

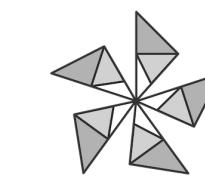
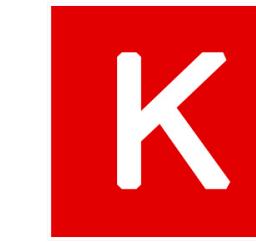
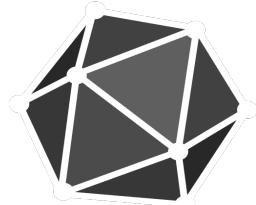
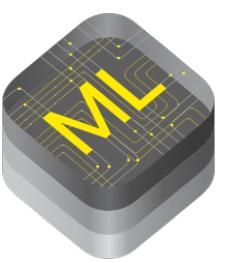
# TVM: Where are we going

Tianqi Chen

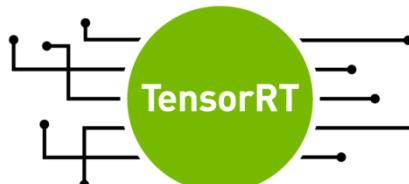


# Current Deep Learning Landscape

Frameworks and  
Inference engines



ONNX  
RUNTIME



DL Compilers



Kernel Libraries

CuDNN

NNPack

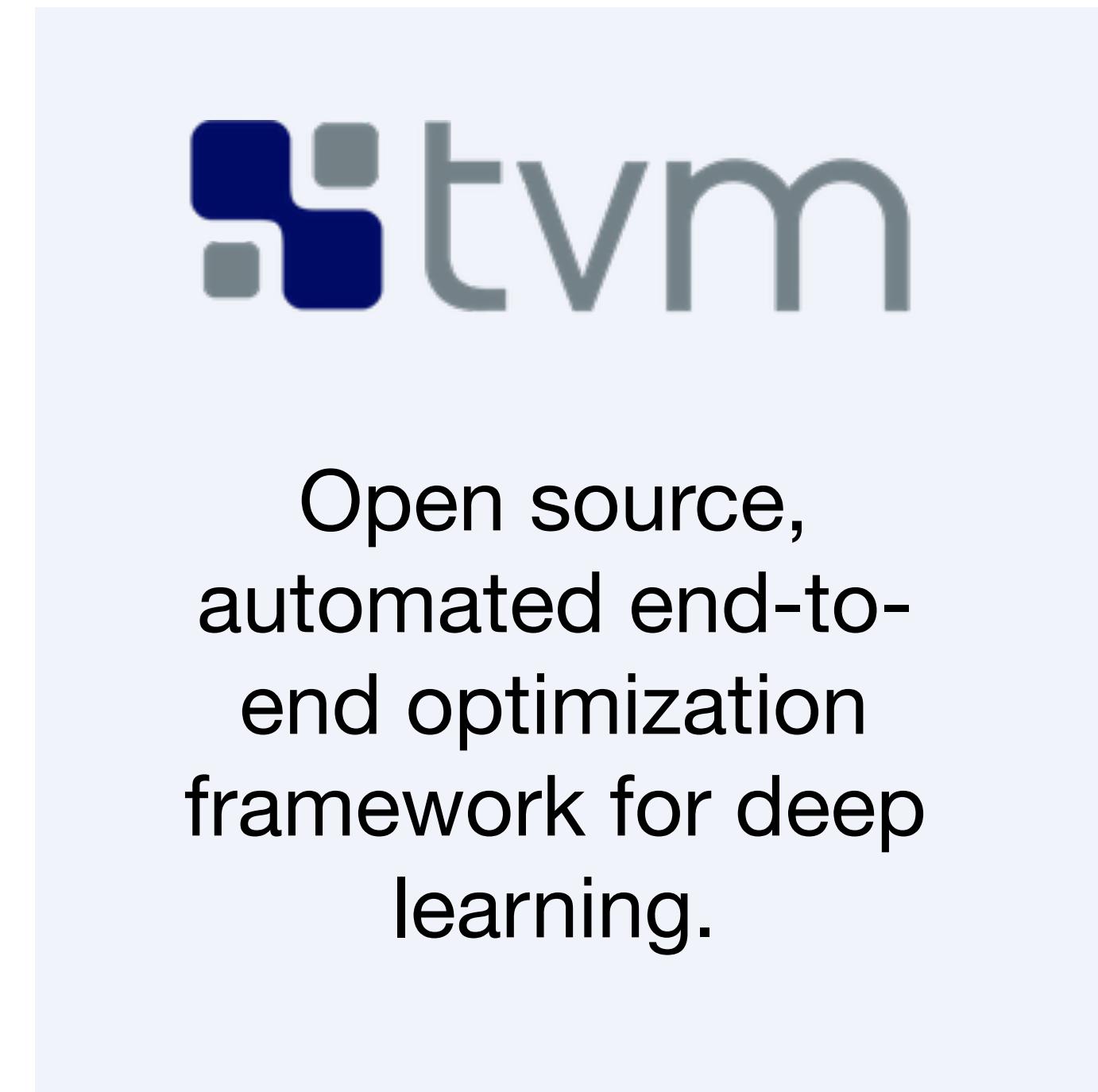
MKL-DNN

Hand optimized

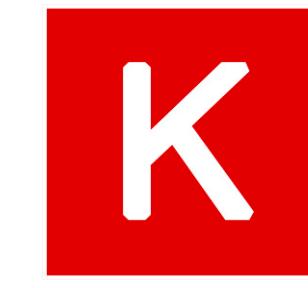
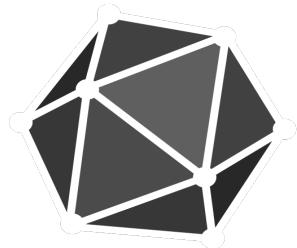
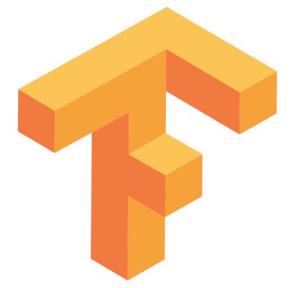
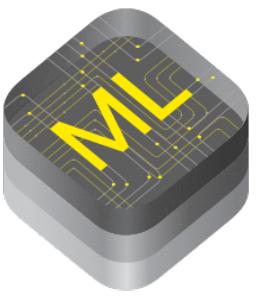
Hardware



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



# TVM Stack



**Optimization**

High-Level Differentiable IR

Tensor Expression and Optimization Search Space

LLVM, CUDA, Metal

VTA

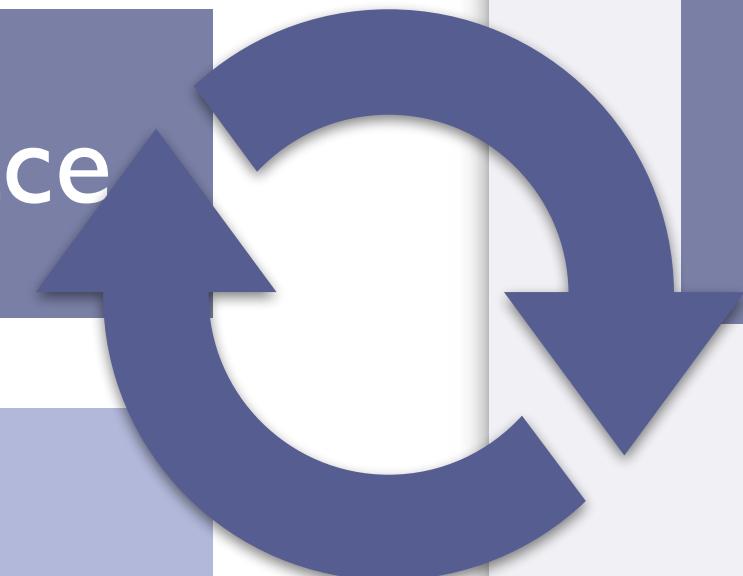


Edge  
FPGA

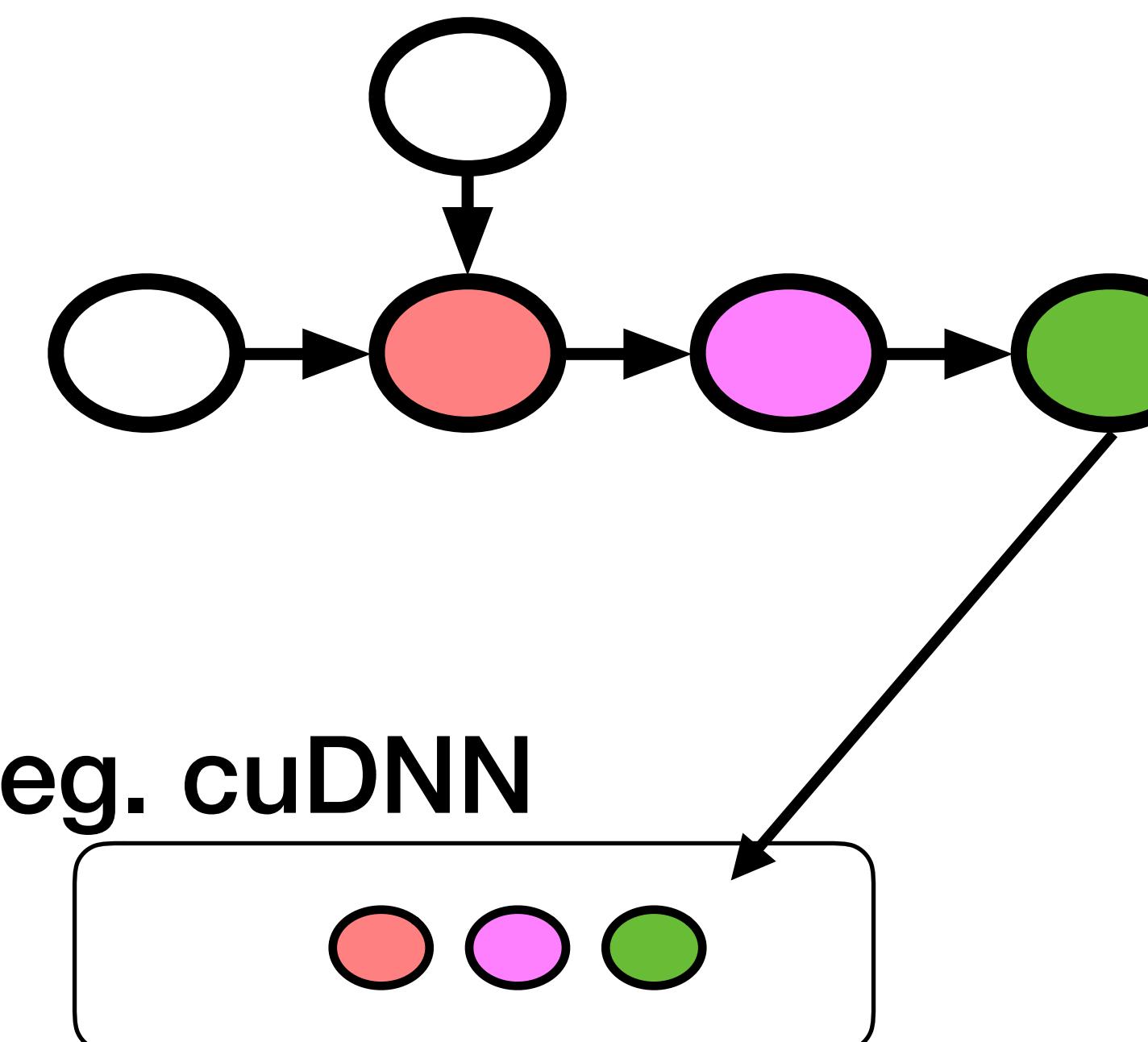
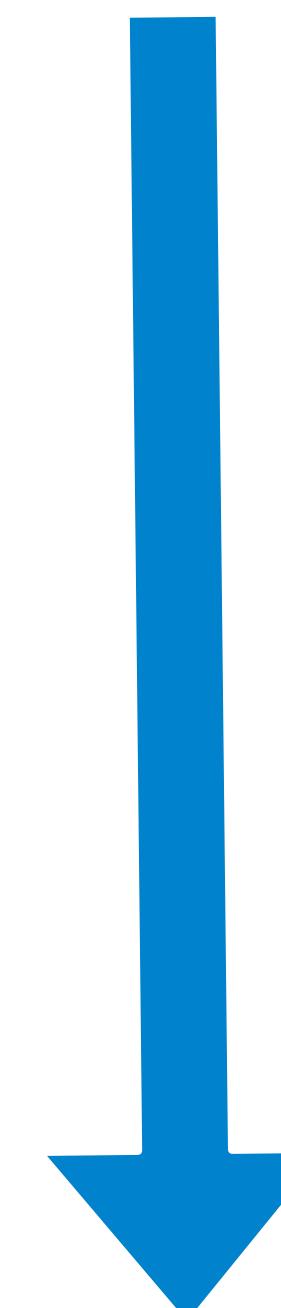
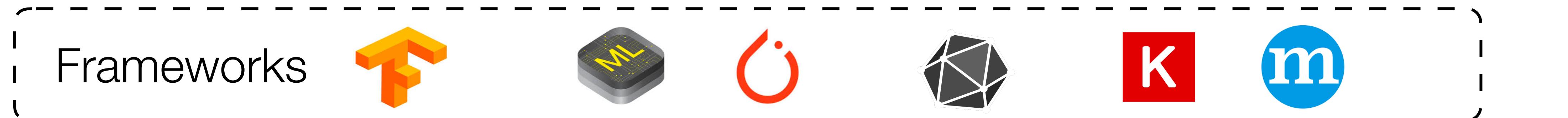
Cloud  
FPGA

ASIC

Device Fleet



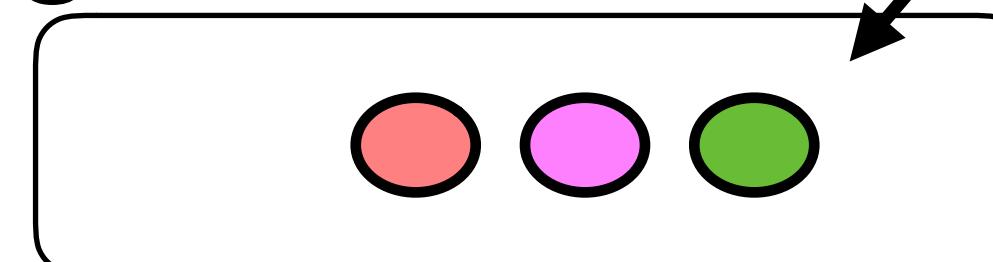
# Existing Deep Learning Frameworks



High-level data flow graph

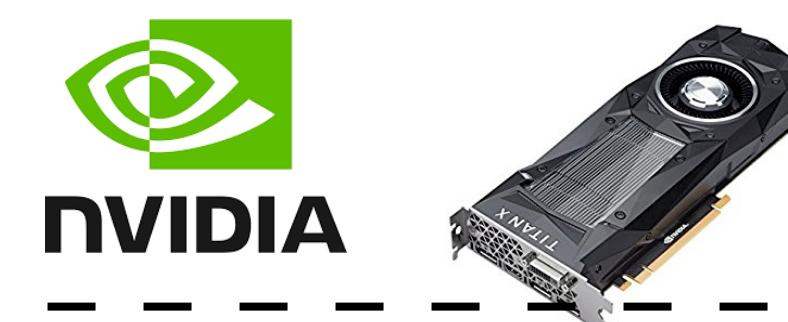
Primitive Tensor operators such as Conv2D

eg. cuDNN

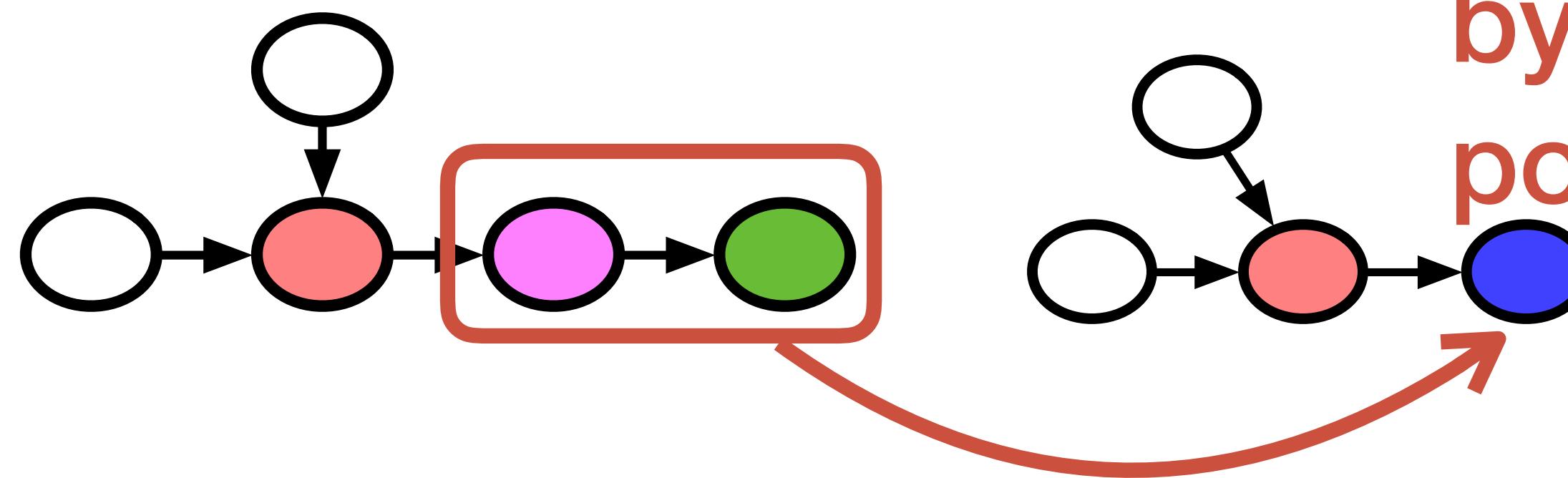
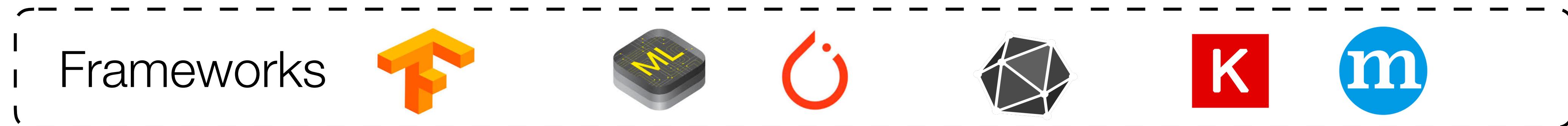


Offload to heavily optimized  
DNN operator library

Hardware



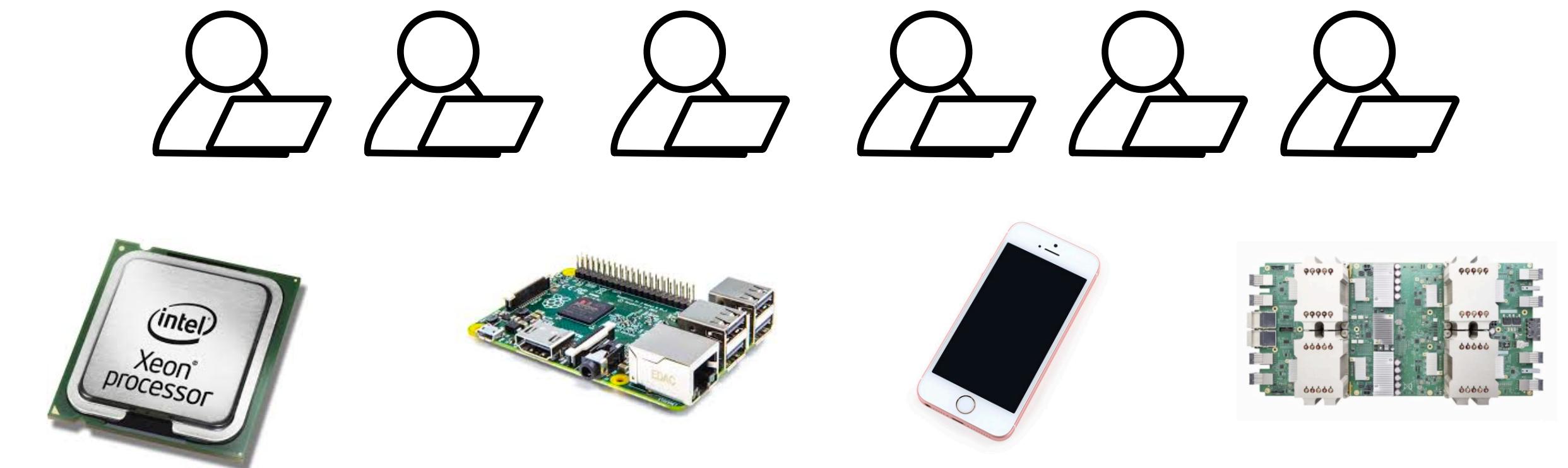
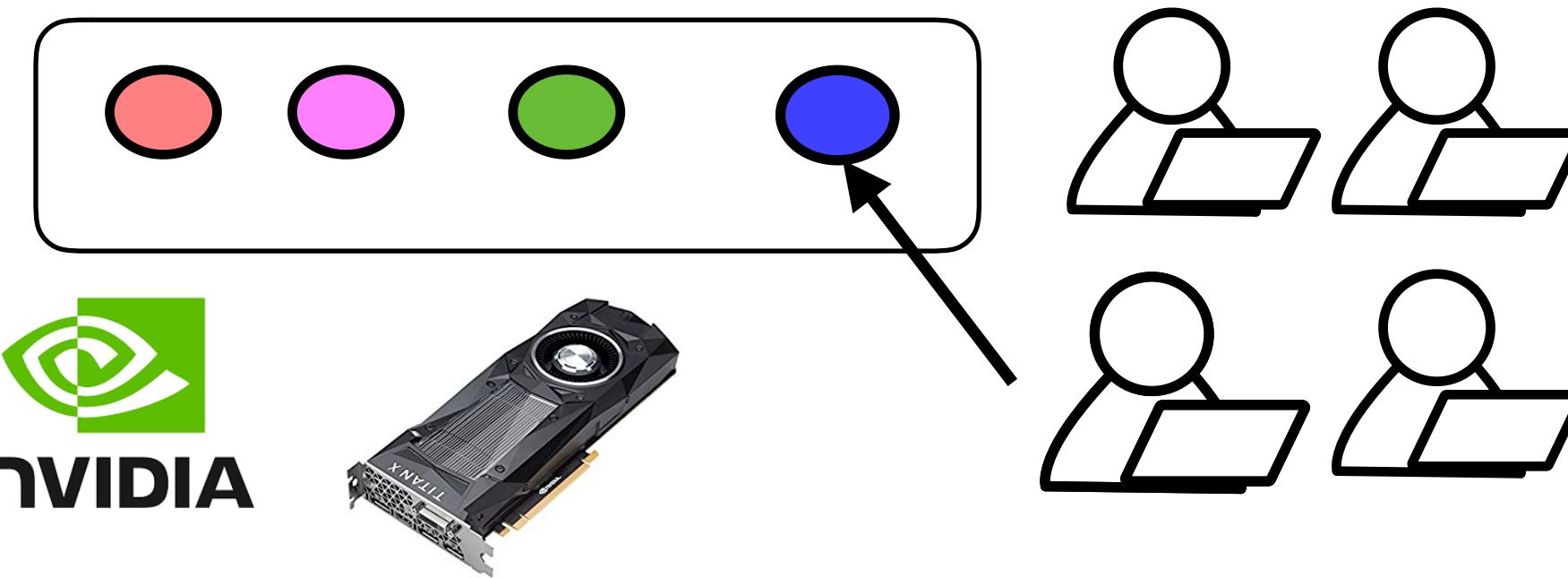
# Limitations of Existing Approach



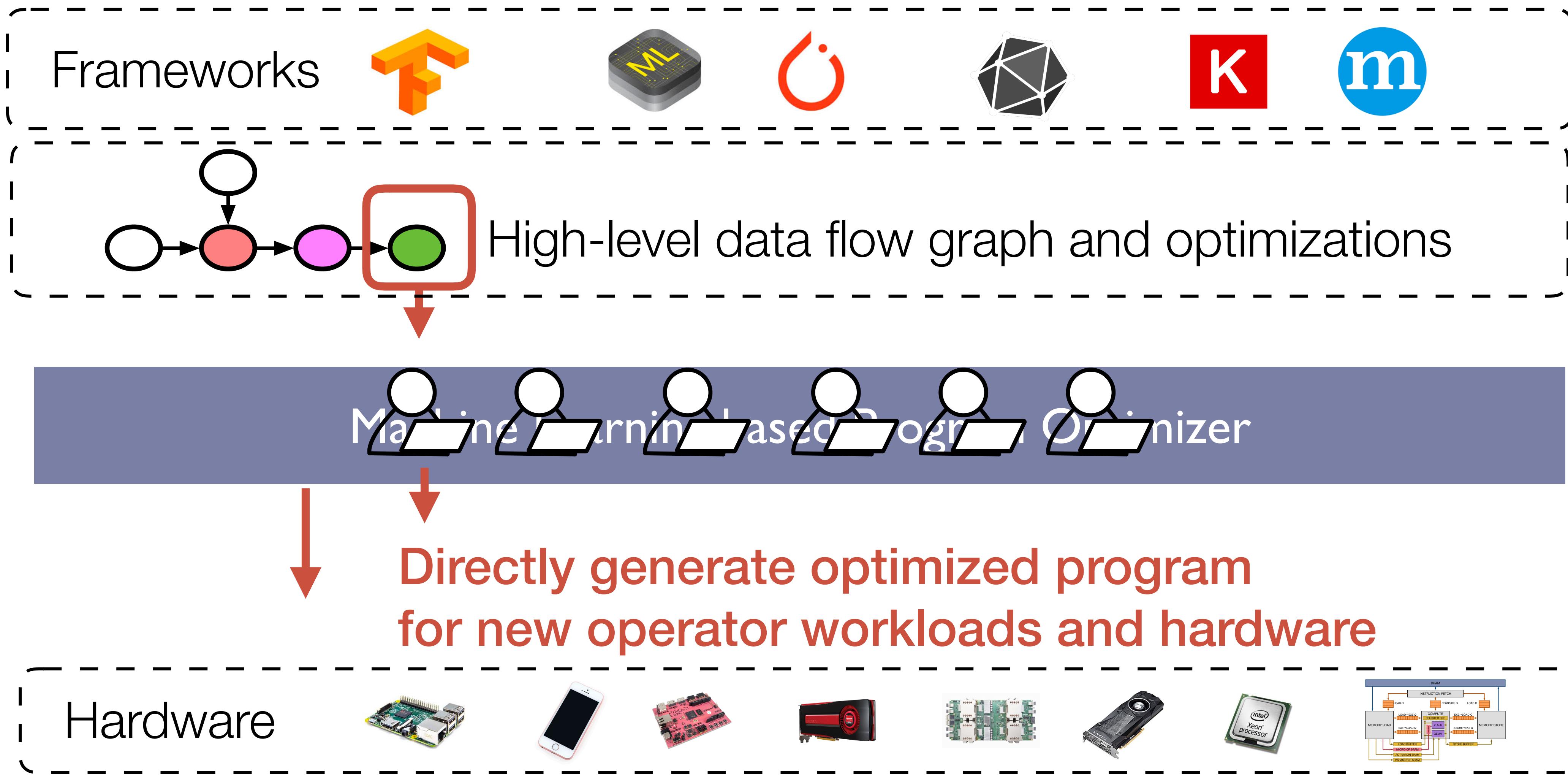
New operator introduced  
by operator fusion optimization  
potential benefit: 1.5x speedup

Engineering intensive

cuDNN



# TVM: Learning-based Learning System



# Why Automation is the Future

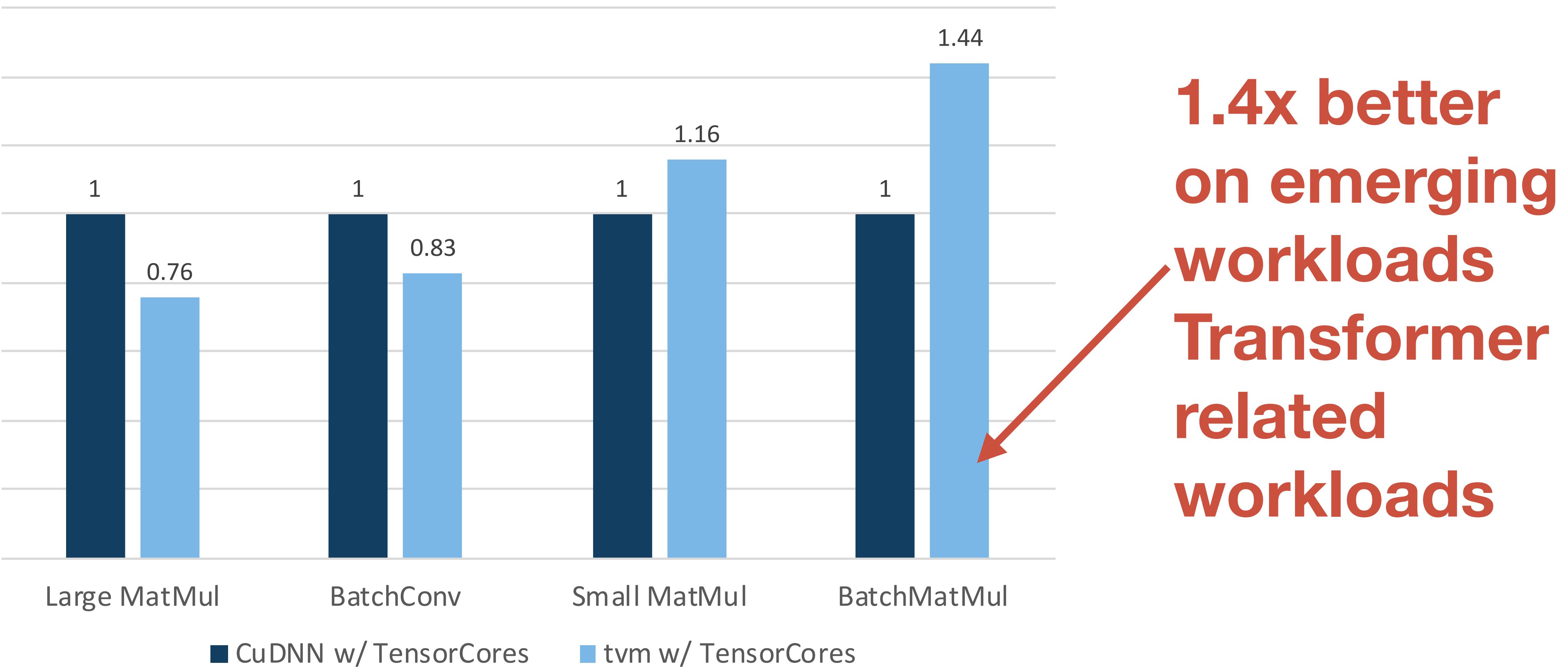
Clear winner on emerging models in product

Competitive on benchmarking type model

Quickly enables other optimizations: fusion, layout, parallelization

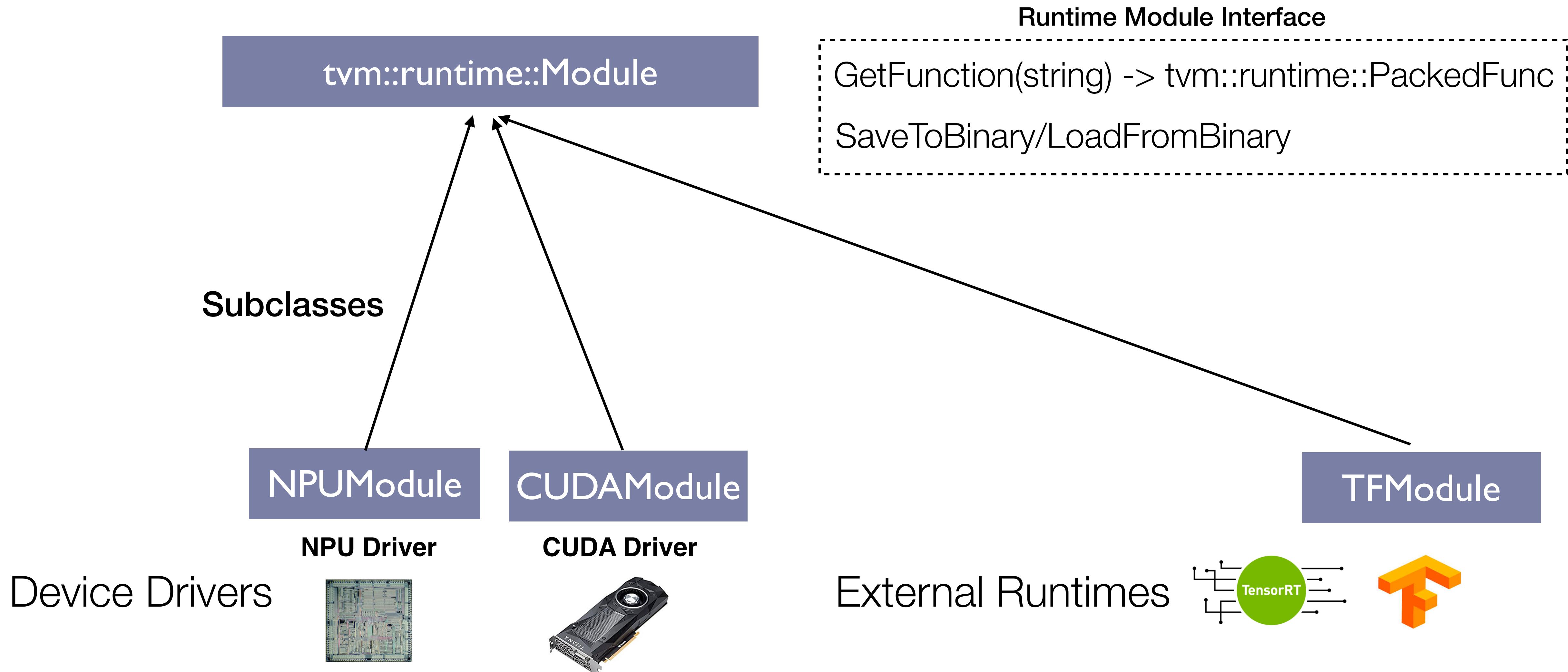
Portable performance across devices

# Why Automation is the Future



Where are we going

# 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)
```

# Virtual Machine: Supporting Dynamic Workload

Dynamic shape workloads

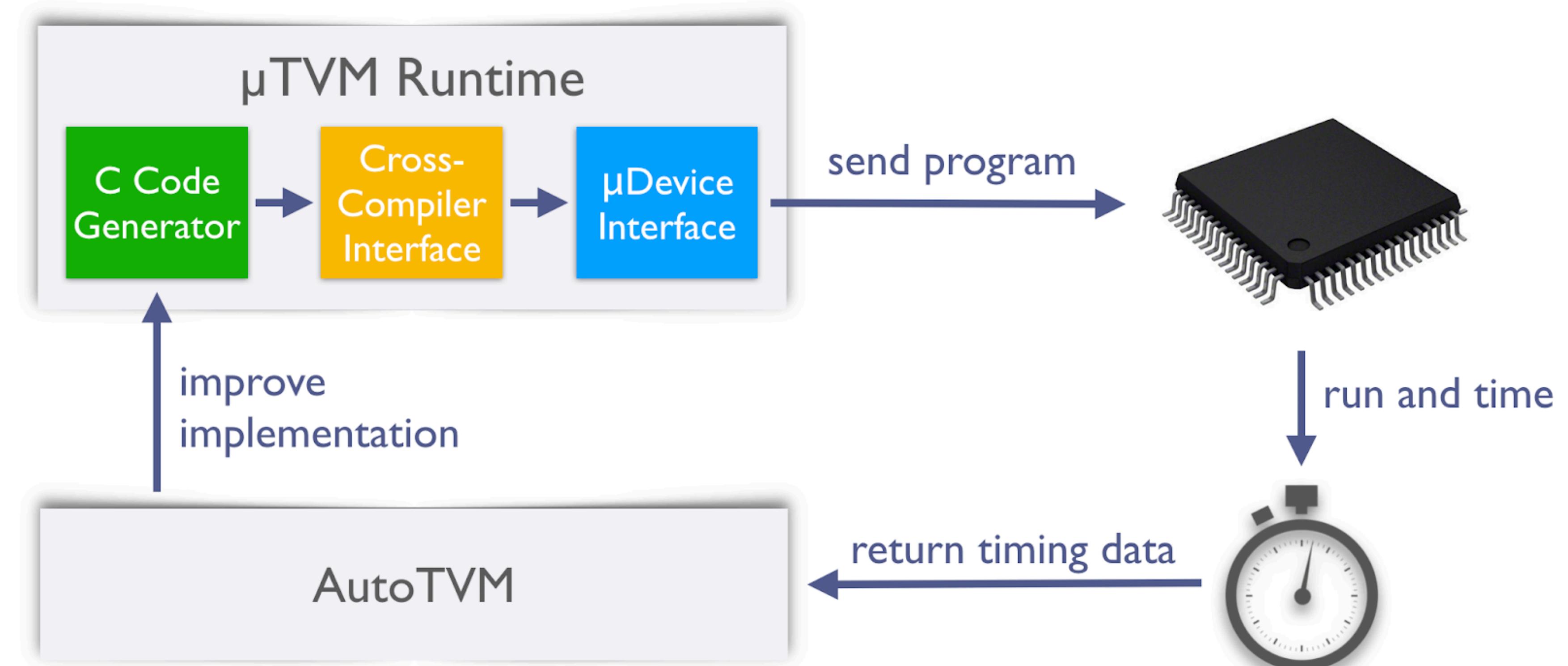
More runtime objects: Arrays, Tuples, Trees, ADTs

Minimum runtime for dynamic models

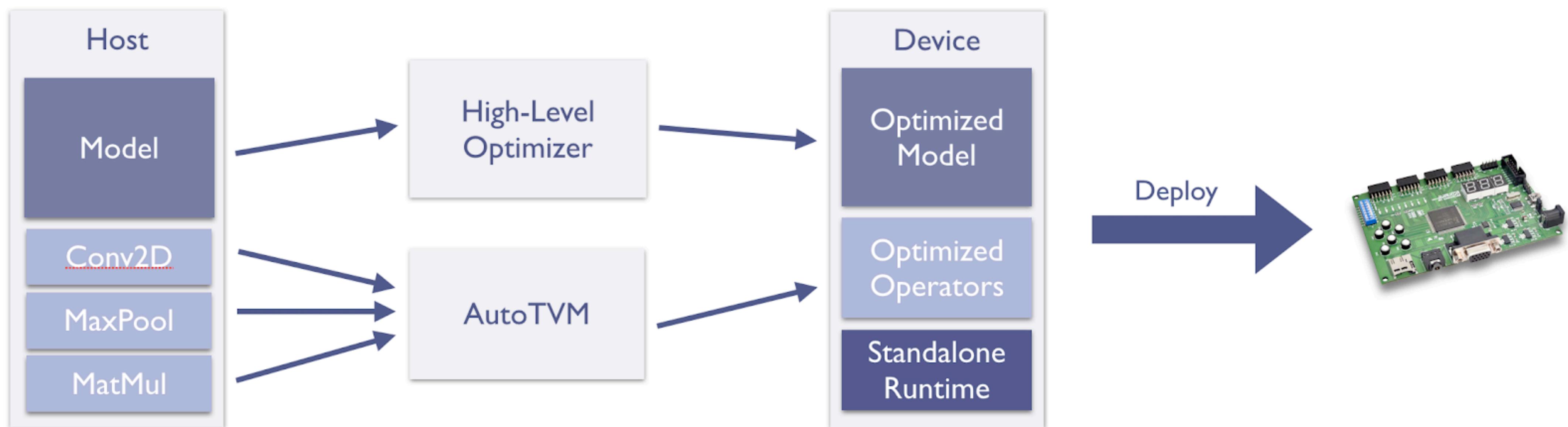
# uTVM: TVM on bare-metal Devices

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

ARM Cortex-M  
RISC-V



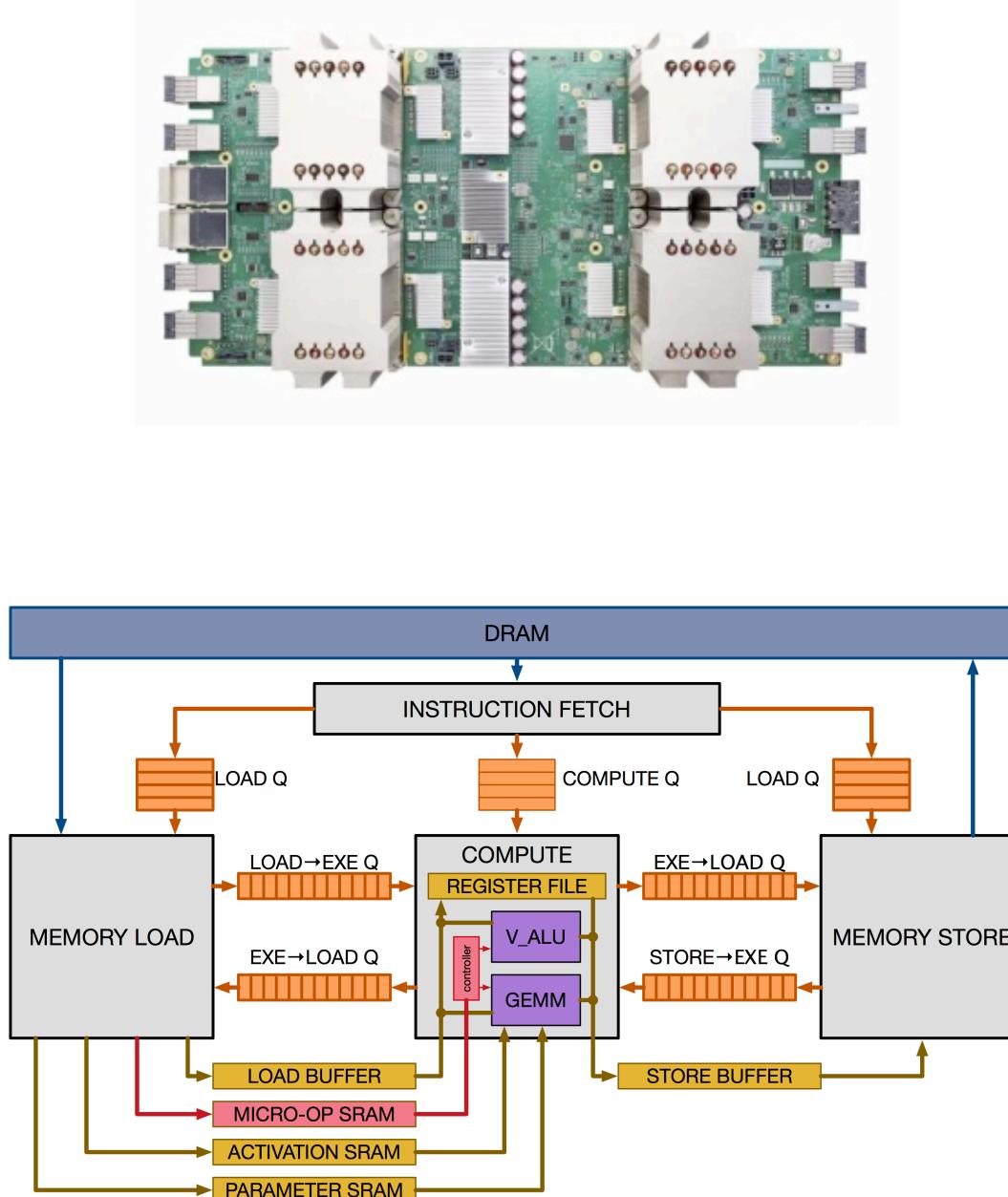
# uTVM upcoming: Self Hosted Runtime



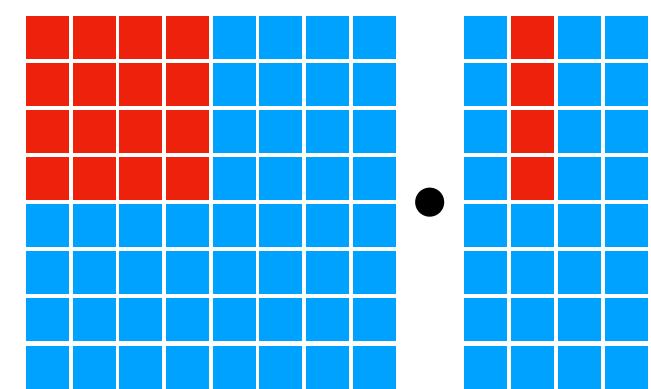
Designed for Accelerators(NPU)

# Search Space for TPU-like Specialized Accelerators

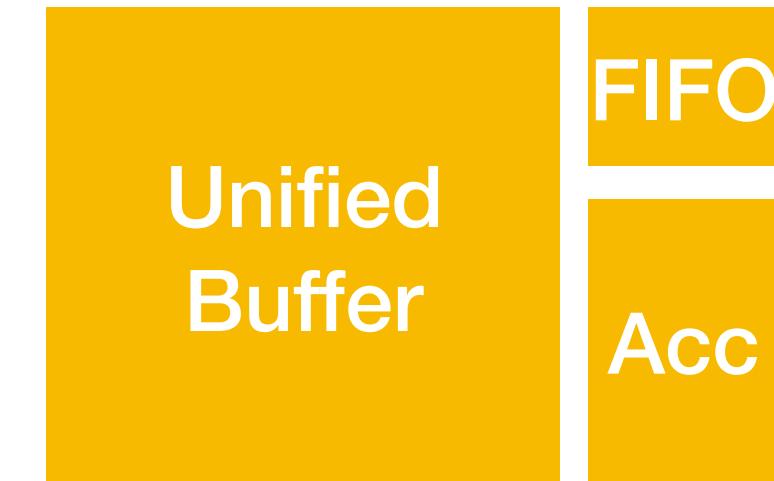
## TPUs



**Tensor  
Compute Primitives**

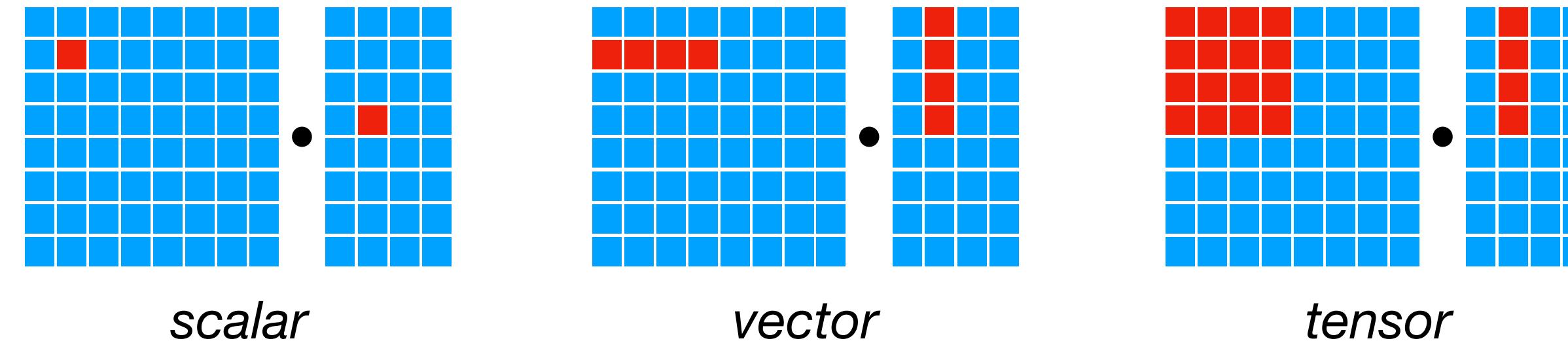


**Explicitly Managed  
Memory Subsystem**



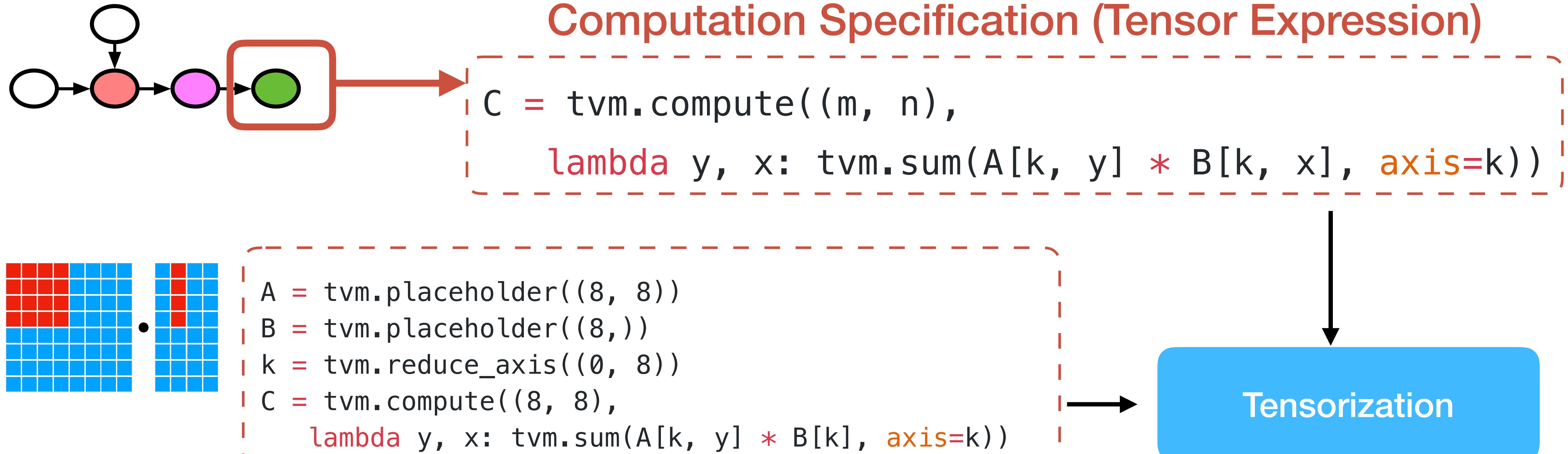
# Tensorization Challenge

**Compute  
primitives**

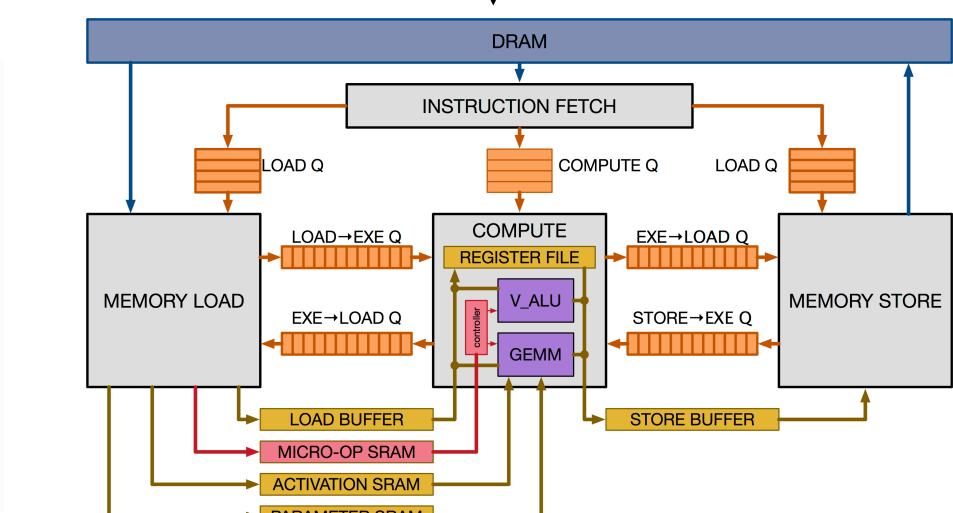
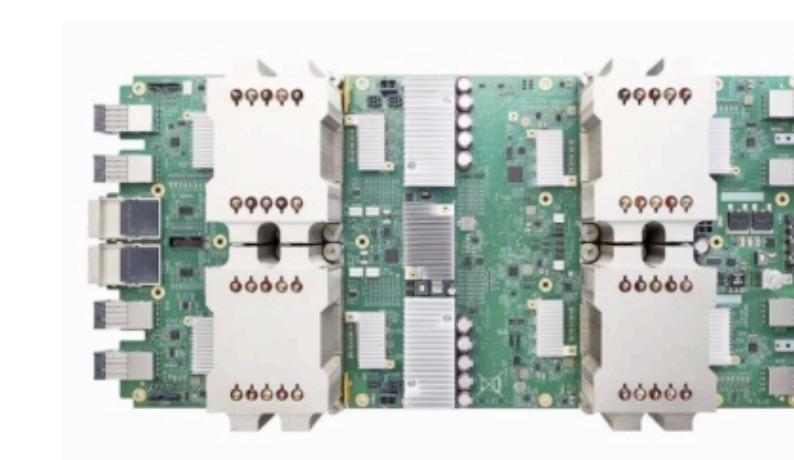


**Challenge: Build systems to support  
emerging tensor instructions**

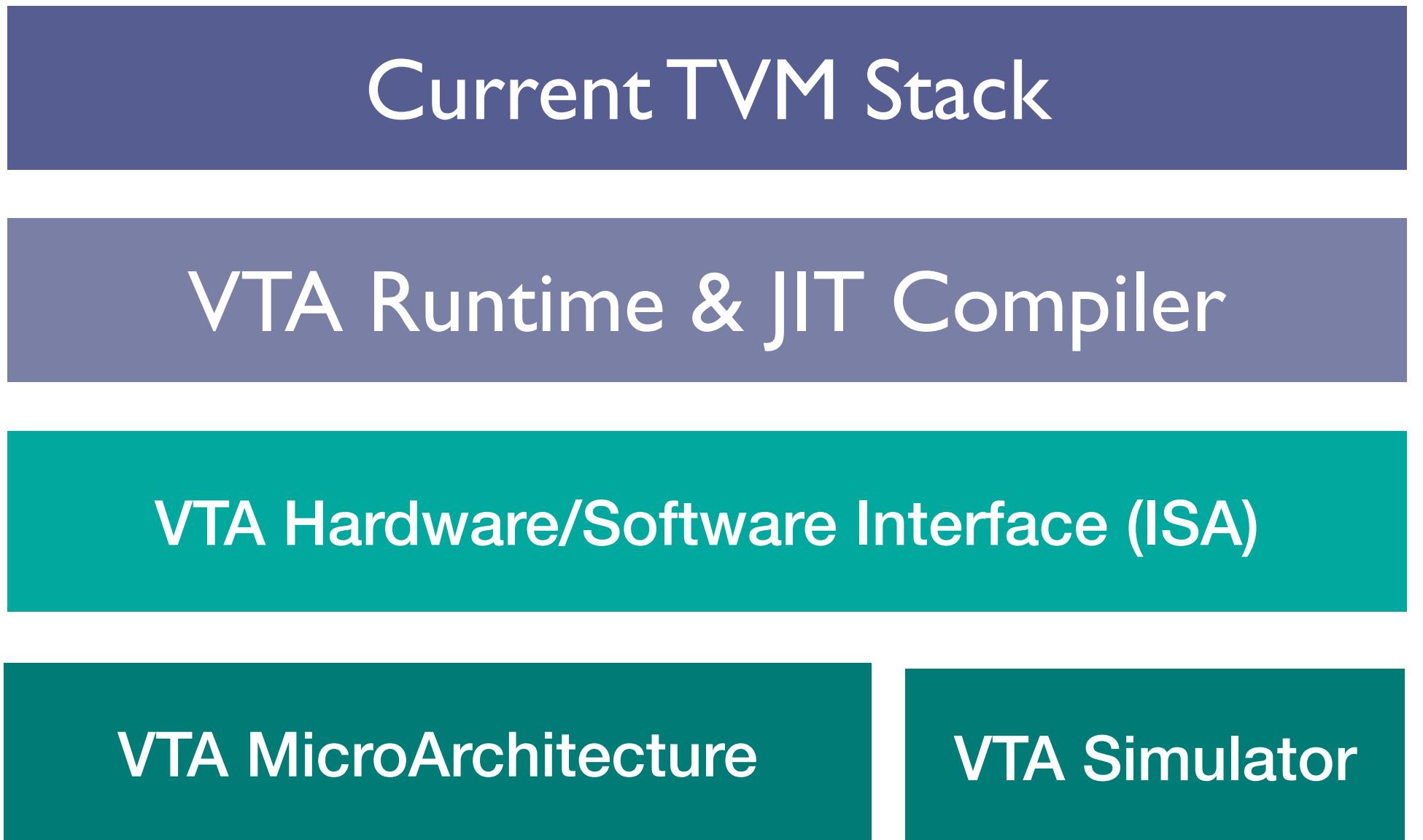
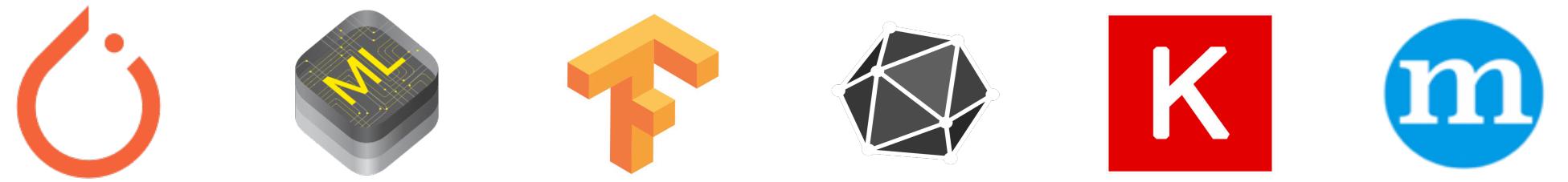
# Tensorization Challenge



HW Interface Specification by Tensor Expression



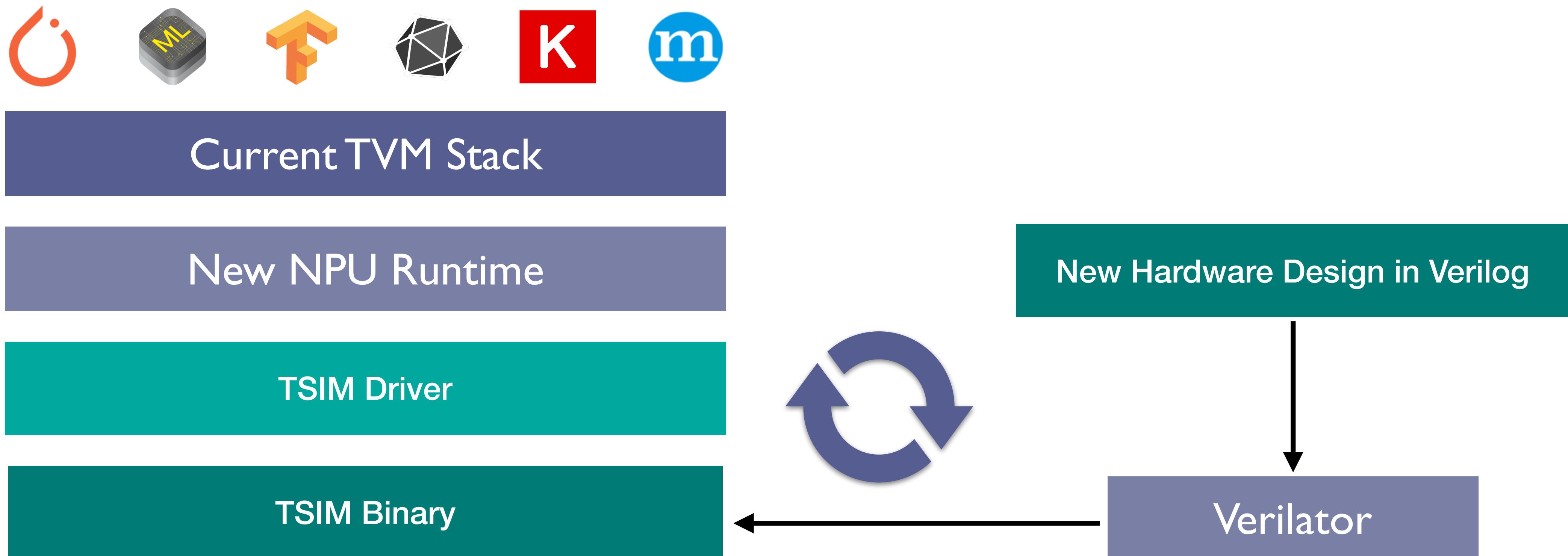
# VTA: Open & Flexible Deep Learning Accelerator



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

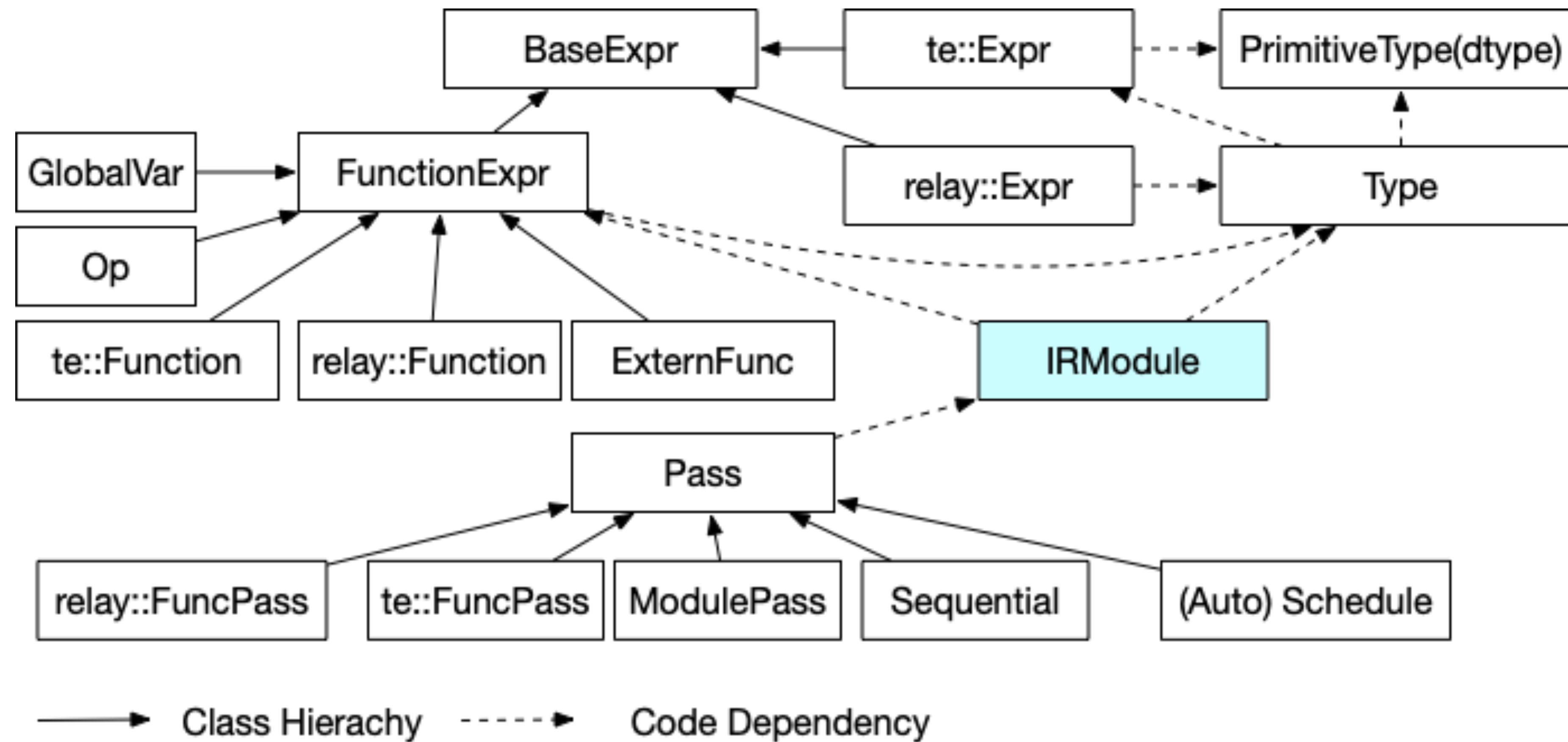
**compiler, driver,  
hardware design  
full stack open source**

# TSIM: Support for Future Hardware



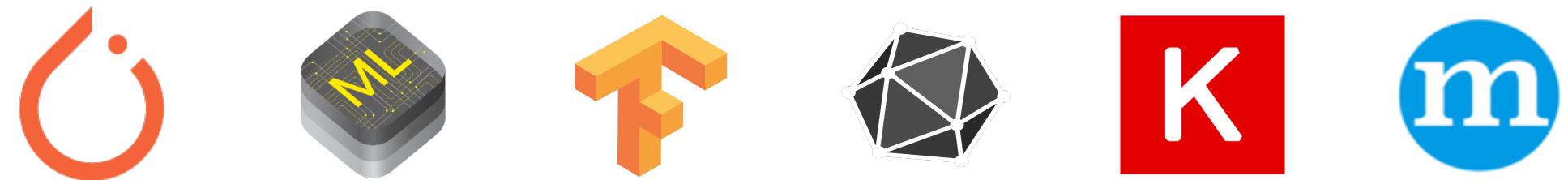
# Toward Unified IR Infra

# Overview of New IR Infra



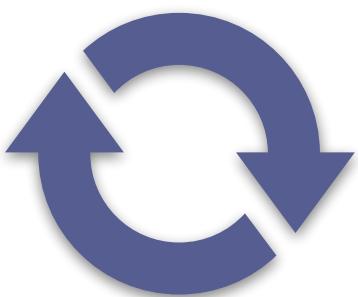
Single unified module/pass, type system, with function variants support

# Compilation Flow under the New Infra



Import

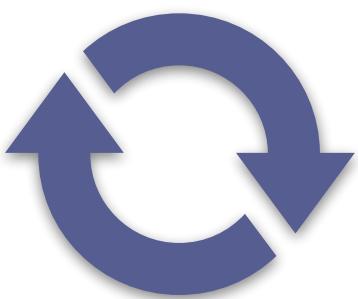
IRModule (relay::Function)



High-level optimizations

Lower

IRModule (te::Function, ExternFunc, ...)



(Auto) Schedules  
Low-level optimizations

Codegen

runtime::Module

# Mixed Function Variants in the Same Module

```
def @relay_add_one(%x : Tensor((10,), f32)) {
    call_destination_passing @te_add_one(%x,  out=%b)
}

def @te_add_one(%a: NDArray, %b: NDArray) {
    var %n
    %A = decl_buffer(shape=[%n], src=%a)
    %B = decl_buffer(shape=[%n], src=%b)
    for %i = 0 to 10 [data_par] {
        %B[%i] = %A[%i] + 1.0
    }
}
```

# First-class Python Support

```
@tvm.hybrid
def te_add_one(a, b):
    n = var("n")
    A = bind_buffer(shape=[n], a)
    B = bind_buffer(shape=[n], b)
    for i in iter_range(n, iter_type="data_par"):
        A[i] = B[i] + 1
```

Use hybrid script as  
an alternative text  
format

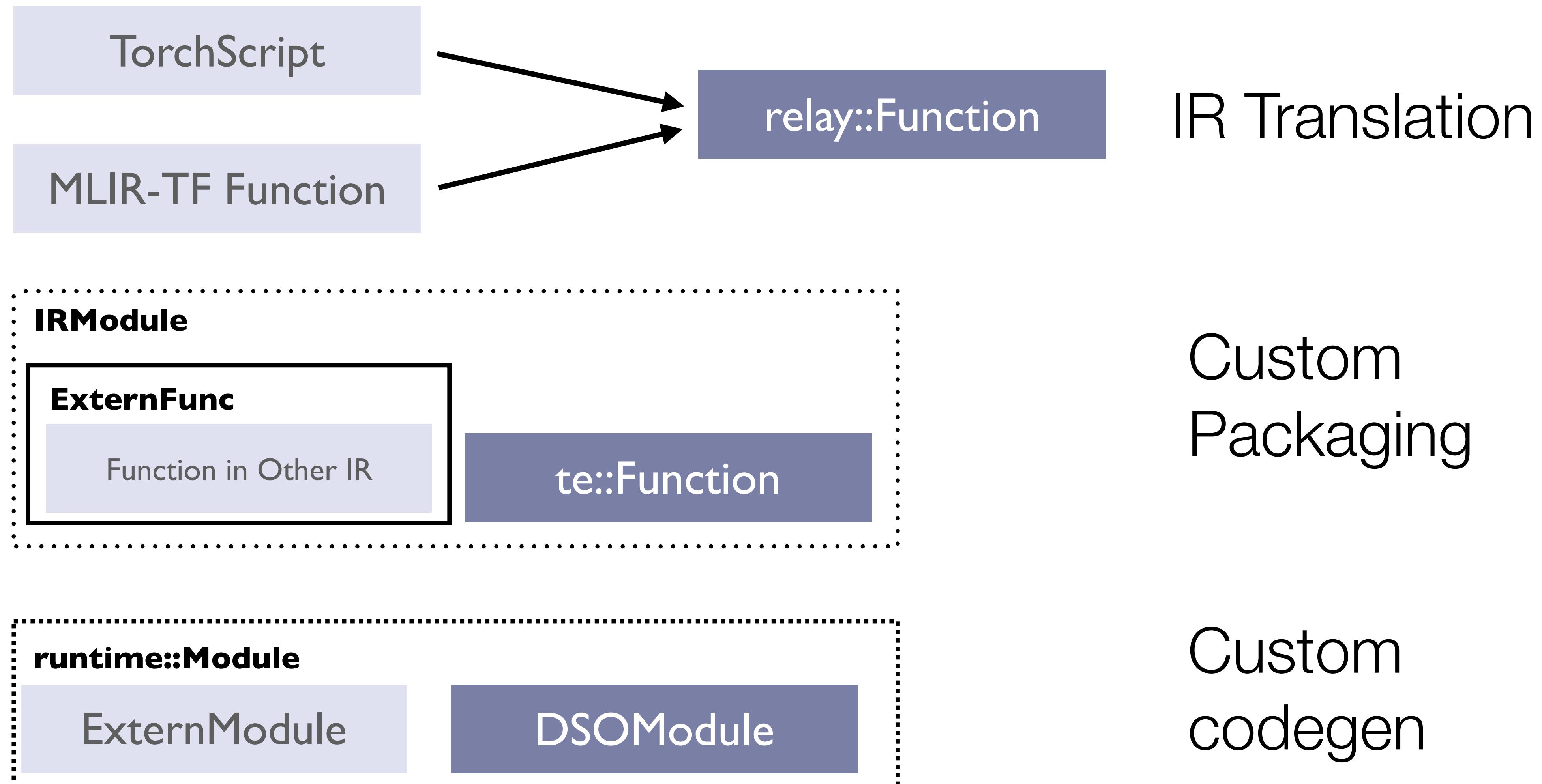
```
mod = tvm.IRModule([te_add_one])
print(mod["te_add_one"].args)
```

Directly write pass,  
manipulate IR structures

Accelerate innovation,  
e.g. use (GA/RL/BayesOpt/your favorite ML method) for AutoSchedule

Easy shift to C++ when product ready

# Interpolate with Other Compilers



# Timeline

RFC: now

Main area of focus in next 3-4 month

More updates about tensor-level IR at the TVM conference

# Community

# Open Source Community

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



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

Open Code

Open Development

Open Governance

# Acknowledgement

Apache (incubating) TVM community

Our awesome community members!

Alibaba for hosting the meetup

# TVM Conference 2019

**Thu, December 5th 2019, Seattle.**

<https://sampl.cs.washington.edu/tvmconf/>

2 days before NeurIPS, Vancouver is 3h away from SEA