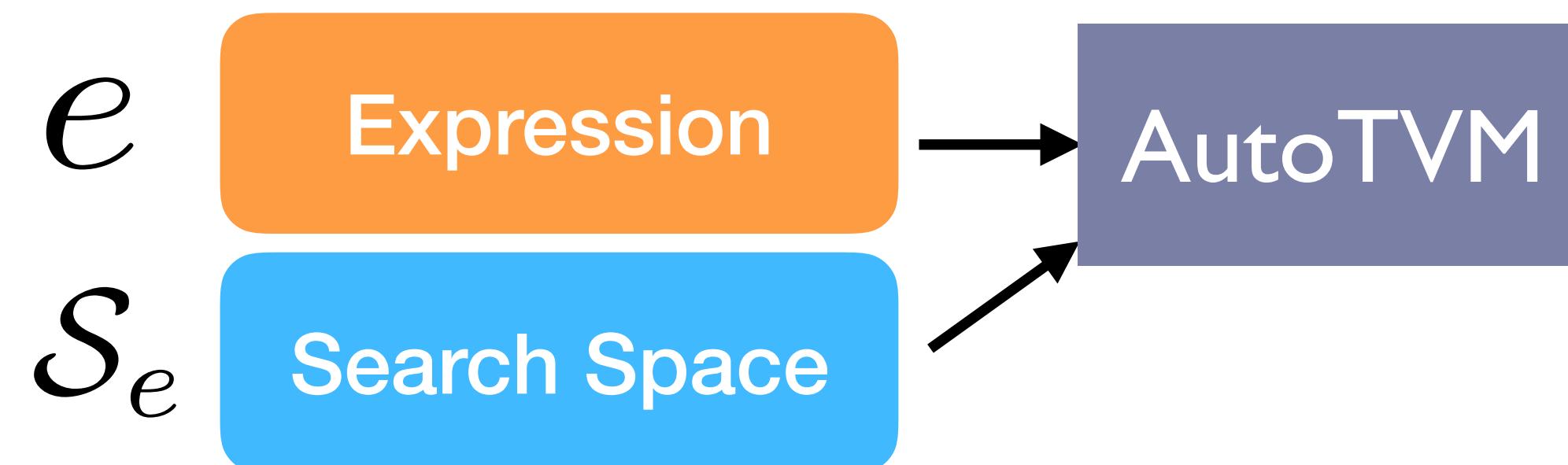
Use cost model to pick configuration



**Challenge:** Need reliable cost model per hardware

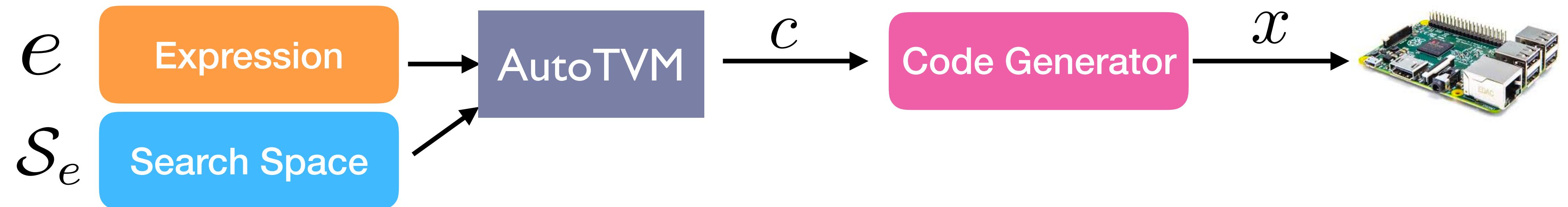
# Statistical Cost Model

Use machine learning to learn a statistical cost model



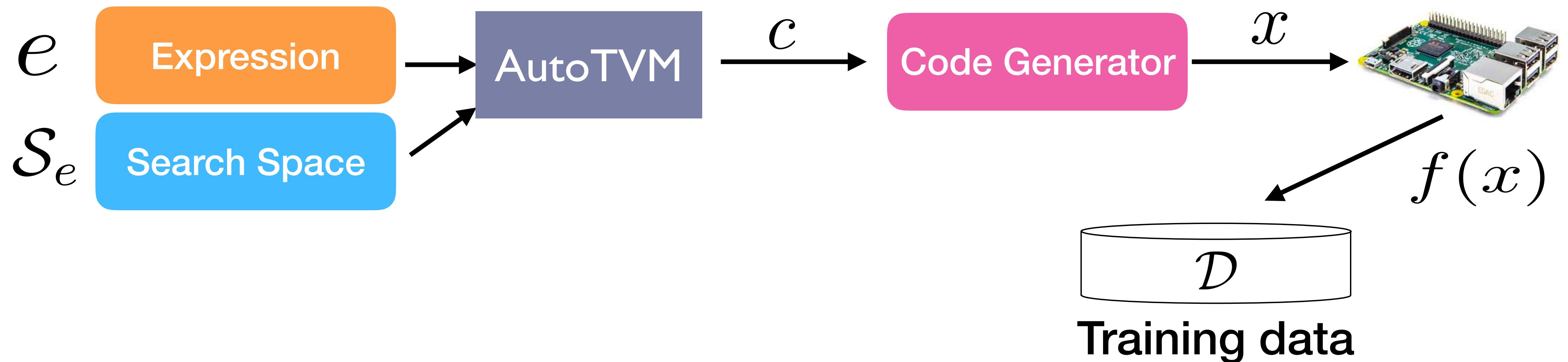
# Statistical Cost Model

Use machine learning to learn a statistical cost model



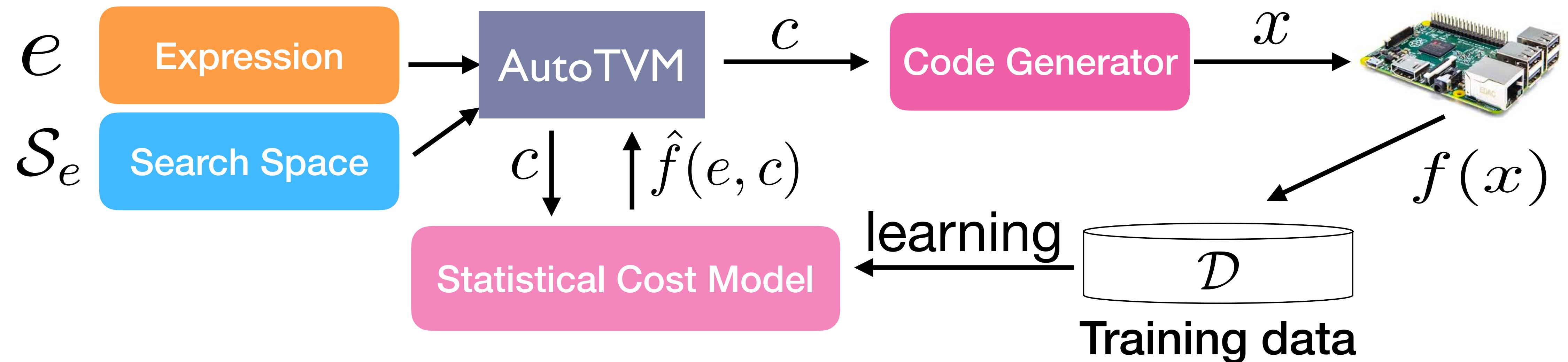
# Statistical Cost Model

Use machine learning to learn a statistical cost model



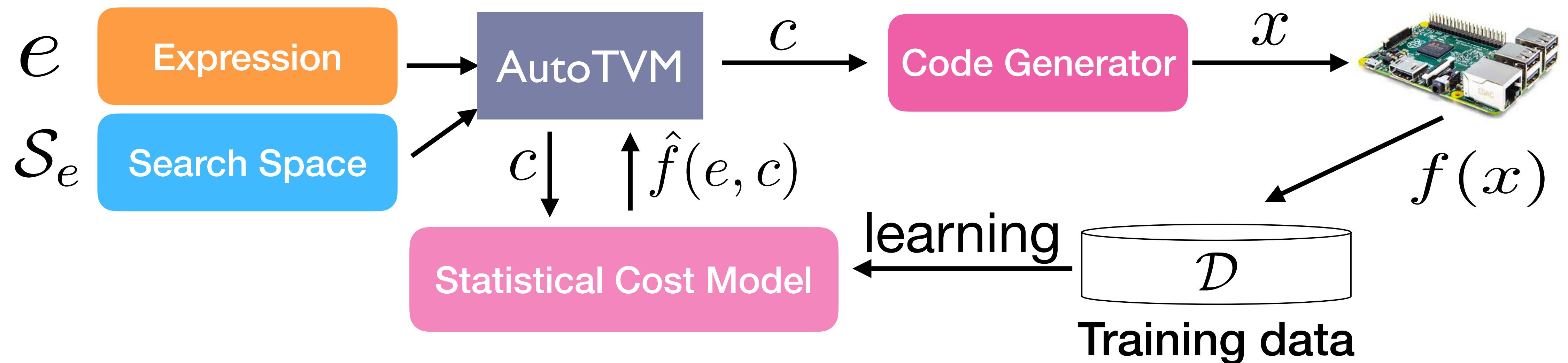
# Statistical Cost Model

Use machine learning to learn a statistical cost model



# Statistical Cost Model

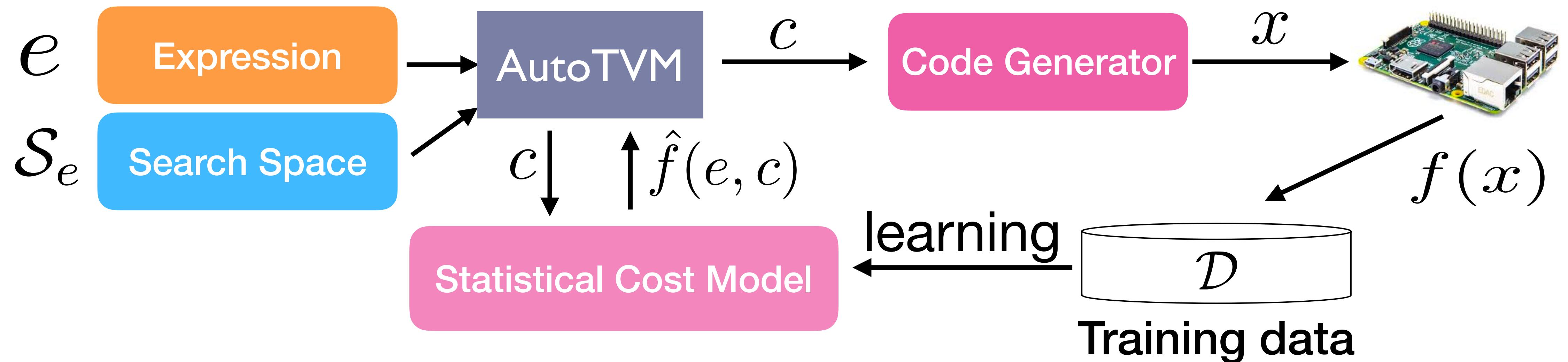
Use machine learning to learn a statistical cost model



**Benefit: Automatically adapt to hardware type**

# Statistical Cost Model

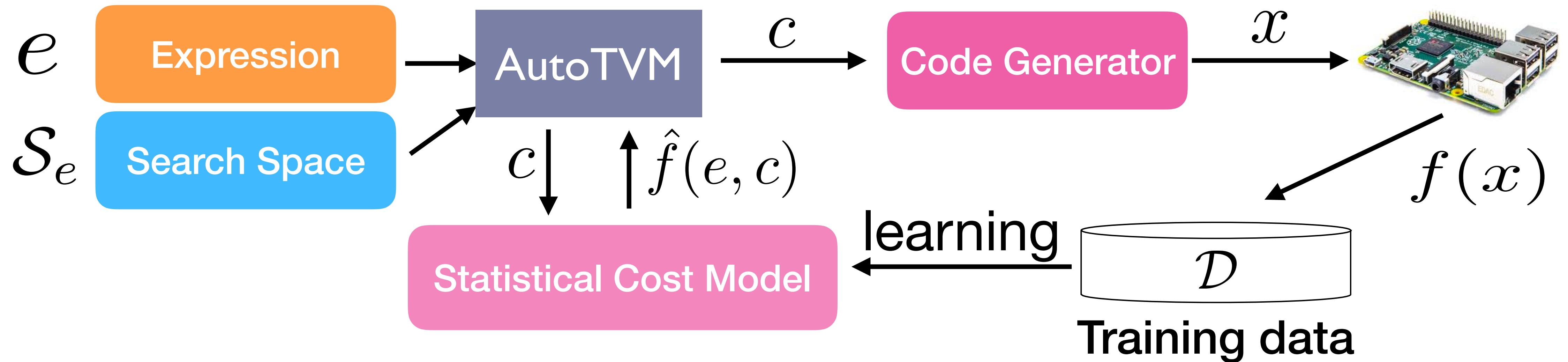
Use machine learning to learn a statistical cost model



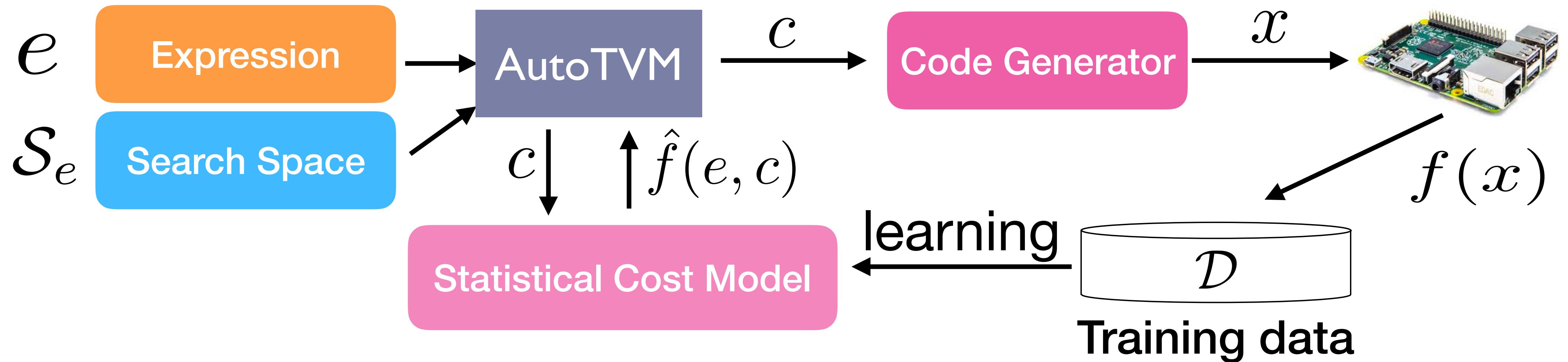
**Benefit: Automatically adapt to hardware type**

**Challenge: How to design the cost model**

# Unique Problem Characteristics

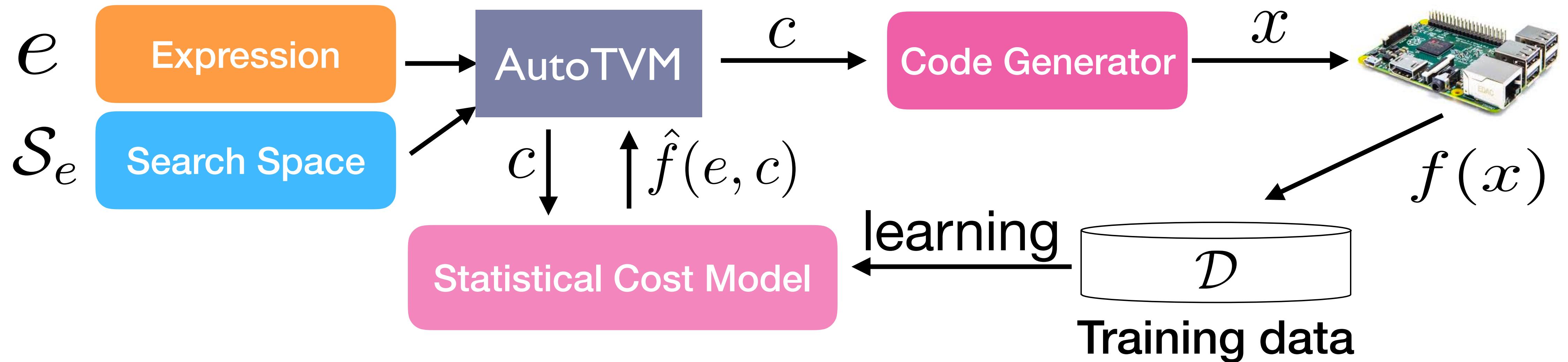


# Unique Problem Characteristics



**Relatively low  
experiment cost**

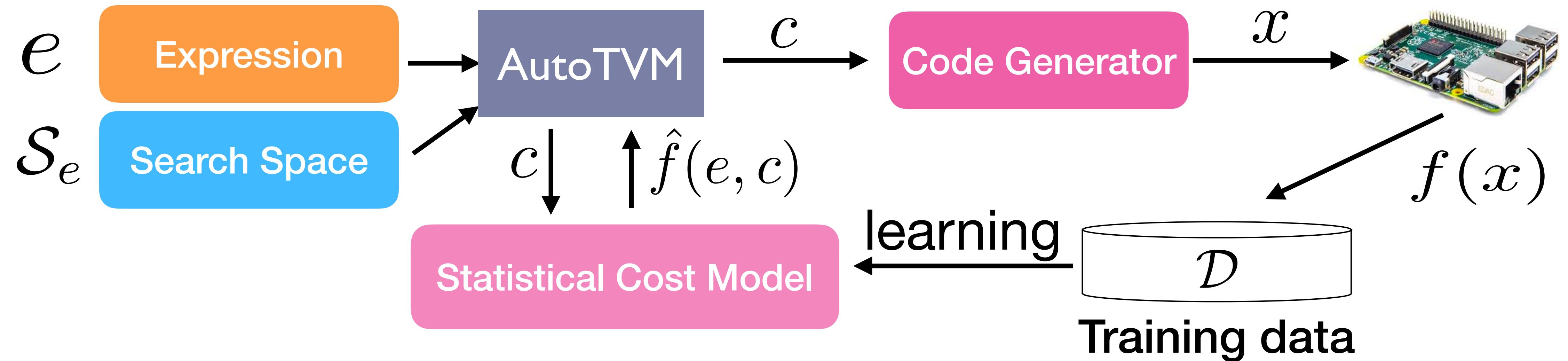
# Unique Problem Characteristics



**Relatively low  
experiment cost**

**Program-aware  
modeling**

# Unique Problem Characteristics

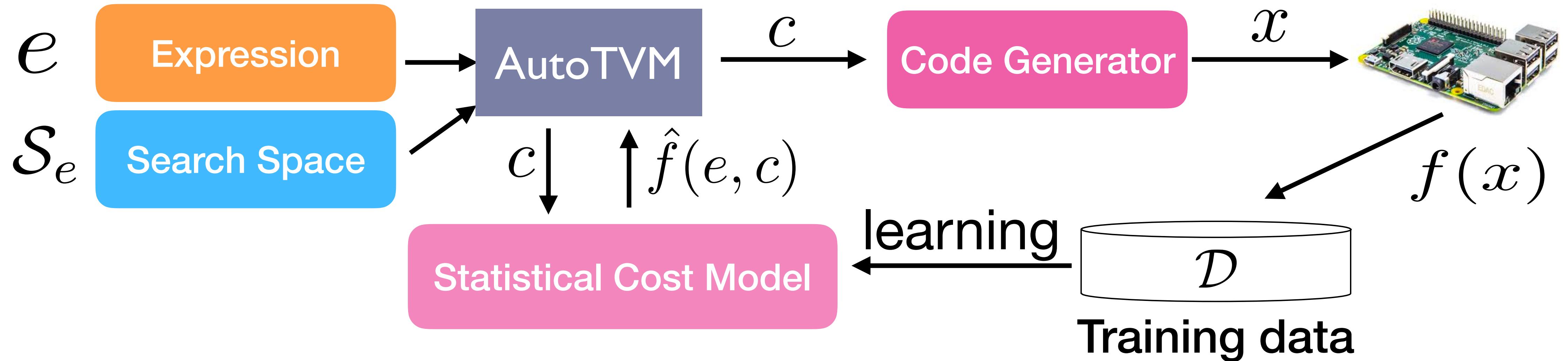


**Relatively low experiment cost**

**Program-aware modeling**

**Large number of similar tasks**

# Unique Problem Characteristics



Relatively low  
experiment cost

**Program-aware  
modeling**

Large number of  
similar tasks

# Vanilla Cost Modeling

*C*

High-level  
Configurations

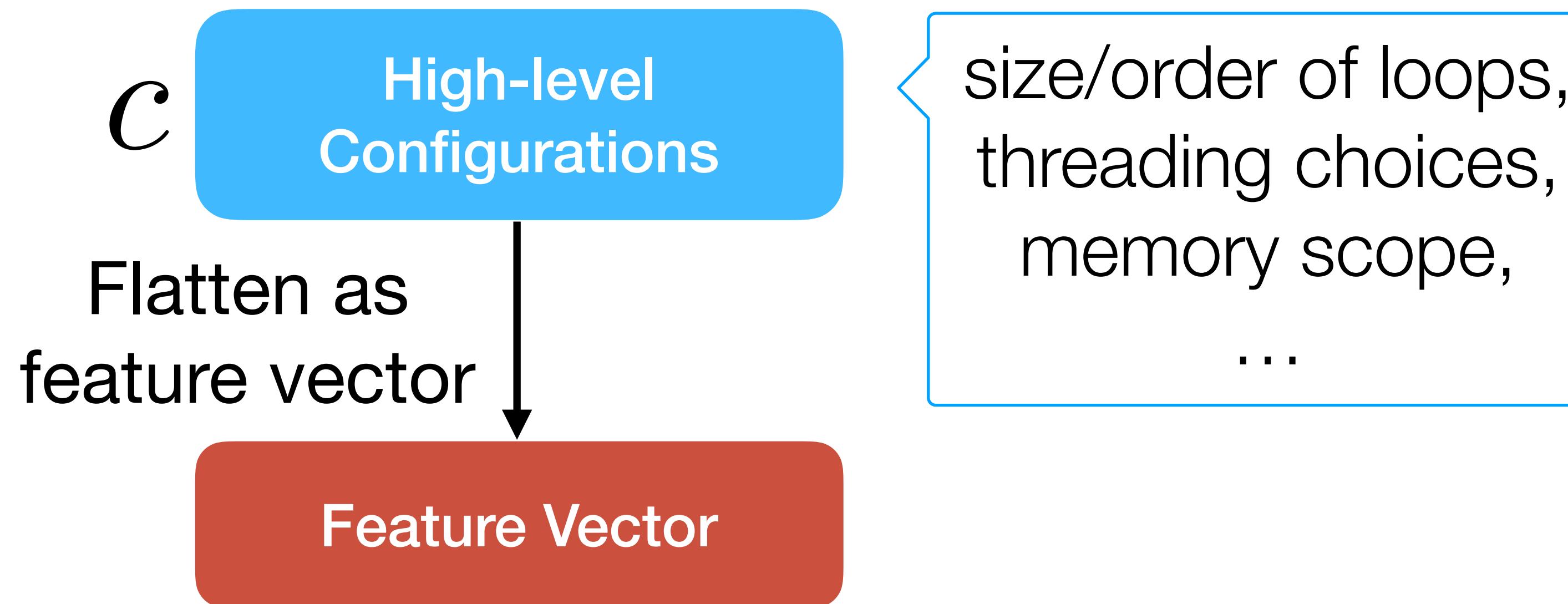
# Vanilla Cost Modeling

*C*

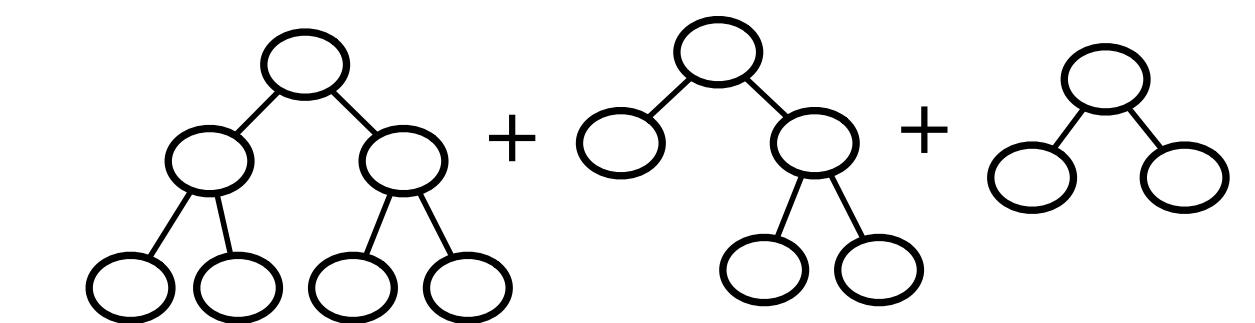
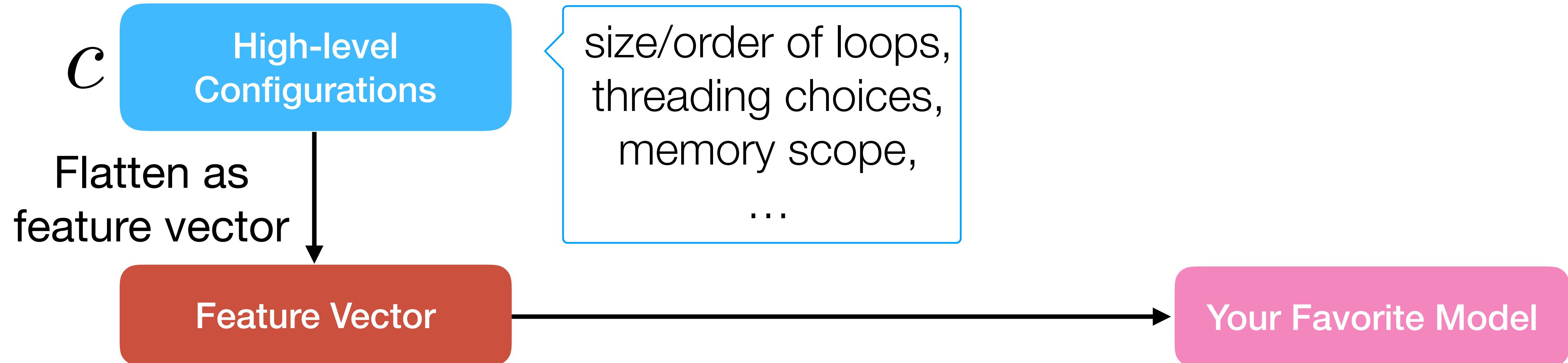
High-level  
Configurations

size/order of loops,  
threading choices,  
memory scope,  
...

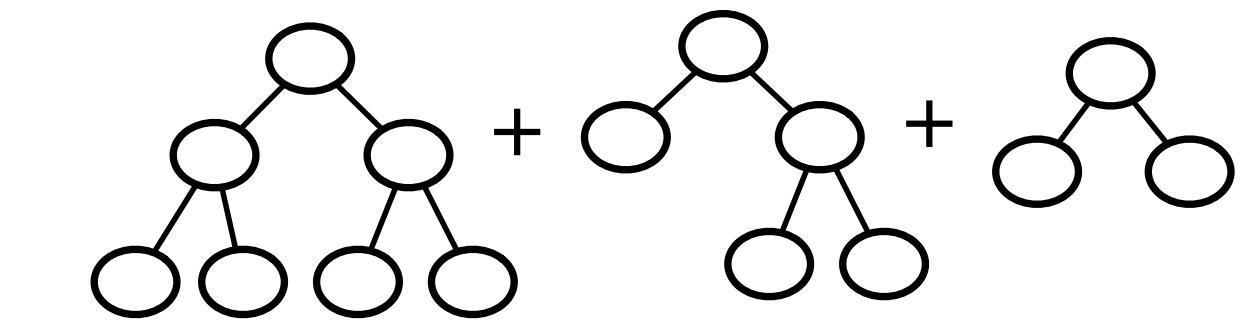
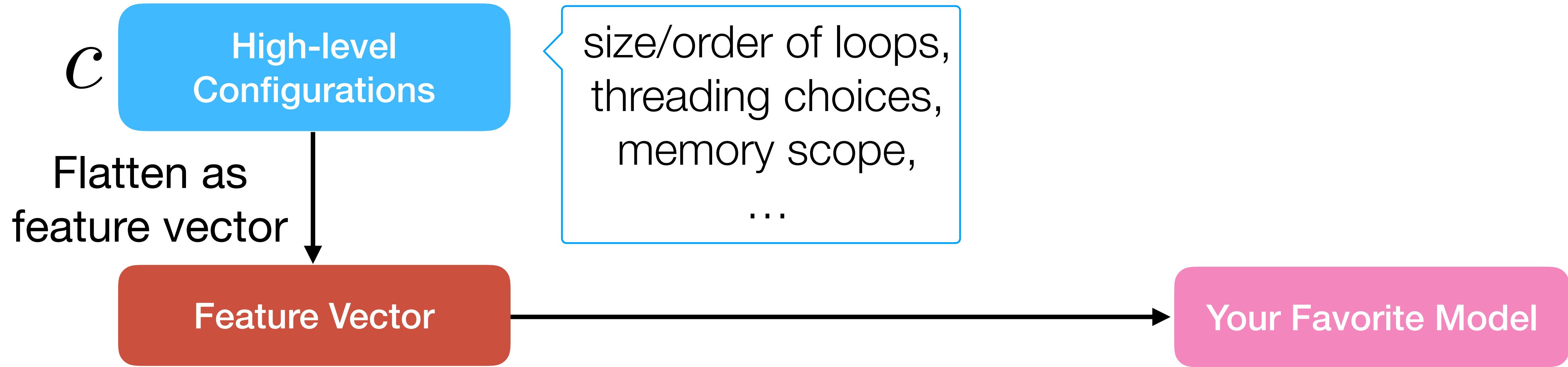
# Vanilla Cost Modeling



# Vanilla Cost Modeling



# Vanilla Cost Modeling



**Drawback:**

Ignores domain knowledge

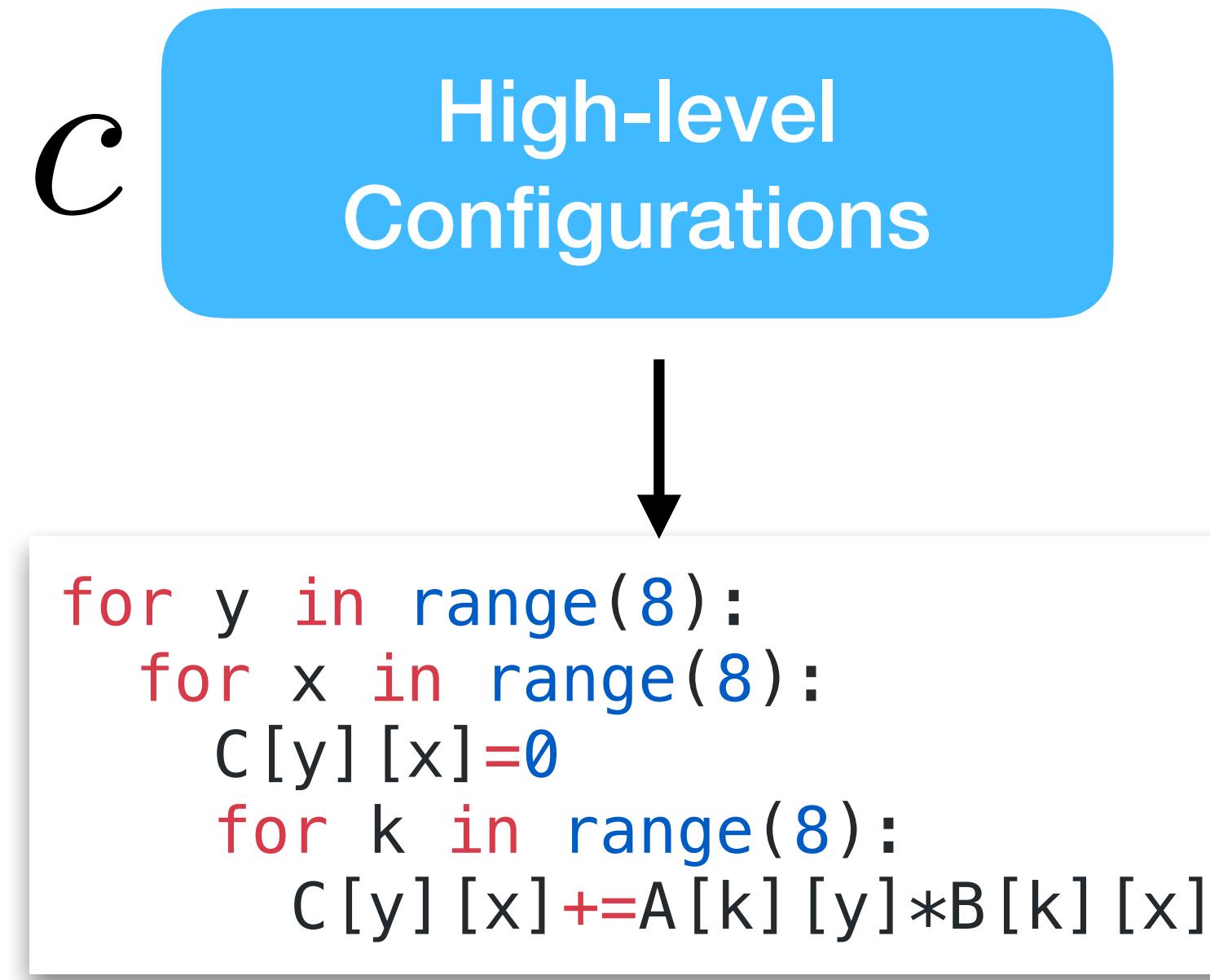
Set of configurations can differ per task (task dependent)

# Program-aware Modeling: Tree-based Approach

$C$

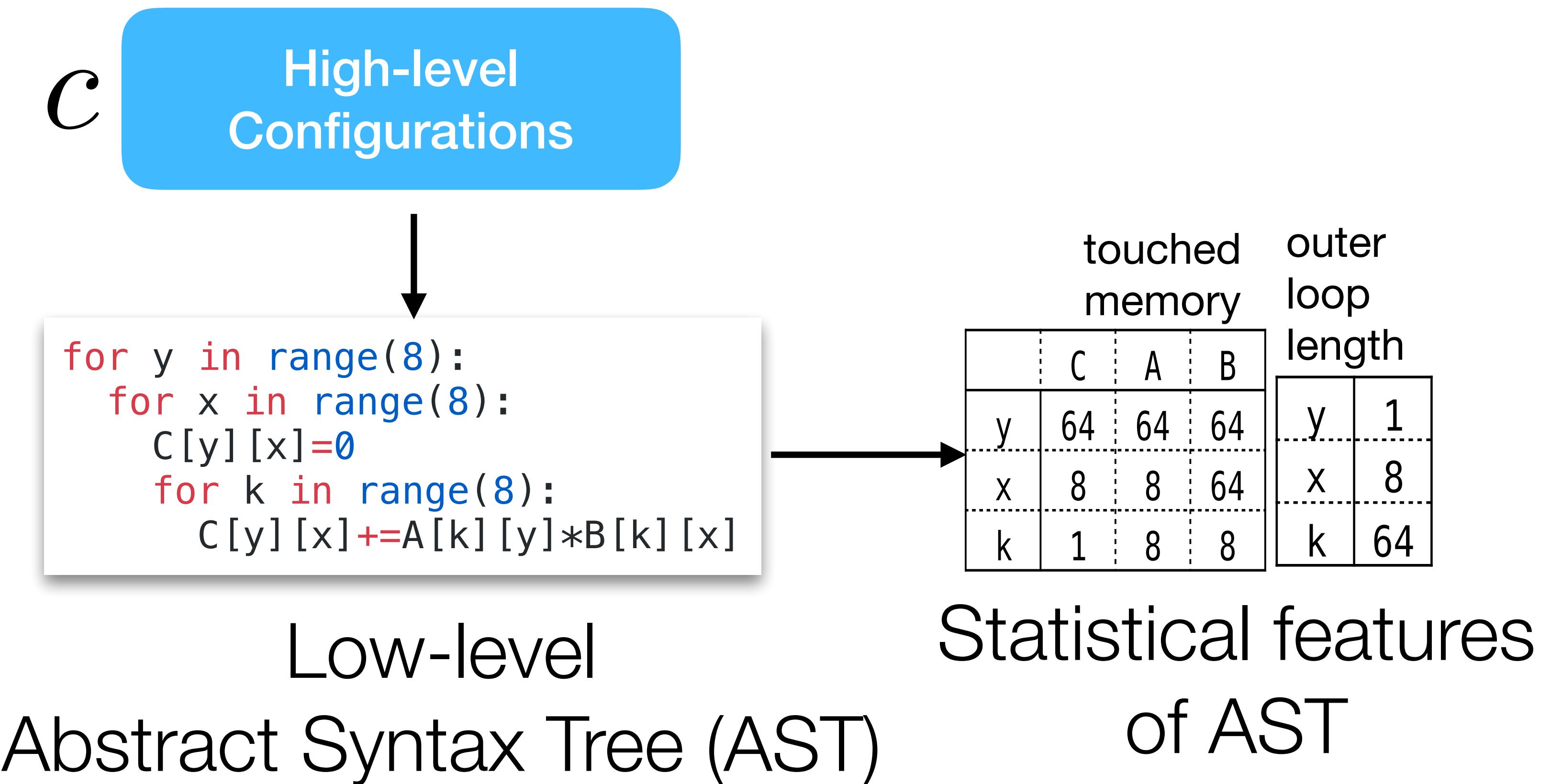
High-level  
Configurations

# Program-aware Modeling: Tree-based Approach

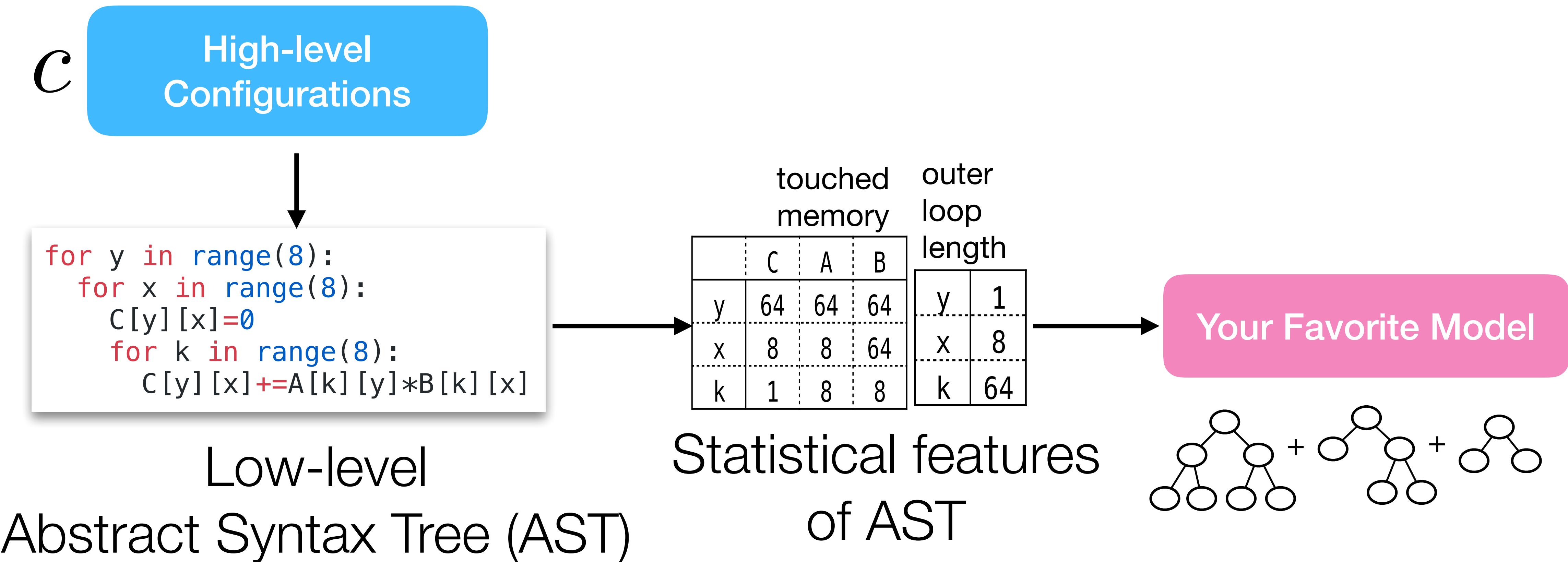


Low-level  
Abstract Syntax Tree (AST)

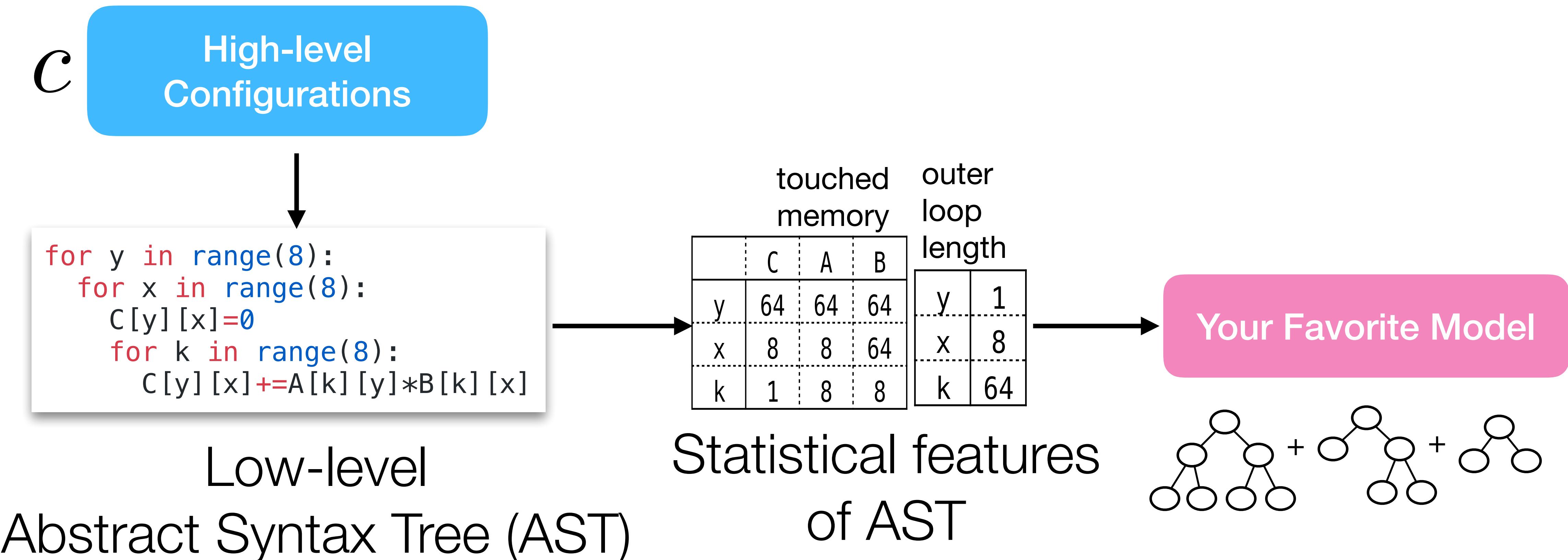
# Program-aware Modeling: Tree-based Approach



# Program-aware Modeling: Tree-based Approach

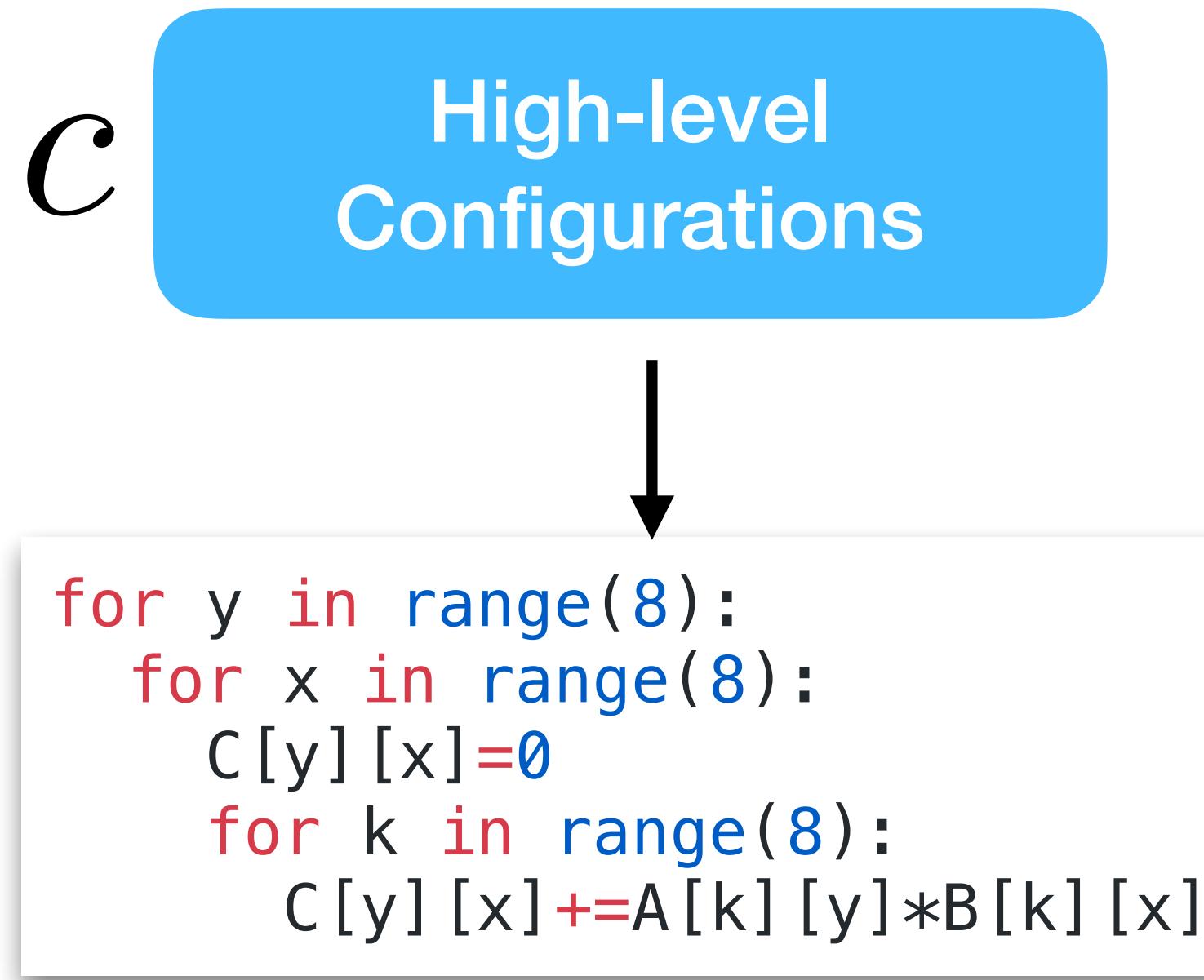


# Program-aware Modeling: Tree-based Approach



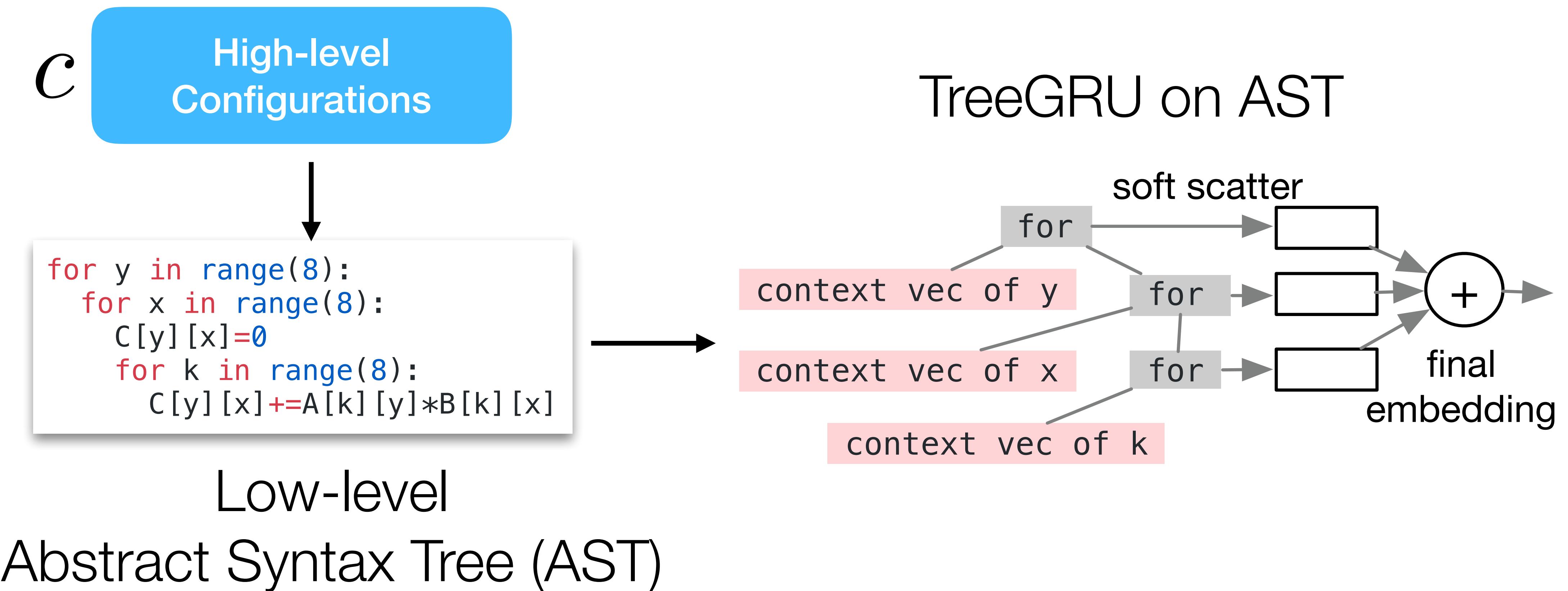
**Benefit: low-level AST is a common representation (task invariant)**

# Program-aware Modeling: Neural Approach



Low-level  
Abstract Syntax Tree (AST)

# Program-aware Modeling: Neural Approach



# Comparisons of Models

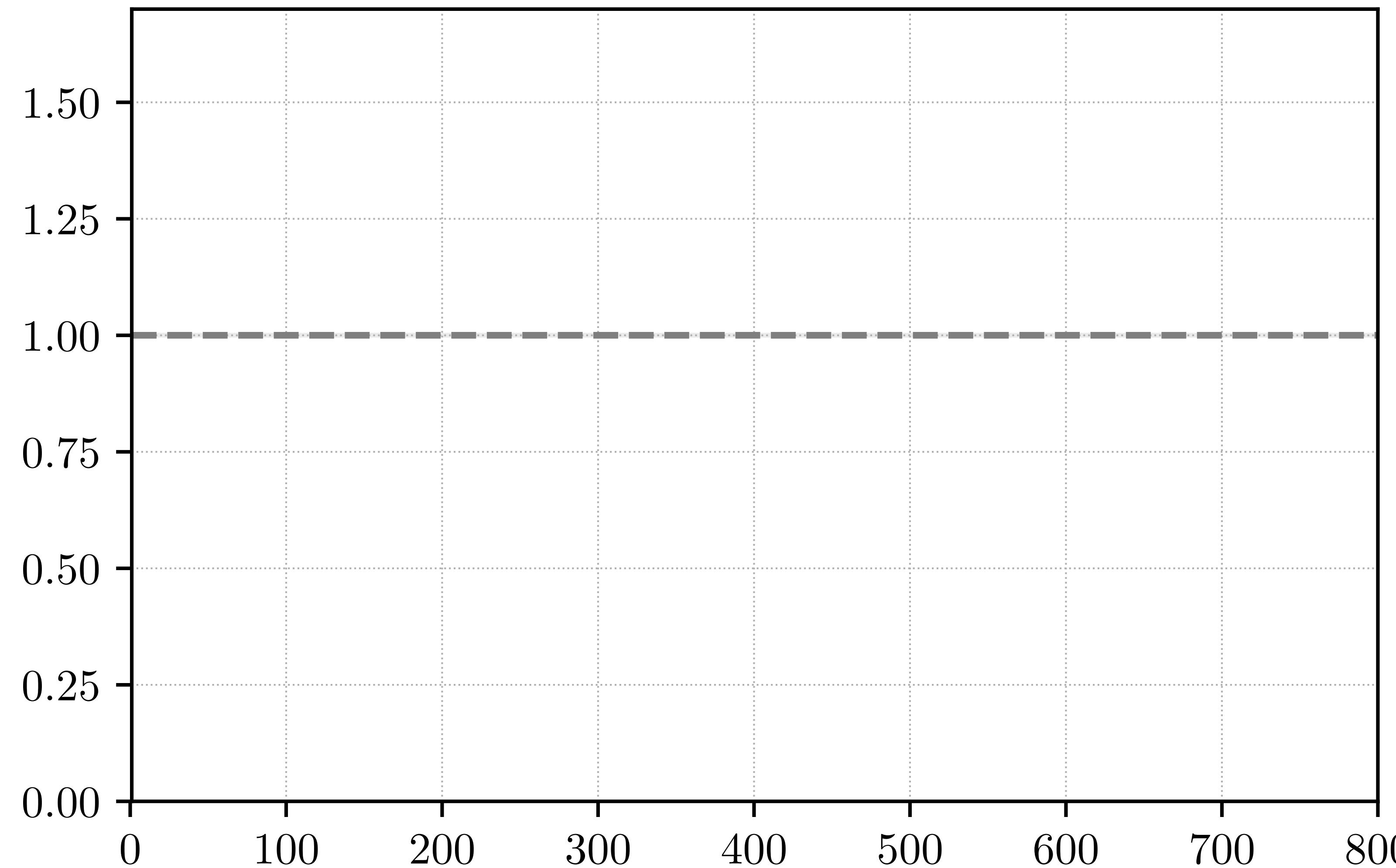
	Task Invariant	Time Cost	Predictive Accuracy
Vanilla Model	No	Low	Medium
Tree-based Model	Yes	Low	Good
Neural Model	Yes	High	Good

# Comparisons of Models

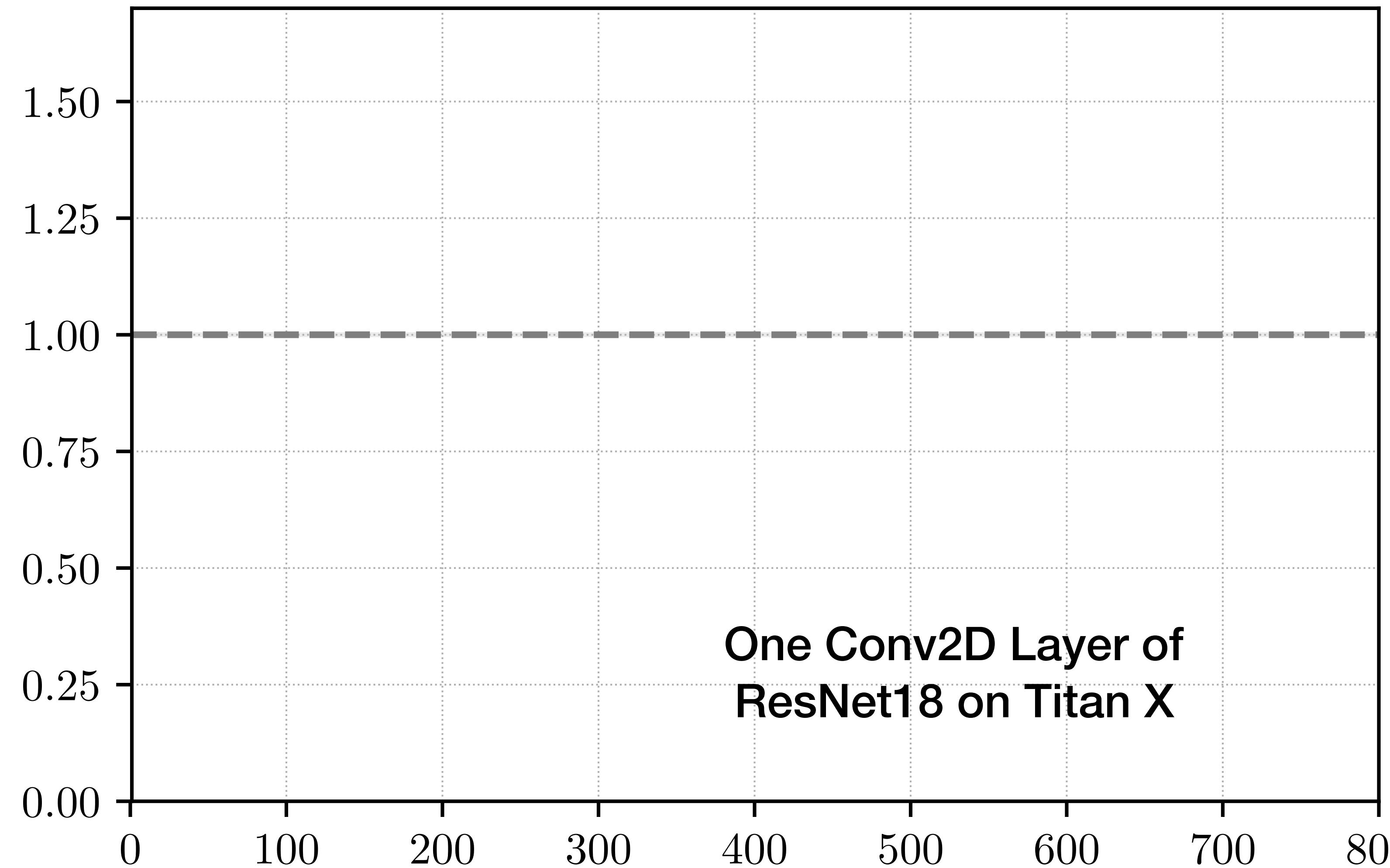
	Task Invariant	Time Cost	Predictive Accuracy
Vanilla Model	No	Low	Medium
Tree-based Model	Yes	Low	Good
Neural Model	Yes	High	Good

Choose tree-based model by default

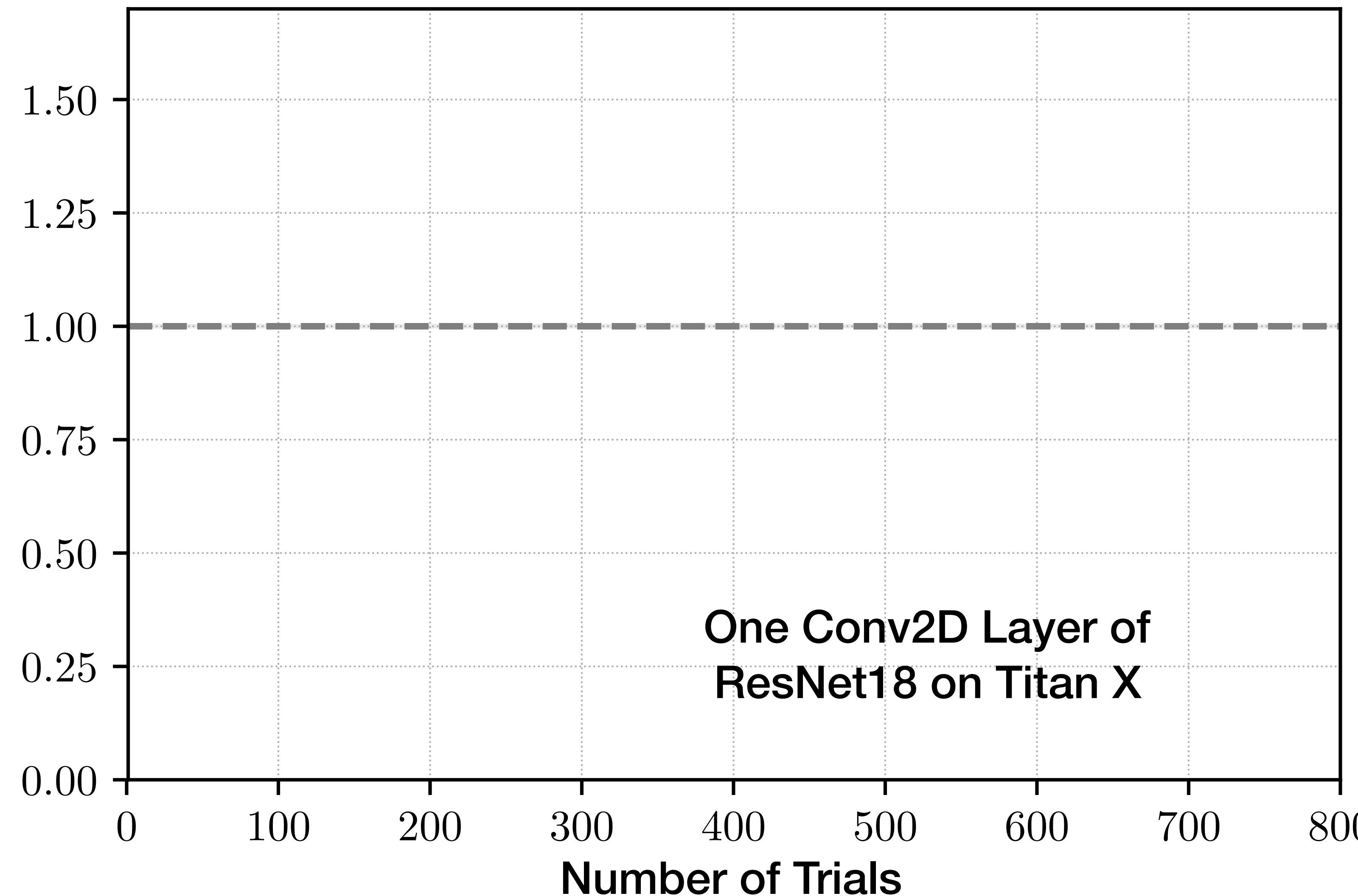
# Effectiveness of ML based Model



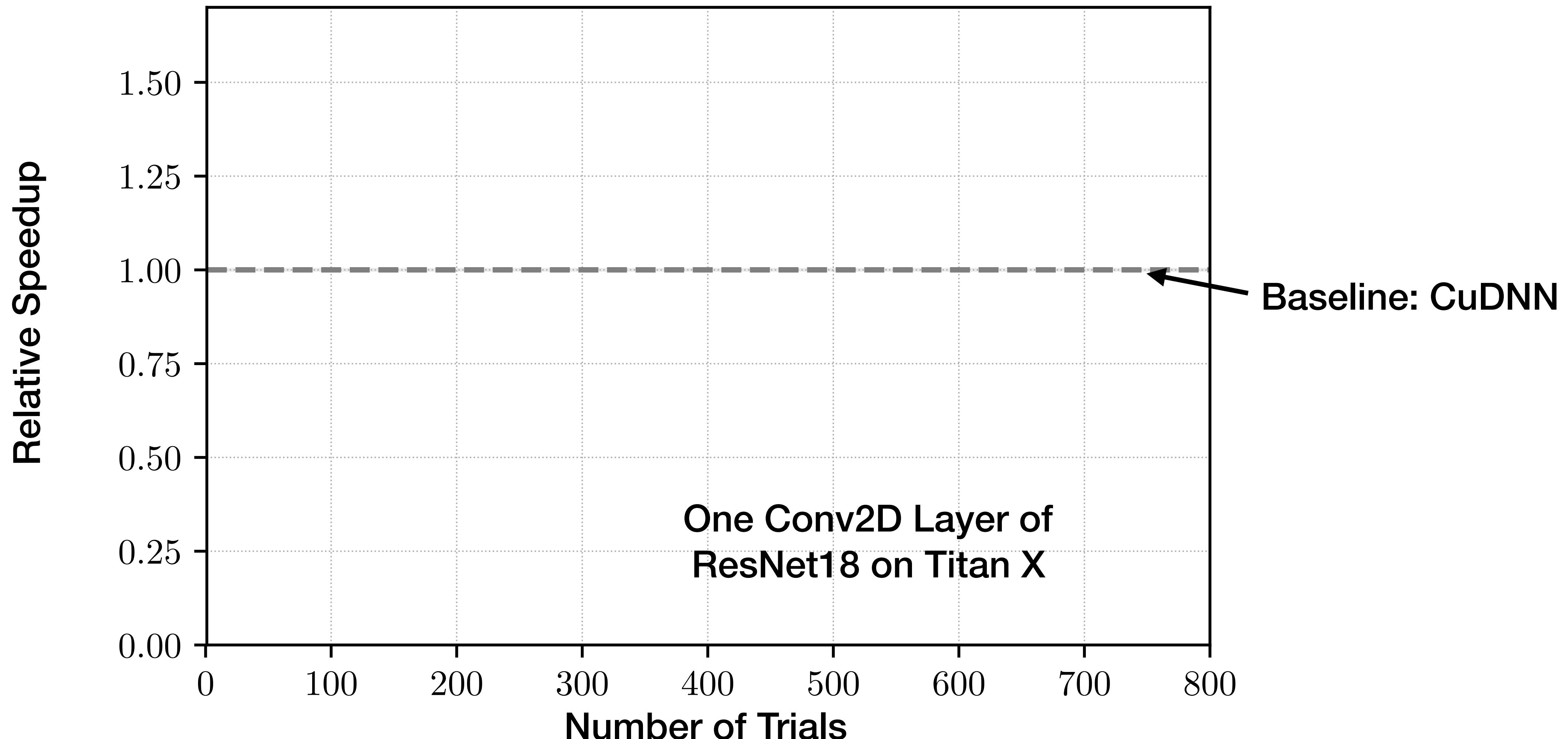
# Effectiveness of ML based Model



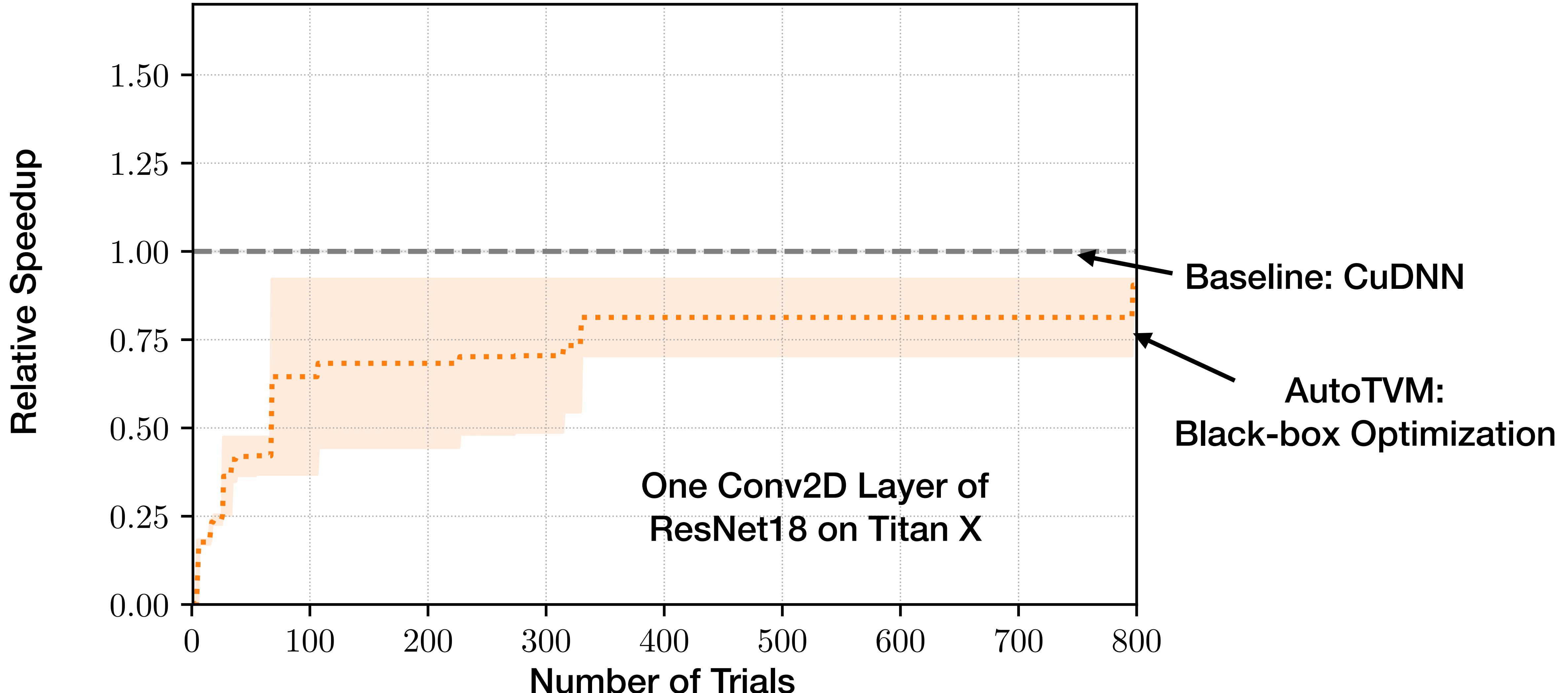
# Effectiveness of ML based Model



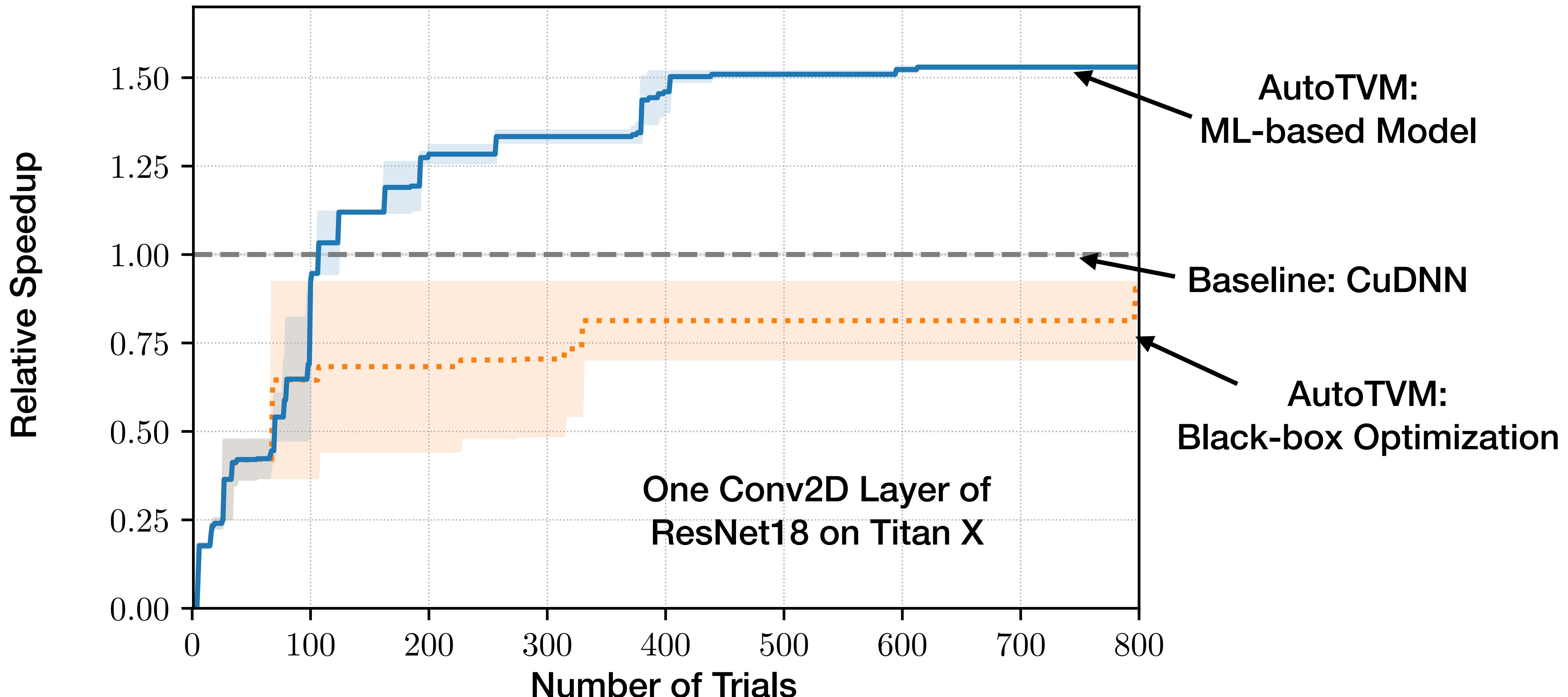
# Effectiveness of ML based Model



# Effectiveness of ML based Model



# Effectiveness of ML based Model



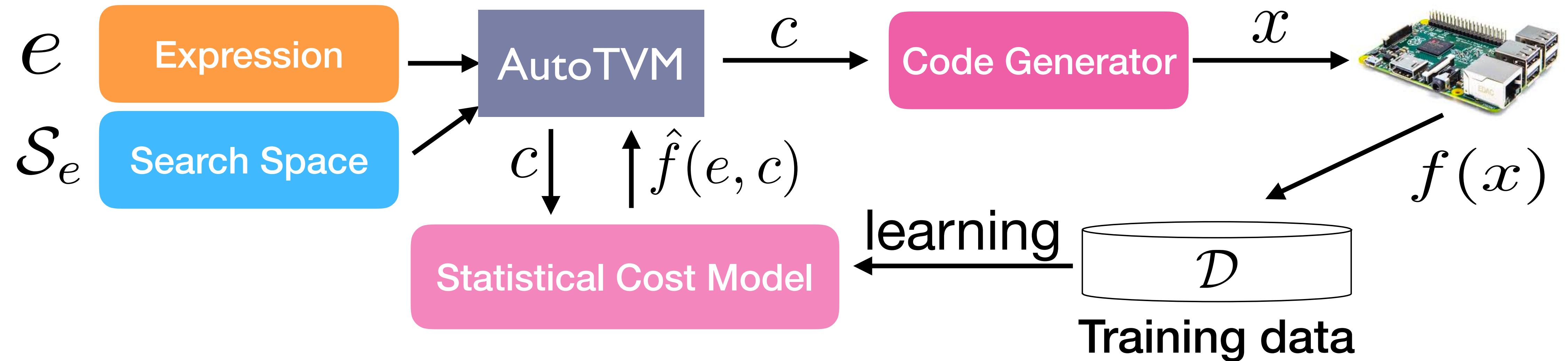
# Comparisons of Models

	Task Invariant	Time Cost	Predictive Accuracy
Vanilla Model	No	Low	Medium
Tree-based Model	Yes	Low	Good
Neural Model	Yes	High	Good

# Comparisons of Models

	Task Invariant	Time Cost	Predictive Accuracy
Vanilla Model	No	Low	Medium
Tree-based Model	Yes	Low	Good
Neural Model	Yes	High	Good

# Unique Problem Characteristics

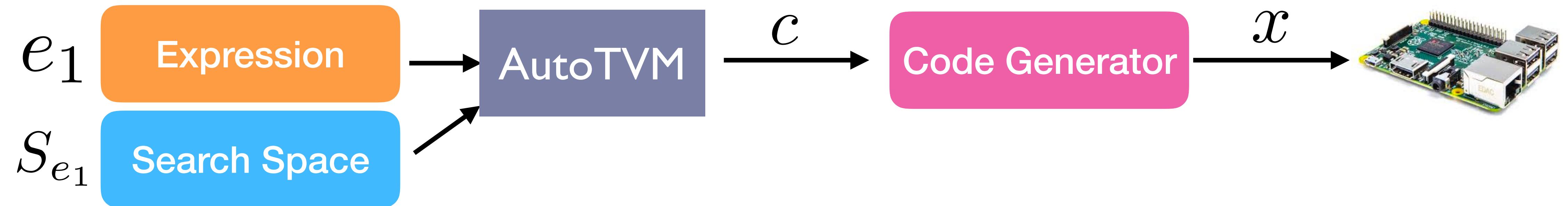


Relatively low  
experiment cost

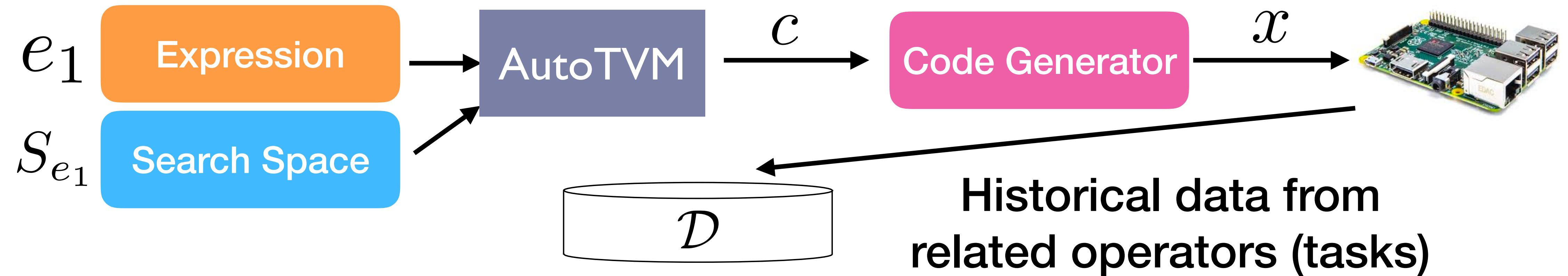
Program-aware  
modeling

**Large number of  
similar tasks**

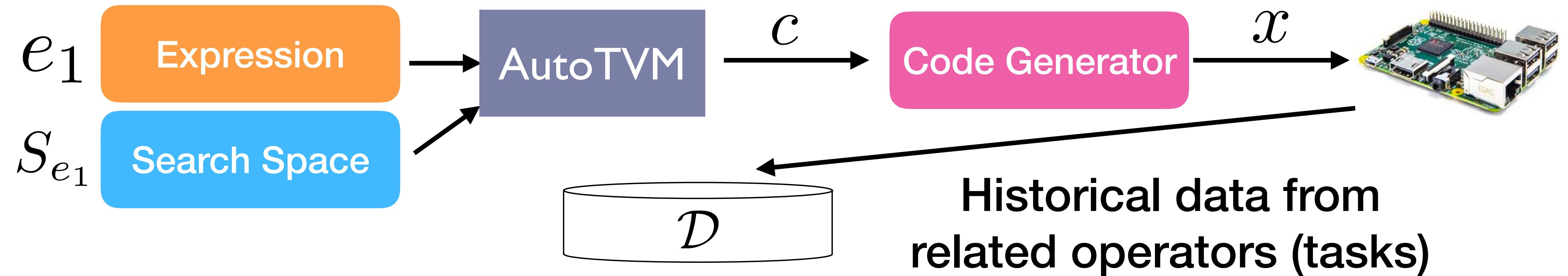
# Transferable Cost Model



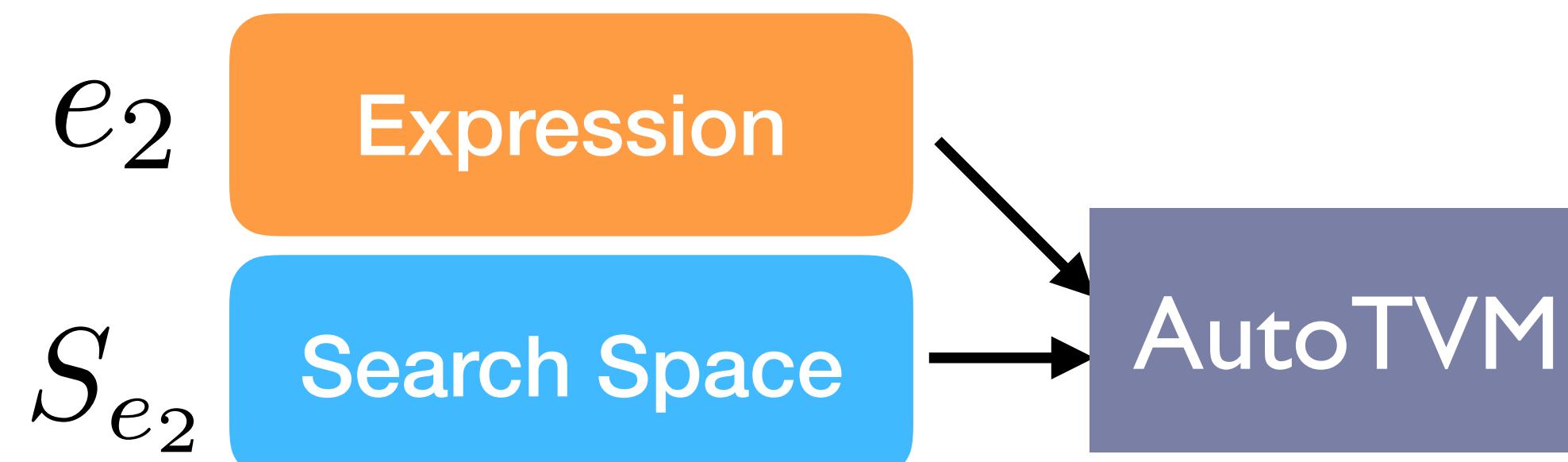
# Transferable Cost Model



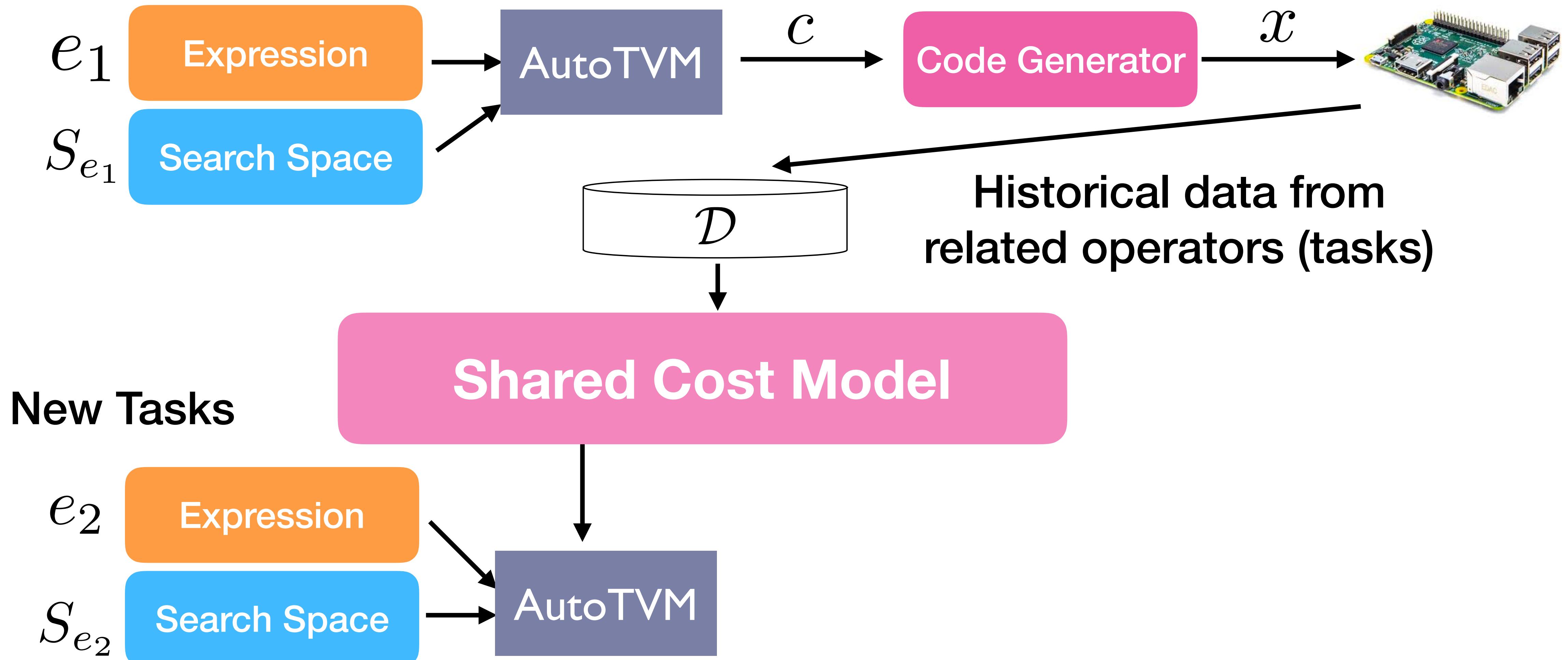
# Transferable Cost Model



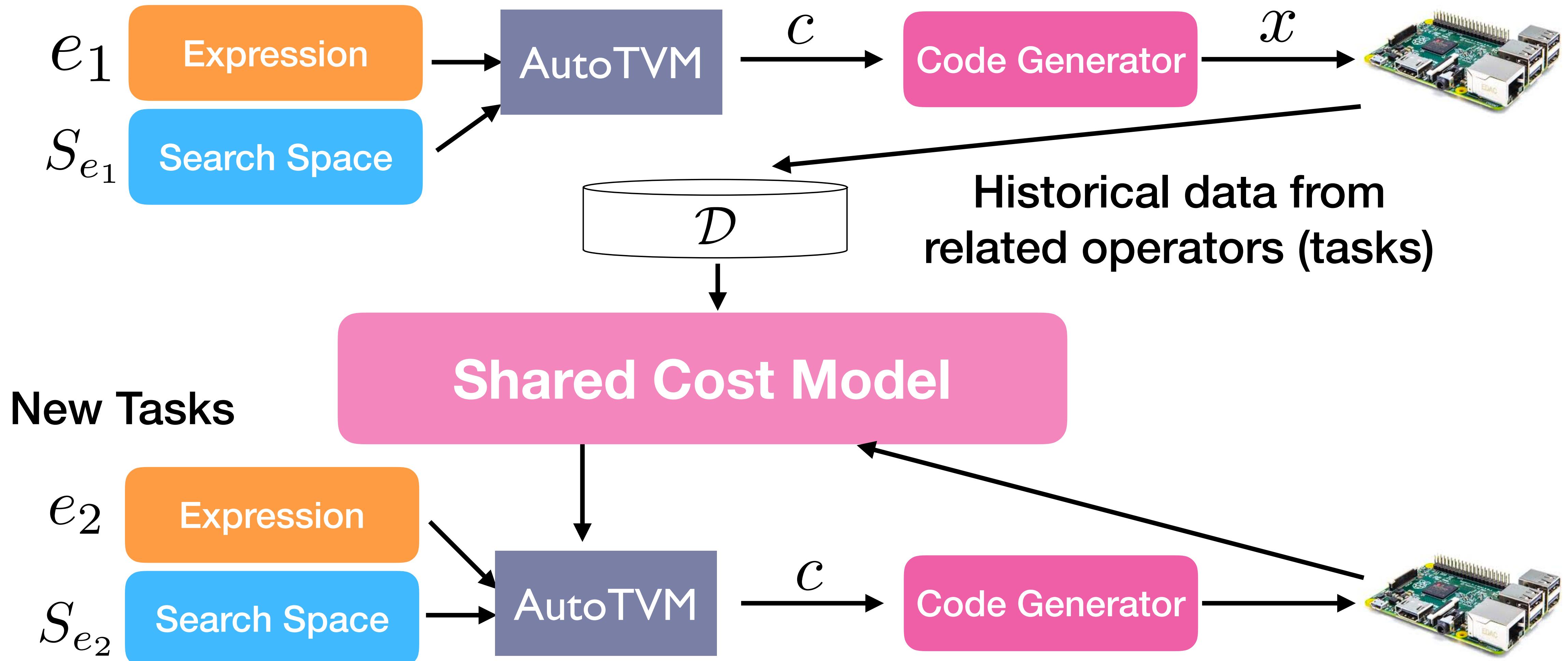
## New Tasks



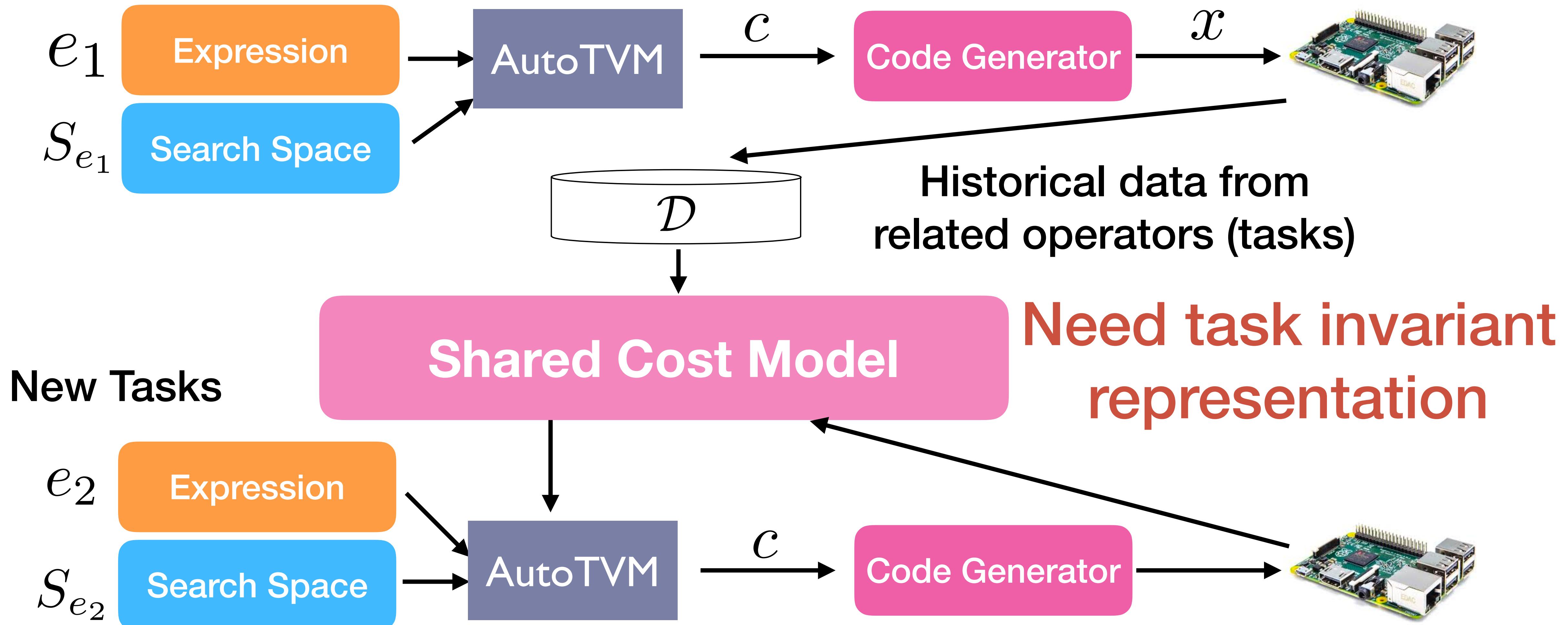
# Transferable Cost Model



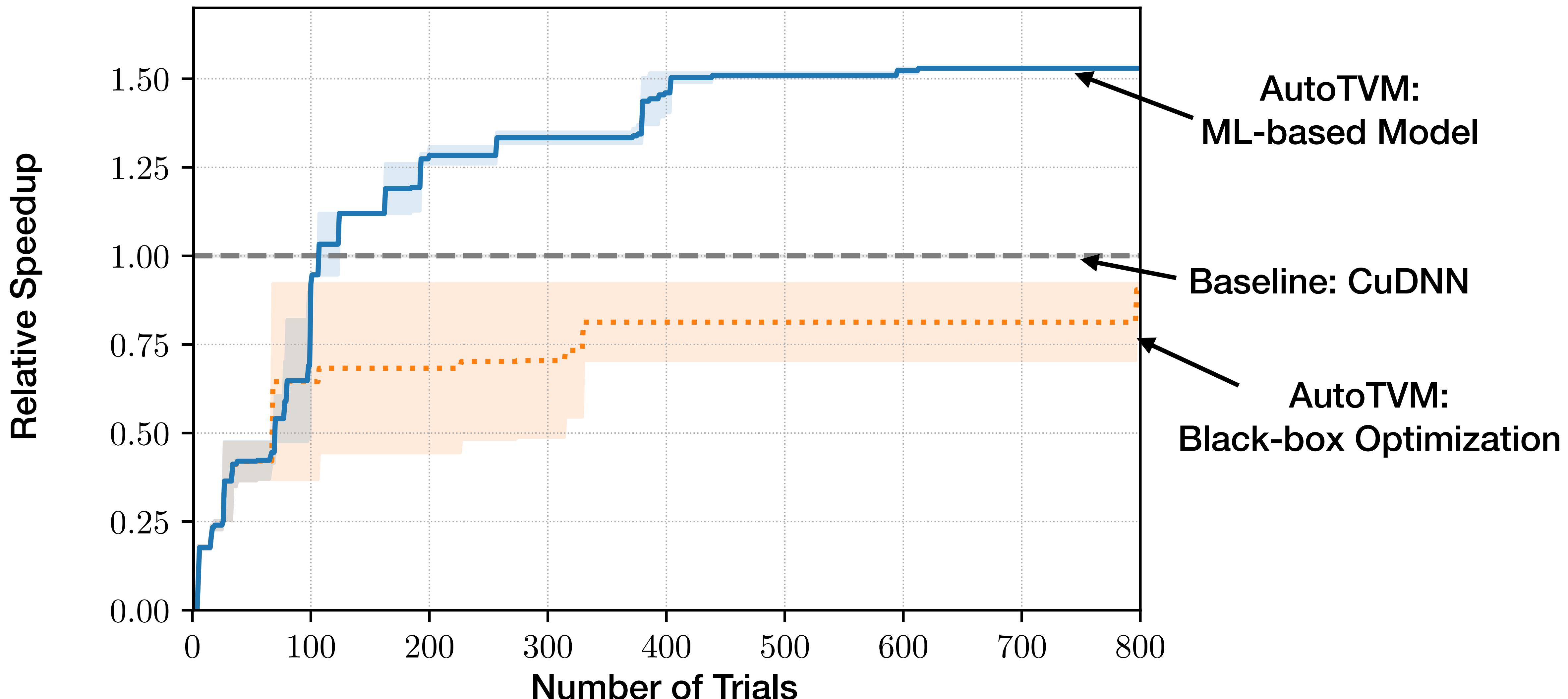
# Transferable Cost Model



# Transferable Cost Model



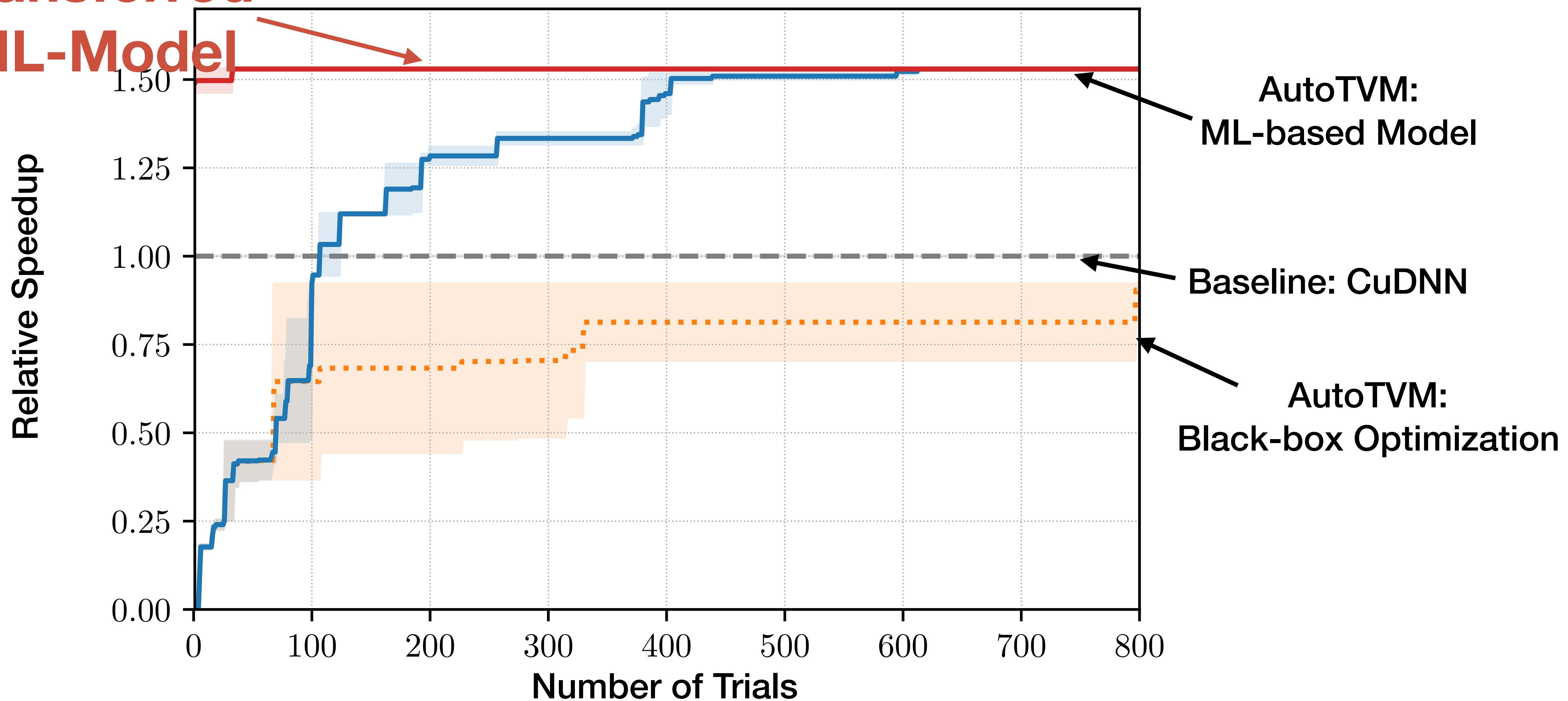
# Impact of Transfer Learning



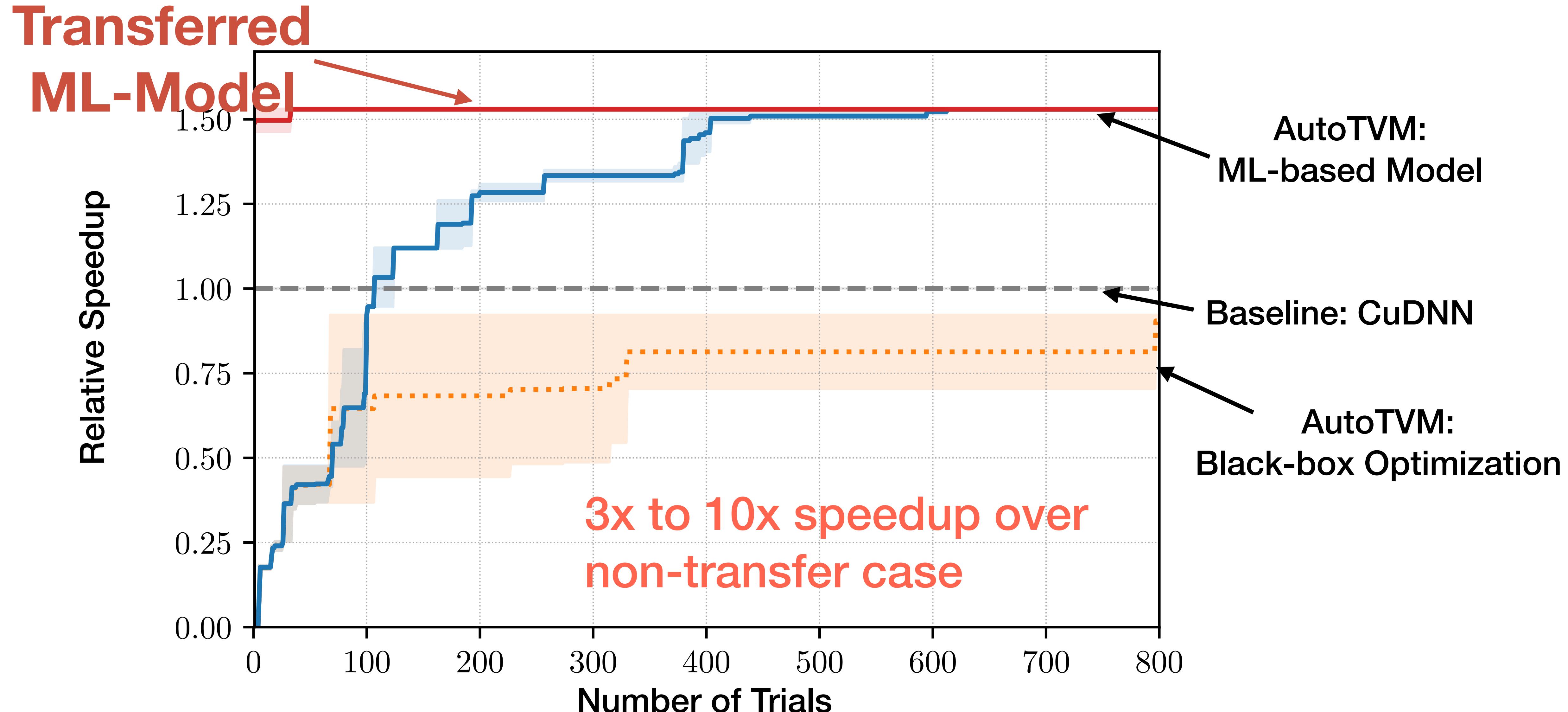
# Impact of Transfer Learning

Transferred

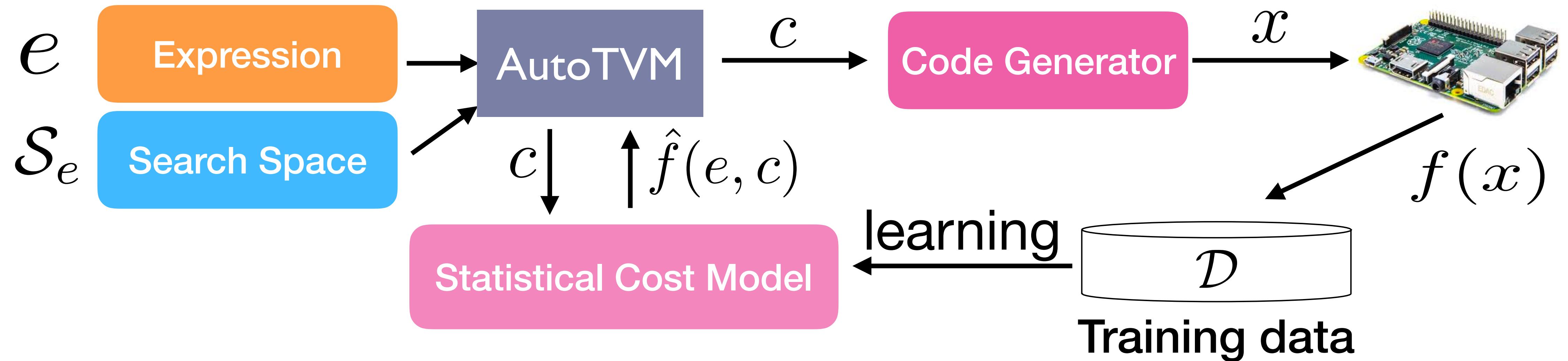
ML-Model



# Impact of Transfer Learning



# Learning to Optimize Tensor Programs

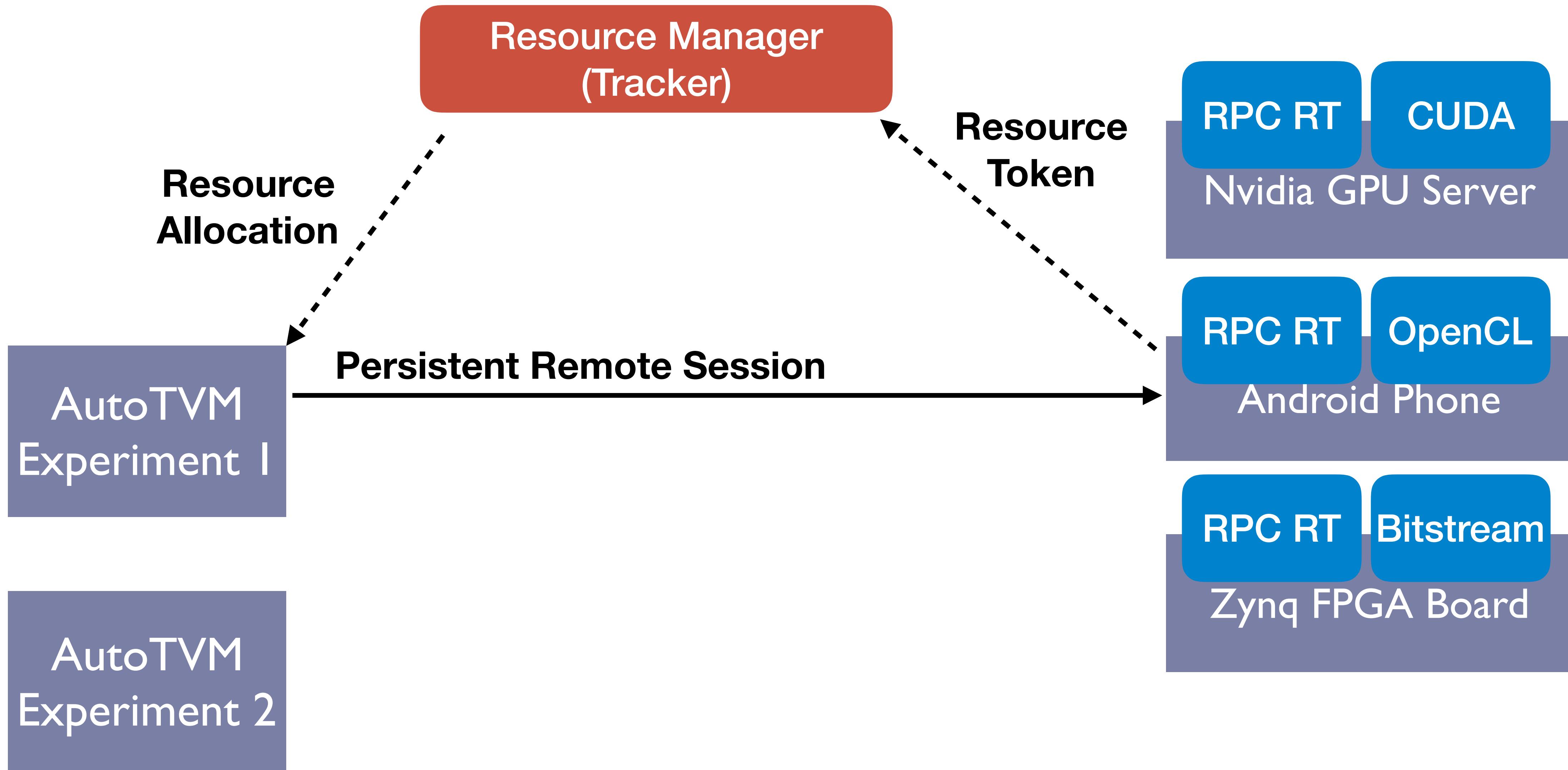


**Relatively low experiment cost**

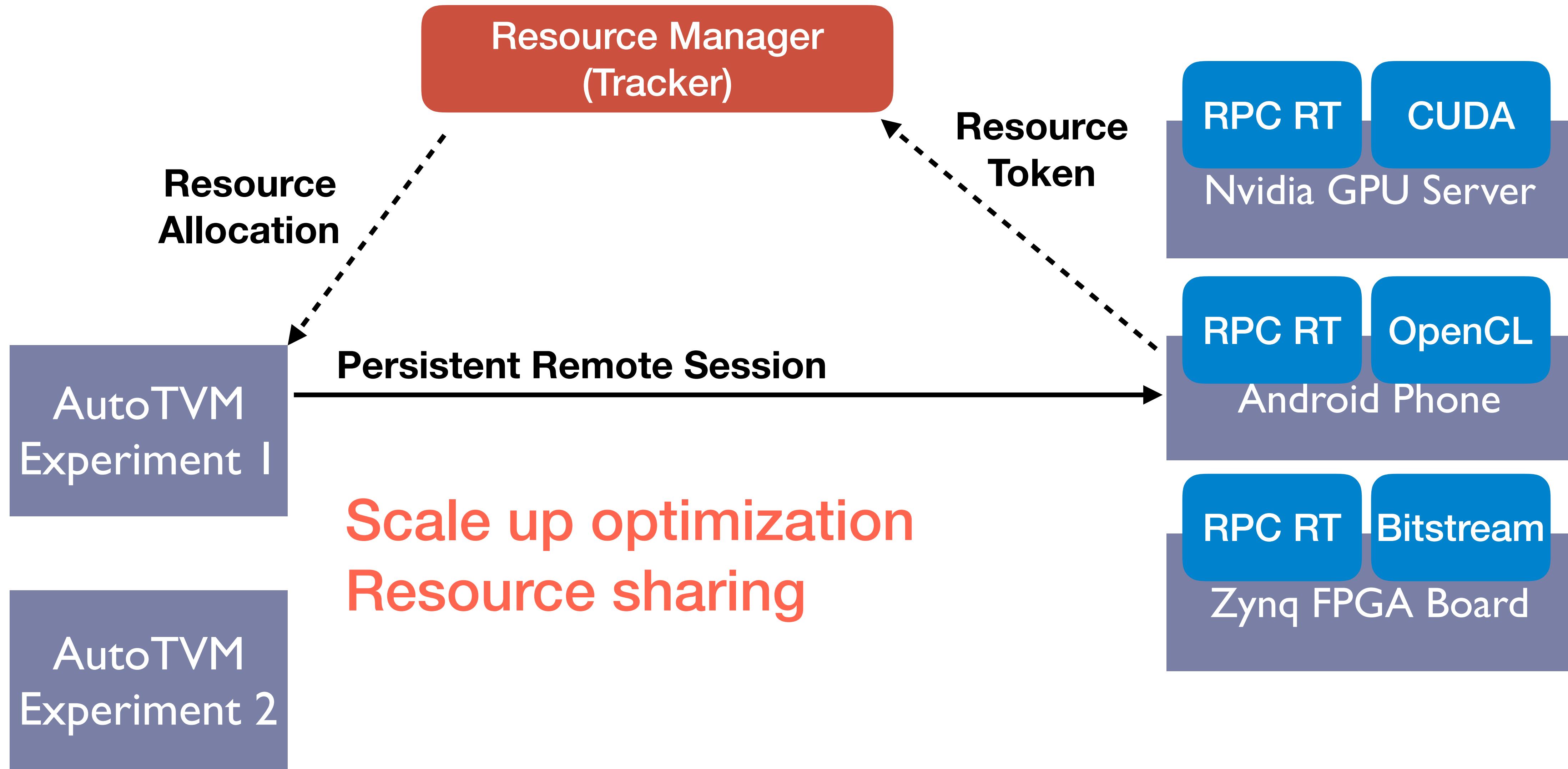
**Program-aware modeling**

**Large number of similar tasks**

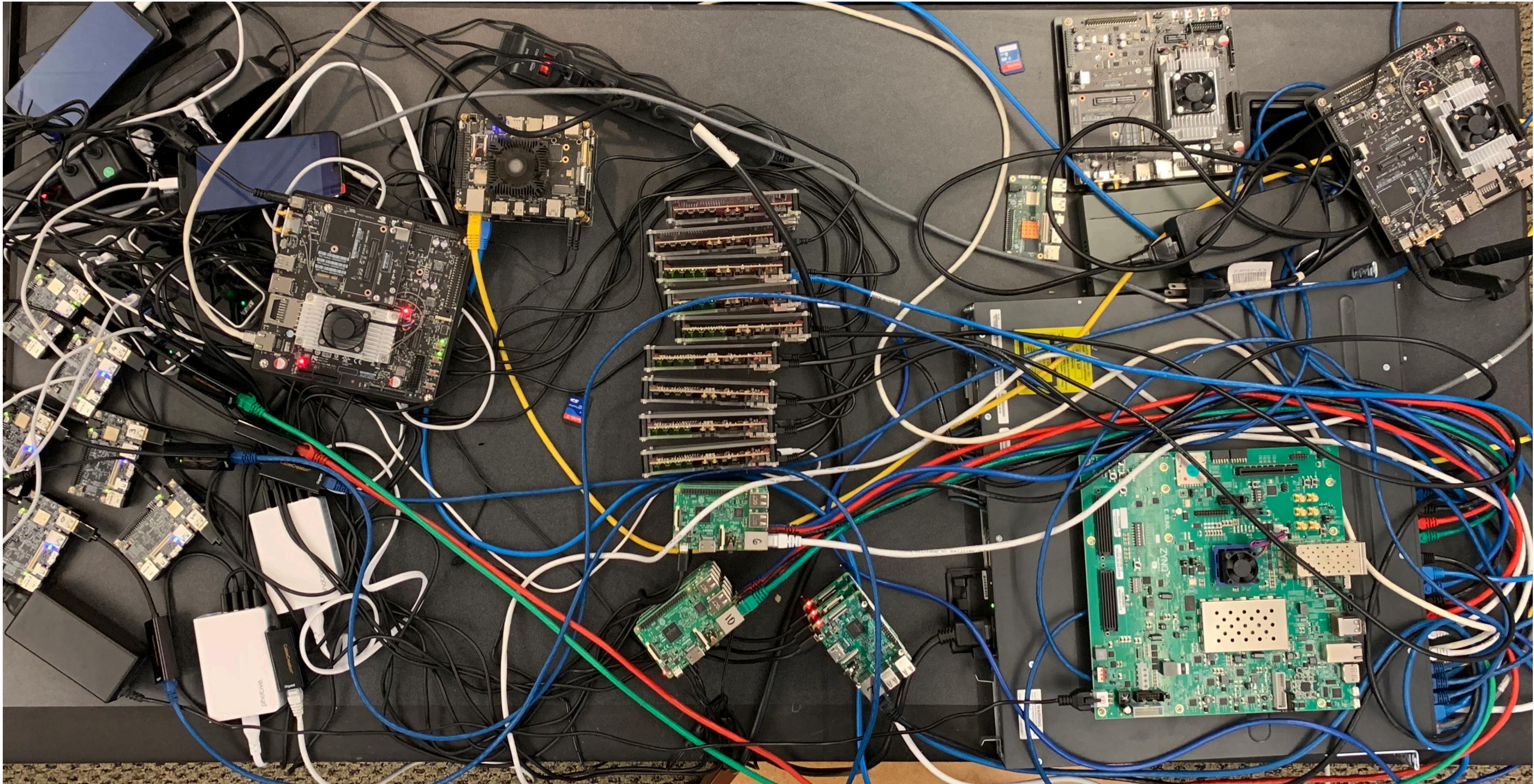
# Device Fleet: Distributed Test Bed for AutoTVM



# Device Fleet: Distributed Test Bed for AutoTVM



# Device Fleet in Action



# TVM: Learning-based Learning System

Why do we need machine learning for systems

How to build intelligent systems with learning

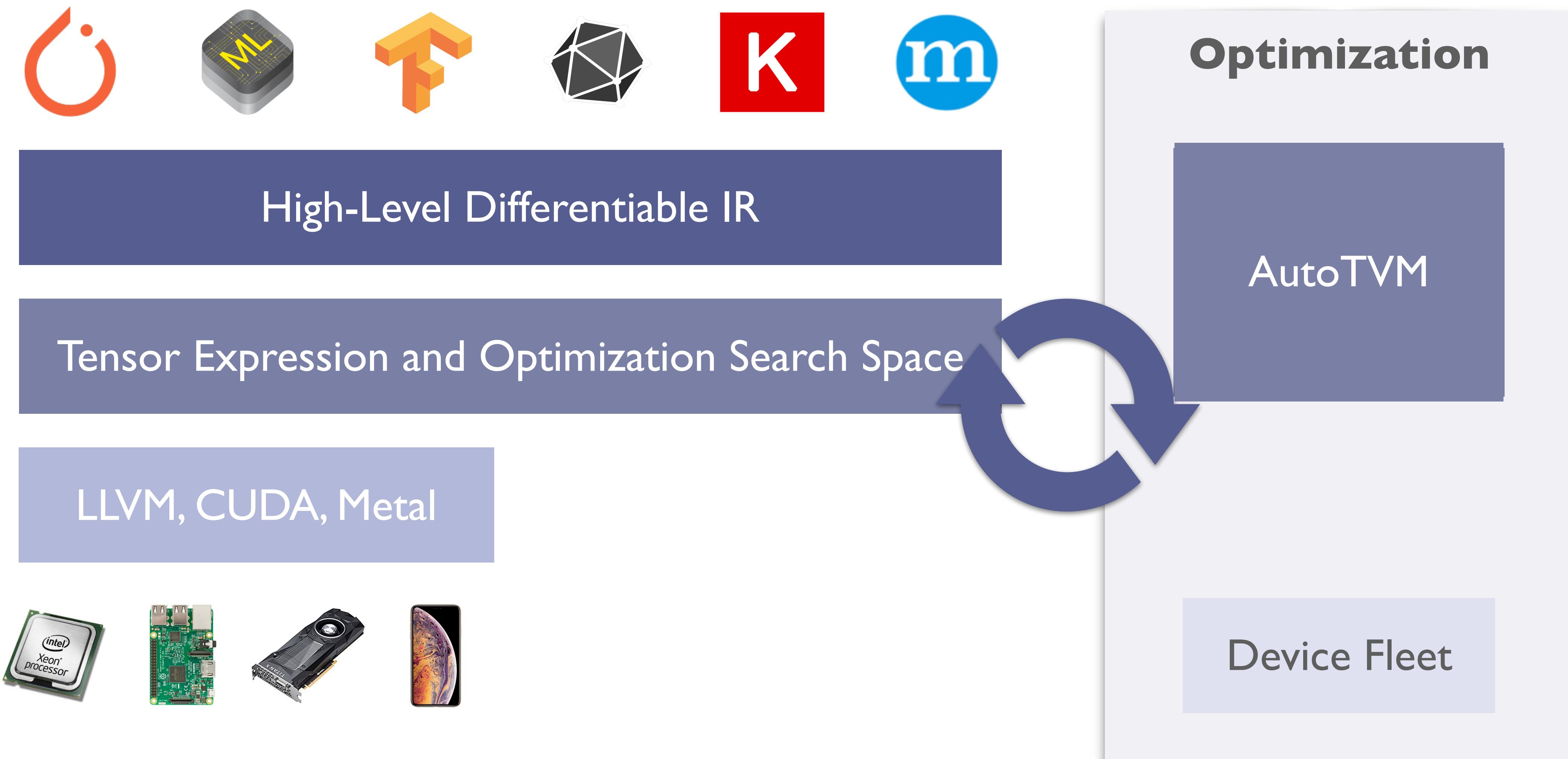
**End to end learning-based learning system stack**

# TVM: End to End Deep Learning Compiler

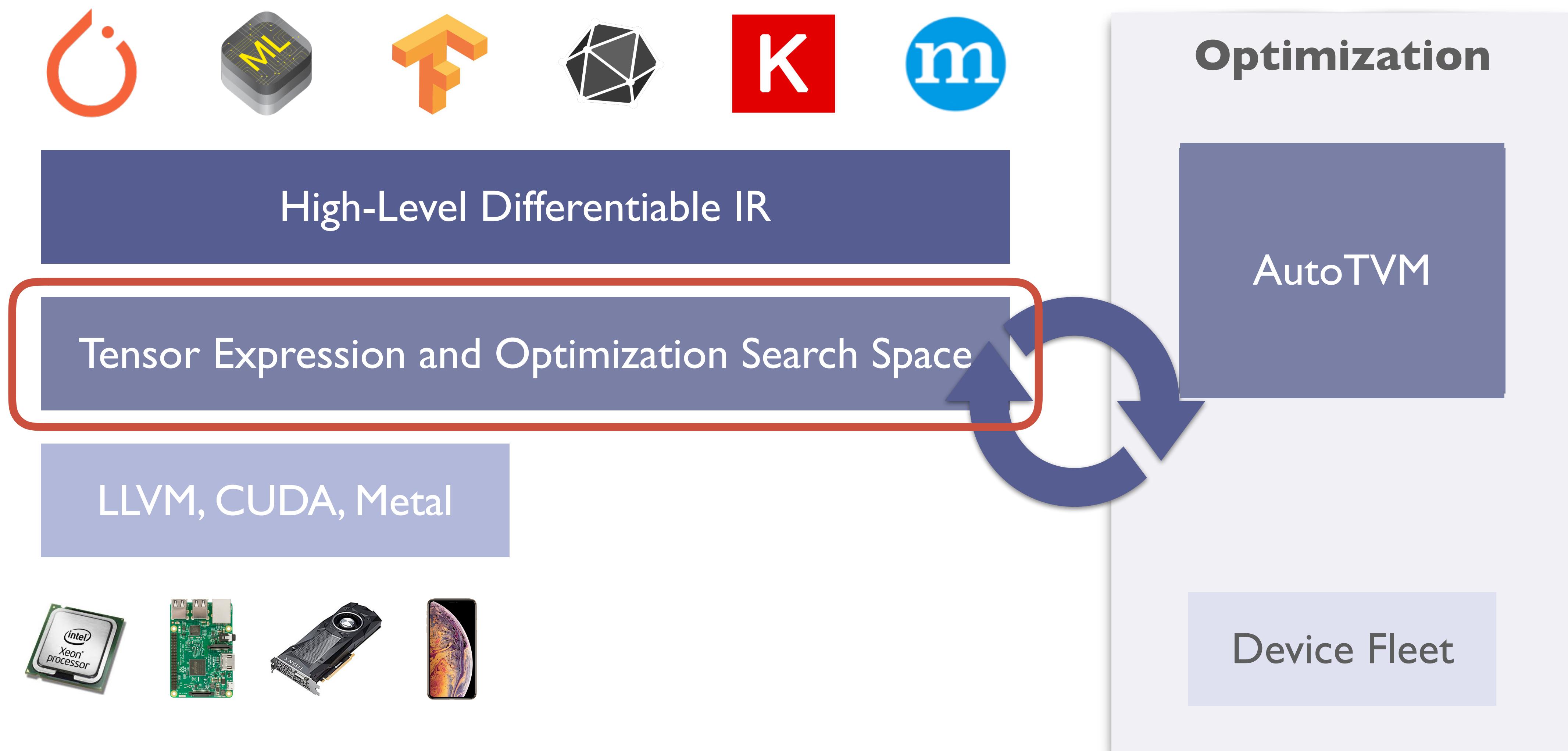


AutoTVM

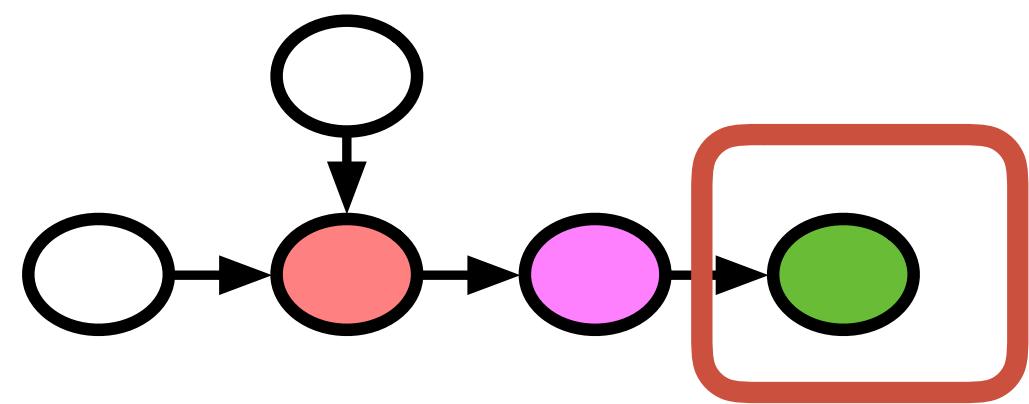
# TVM: End to End Deep Learning Compiler



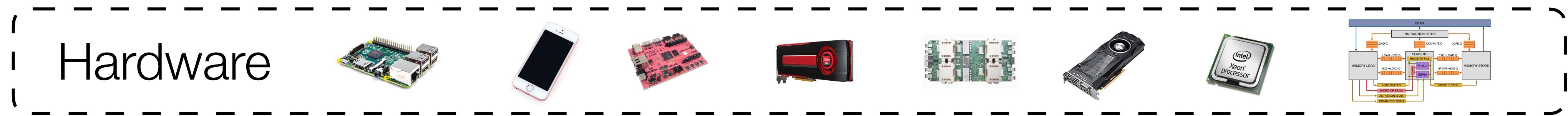
# TVM: End to End Deep Learning Compiler



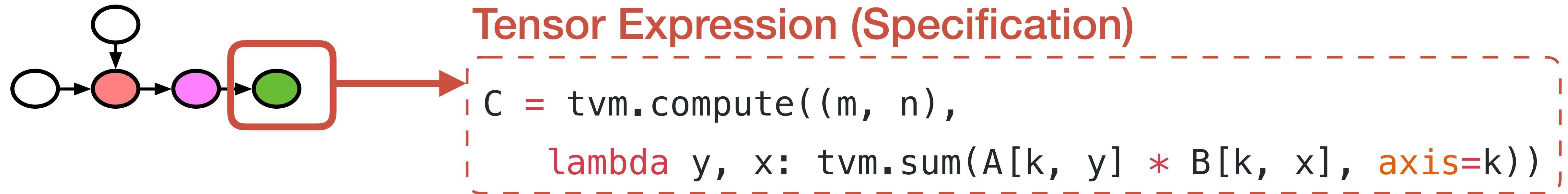
# Tensor Expression and Optimization Search Space



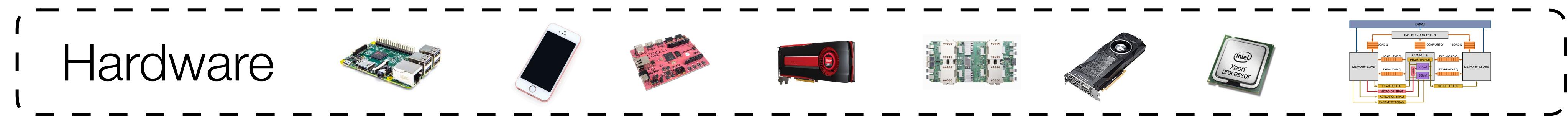
Based on Halide's  
compute/schedule  
separation



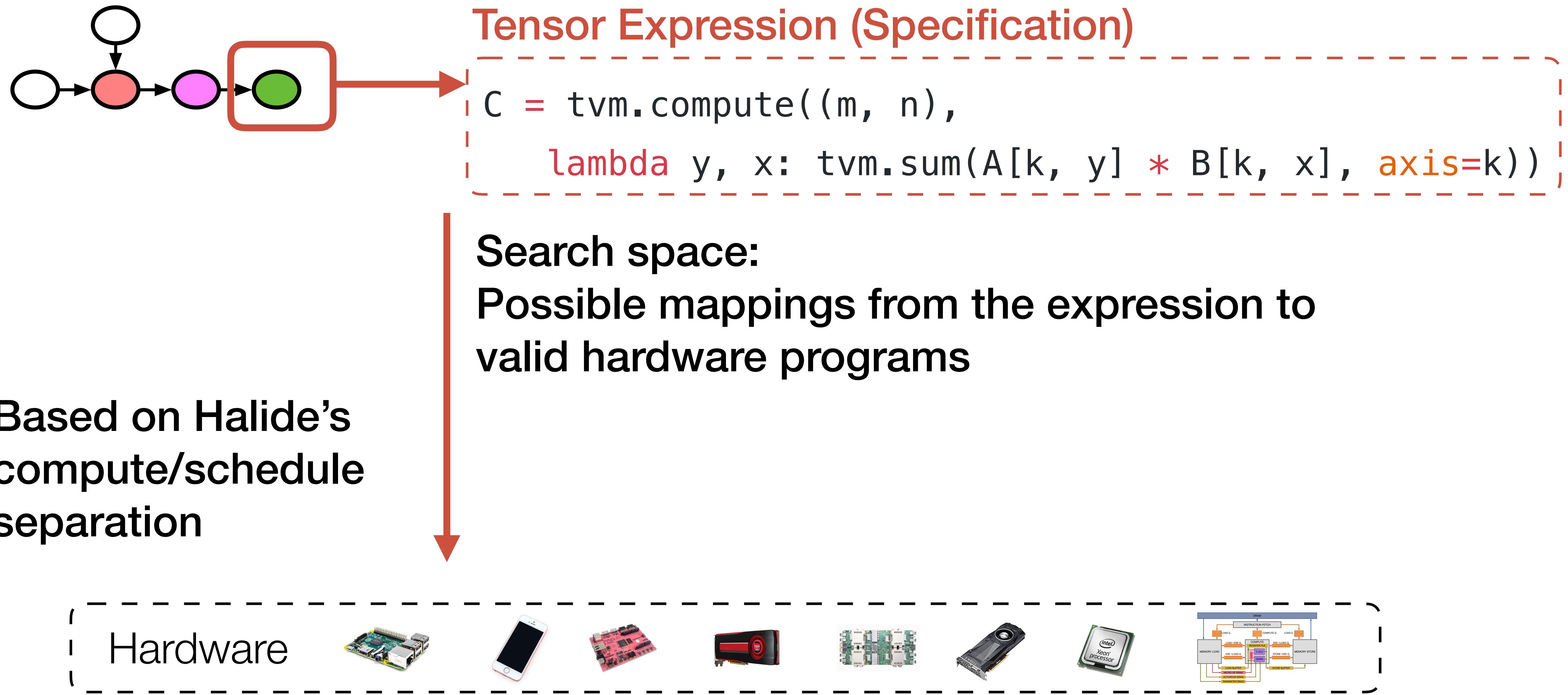
# Tensor Expression and Optimization Search Space



Based on Halide's  
compute/schedule  
separation



# Tensor Expression and Optimization Search Space



# Tensor Expression and Optimization Search Space

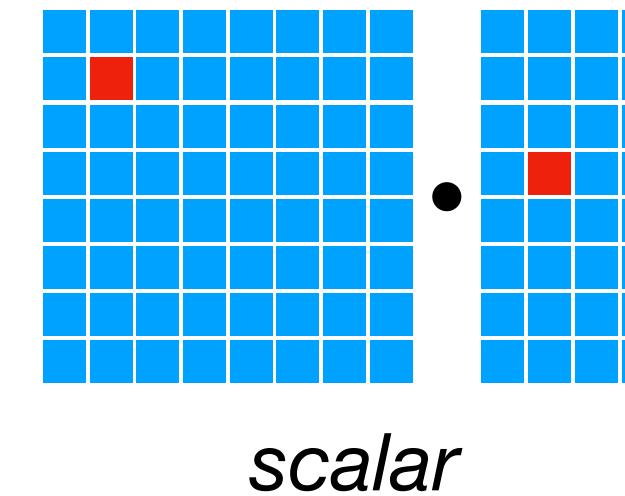


# Search Space for CPUs

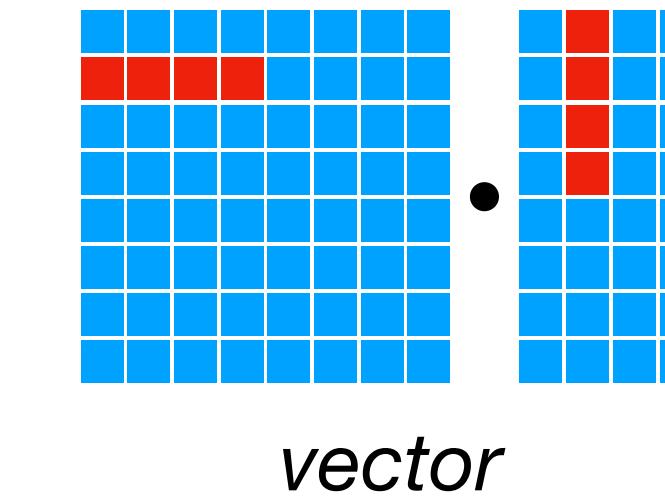
CPUs



Compute Primitives

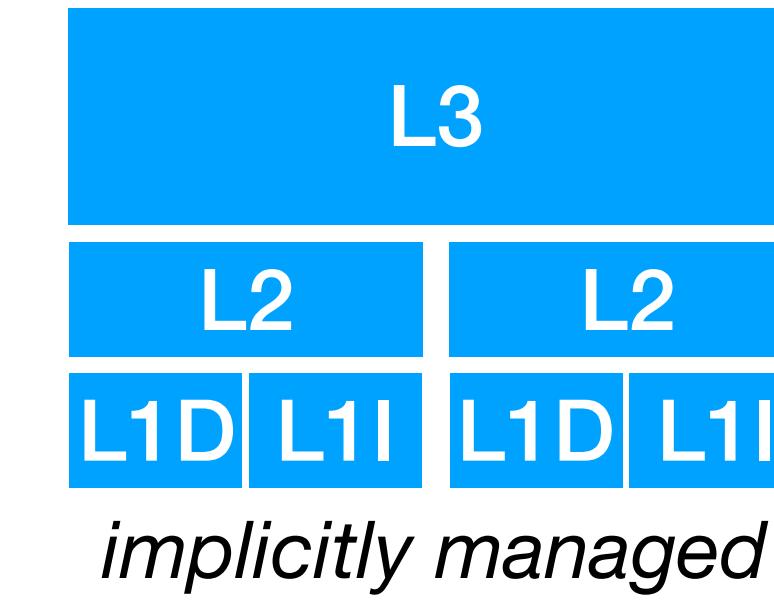


scalar



vector

Memory Subsystem



implicitly managed

Loop  
Transformations

Cache  
Locality

Vectorization

Reuse primitives from prior work:  
Halide, Loopy

# Hardware-aware Search Space

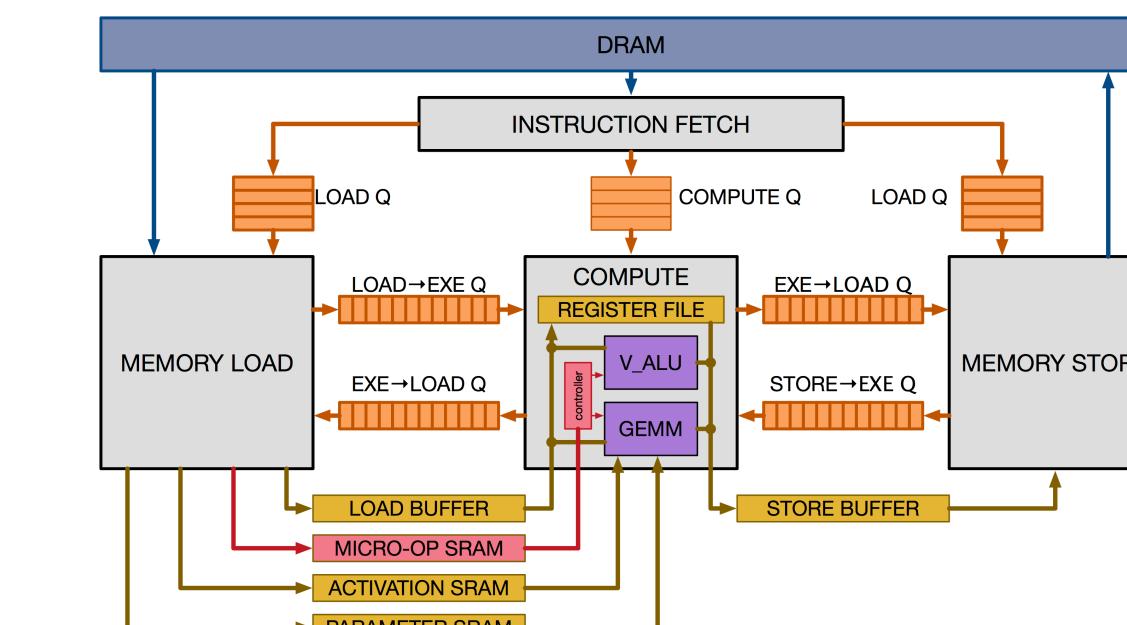
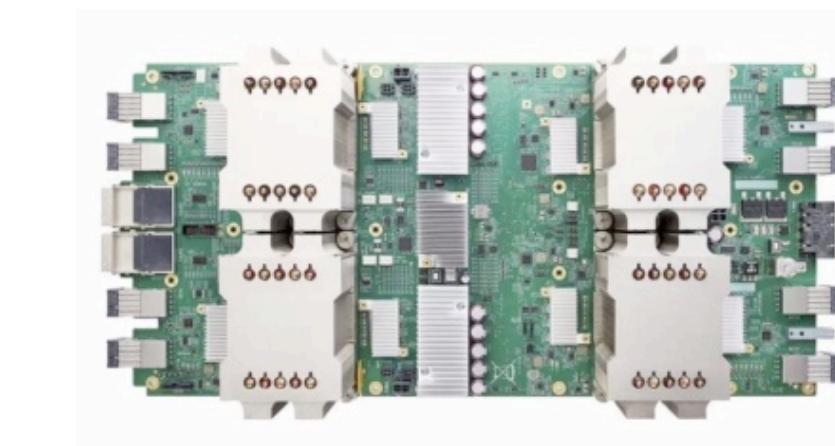
CPUs



GPUs



TPU-like specialized  
Accelerators

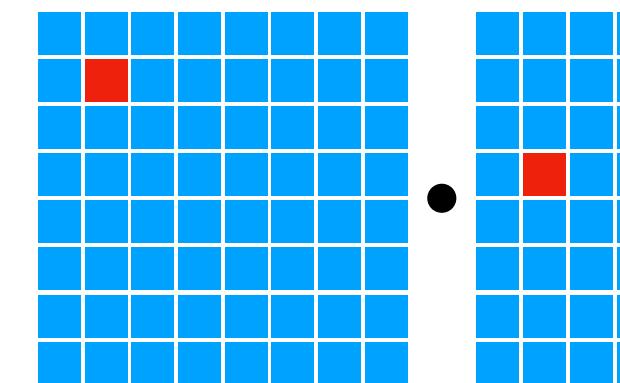


# Search Space for GPUs

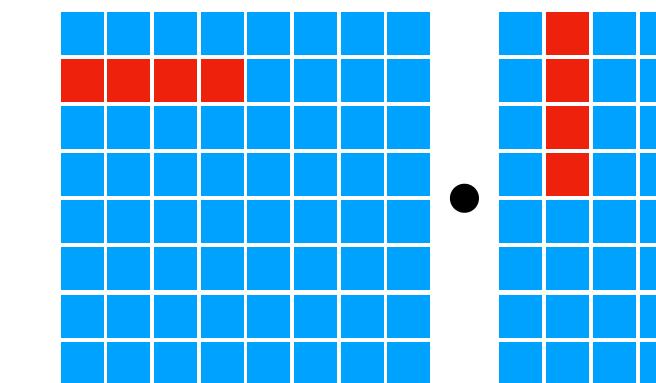
GPUs



Compute Primitives



*scalar*



*vector*

Memory Subsystem



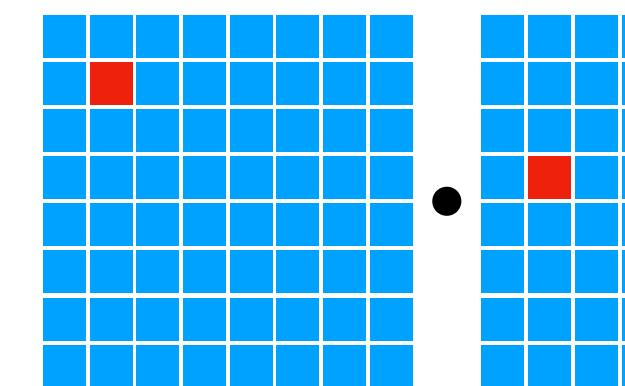
*mixed*

# Search Space for GPUs

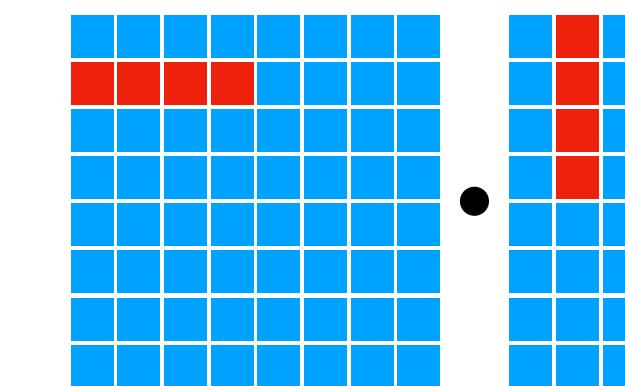
GPUs



Compute Primitives

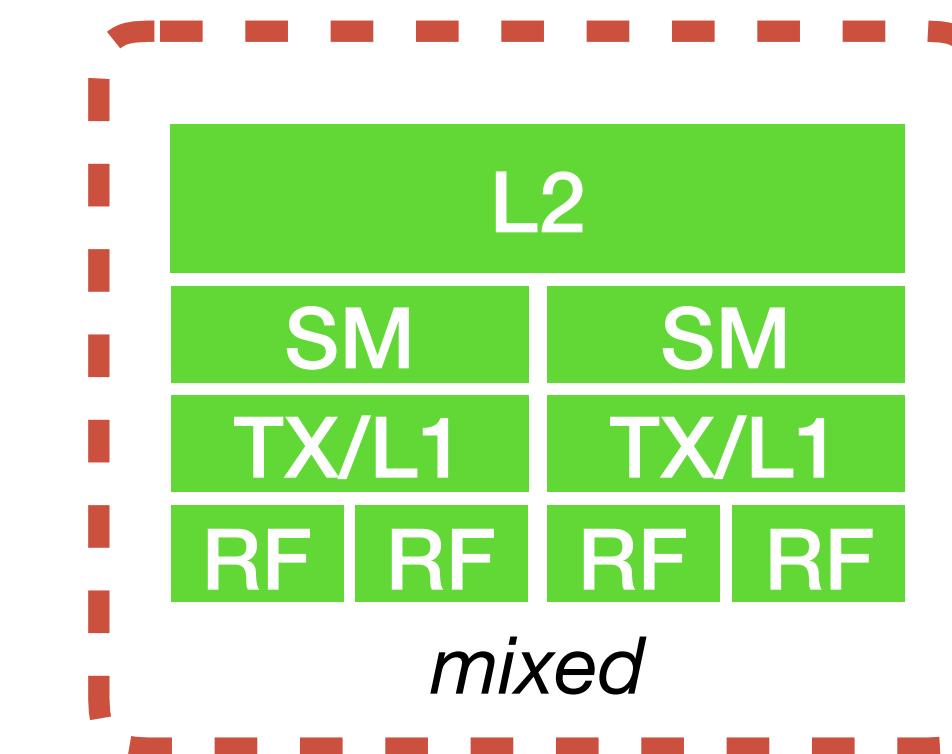


scalar



vector

Memory Subsystem



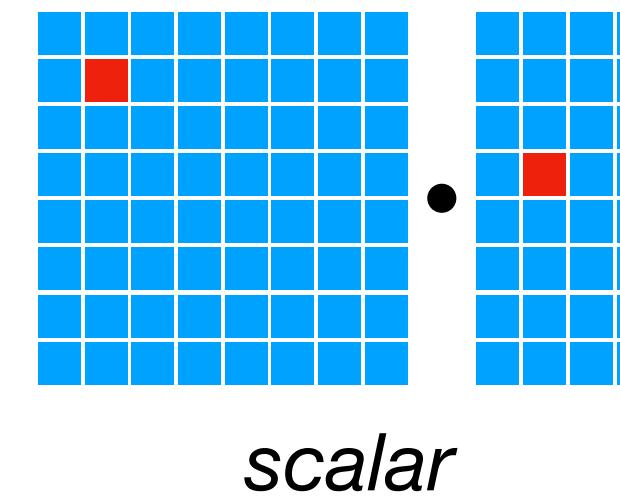
Shared memory among  
compute cores

# Search Space for GPUs

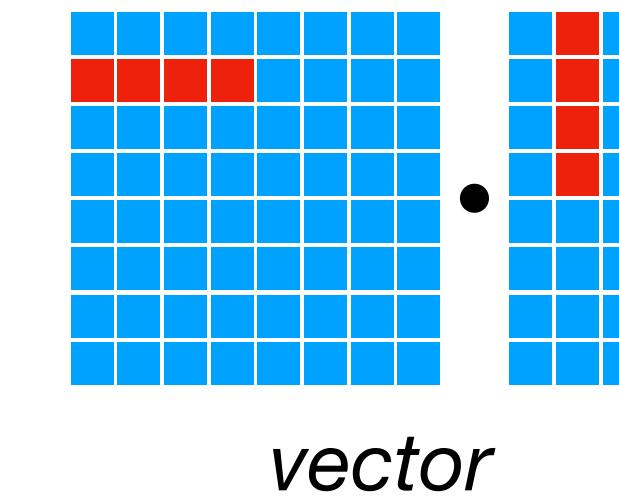
GPUs



Compute Primitives

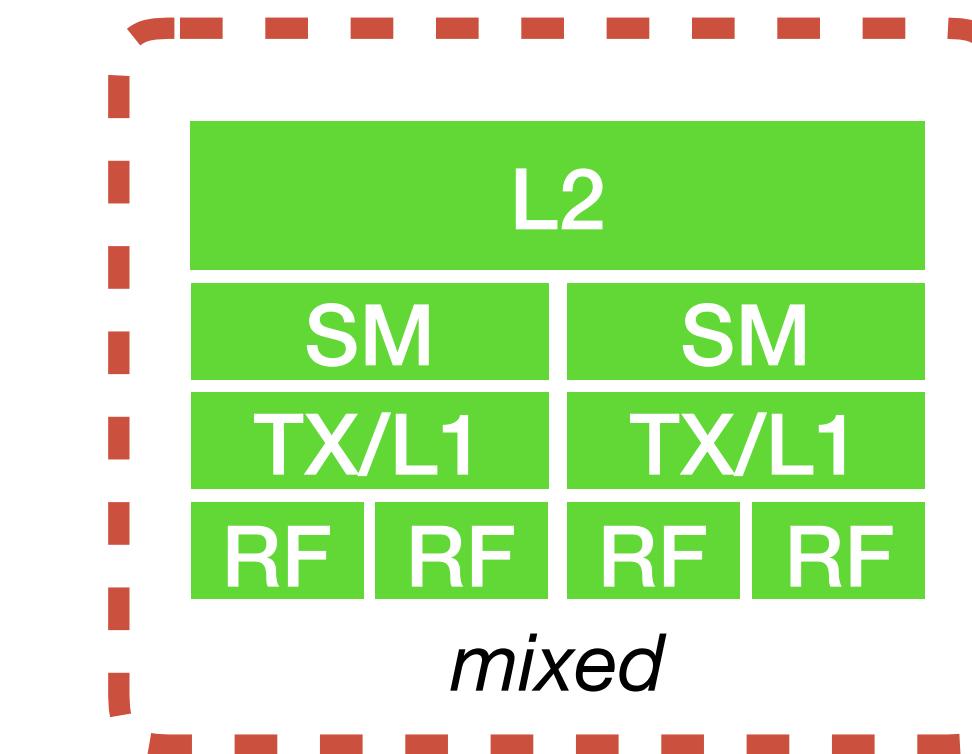


scalar



vector

Memory Subsystem



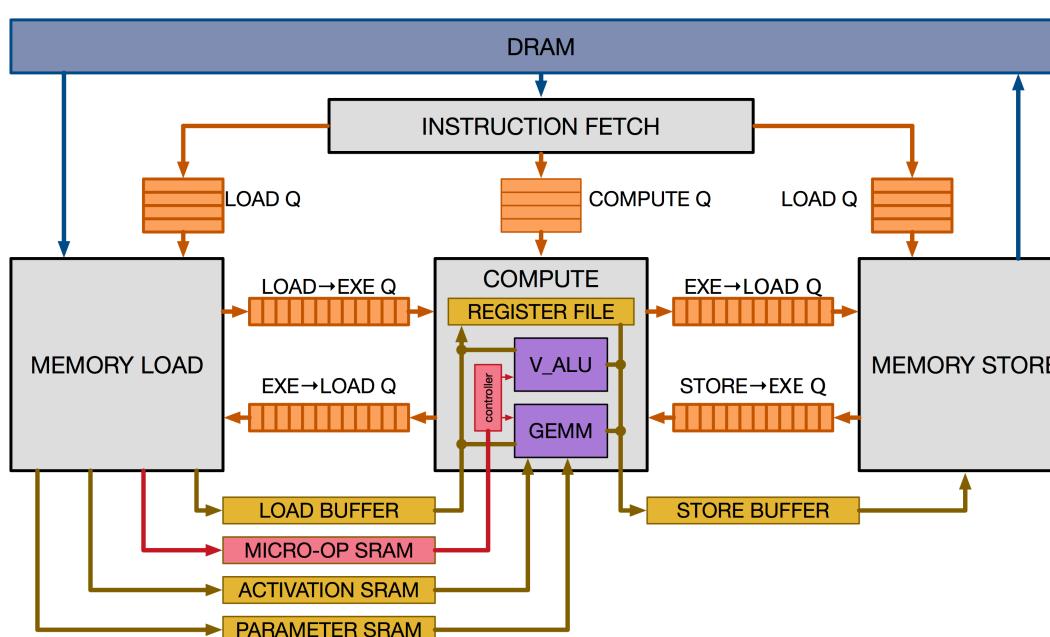
Shared memory among  
compute cores

Use of Shared  
Memory

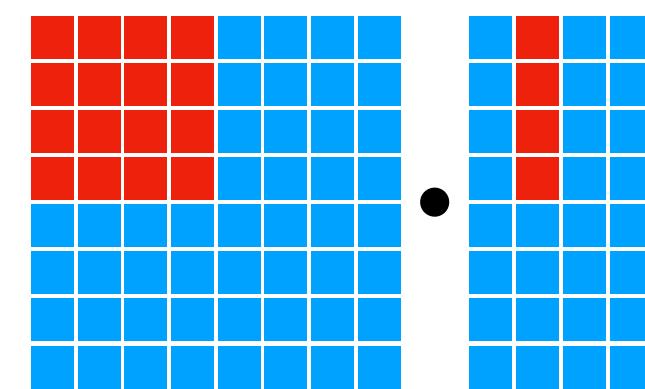
Thread  
Cooperation

# Search Space for TPU-like Specialized Accelerators

## TPUs



## Tensor Compute Primitives

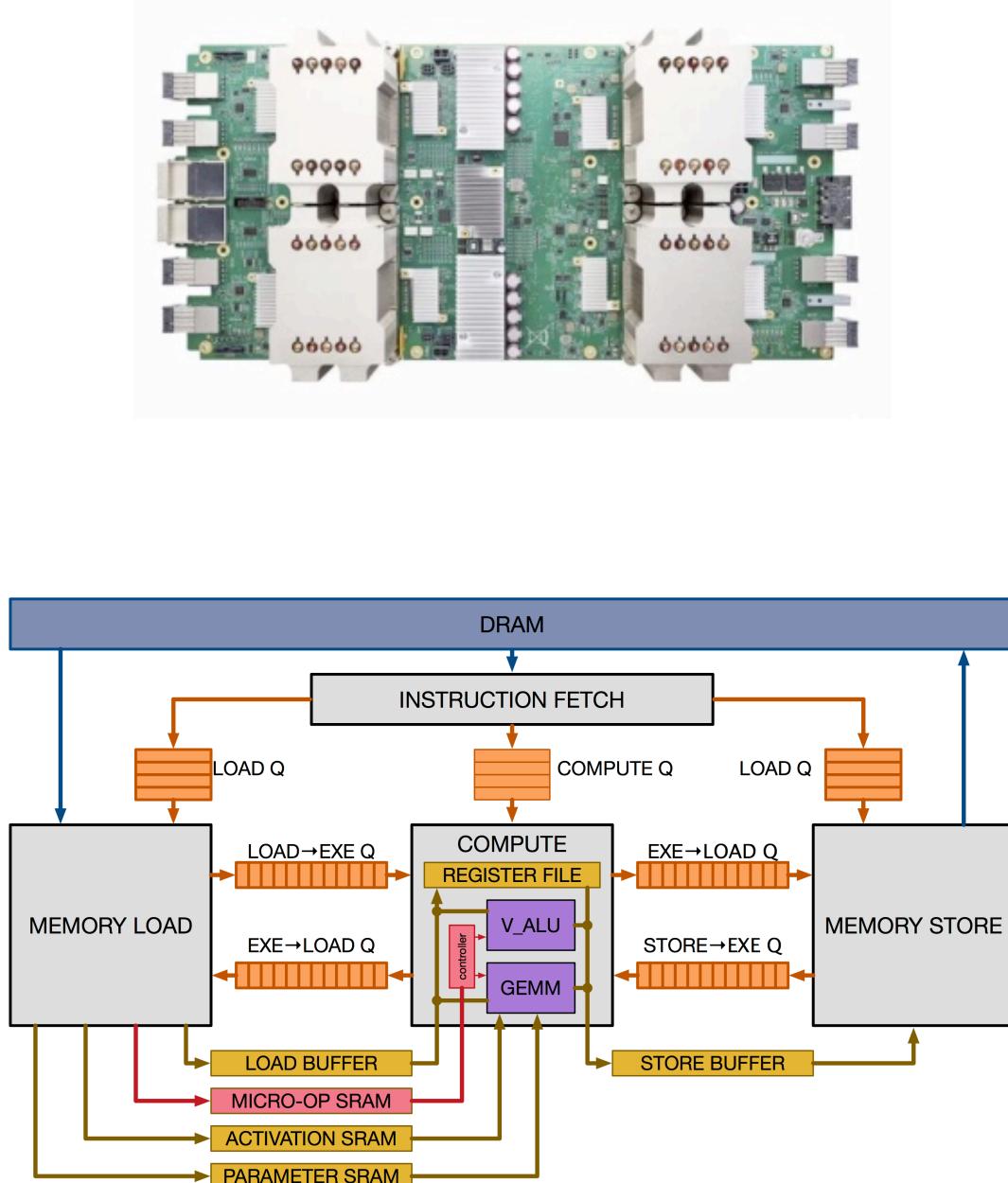


## Explicitly Managed Memory Subsystem

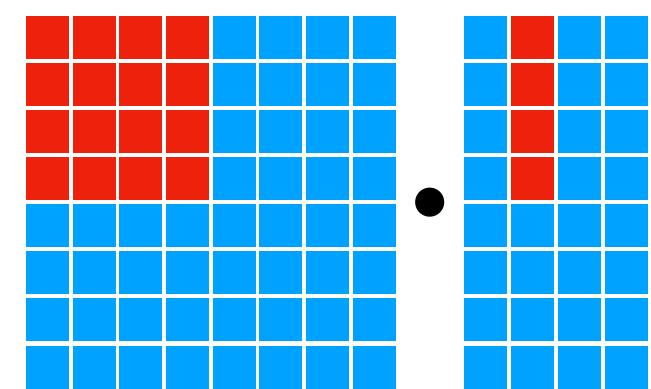


# Search Space for TPU-like Specialized Accelerators

## TPUs



**Tensor  
Compute Primitives**



**Explicitly Managed  
Memory Subsystem**

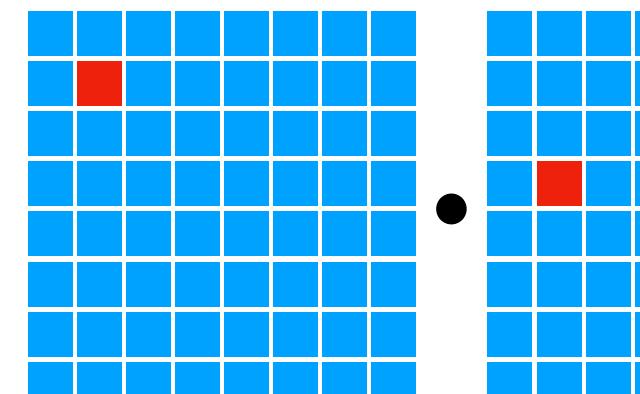


# Tensorization Challenge

**Compute  
primitives**

# Tensorization Challenge

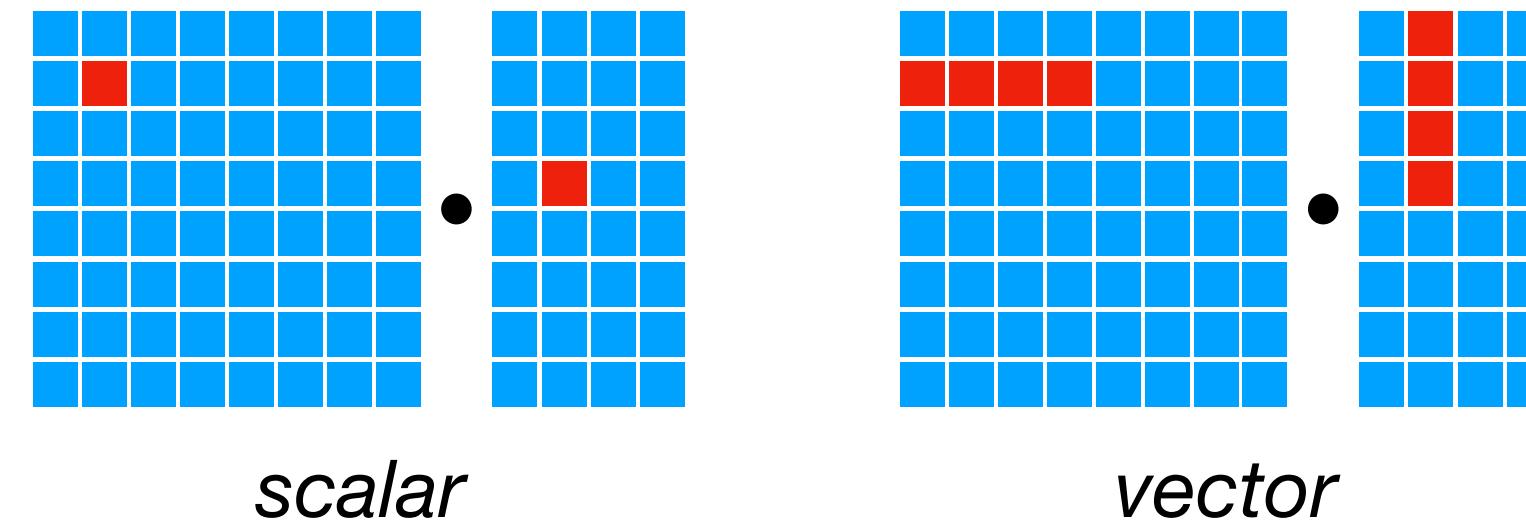
**Compute  
primitives**



*scalar*

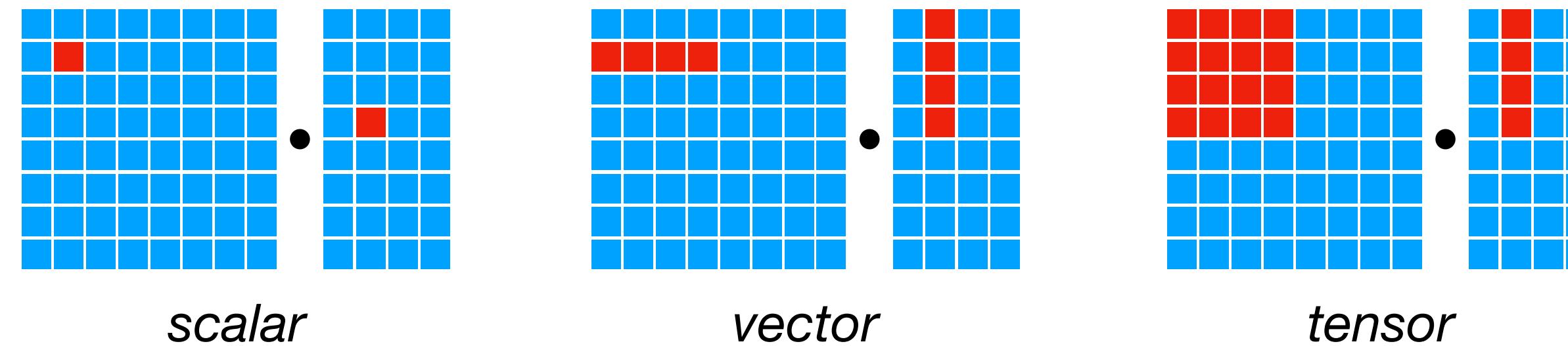
# Tensorization Challenge

**Compute  
primitives**



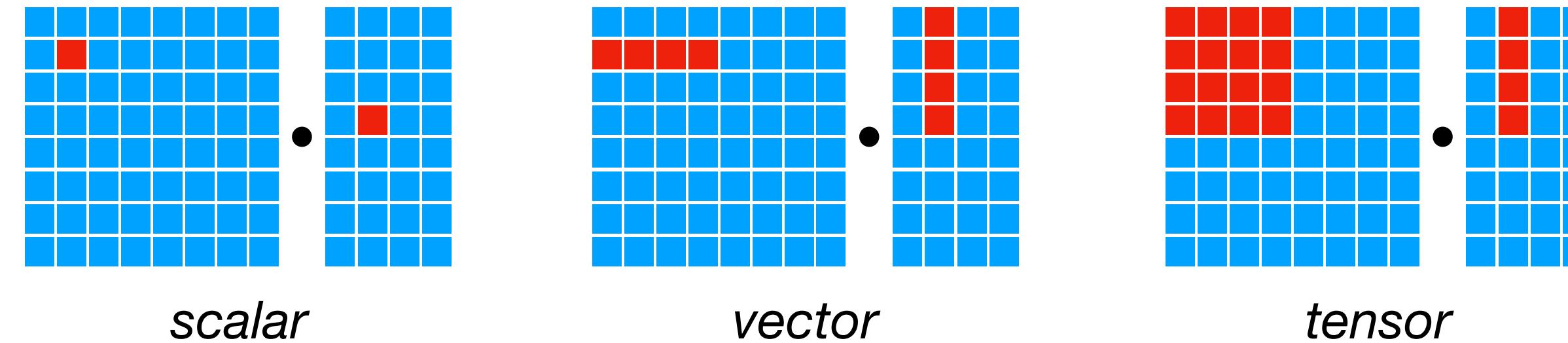
# Tensorization Challenge

**Compute  
primitives**



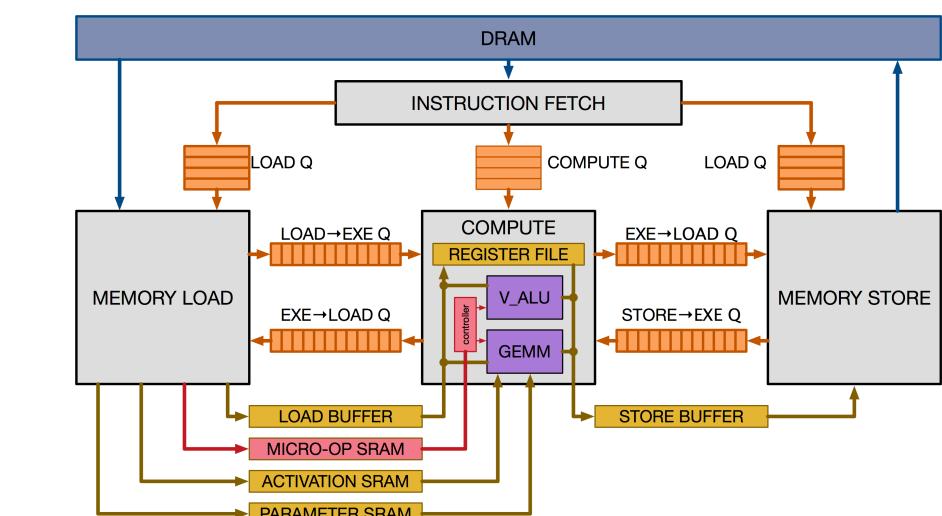
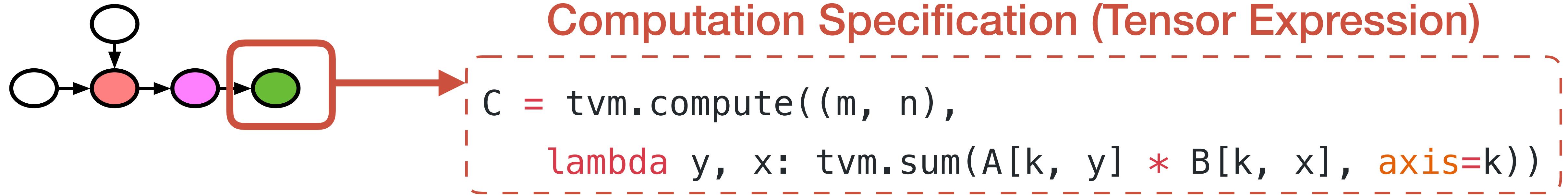
# Tensorization Challenge

**Compute  
primitives**

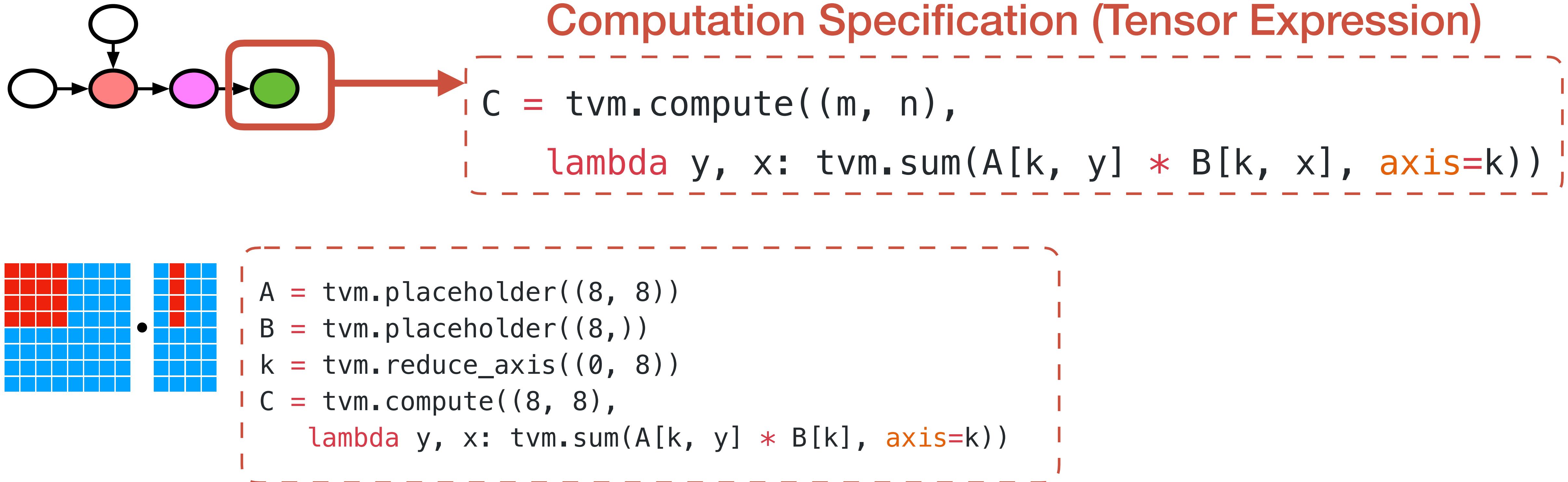


**Challenge: Build systems to support  
emerging tensor instructions**

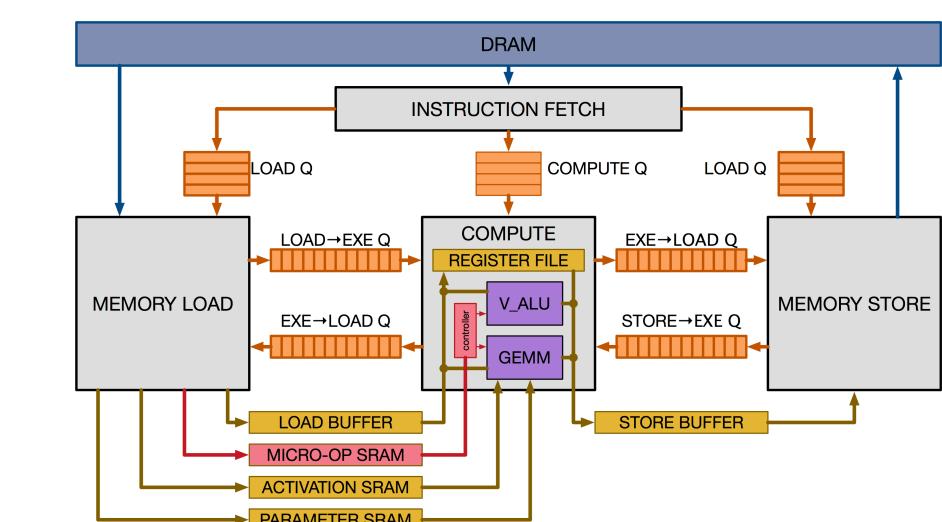
# Tensorization Challenge



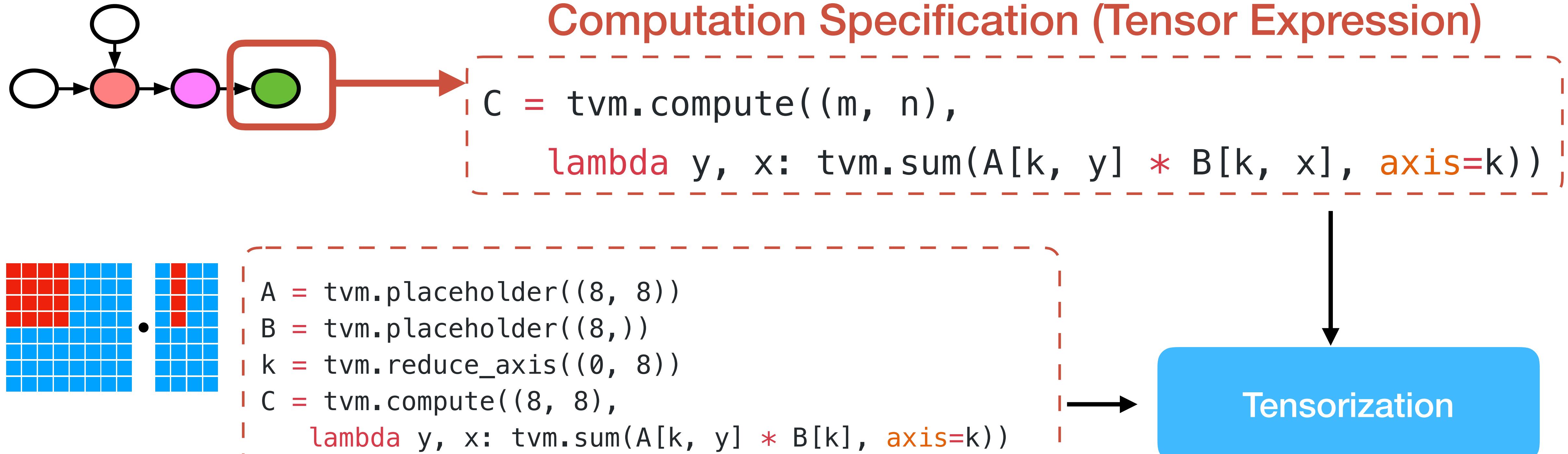
# Tensorization Challenge



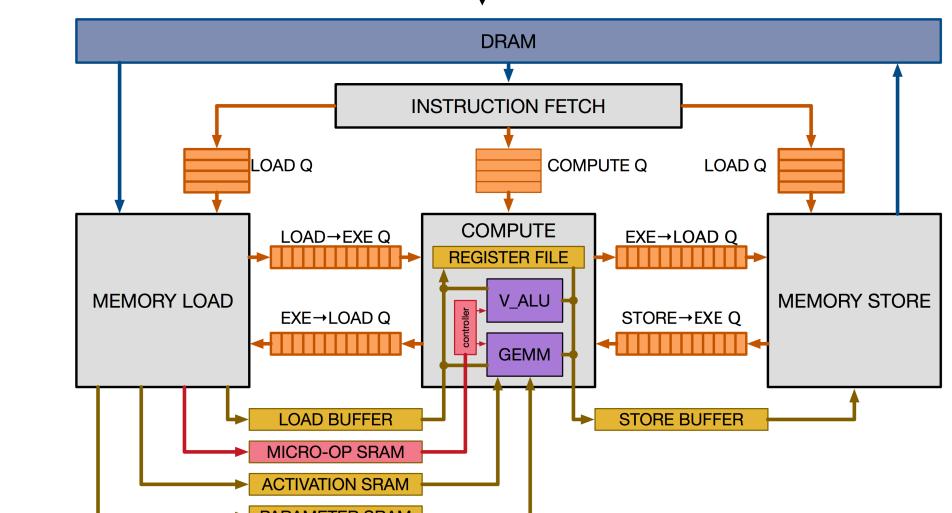
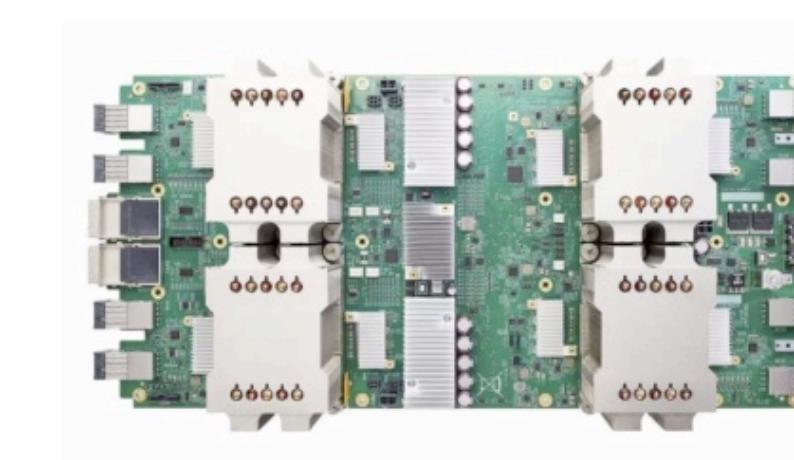
## HW Interface Specification by Tensor Expression



# Tensorization Challenge

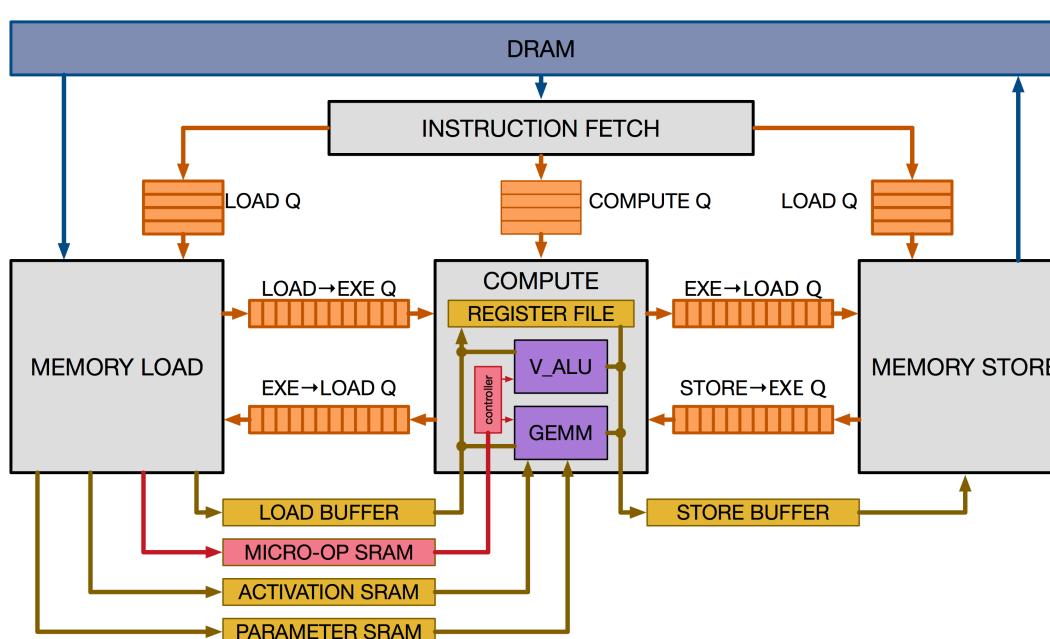


HW Interface Specification by Tensor Expression

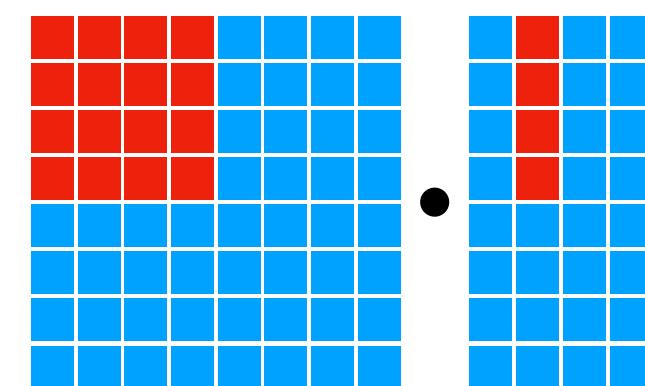


# Search Space for TPU-like Specialized Accelerators

## TPUs



## Tensor Compute Primitives

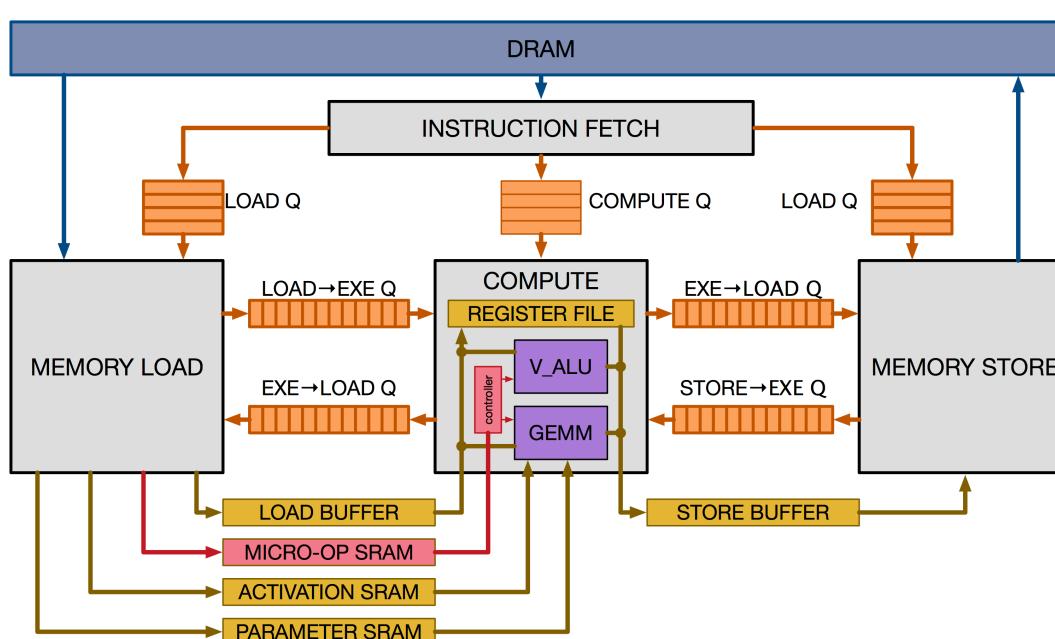


## Explicitly Managed Memory Subsystem

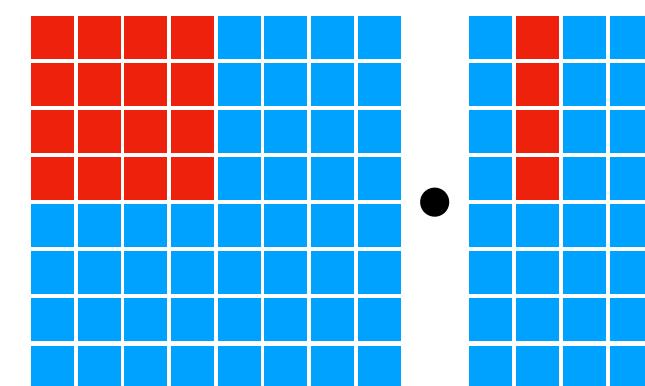


# Search Space for TPU-like Specialized Accelerators

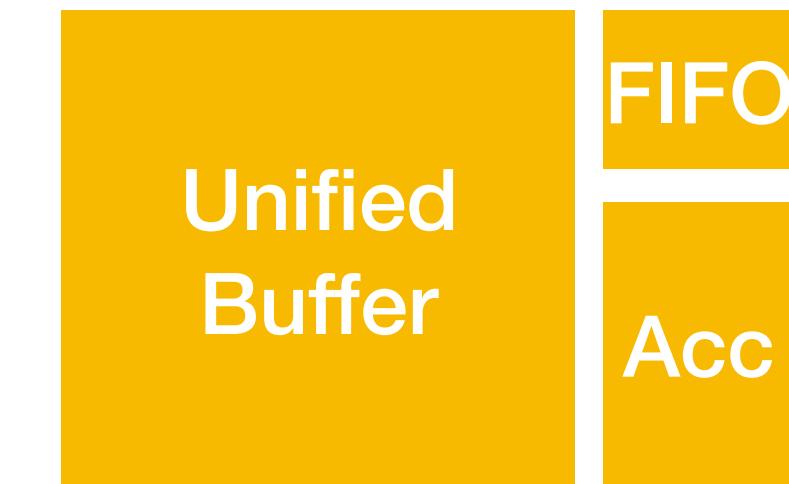
TPUs



Tensor  
Compute Primitives



Explicitly Managed  
Memory Subsystem



# Software Support for Latency Hiding

**Single Module  
No Task-Pipelining**

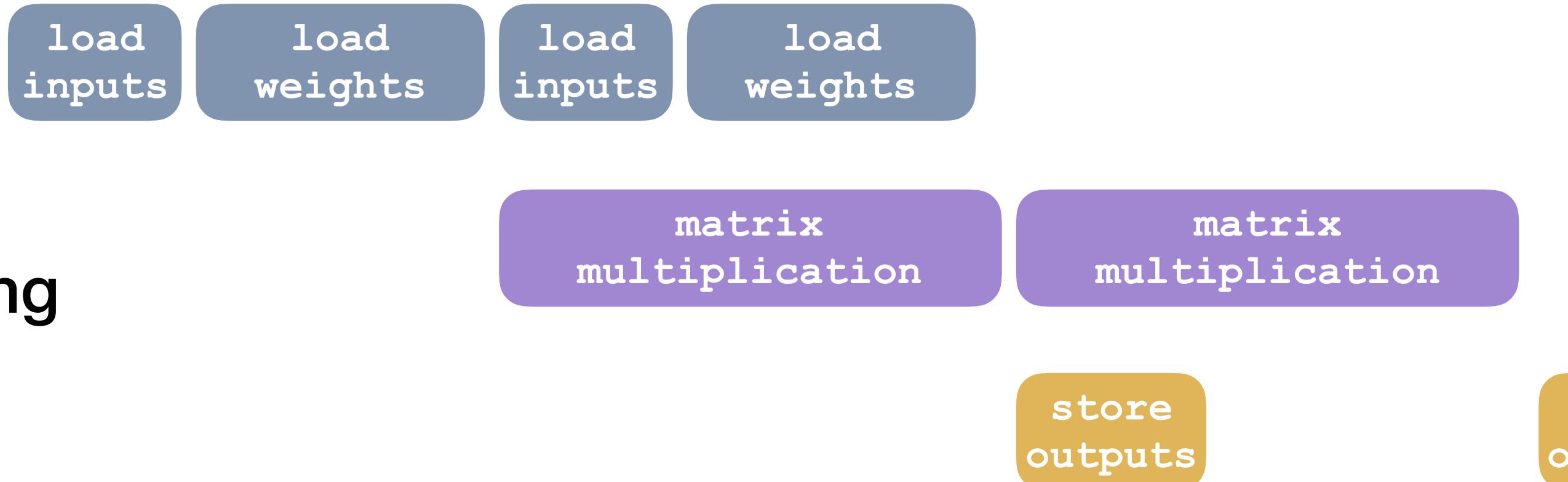


# Software Support for Latency Hiding

**Single Module  
No Task-Pipelining**



**Multiple-Module  
Task-Level Pipelining**

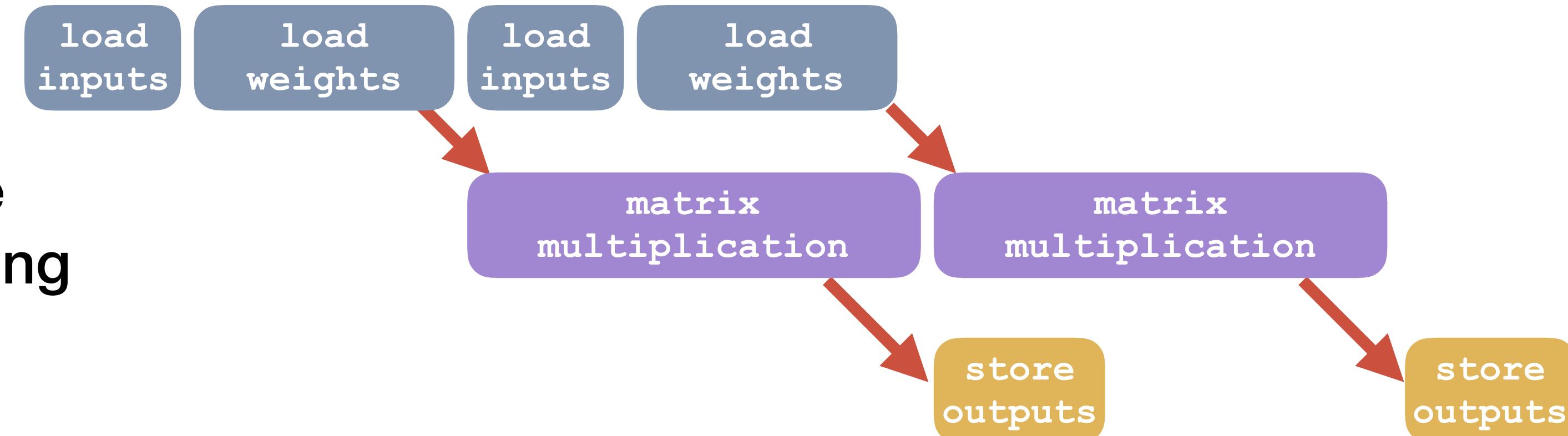


# Software Support for Latency Hiding

**Single Module  
No Task-Pipelining**



**Multiple-Module  
Task-Level Pipelining**

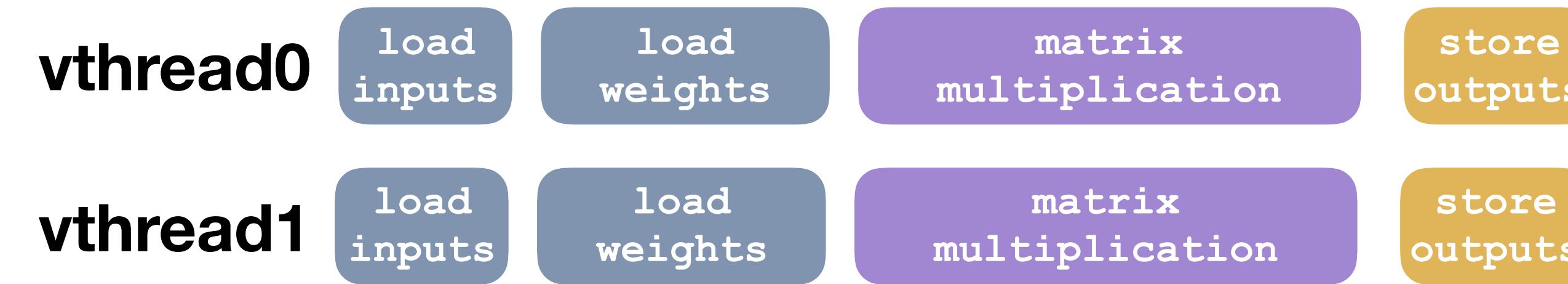


**Explicit dependencies  
managed by software to hide memory latency**

# Software Support for Latency Hiding

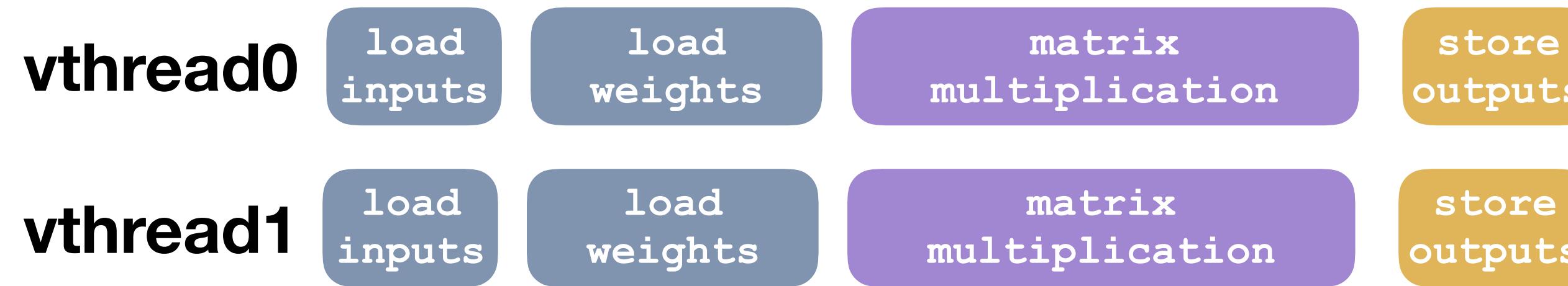
# Software Support for Latency Hiding

Multi-threaded  
Program

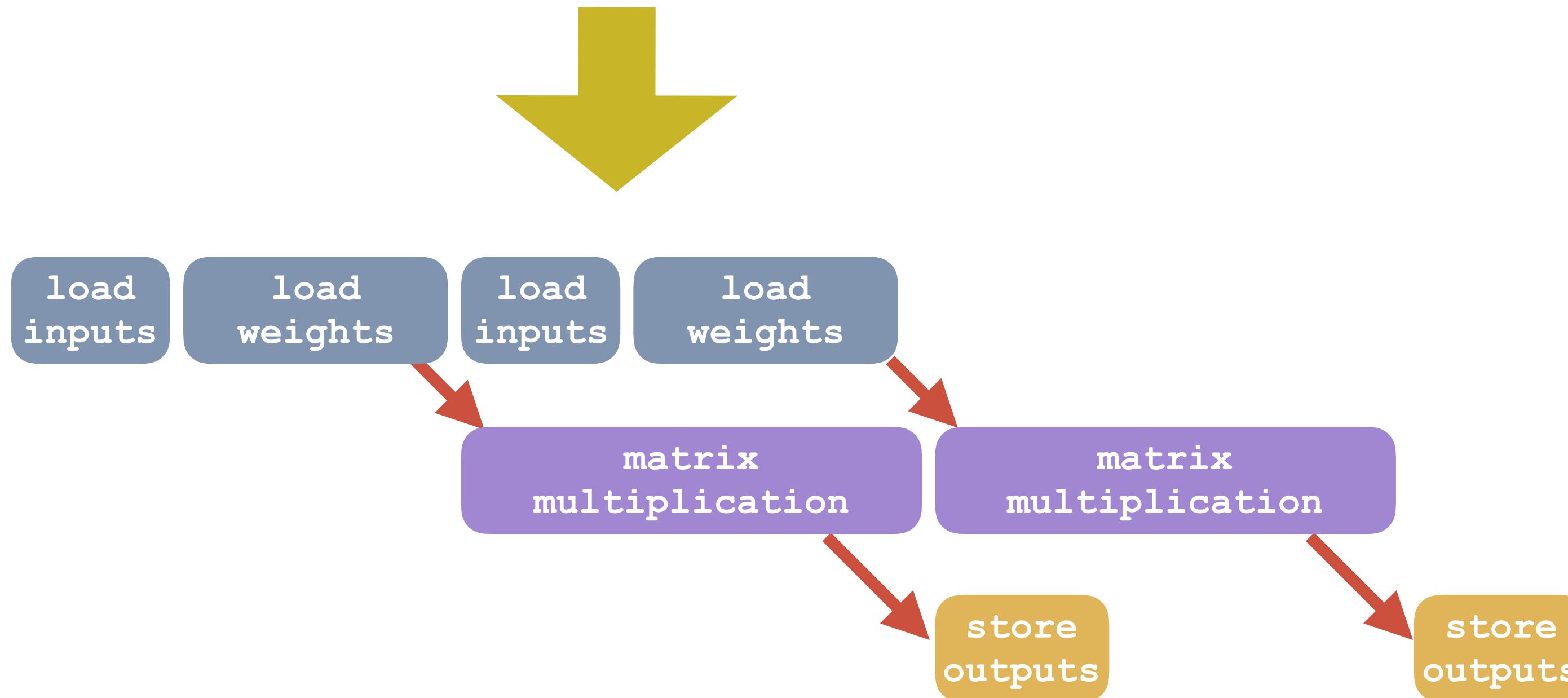


# Software Support for Latency Hiding

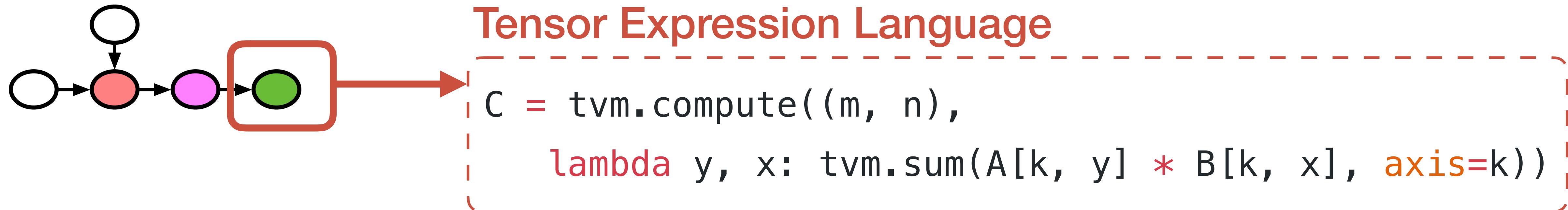
Multi-threaded  
Program



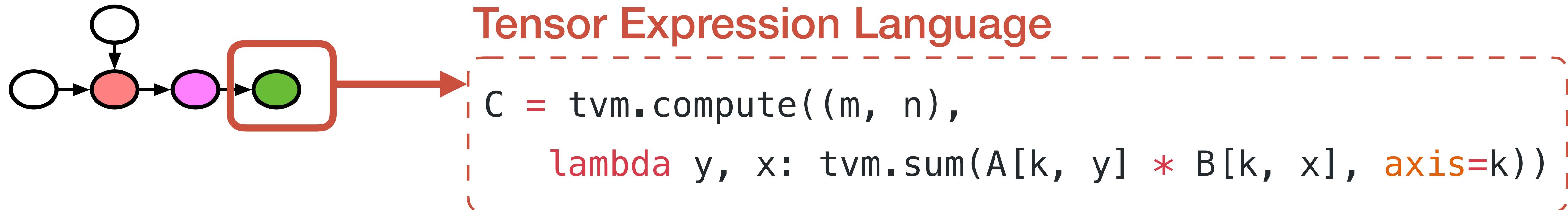
Program with  
Task-level Pipeline  
Instructions



# Summary: Hardware-aware Search Space



# Summary: Hardware-aware Search Space



Primitives in prior work:

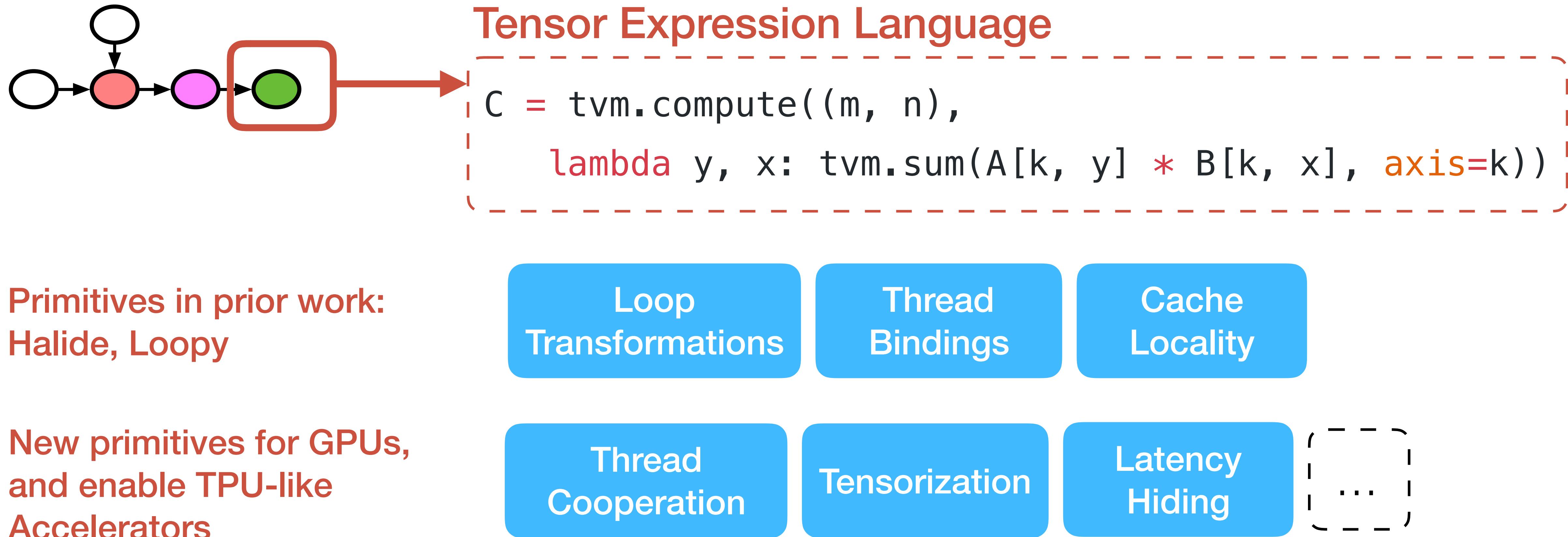
Halide, Loopy

Loop  
Transformations

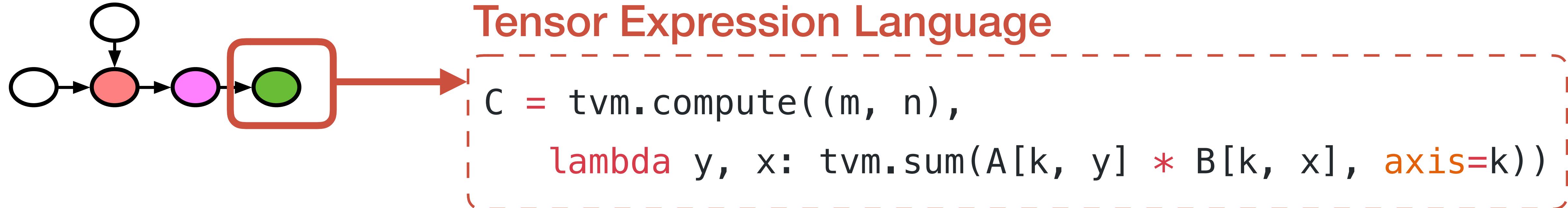
Thread  
Bindings

Cache  
Locality

# Summary: Hardware-aware Search Space



# Summary: Hardware-aware Search Space



Primitives in prior work:  
Halide, Loopy

Loop  
Transformations

Thread  
Bindings

Cache  
Locality

New primitives for GPUs,  
and enable TPU-like  
Accelerators

Thread  
Cooperation

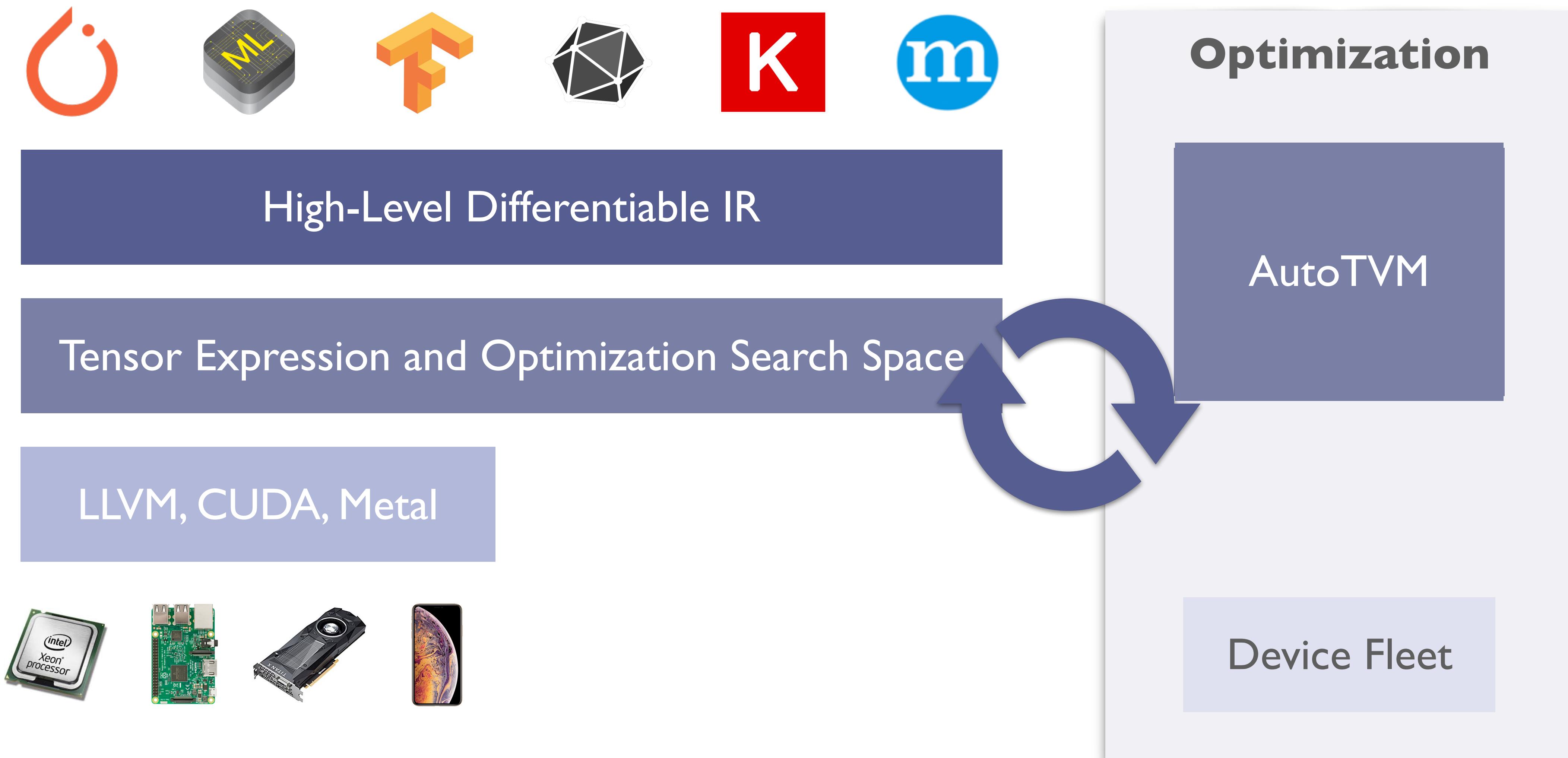
Tensorization

Latency  
Hiding

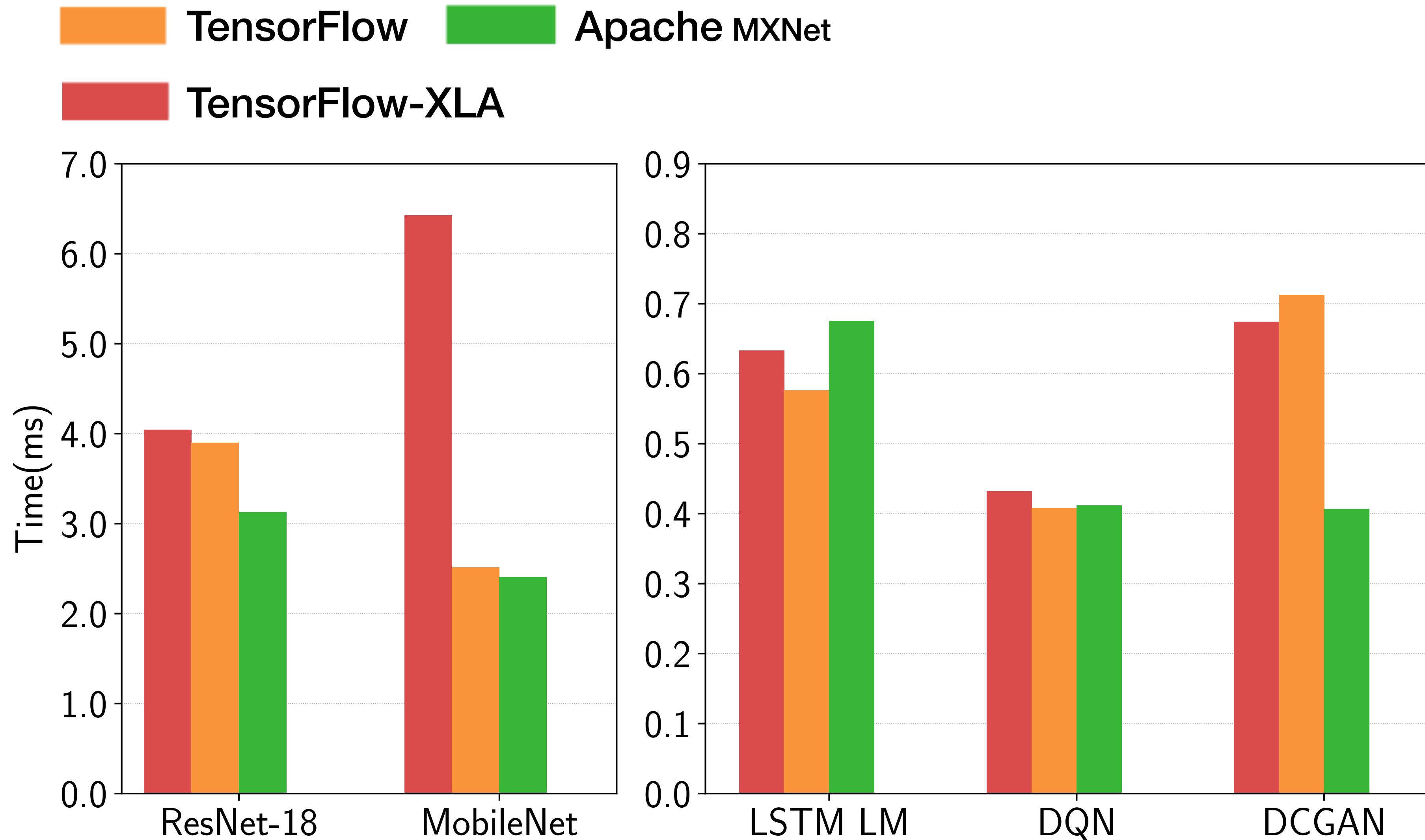
Hardware



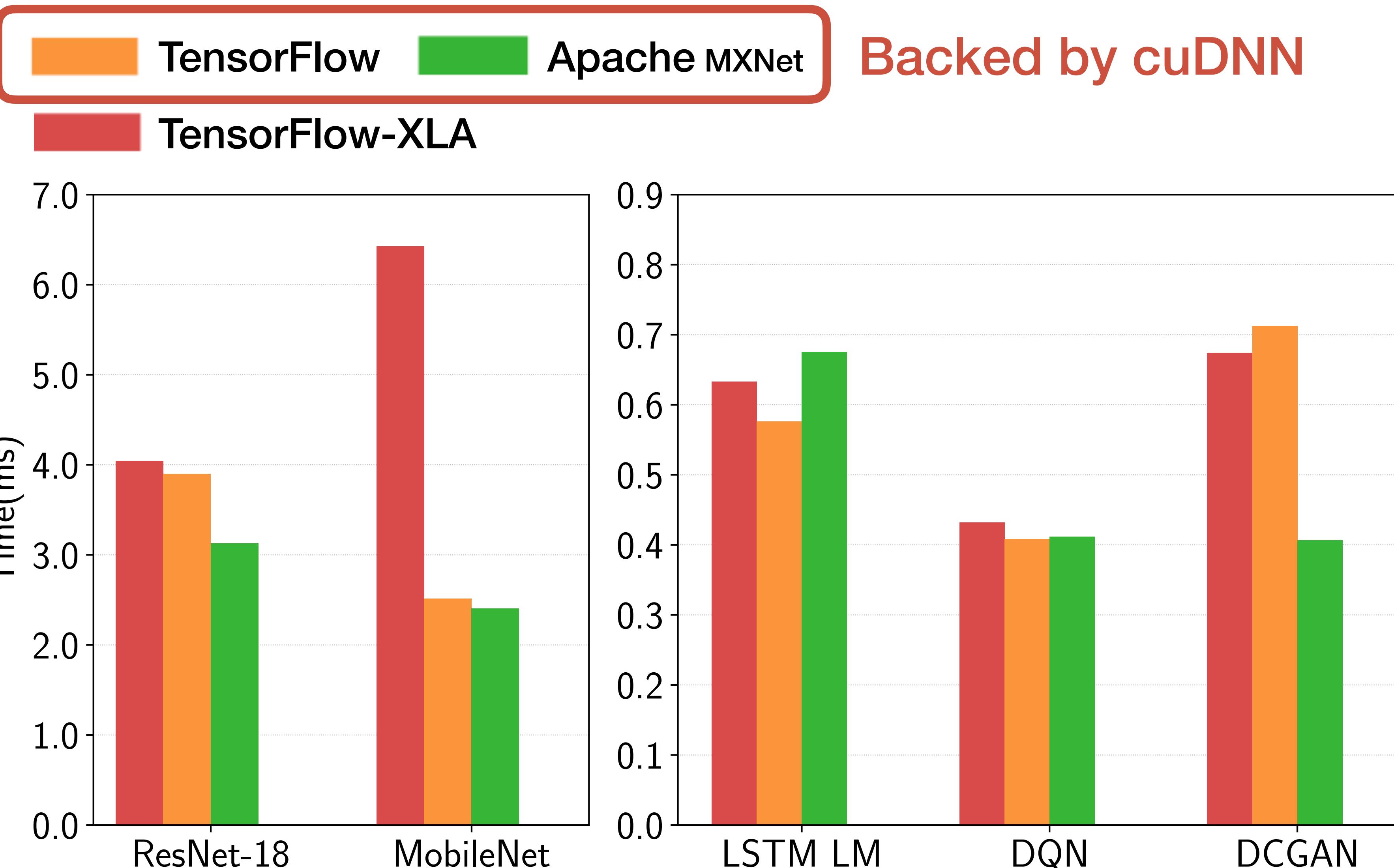
# TVM: End to End Deep Learning Compiler



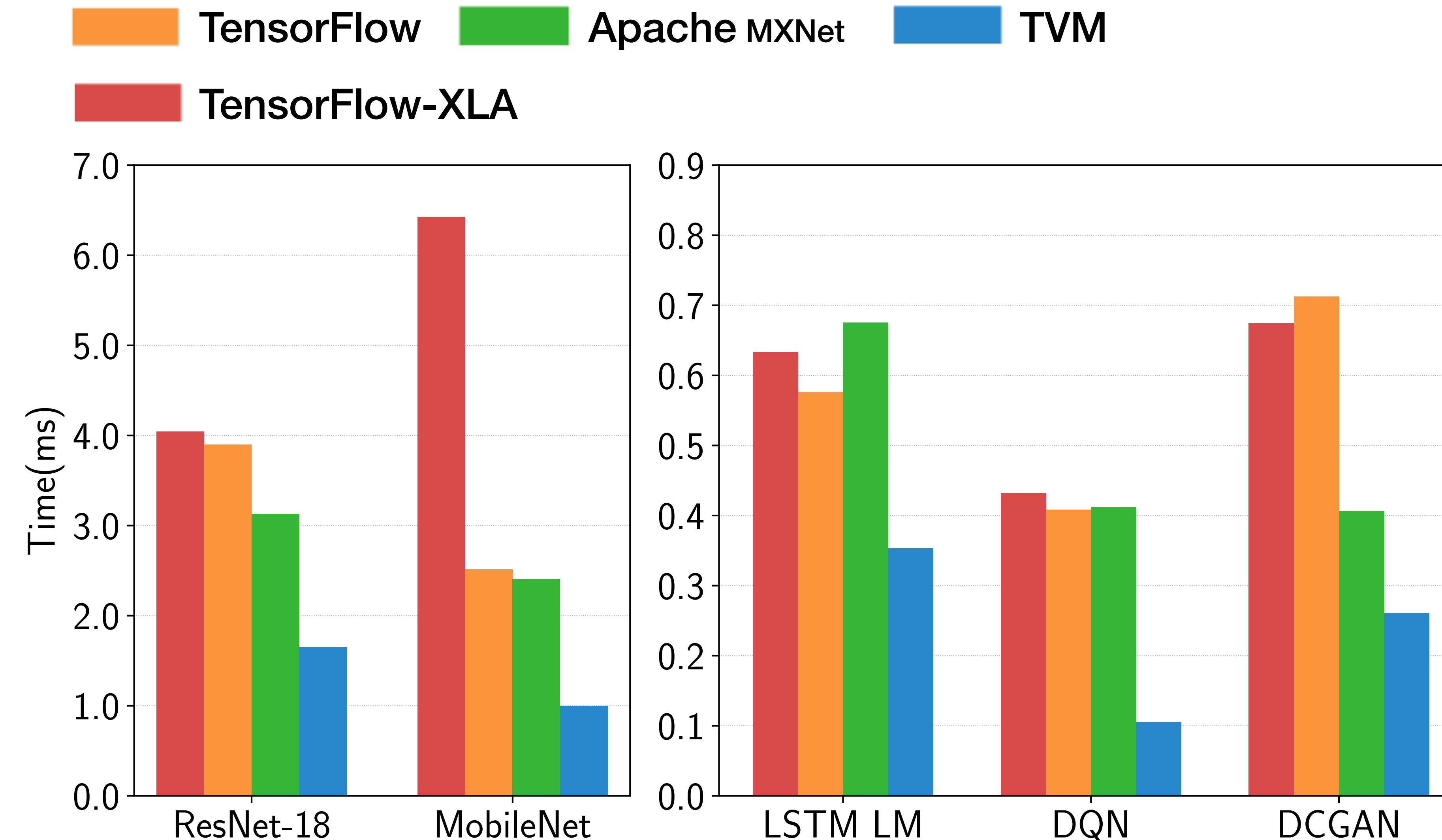
# End to End Inference Performance (Nvidia Titan X)



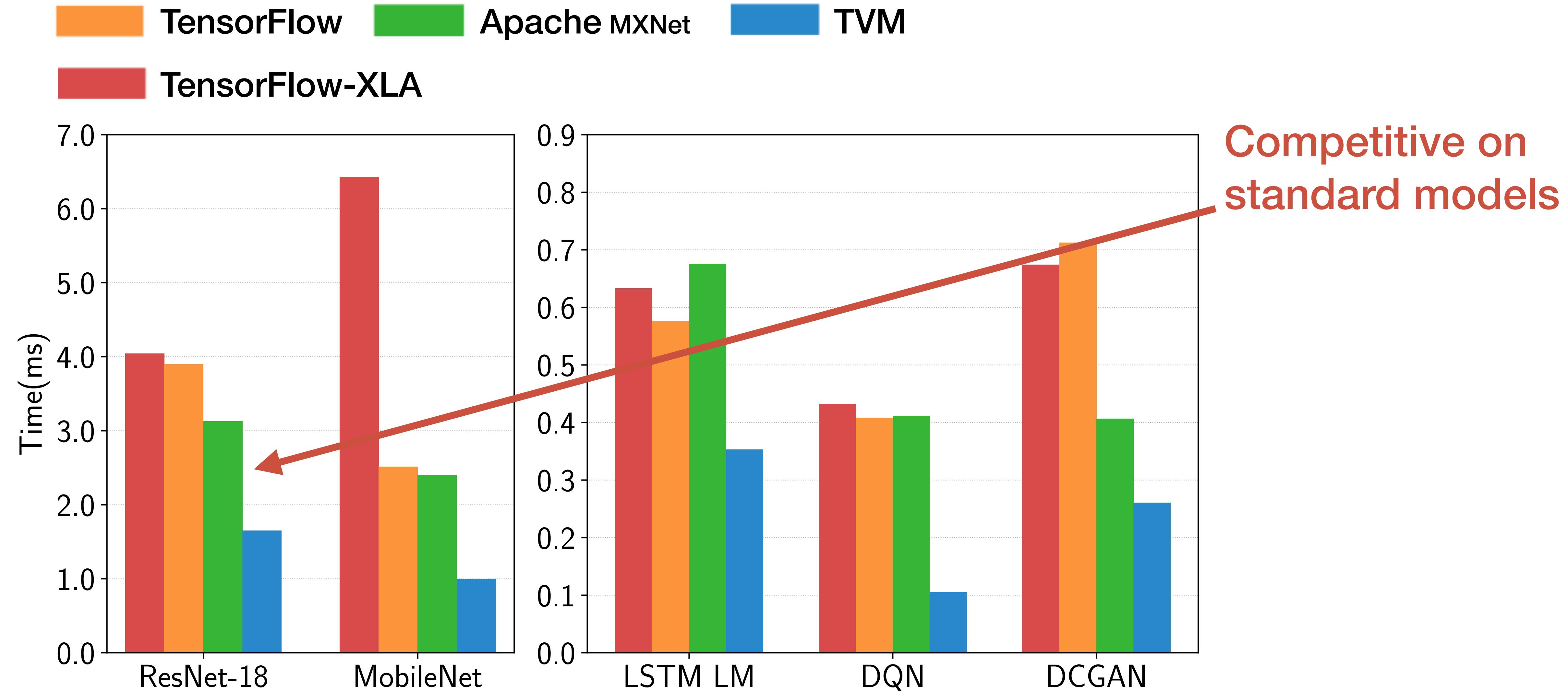
# End to End Inference Performance (Nvidia Titan X)



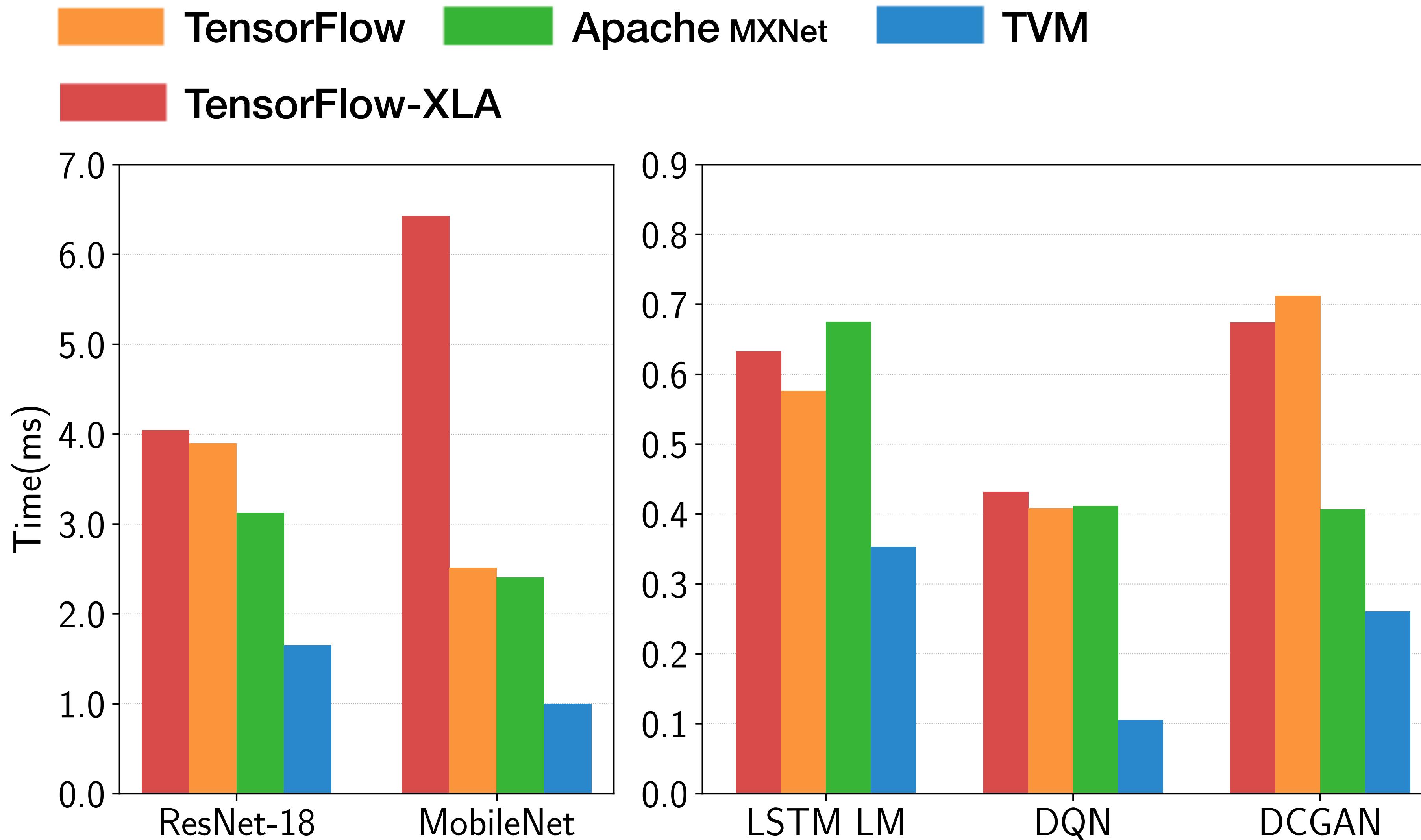
# End to End Inference Performance (Nvidia Titan X)



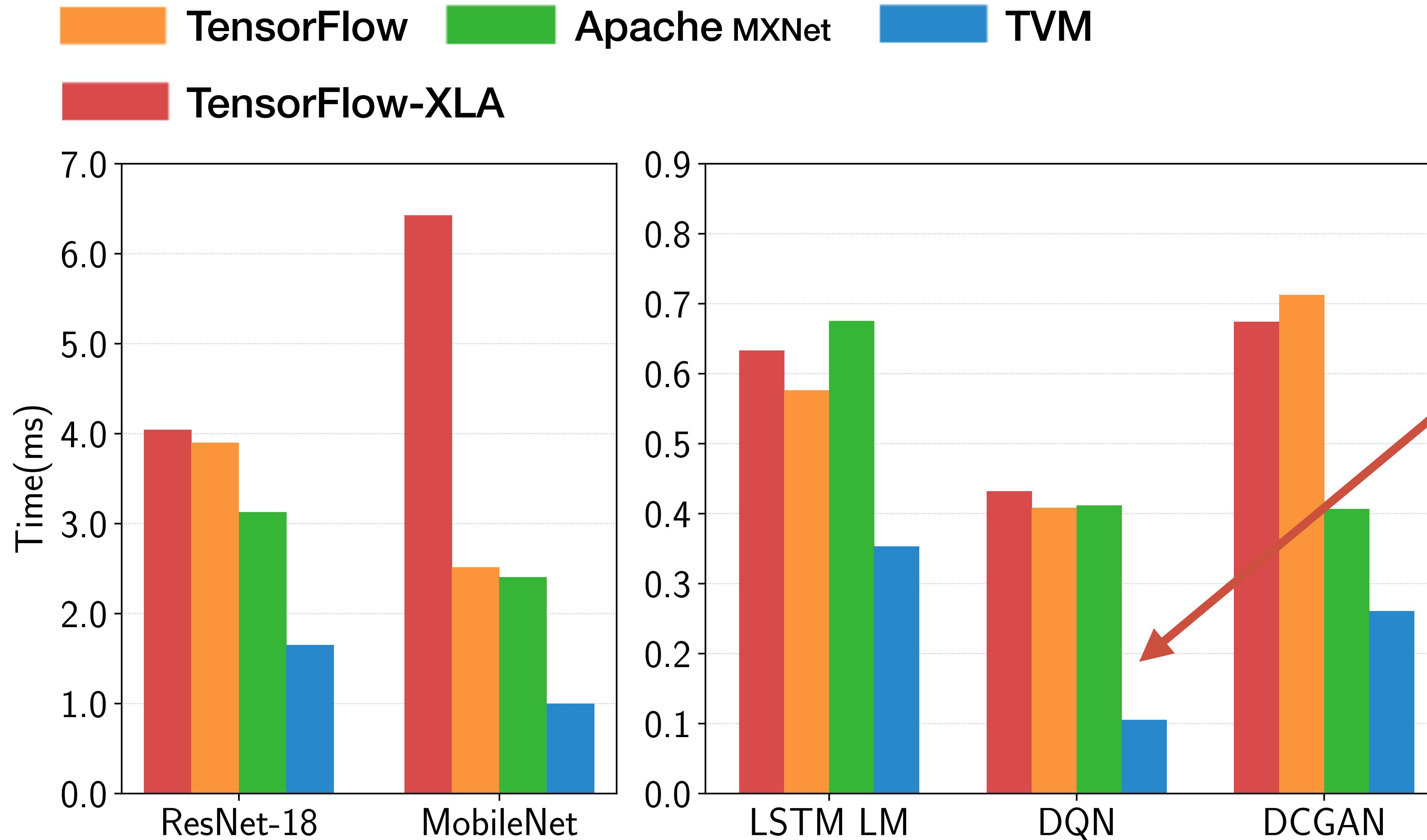
# End to End Inference Performance (Nvidia Titan X)



# End to End Inference Performance (Nvidia Titan X)

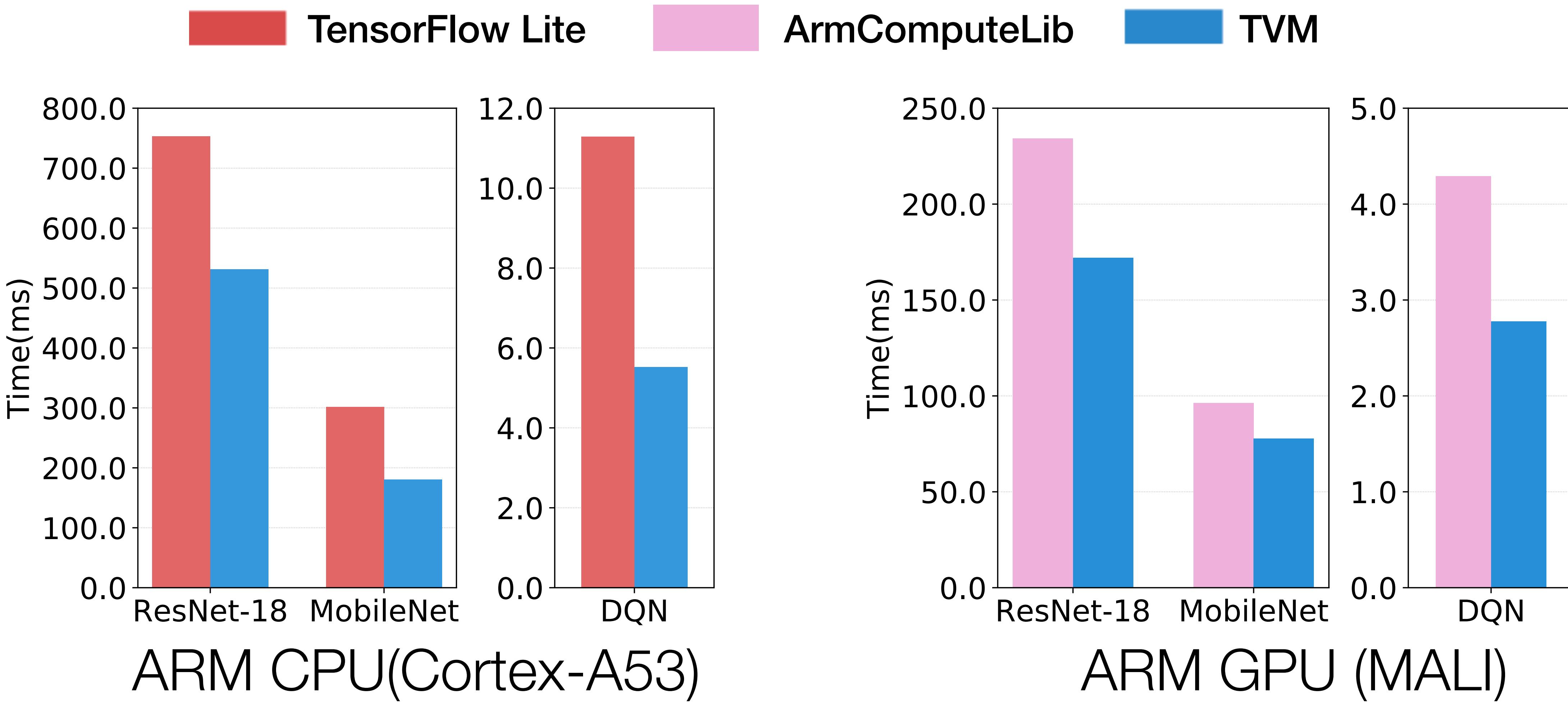


# End to End Inference Performance (Nvidia Titan X)



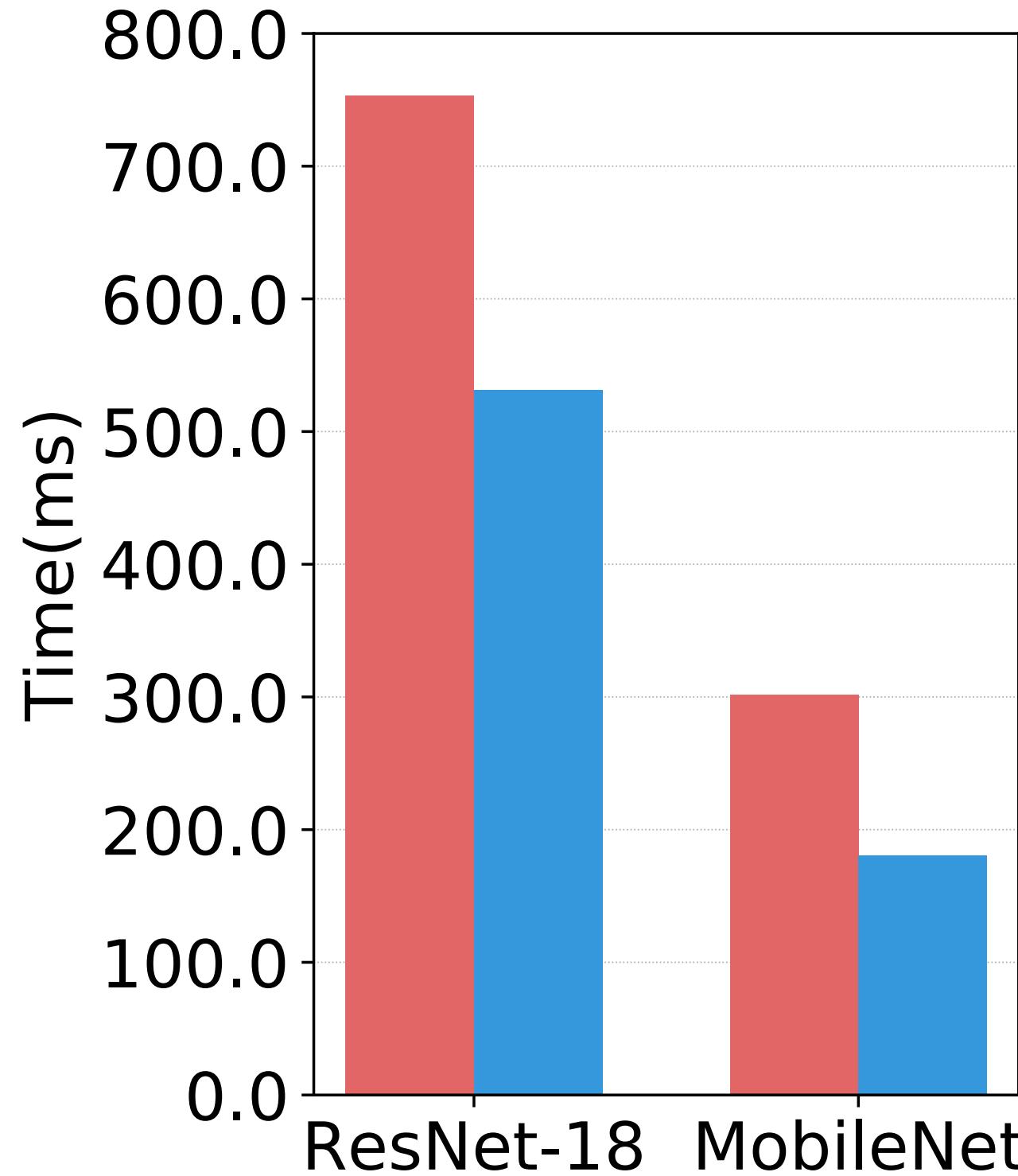
3x better  
on  
emerging  
models

# Portable Performance Across Hardware Platforms

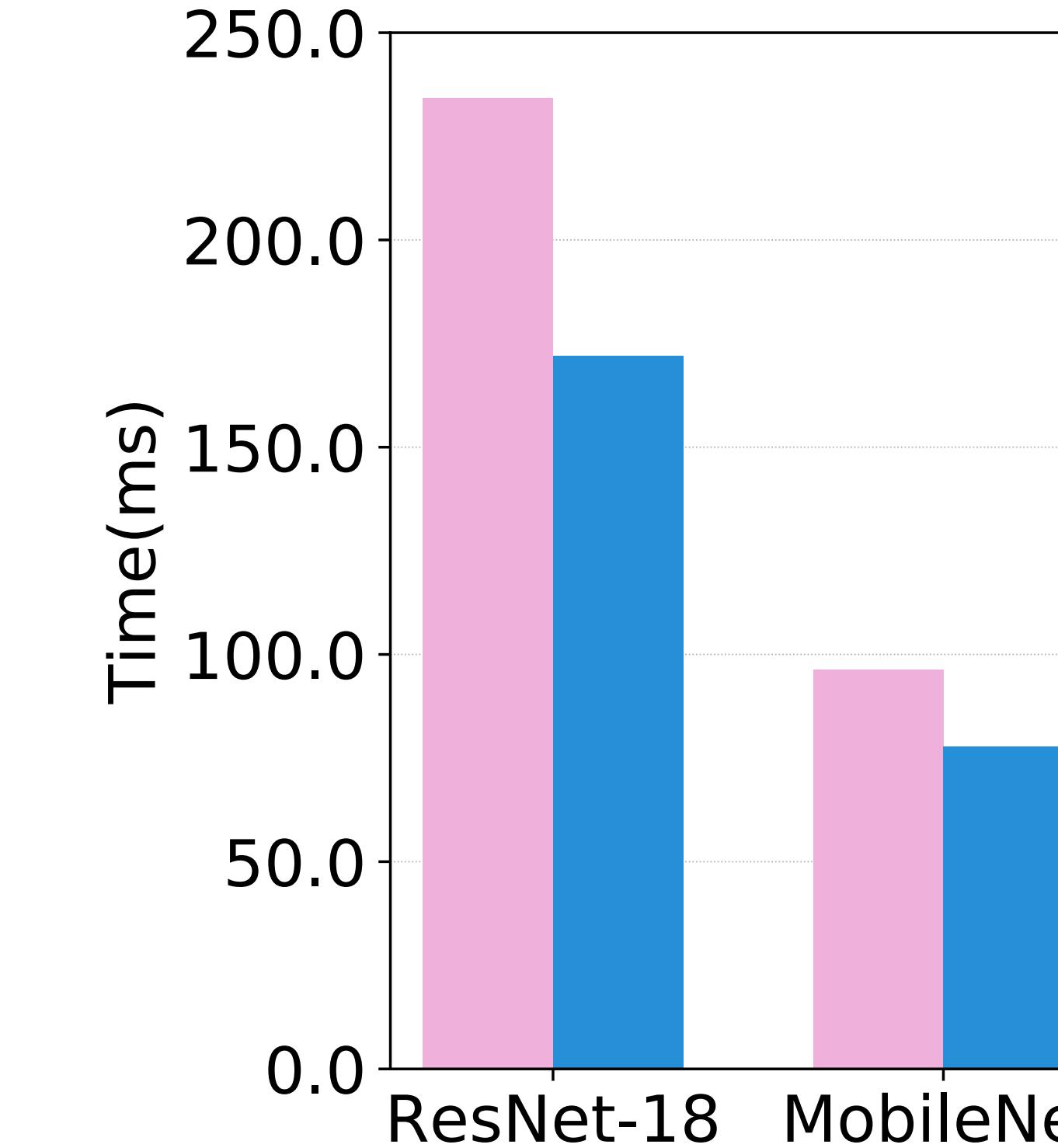
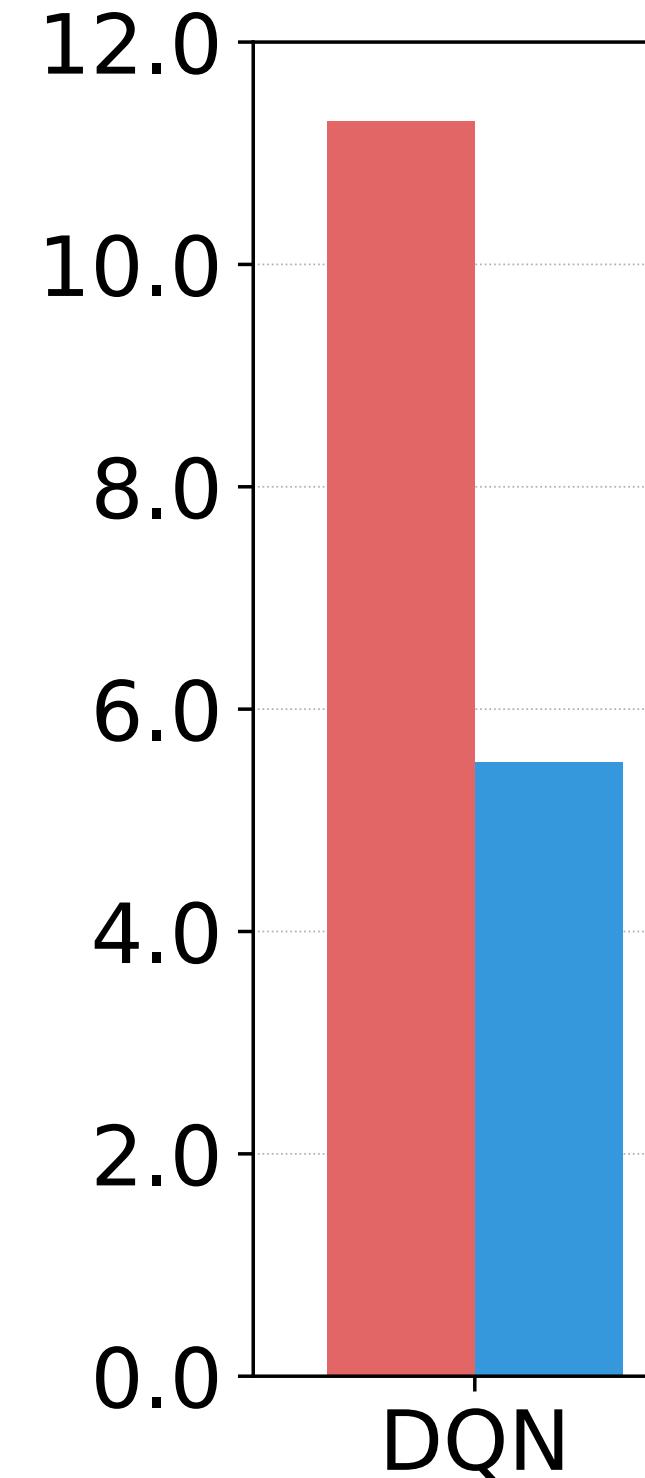


# Portable Performance Across Hardware Platforms

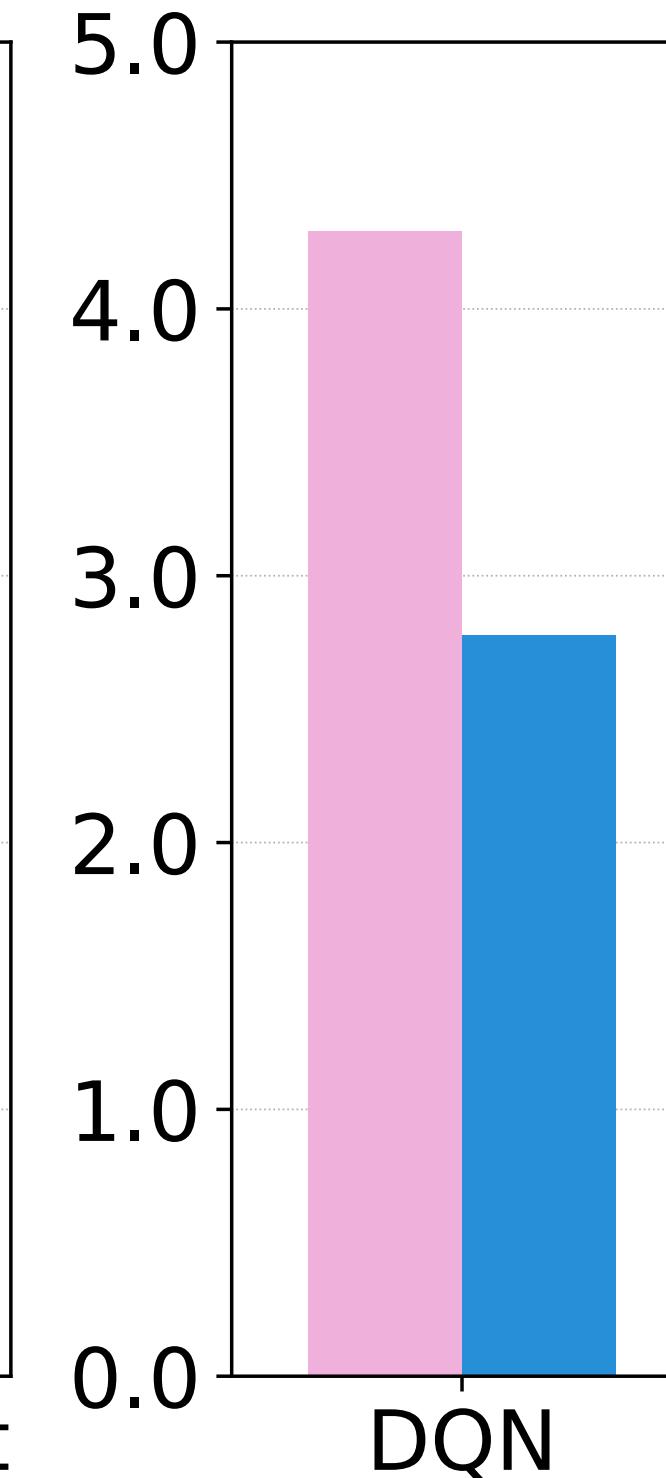
**Special frameworks for the particular hardware platform**



ARM CPU(Cortex-A53)



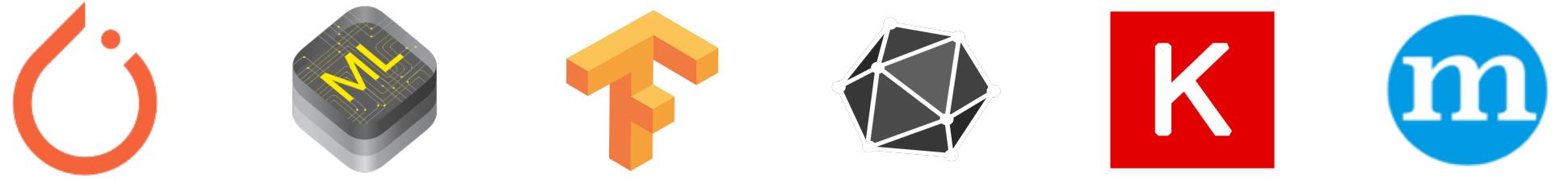
ARM GPU (MALI)



# TVM: End to End Deep Learning Compiler

**What about Accelerator Support?**

# VTA: Open & Flexible Deep Learning Accelerator

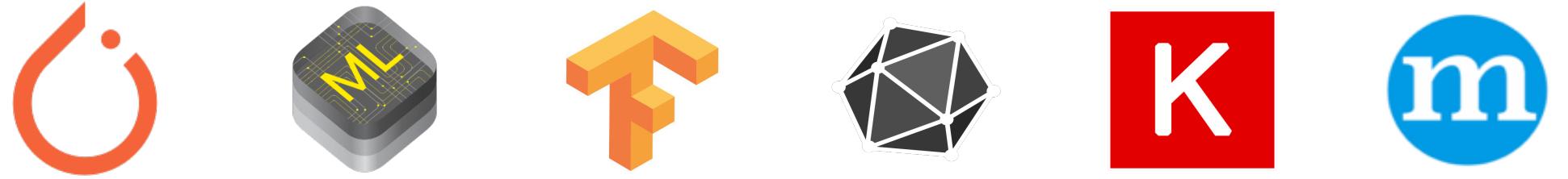


Current TVM Stack



**Moreau, Chen, et al. work in progress**

# VTA: Open & Flexible Deep Learning Accelerator



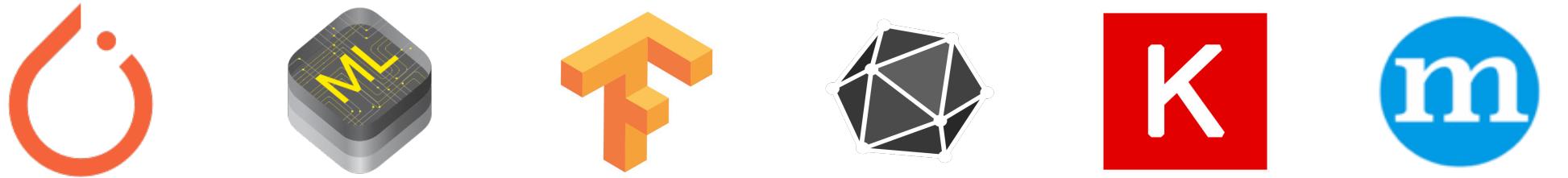
Current TVM Stack

VTA MicroArchitecture



**Moreau, Chen, et al. work in progress**

# VTA: Open & Flexible Deep Learning Accelerator



Current TVM Stack

VTA Hardware/Software Interface (ISA)

VTA MicroArchitecture



**Moreau, Chen, et al. work in progress**

# VTA: Open & Flexible Deep Learning Accelerator



Current TVM Stack

VTA Runtime & JIT Compiler

VTA Hardware/Software Interface (ISA)

VTA MicroArchitecture



**Moreau, Chen, et al. work in progress**

# VTA: Open & Flexible Deep Learning Accelerator



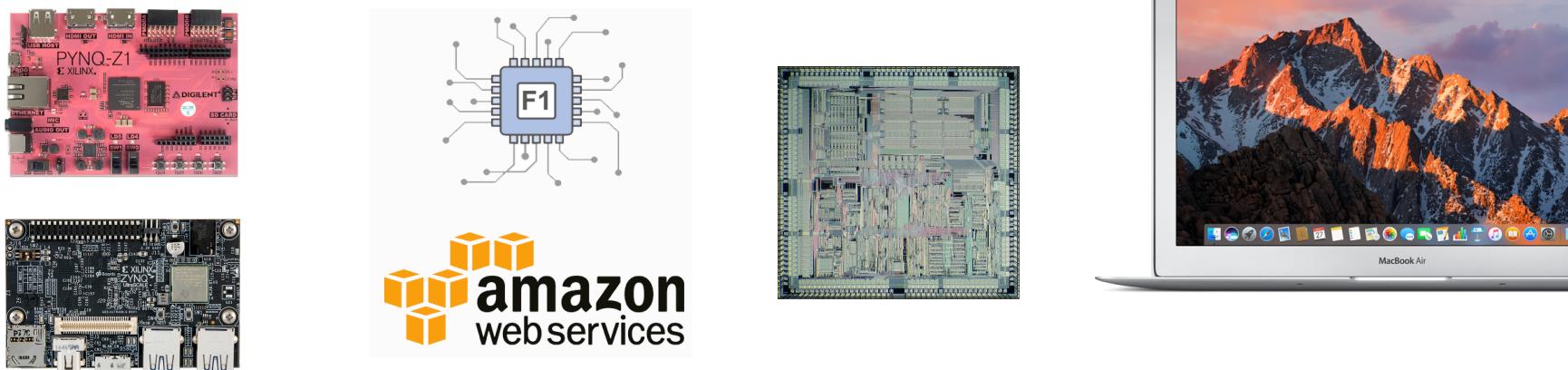
Current TVM Stack

VTA Runtime & JIT Compiler

VTA Hardware/Software Interface (ISA)

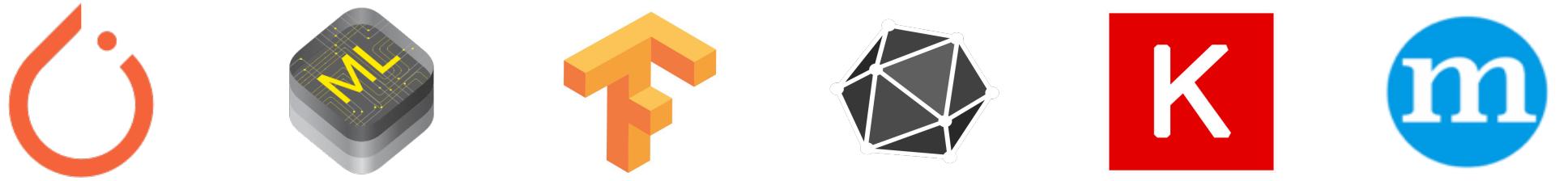
VTA MicroArchitecture

VTA Simulator



**Moreau, Chen, et al. work in progress**

# VTA: Open & Flexible Deep Learning Accelerator



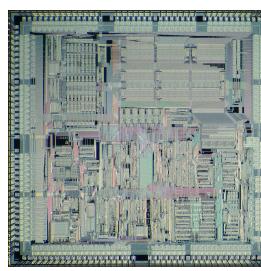
Current TVM Stack

VTA Runtime & JIT Compiler

VTA Hardware/Software Interface (ISA)

VTA MicroArchitecture

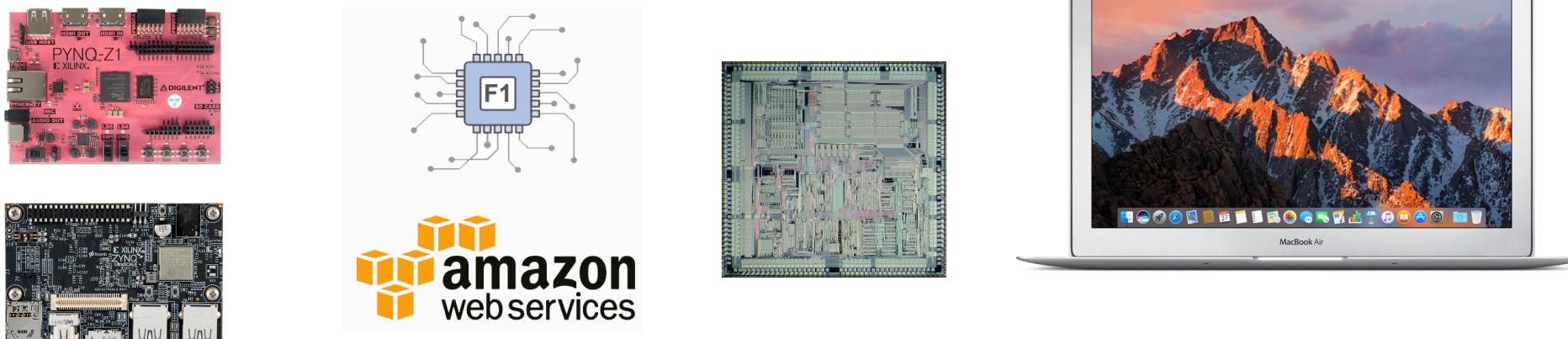
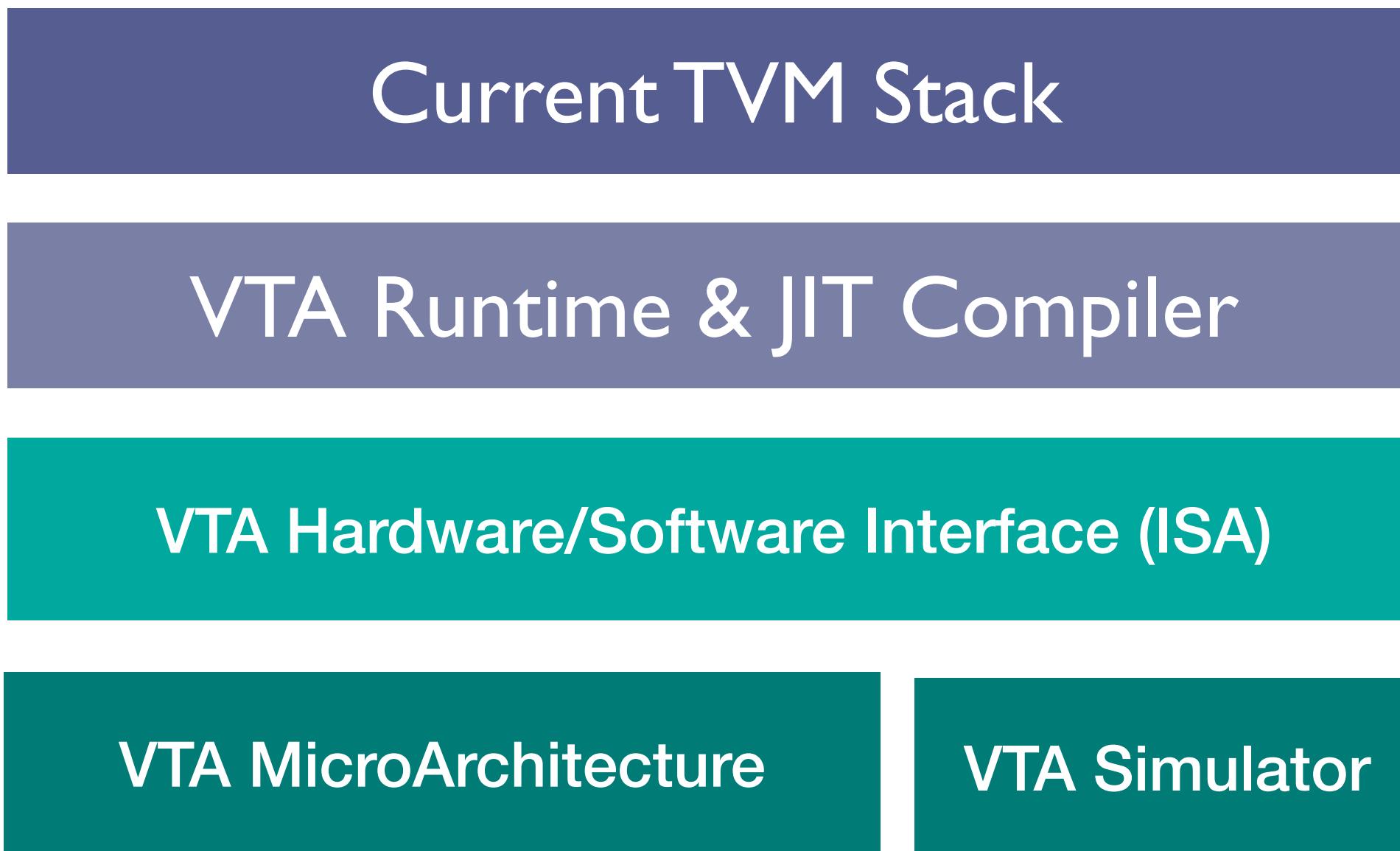
VTA Simulator



- Runtime JIT compile accelerator micro code
- Support heterogenous devices, 10x better than CPU on the same board.
- Move hardware complexity to software

**Moreau, Chen, et al. work in progress**

# VTA: Open & Flexible Deep Learning Accelerator

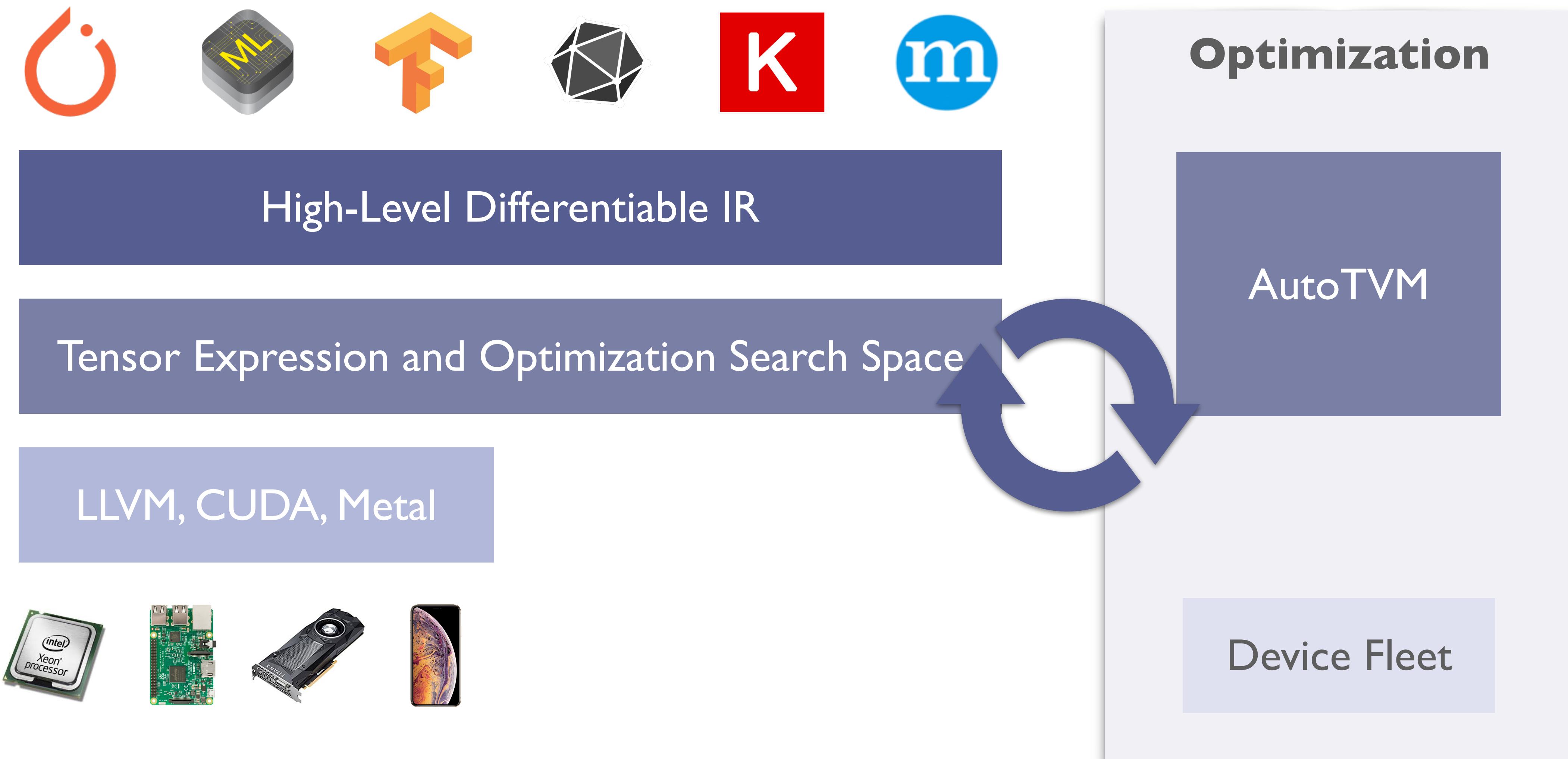


**Moreau, Chen, et al. work in progress**

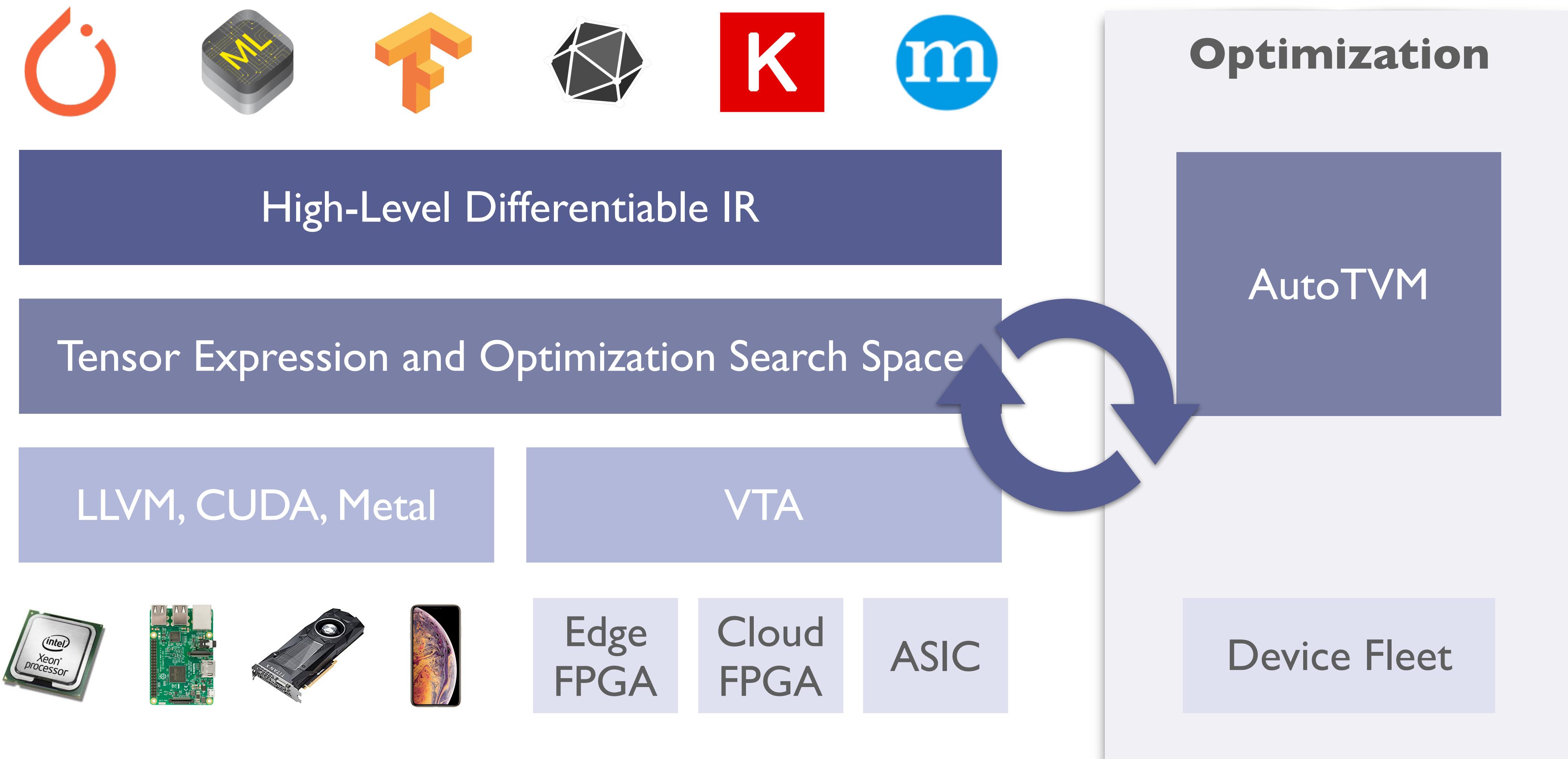
- Runtime JIT compile accelerator micro code
- Support heterogenous devices, 10x better than CPU on the same board.
- Move hardware complexity to software

**compiler, driver,  
hardware design  
full stack open source**

# TVM: End to End Deep Learning Compiler



# TVM: End to End Deep Learning Compiler



# Recent Developments

# Community Highlights

More **Dynamism**

**Tiny** machine learning

Better core **Infra**

More Specialized **Accelerator Support**

# Need for More Dynamism

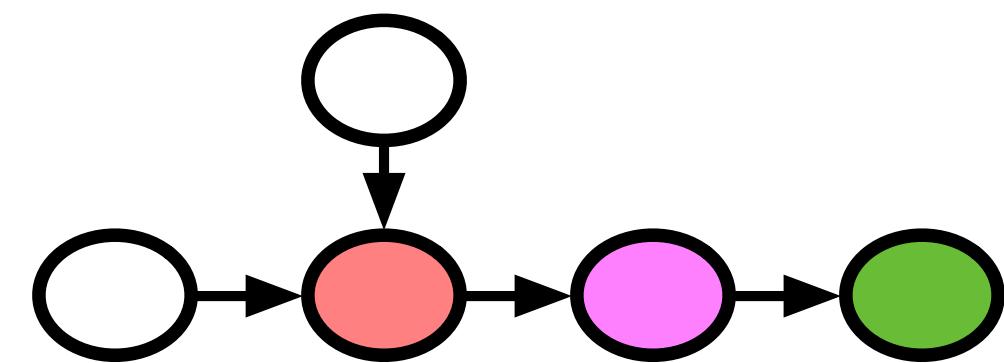
**Model**

**Data**

# Need for More Dynamism

static  
computacional graph

Model

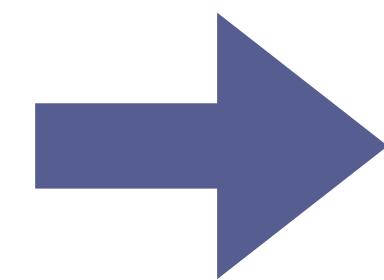
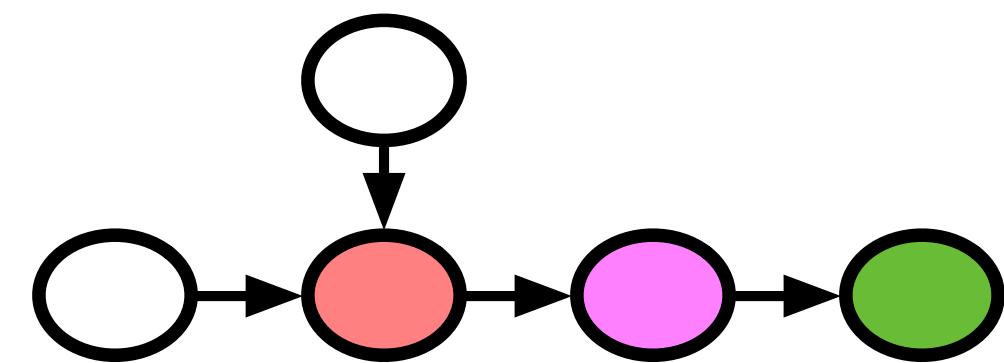


Data

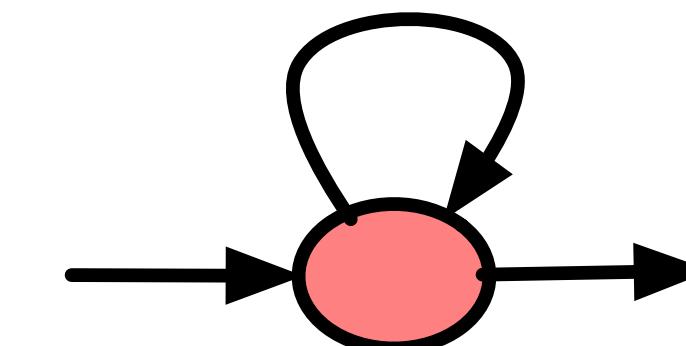
# Need for More Dynamism

**Model**

static  
computacional graph



program with  
loops and recursions

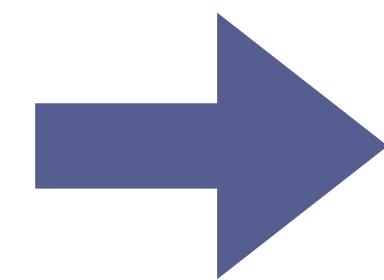
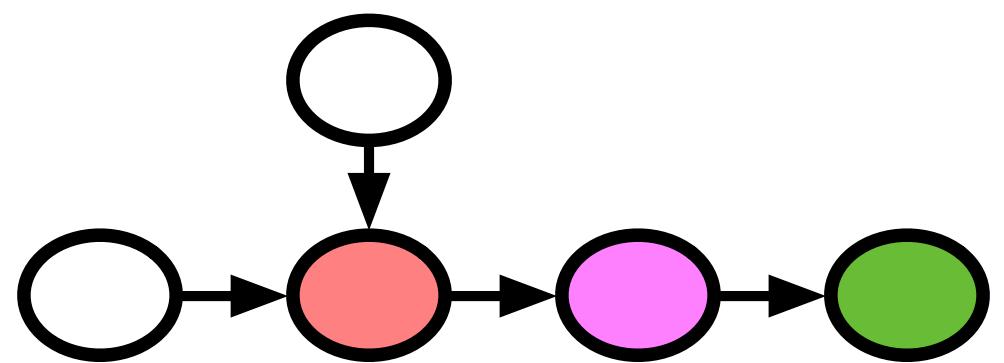


**Data**

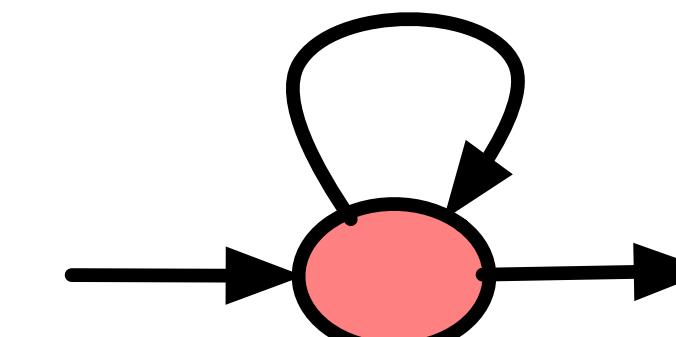
# Need for More Dynamism

Model

static  
computacional graph



program with  
loops and recursions



Data

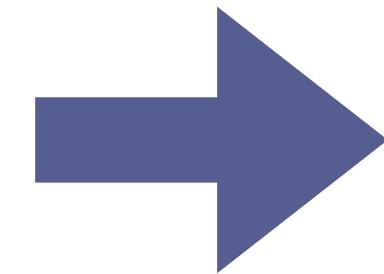
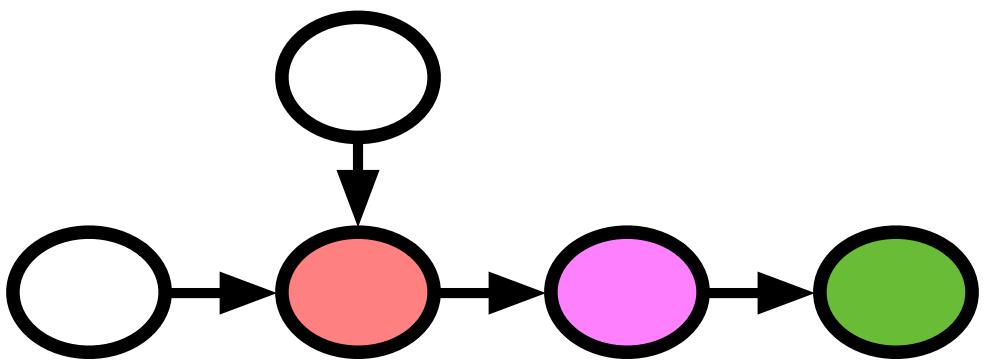
single tensor  
with known shape



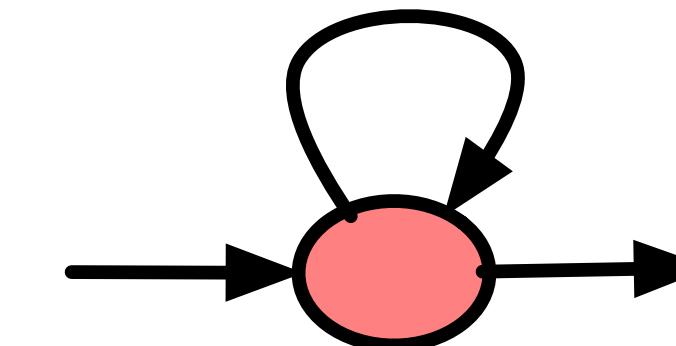
# Need for More Dynamism

Model

static  
computacional graph

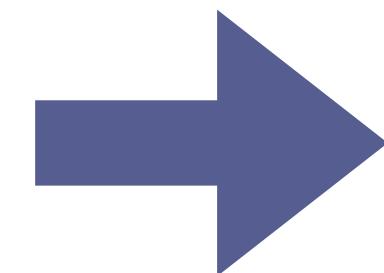


program with  
loops and recursions

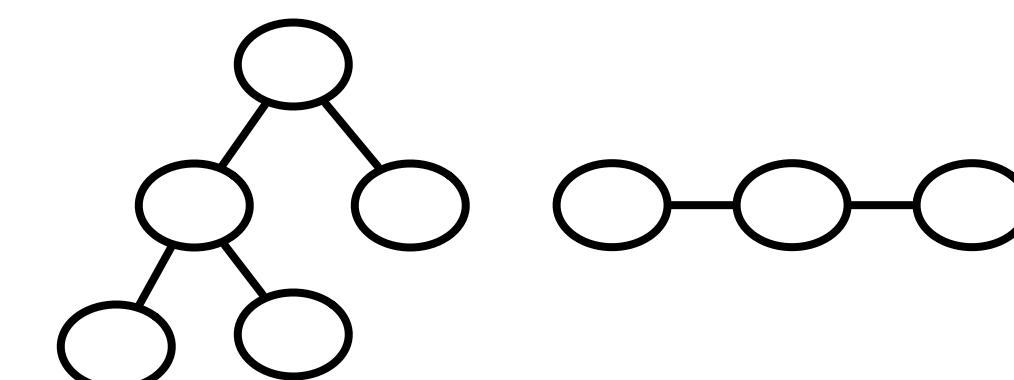


Data

single tensor  
with known shape

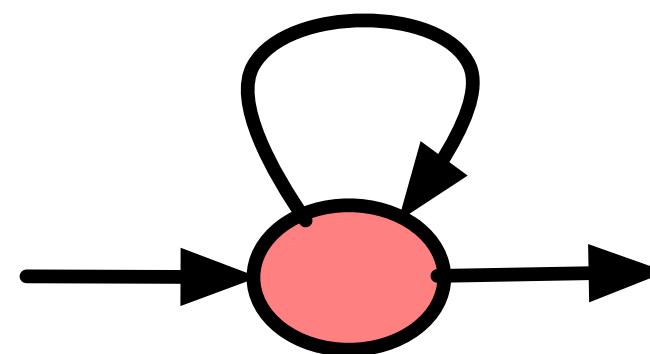


sequence, trees,  
nested data structure

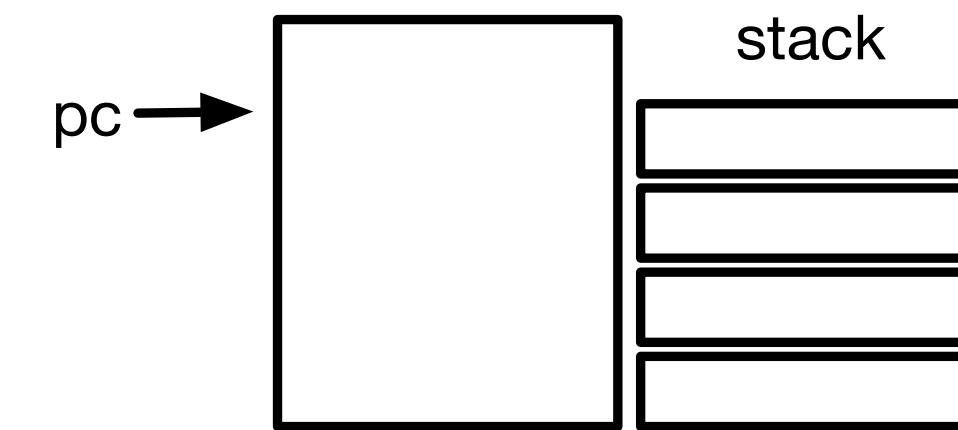


# Relay Virtual Machine

source program



VM bytecode and runtime



Dynamic shape workloads

More runtime objects: Arrays, Tuples, Trees, ADTs

Minimum runtime for dynamic models

# Community Highlights

More **Dynamism**

**Tiny** machine learning

Better core **Infra**

More Specialized **Accelerator Support**

# Machine Learning is Getting into Tiny Devices

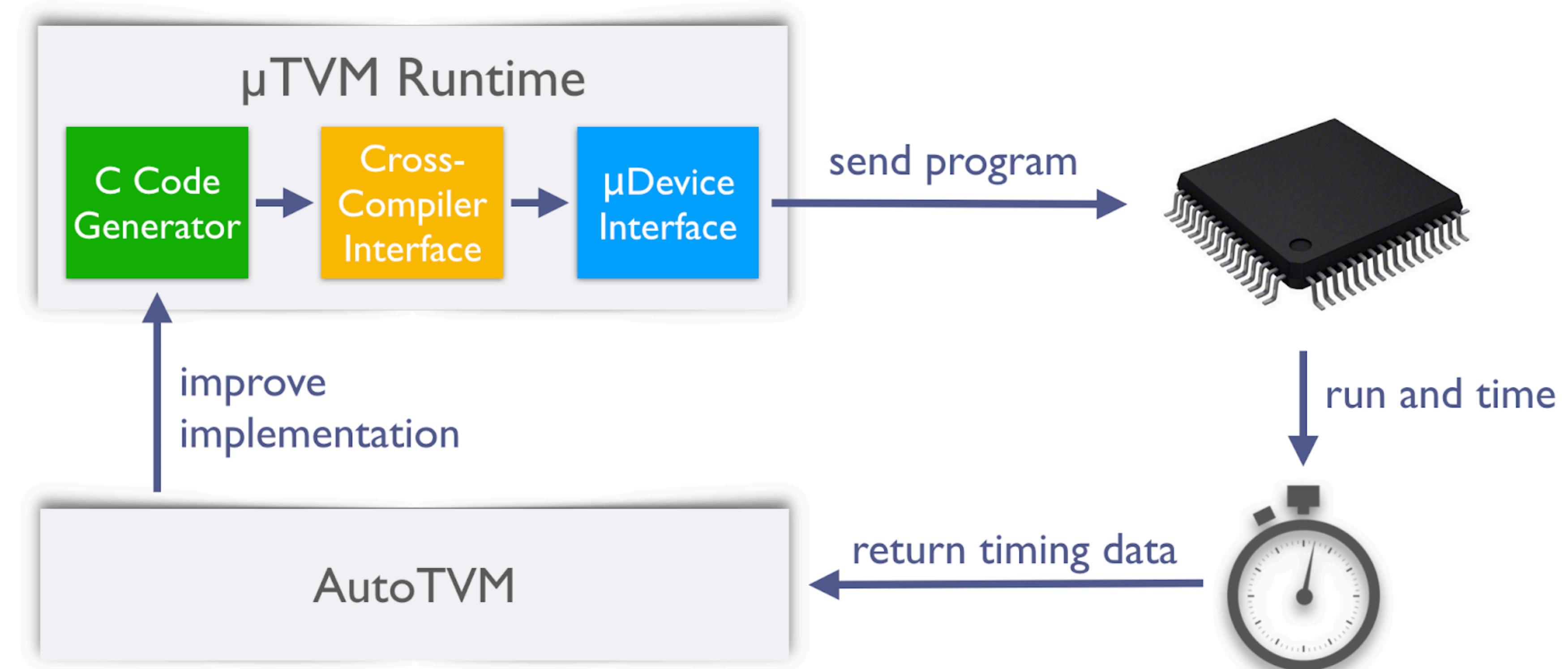
**Challenges: limited resources, OS support**



# uTVM: TVM on bare-metal Devices

Support bare-metal J-TAG devices, **no OS is needed**

ARM Cortex-M  
RISC-V



# Community Highlights

More **Dynamism**

**Tiny** machine learning

Better core **Infra**

More Specialized **Accelerator Support**

# Core Infrastructure

New integer simplification and analysis

Unified runtime object protocol

# Core Infrastructure

New integer simplification and analysis

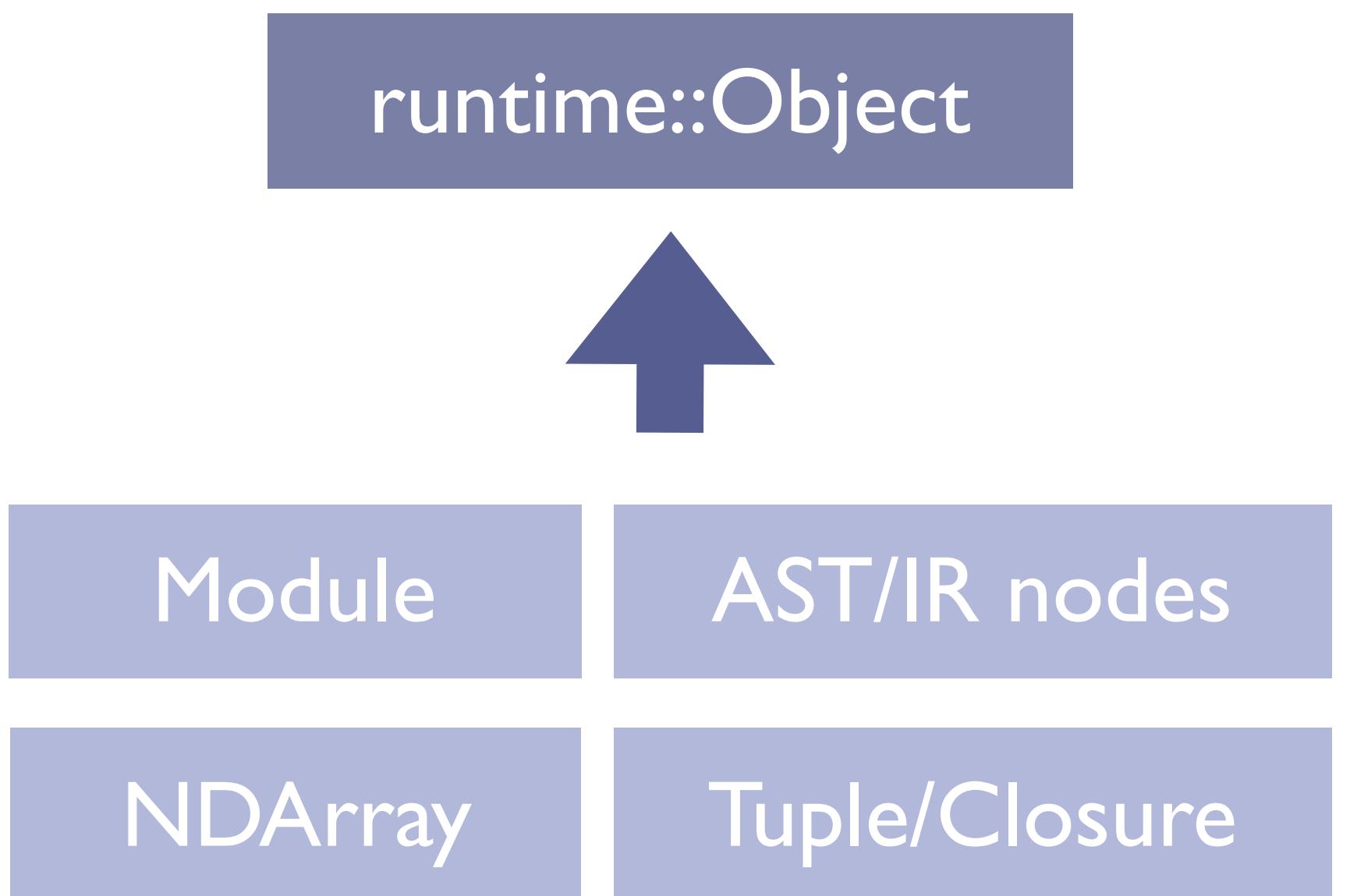
Unified runtime object protocol

Module	AST/IR nodes
NDArray	Tuple/Closure

# Core Infrastructure

New integer simplification and analysis

Unified runtime object protocol



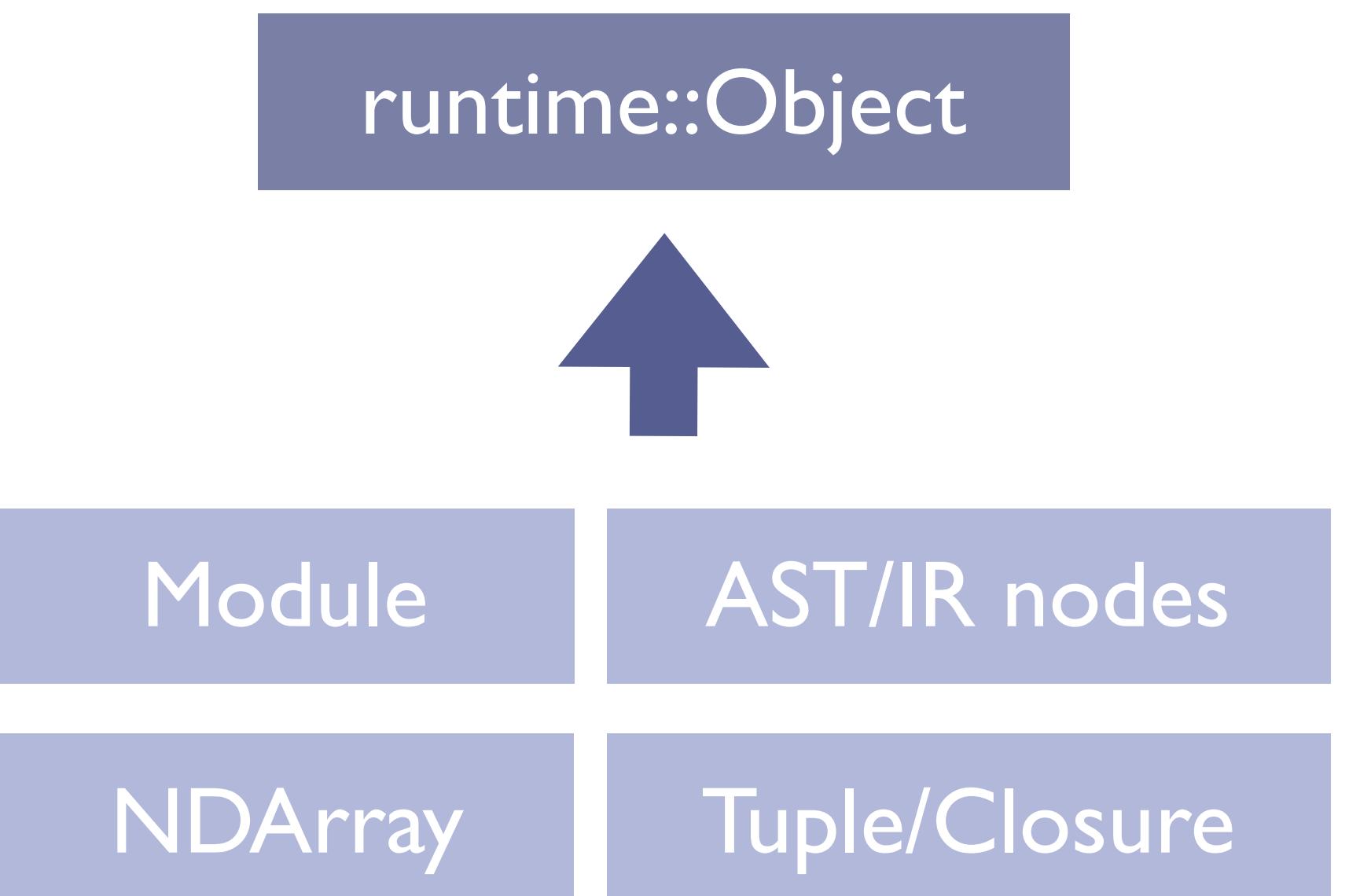
# Core Infrastructure

New integer simplification and analysis

Unified runtime object protocol

Easy to add new objects (trees, graphs)

Cross language support



# Community Highlights

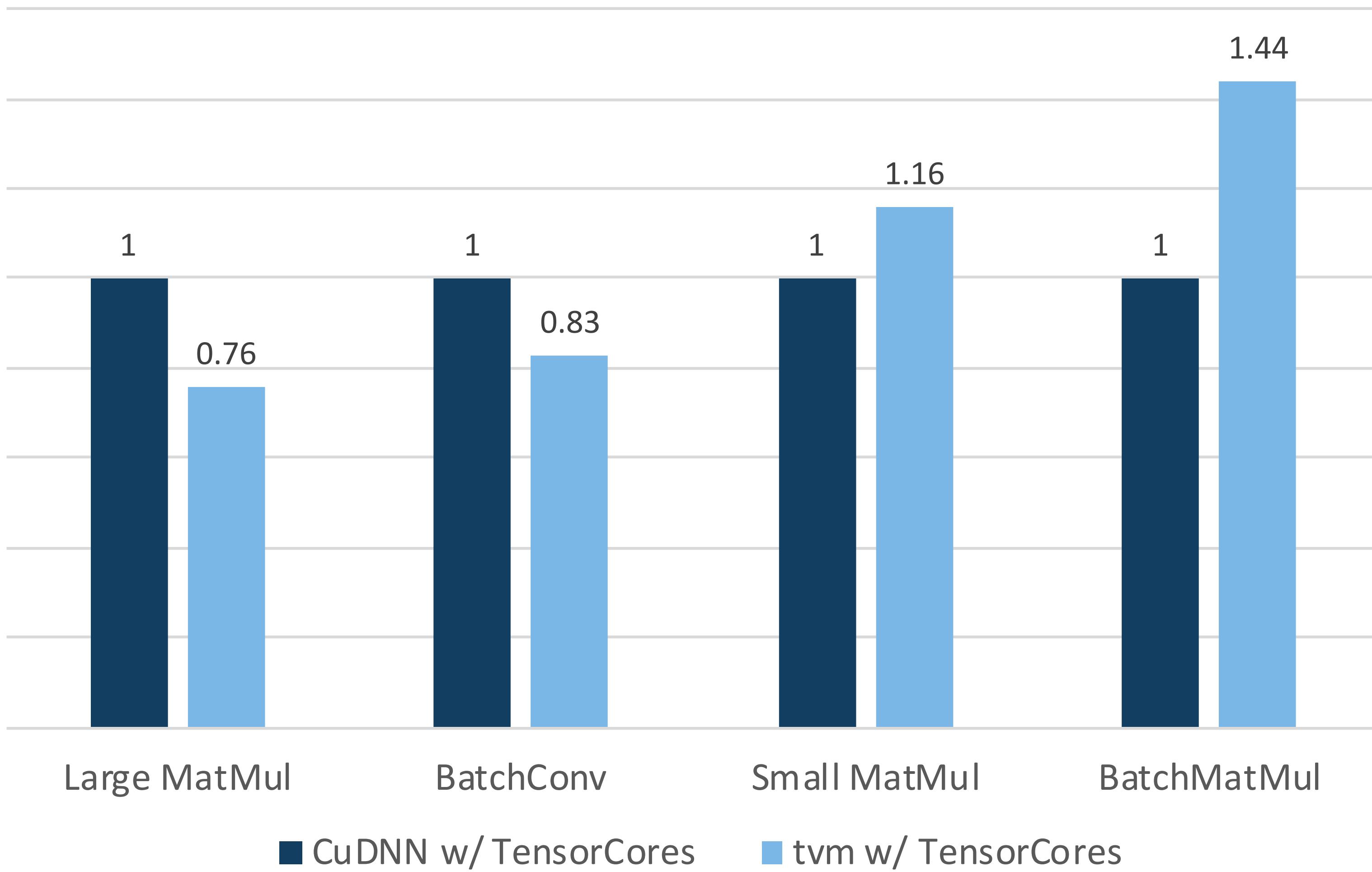
More **Dynamism**

**Tiny** machine learning

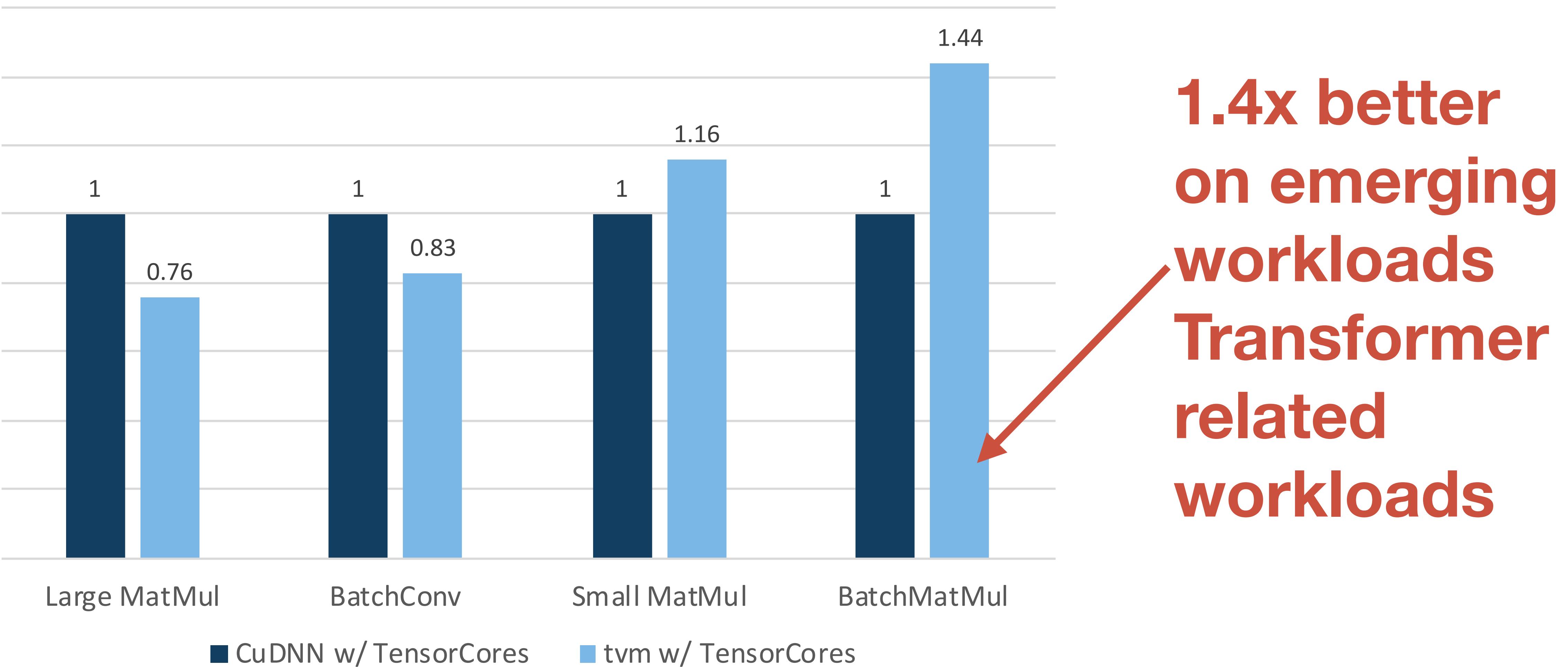
Better core **Infra**

More Specialized **Accelerator Support**

# TVM for TensorCore

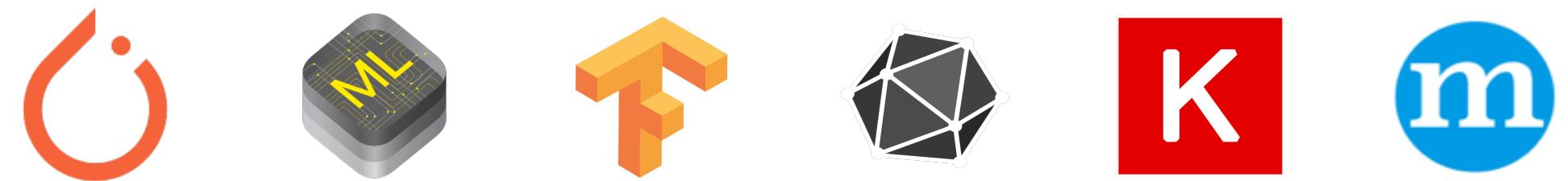


# TVM for TensorCore



Credit: Siyuan Feng

# TSIM: Support for Future Hardware



Current TVM Stack

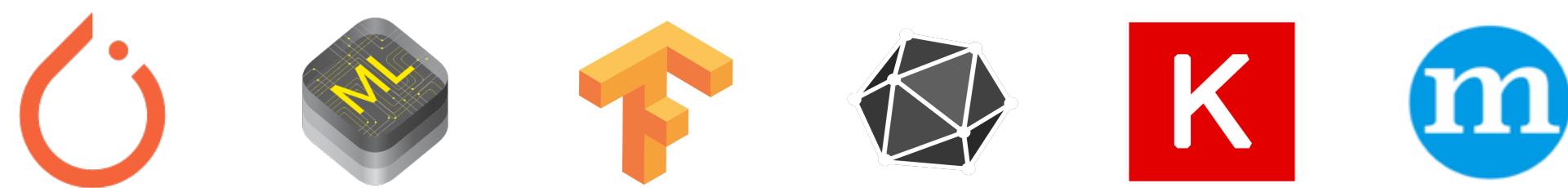
New NPU Runtime

TSIM Driver



Credit: Luis Vega, Thierry Moureau

# TSIM: Support for Future Hardware



Current TVM Stack

New NPU Runtime

New Hardware Design in Verilog

TSIM Driver

TSIM Binary

Verilator



Credit: Luis Vega, Thierry Moureau

# TSIM: Support for Future Hardware



Current TVM Stack

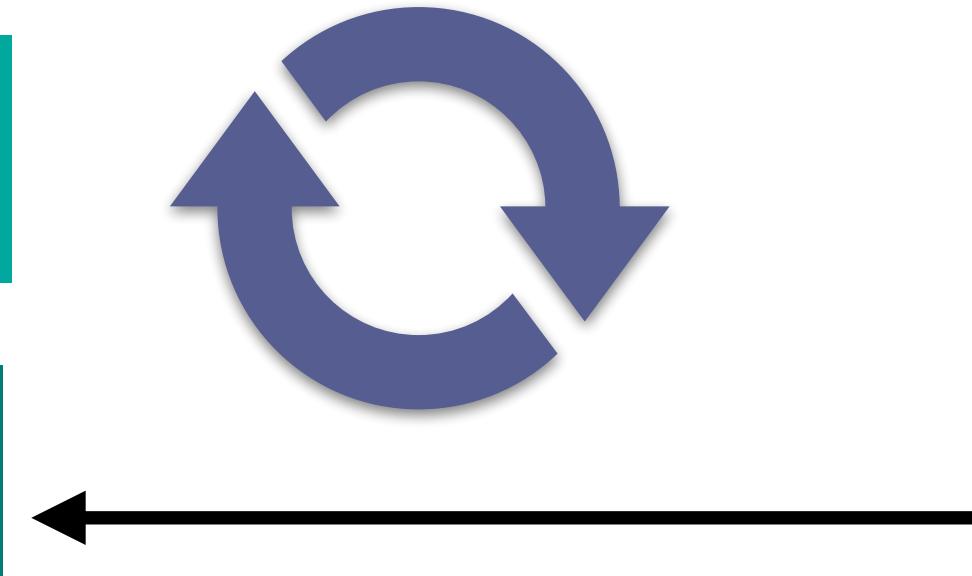
New NPU Runtime

TSIM Driver

TSIM Binary

New Hardware Design in Verilog

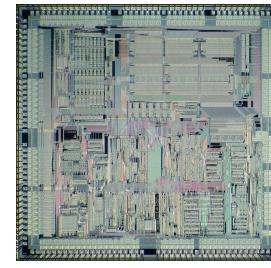
Verilator



# Unified Runtime For Heterogeneous Devices

Device Drivers

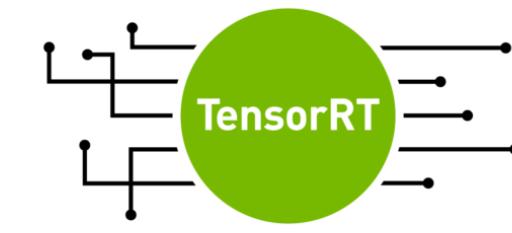
**NPU Driver**



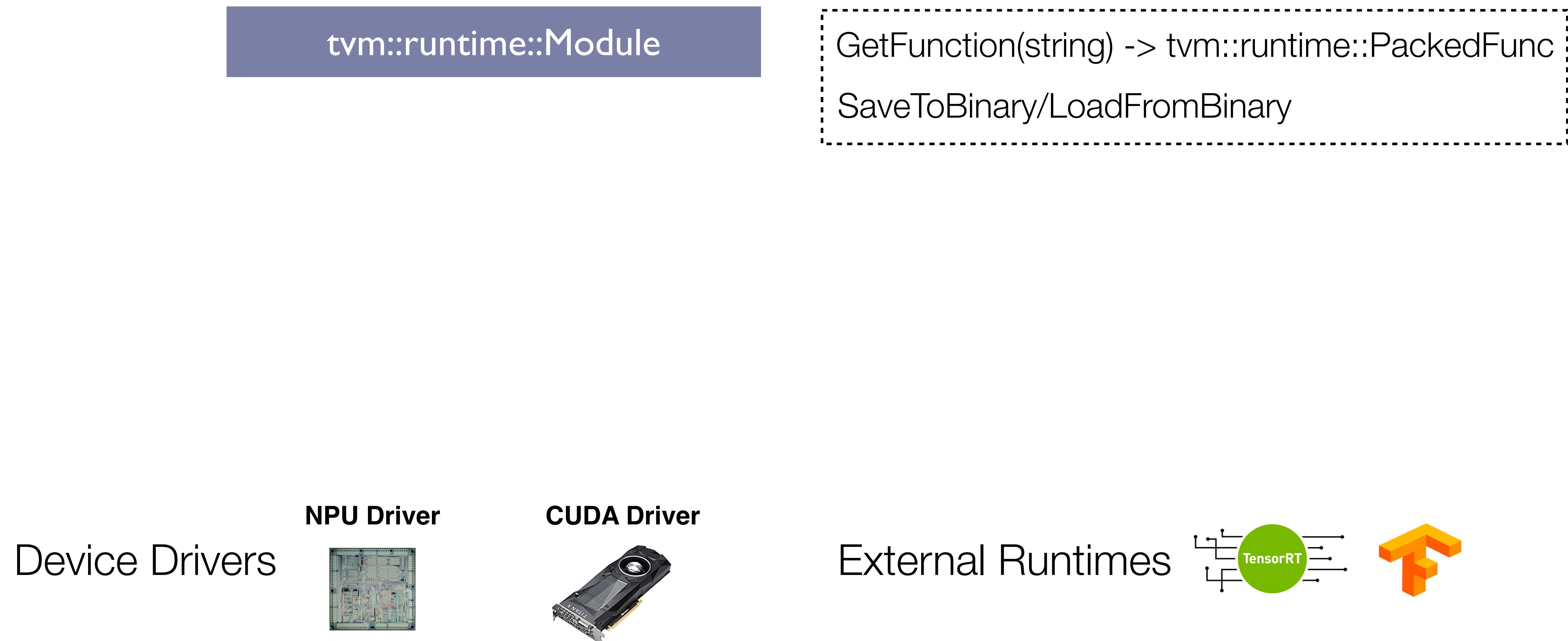
**CUDA Driver**



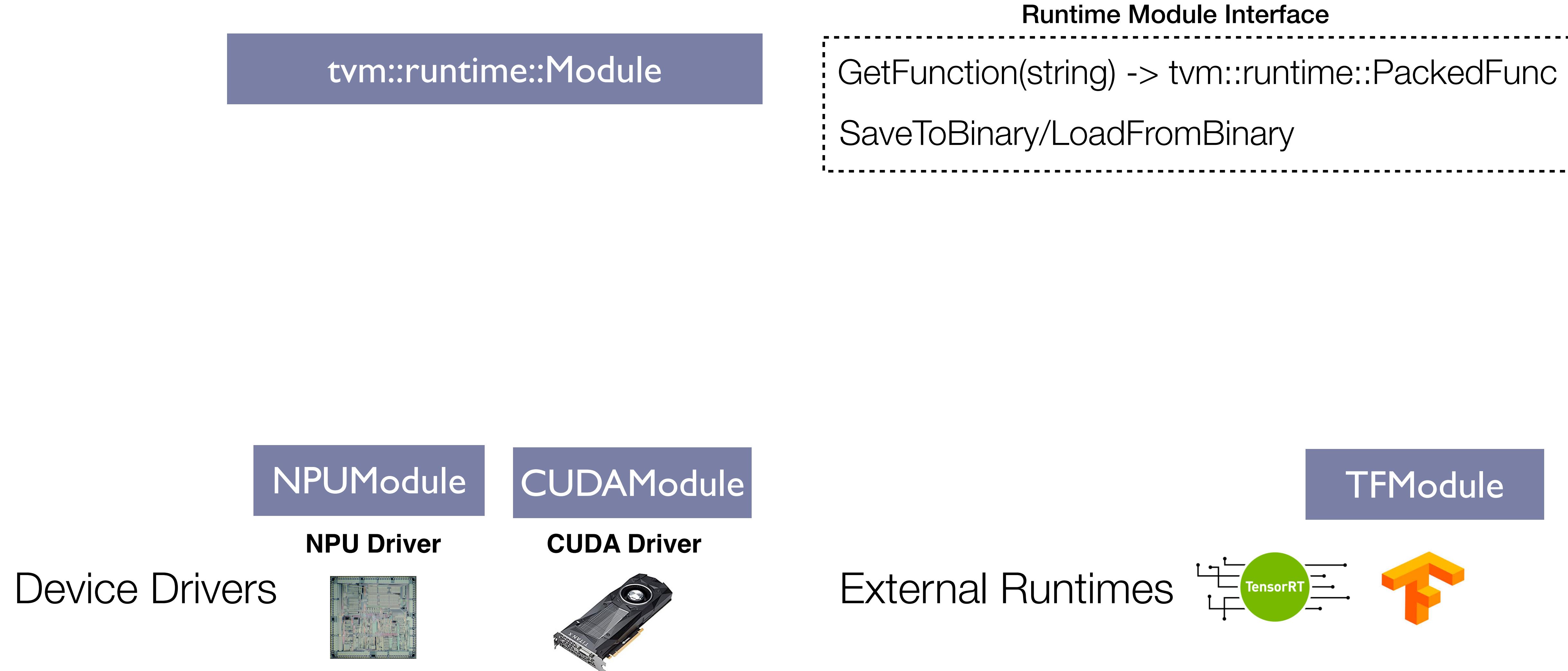
External Runtimes



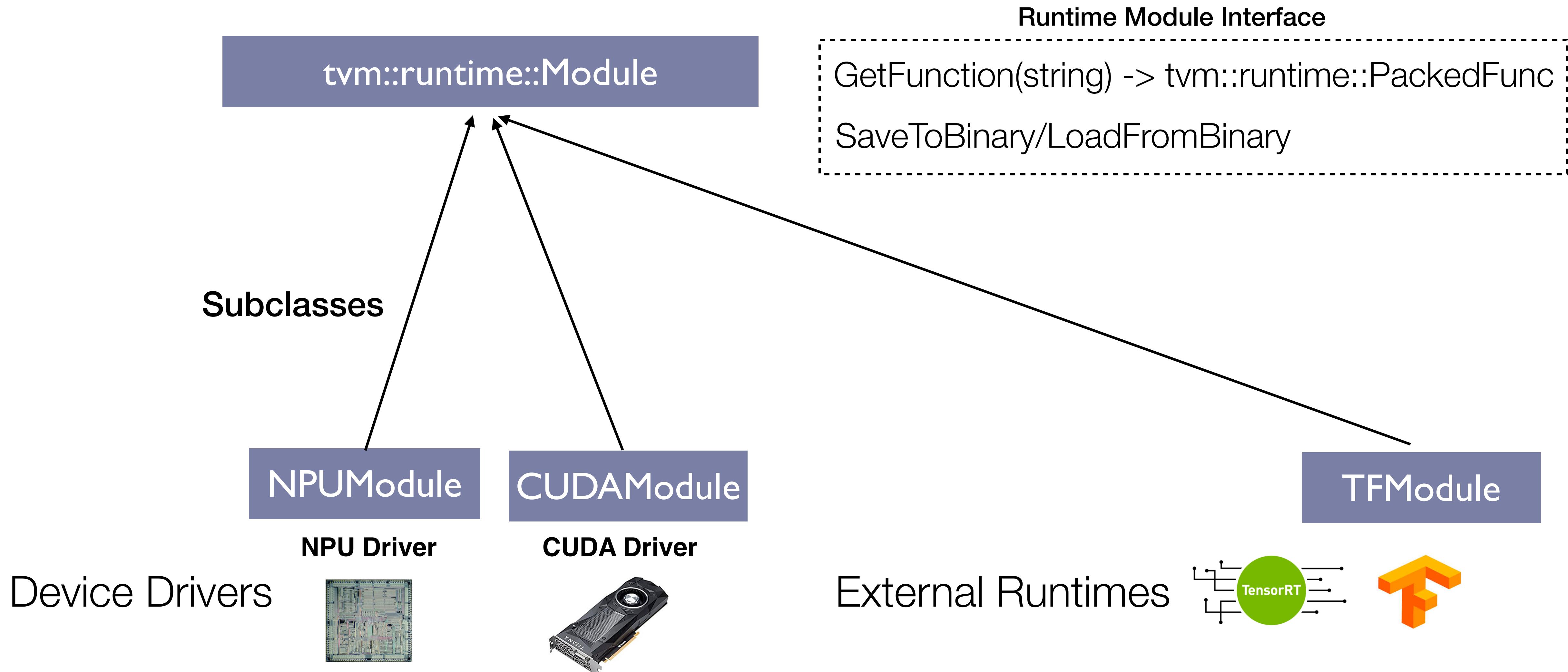
# Unified Runtime For Heterogeneous Devices



# Unified Runtime For Heterogeneous Devices



# Unified Runtime For Heterogeneous Devices



# Unified Runtime Benefit

Unified library packaging

```
mod.export_library("mylib.so")
```

Free API (Py/Java/Go)

```
lib = tvm.module.load("mylib.so")
func = lib["npufunction0"]
func(a, b)
```

Automatic RPC Support

```
remote = tvm.rpc.connect(board_url, port)
remote.upload("mylib.so")
remote_mod = remote.load_module("mylib.so")
func = remote_mod["npufunction0"]
func(remote_a, remote_b)
```

# Community

# Open Source Community



Incubated as Apache TVM. Independent governance,  
allowing competitors to collaborate.

# Open Source Community



Incubated as Apache TVM. Independent governance,  
allowing competitors to collaborate.

Open Source Code

Open Development

Open Governance

# Open Source Community



Incubated as Apache TVM. Independent governance,  
allowing competitors to collaborate.

# Open Source Community



Incubated as Apache TVM. Independent governance, allowing competitors to collaborate.

## Growing Developer Community

22 committers, 47 reviewers, 295 contributors

# Open Source Community



Incubated as Apache TVM. Independent governance, allowing competitors to collaborate.

## Growing Developer Community

22 committers, 47 reviewers, 295 contributors

**~70% growth since TVM Conf 2018**

# Open Source Community



Incubated as Apache TVM. Independent governance, allowing competitors to collaborate.

## Growing Developer Community

22 committers, 47 reviewers, 295 contributors

**~70% growth since TVM Conf 2018**

## Monthly Statistics

~50 authors, ~140 PRs, ~1000 discuss forum posts

# More on TVM Community Conference

TensorCore

ARM

Automation

Driving  
Sparse

DSPs

Quantization  
Dynamic  
Workloads

Binary Neural  
Networks

RISCV

...

**Videos and slides available at <http://tvmconf.org/>**

# Open Source Community

Open source: ~280 contributors from UW, Berkeley, Cornell, UCLA, Amazon, Huawei, NTT, Facebook, Microsoft, Qualcomm, Alibaba, Intel, ...

# Open Source Community

Open source: ~280 contributors from UW, Berkeley, Cornell, UCLA, Amazon, Huawei, NTT, Facebook, Microsoft, Qualcomm, Alibaba, Intel, ...

Used in production

# Open Source Community

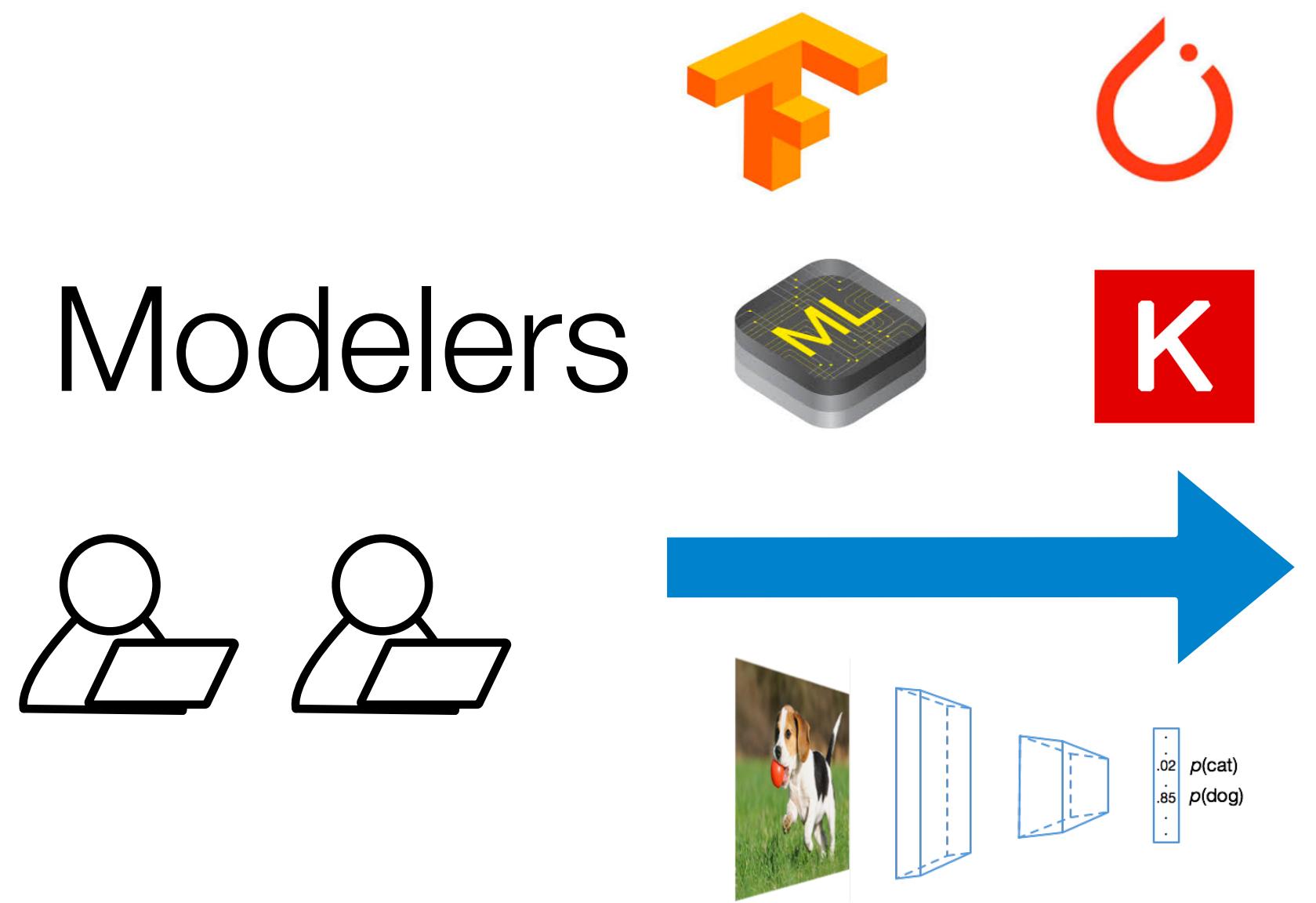
Open source: ~280 contributors from UW, Berkeley, Cornell, UCLA, Amazon, Huawei, NTT, Facebook, Microsoft, Qualcomm, Alibaba, Intel, ...

Used in production

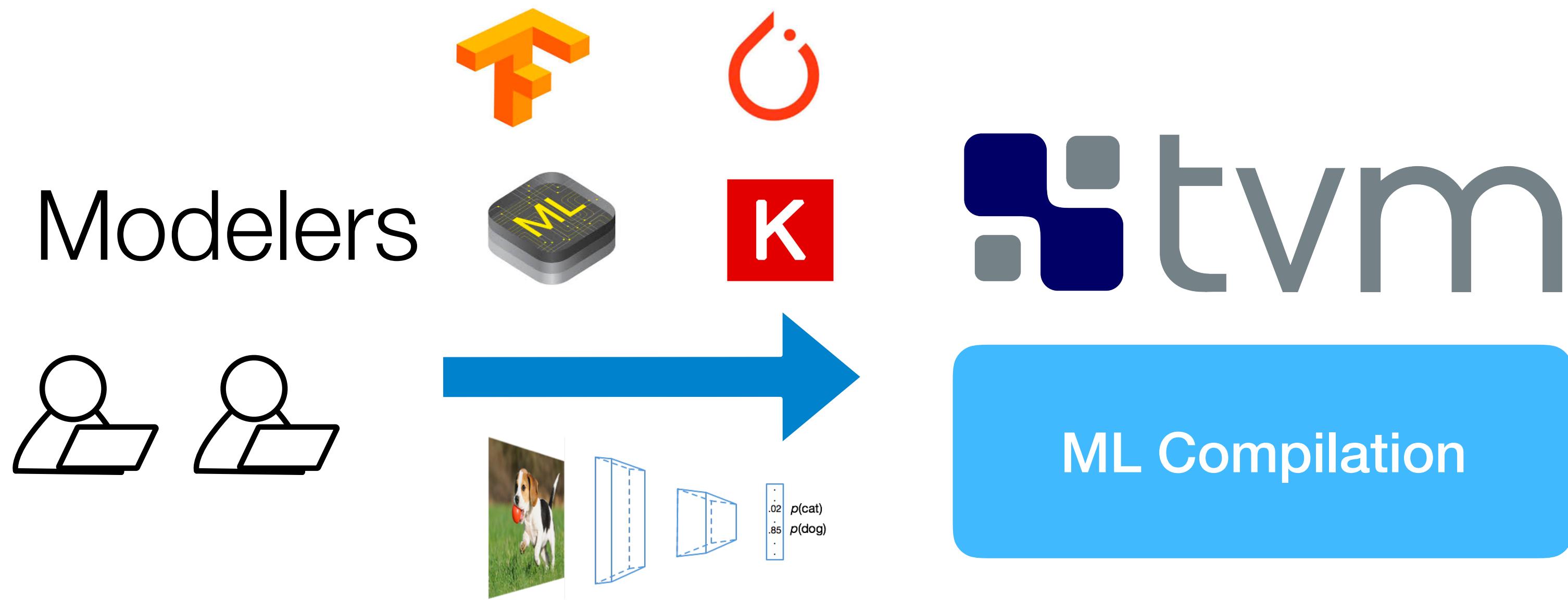


Incubated as Apache TVM recently. Independent governance, allowing competitors to collaborate.

# ML Deployment Flow with Automated Compiler



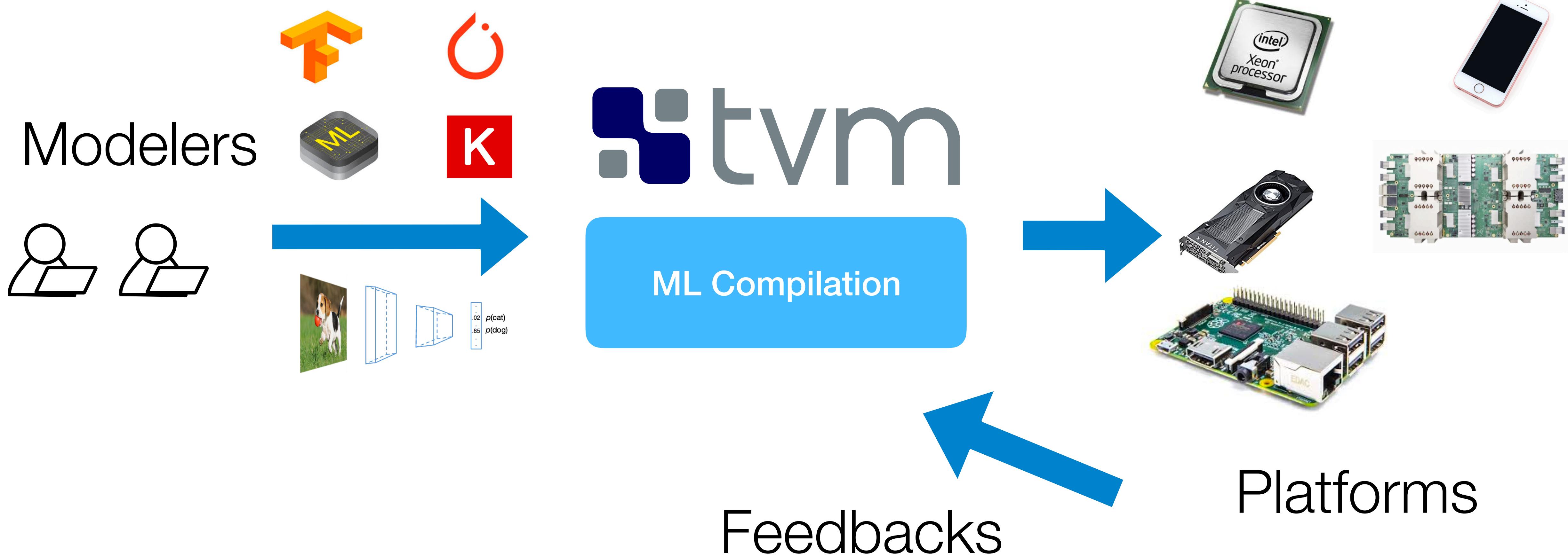
# ML Deployment Flow with Automated Compiler



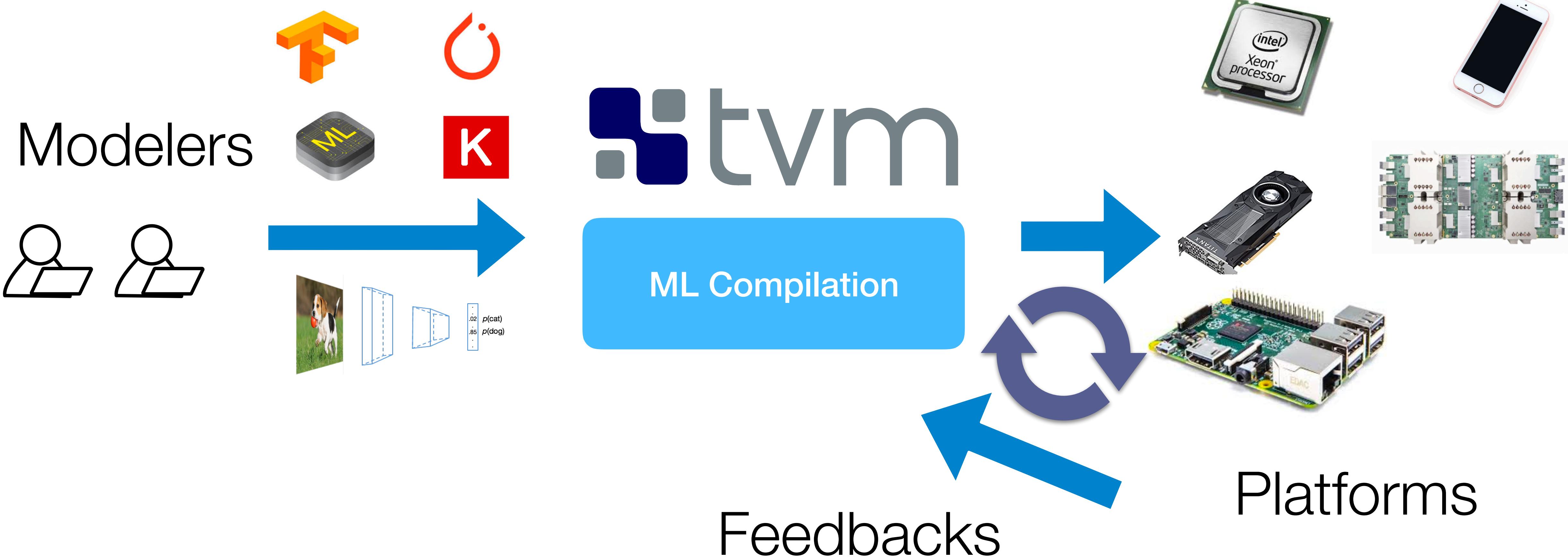
# ML Deployment Flow with Automated Compiler



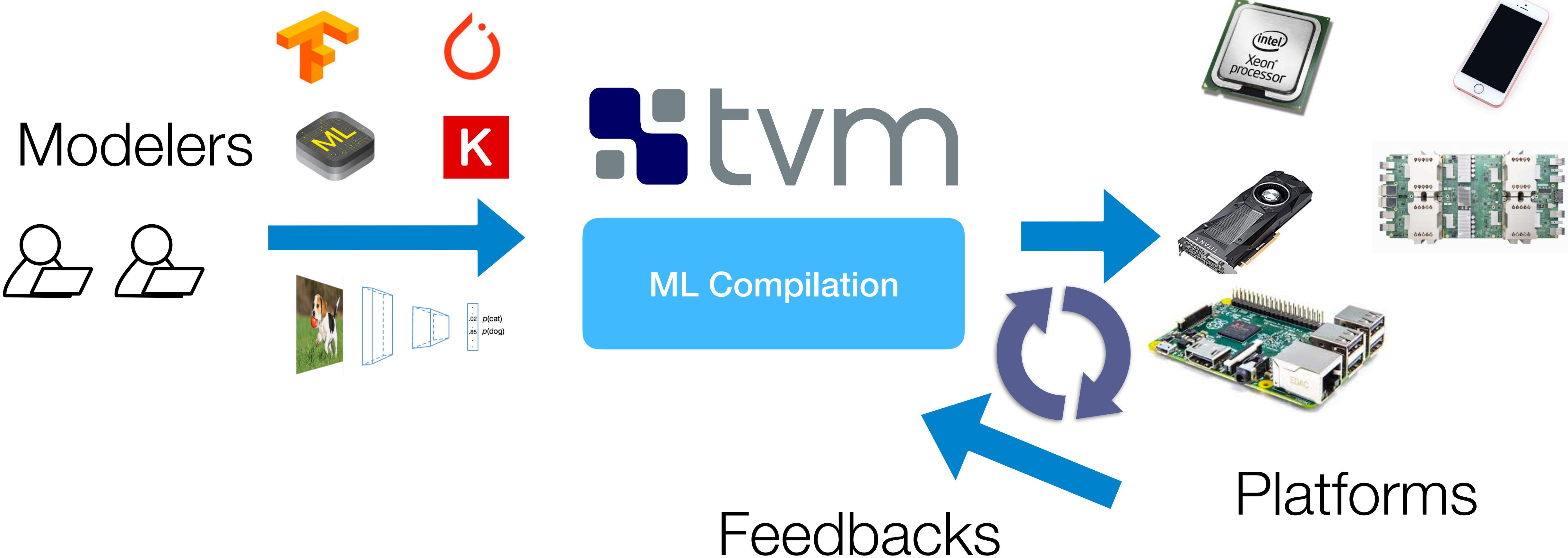
# ML Deployment Flow with Automated Compiler



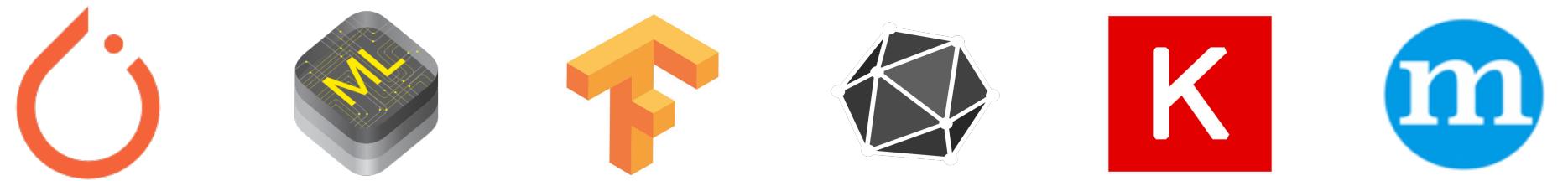
# ML Deployment Flow with Automated Compiler



# ML Deployment Flow with Automated Compiler



# Full Stack Learning-based Learning System



High-Level Differentiable IR

Tensor Expression IR

LLVM, CUDA

VTA



# Full Stack Learning-based Learning System

