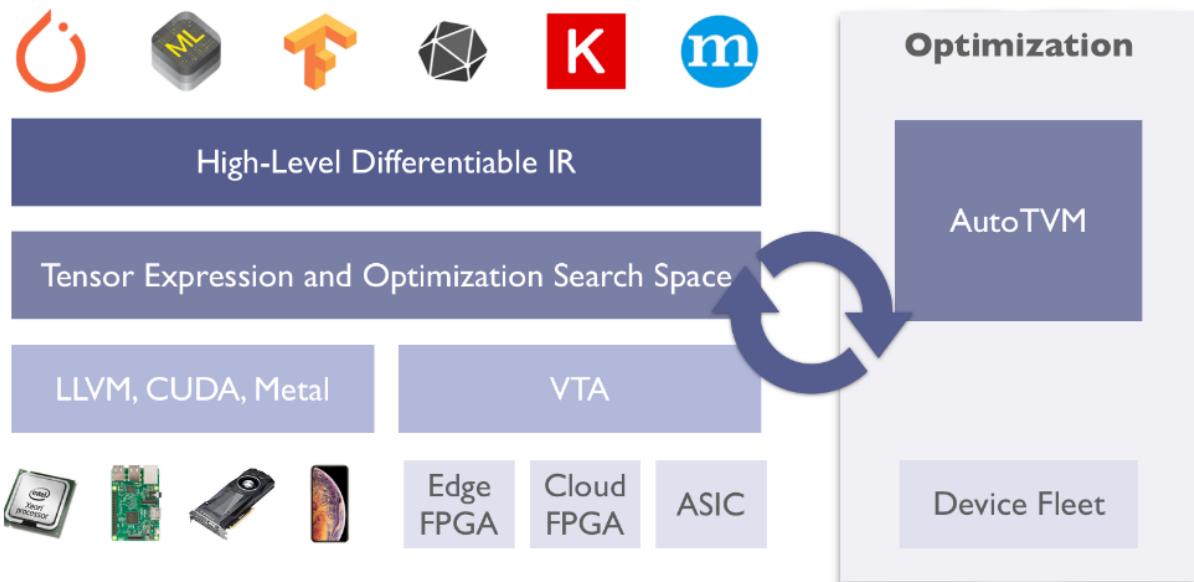


End to End Deep Learning Compiler Stack

for CPUs, GPUs and specialized accelerators



Apache TVM (incubating) is a compiler stack for deep learning systems. It is designed to close the gap between the productivity-focused deep learning frameworks, and the performance- and efficiency-focused hardware backends. TVM works with deep learning frameworks to provide end to end compilation to different backends.

A screenshot of the Apache incubator-tvm GitHub repository page. The header includes links for Pull requests, Issues, Marketplace, and Explore. The main navigation bar shows the repository path apache/incubator-tvm, with options to Unwatch releases (360), Unstar (4.9k), Fork (1.3k), and Insights. Below the navigation are sections for Code (Issues 83, Pull requests 57, Actions, Projects 1, Wiki, Security, Insights), and a list of tags: compiler, tensor, deep-learning, gpu, opencl, metal, performance, javascript, rocm, tvm, vulkan, spirv, machine-learning. At the bottom, there are statistics: 3,530 commits, 3 branches, 0 packages, 6 releases, 308 contributors, and Apache-2.0 license. Buttons for Branch: master, New pull request, Create new file, Upload files, Find file, and Clone or download are also visible.

[HTML] TVM: end-to-end optimization stack for deep learning

T Chen, T Moreau, Z Jiang, H Shen, E Yan... - arXiv preprint arXiv ..., 2018 - arxiv-vanity.com
Scalable frameworks, such as TensorFlow, MXNet, Caffe, and PyTorch drive the current popularity and utility of deep learning. However, these frameworks are optimized for a narrow range of server-class GPUs and deploying workloads to other platforms such as ...

☆ 82 被引用次数: 82 相关文章 所有 6 个版本 图书馆搜索 >>

Learning to optimize tensor programs

T Chen, L Zheng, E Yan, Z Jiang, T Moreau... - Advances in Neural ..., 2018 - papers.nips.cc

Paper accepted and presented at the Neural Information Processing Systems Conference (<http://nips.cc/>).

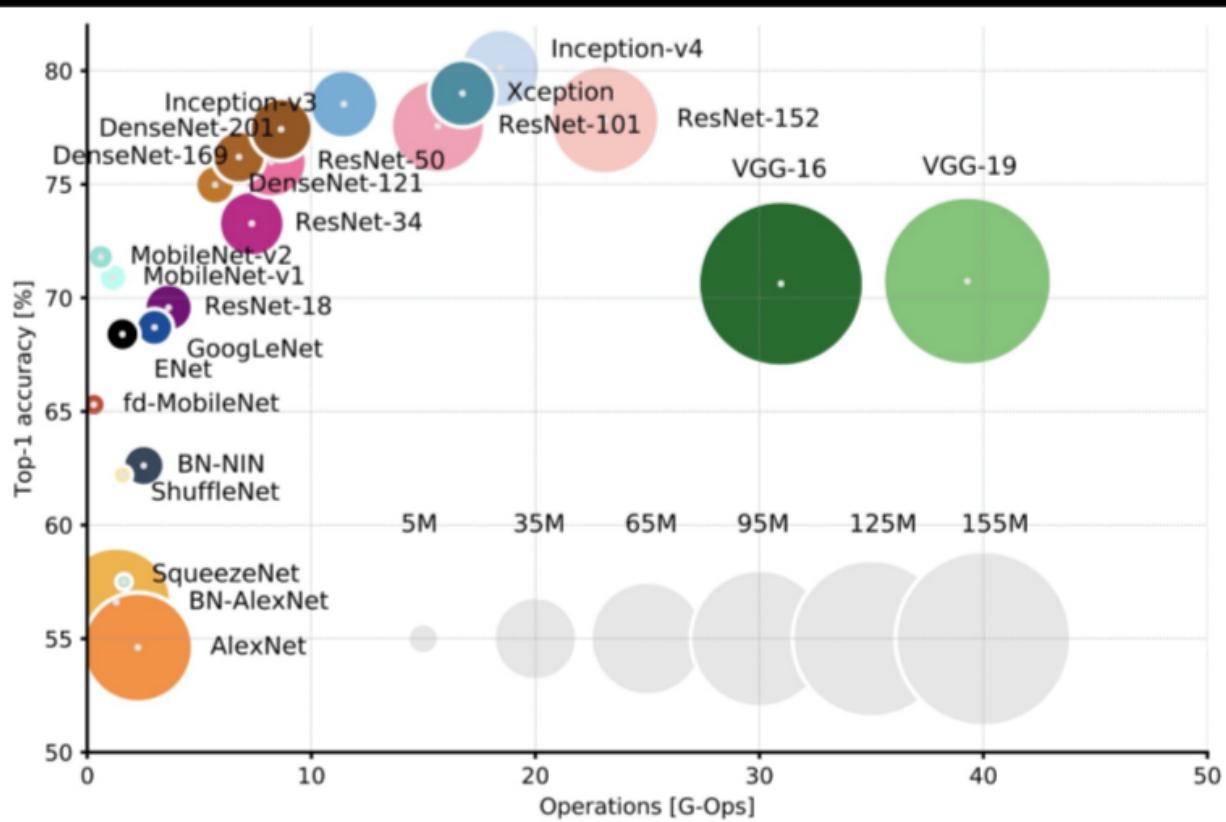
☆ 99 被引用次数: 30 相关文章 所有 10 个版本 图书馆搜索 ≈

We learned a lot from the following projects when building TVM.

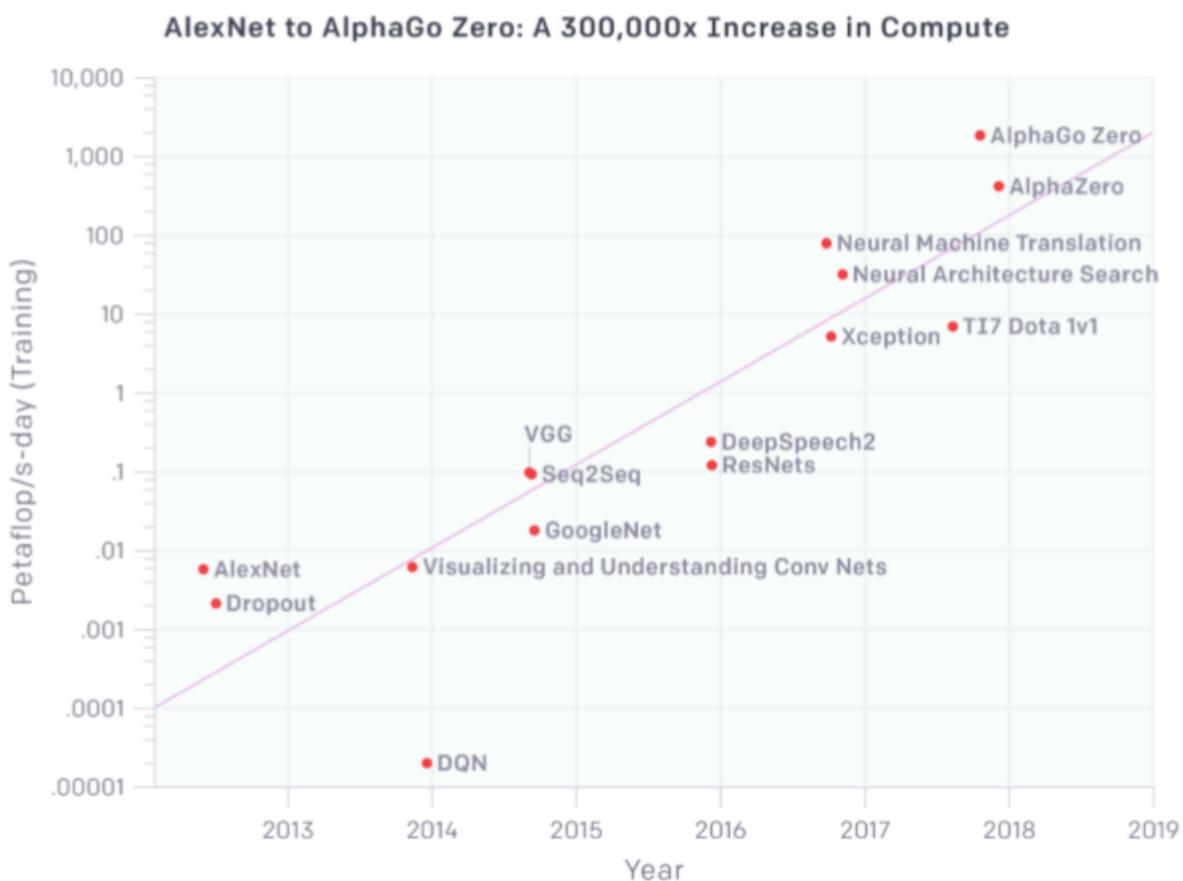
- [Halide](#): TVM uses [HalideIR](#) as data structure for arithmetic simplification and low level lowering. We also learned and adapted some part of lowering pipeline from Halide.
- [Loopy](#): use of integer set analysis and its loop transformation primitives.
- [Theano](#): the design inspiration of symbolic scan operator for recurrence.



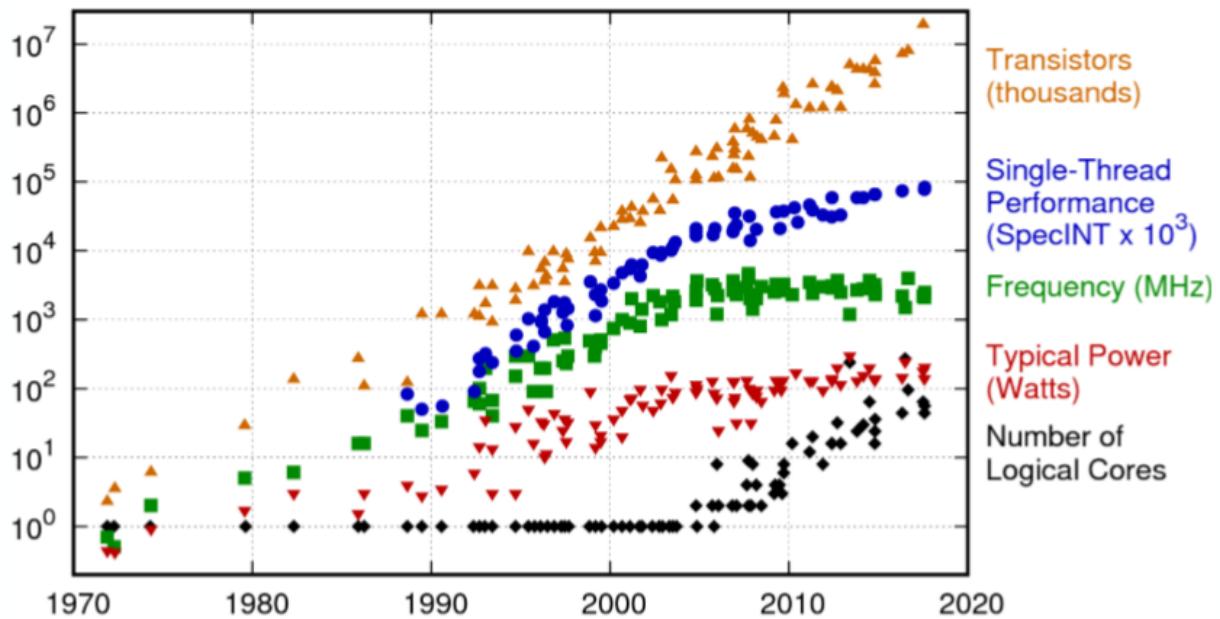
Model size and compute cost growing fast



Training costs growing exponentially



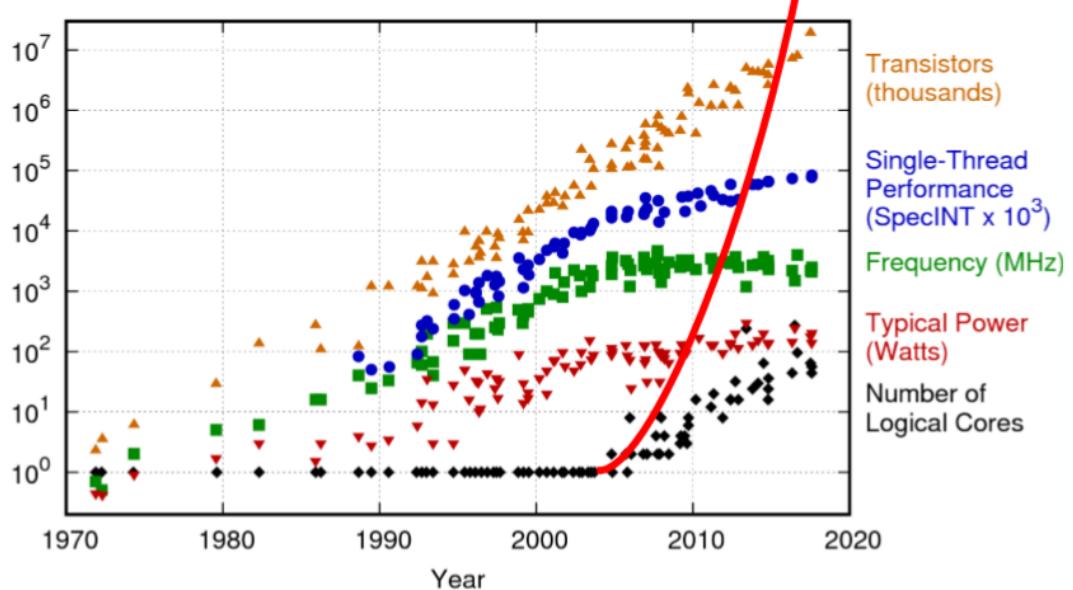
42 Years of Microprocessor Trend Data



'e serious...

Computational cost of
ML. Oops. :)

42 Years of Microprocessor Trend Data



A perfect storm

Growing set of requirements: **cost, latency, power, security & privacy**

Cambrian explosion of models, workloads, and use cases.

CNN GAN RNN MLP DQNN

Rapidly evolving ML software ecosystem quickly fragmenting

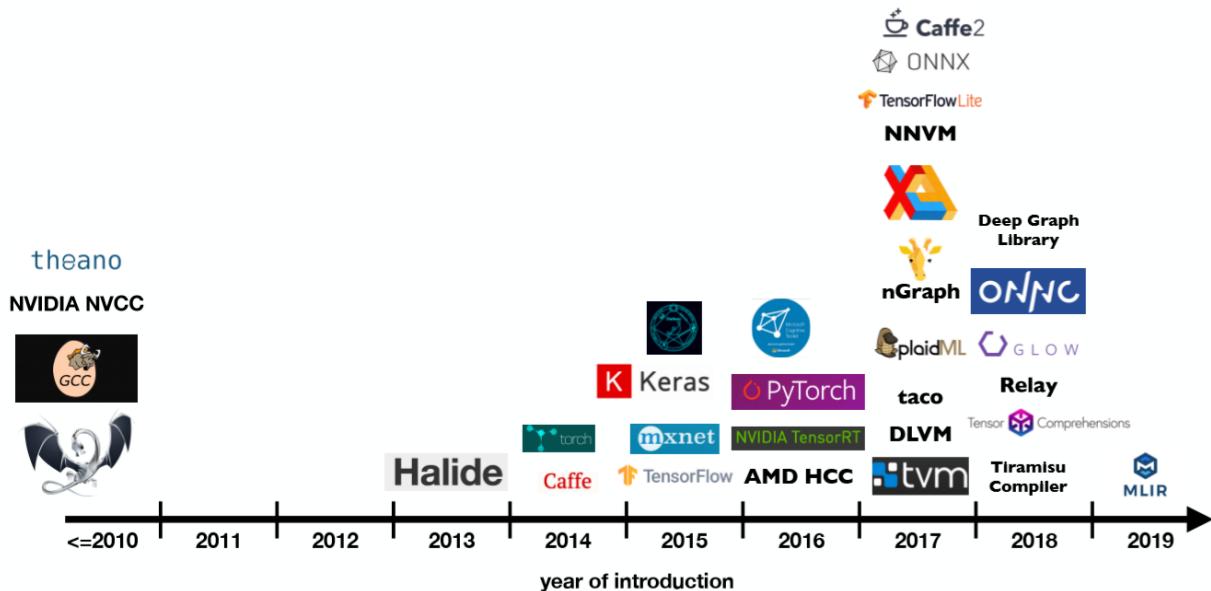


Silicon scaling limitations (Dennard and Moore):

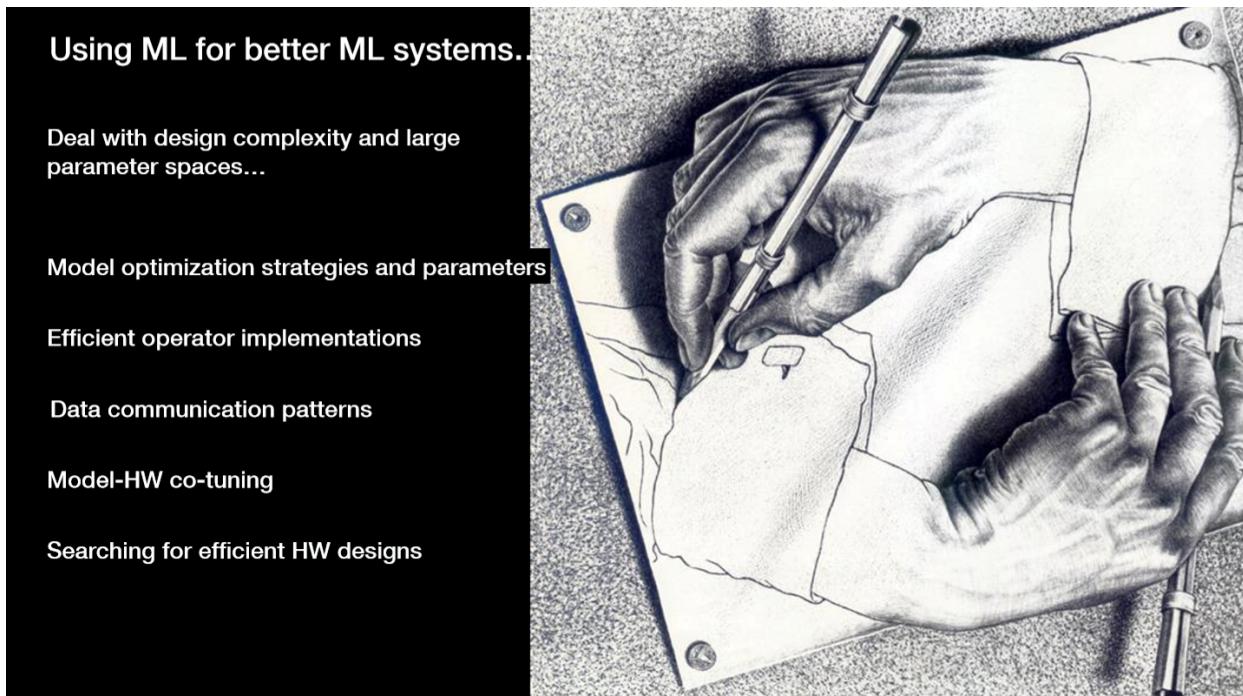
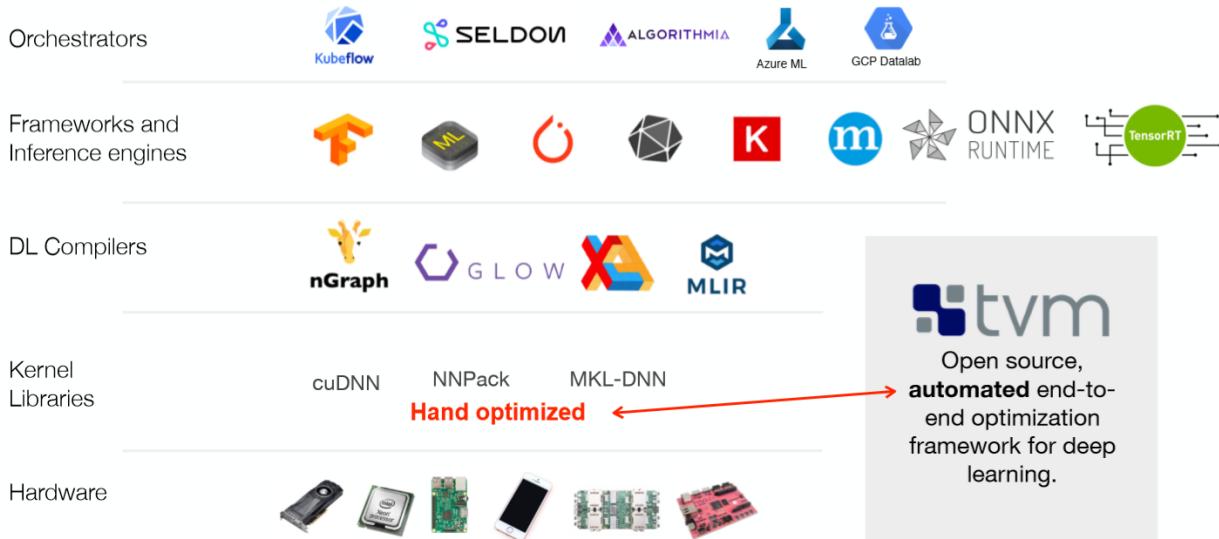


Cambrian explosion of HW backends. Heterogeneous HW.

Deep learning “stack” (r?)evolution



Current Dominant Deep Learning Systems Landscape



tvm Open Source Community Growth and Impact

70% growth from Dec 2018 to **295 contributors** from UW, Berkeley, Cornell, UCLA, Amazon, Huawei, NTT, Facebook, Microsoft, Qualcomm, Alibaba, Intel, ...

Used in production at leading vendors:



Deep Learning Compiler Service



Tensor Engine for mobile ASIC

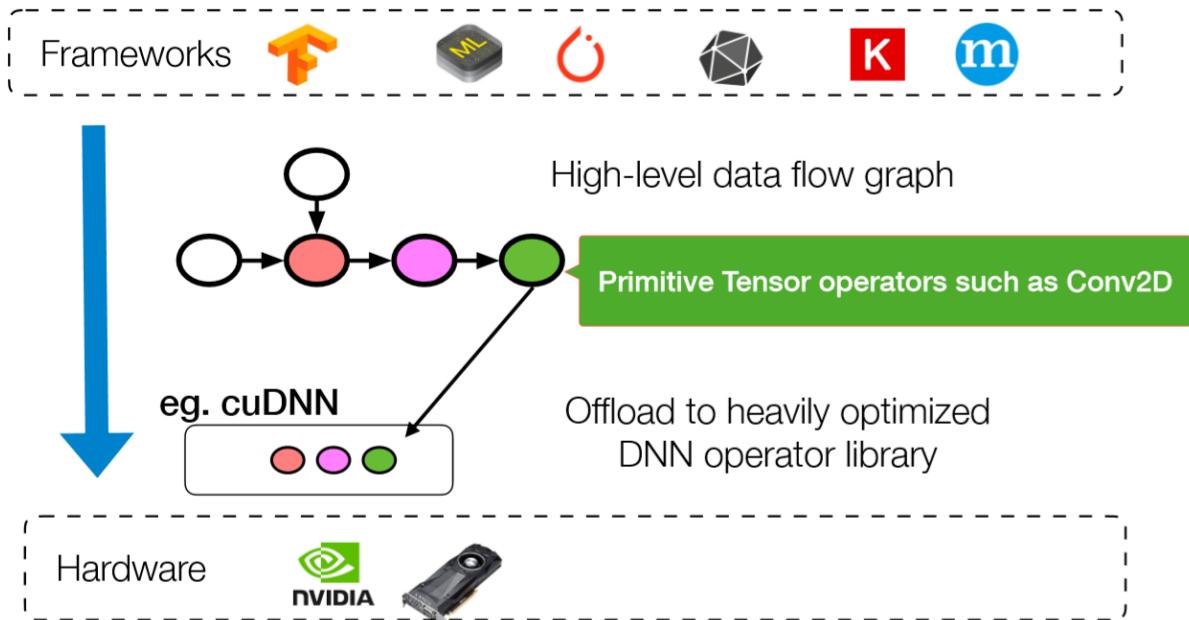


Mobile and Server Optimizations

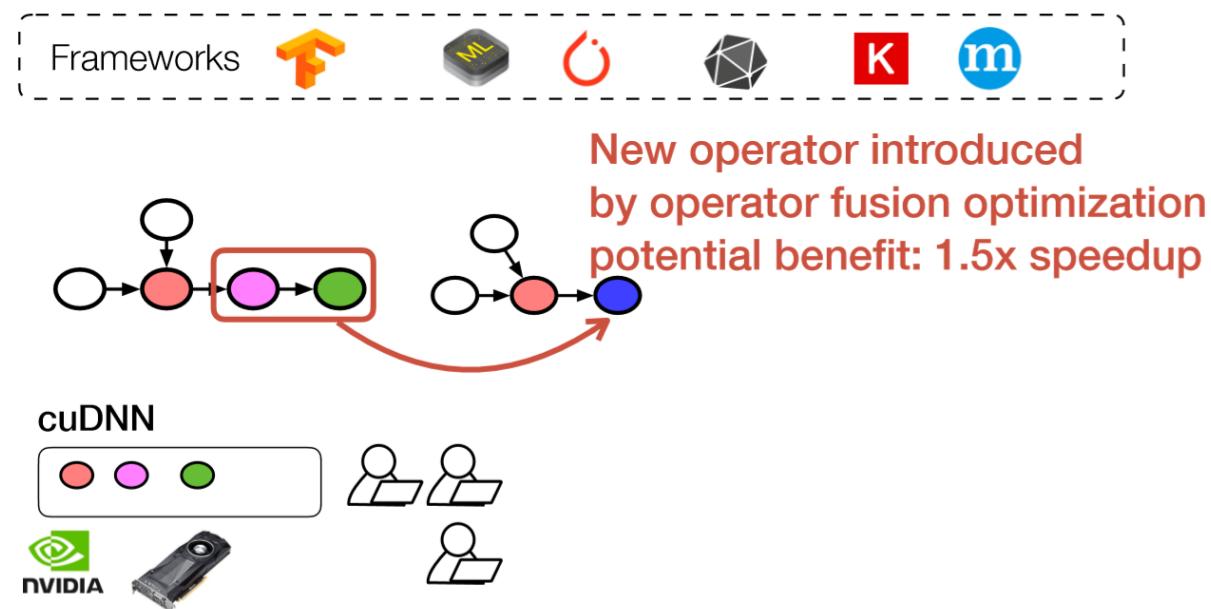


Cloud-side model optimization

Existing Deep Learning Frameworks



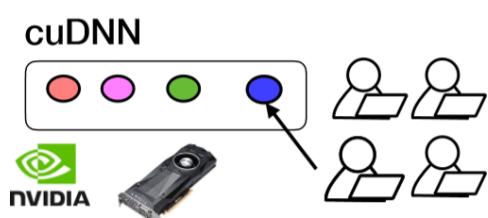
Limitations of Existing Approach



Limitations of Existing Approach



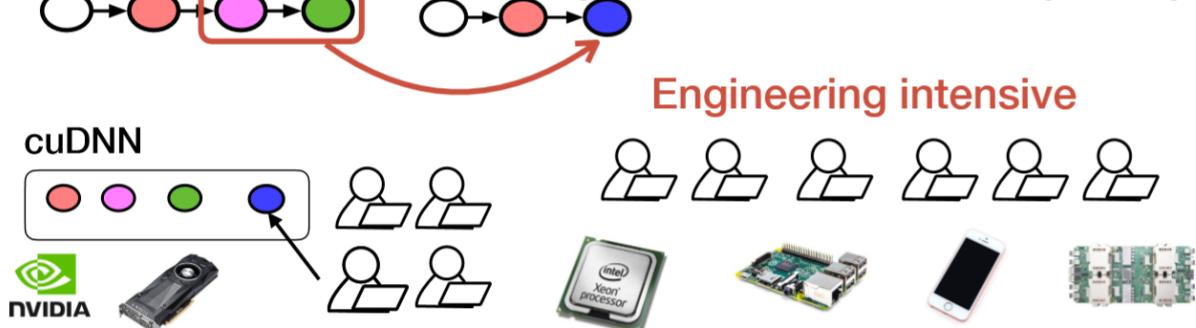
New operator introduced
by operator fusion optimization
potential benefit: 1.5x speedup



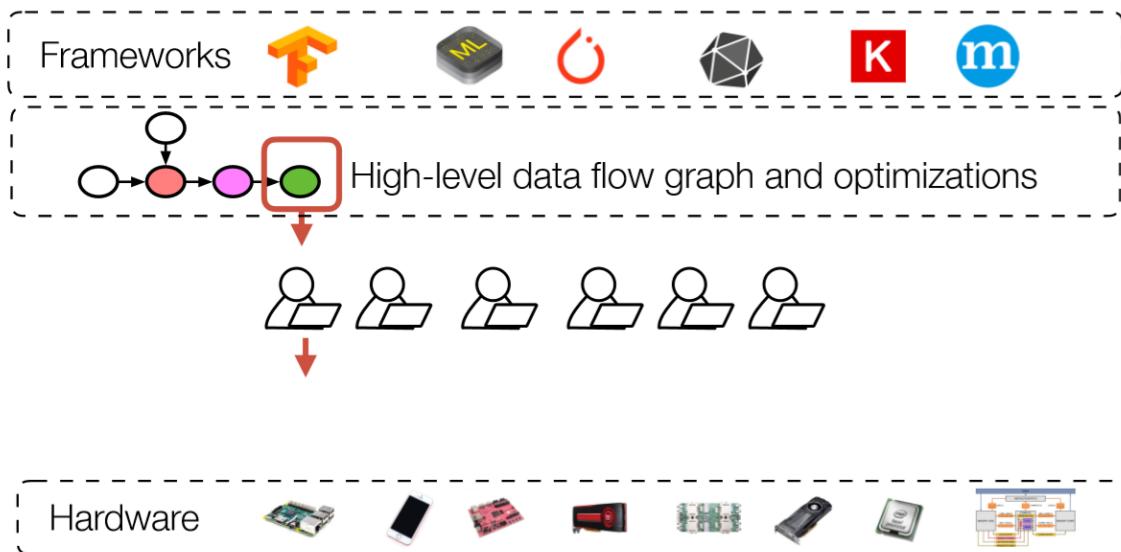
Limitations of Existing Approach



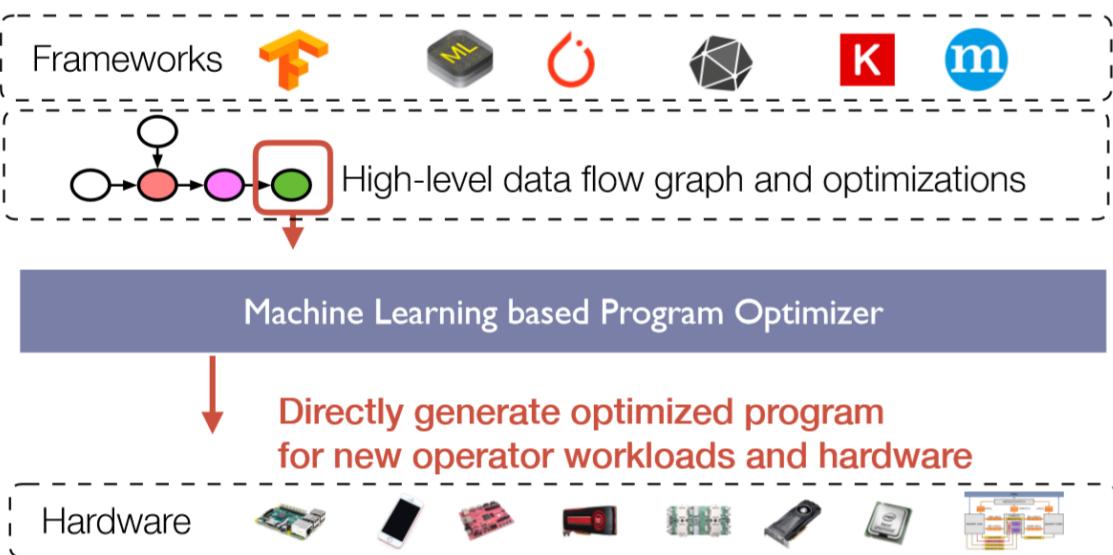
New operator introduced
by operator fusion optimization
potential benefit: 1.5x speedup



TVM: Learning-based Learning System

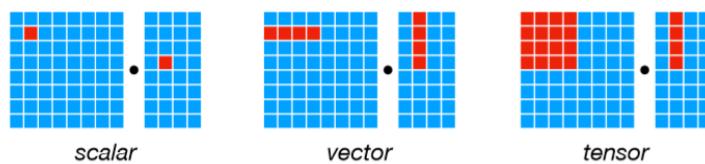


TVM: Learning-based Learning System



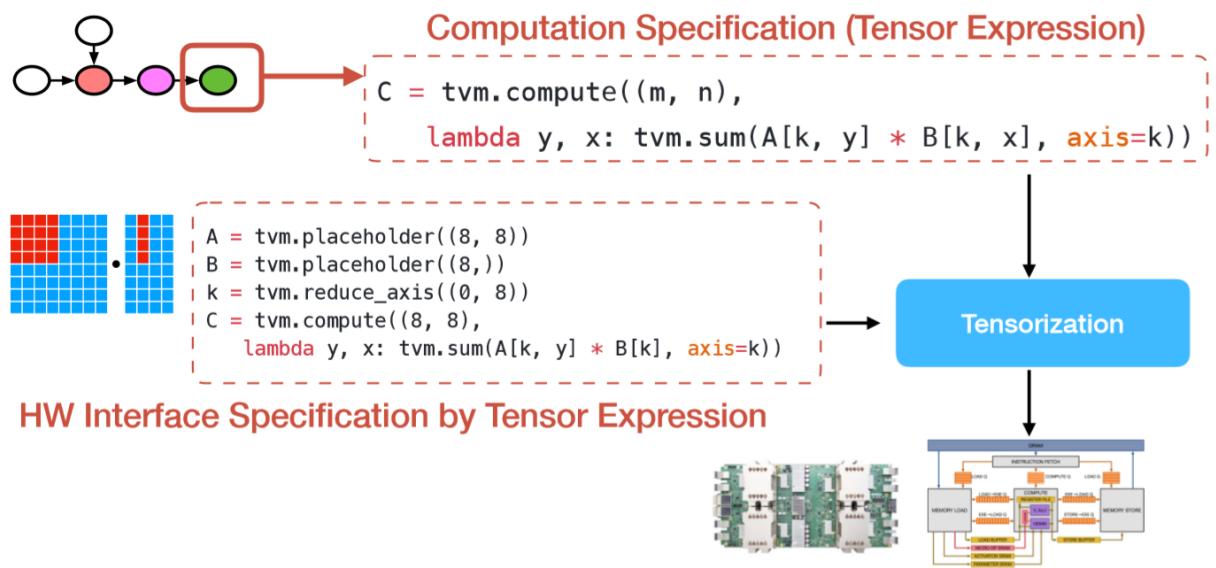
Tensorization Challenge

Compute primitives

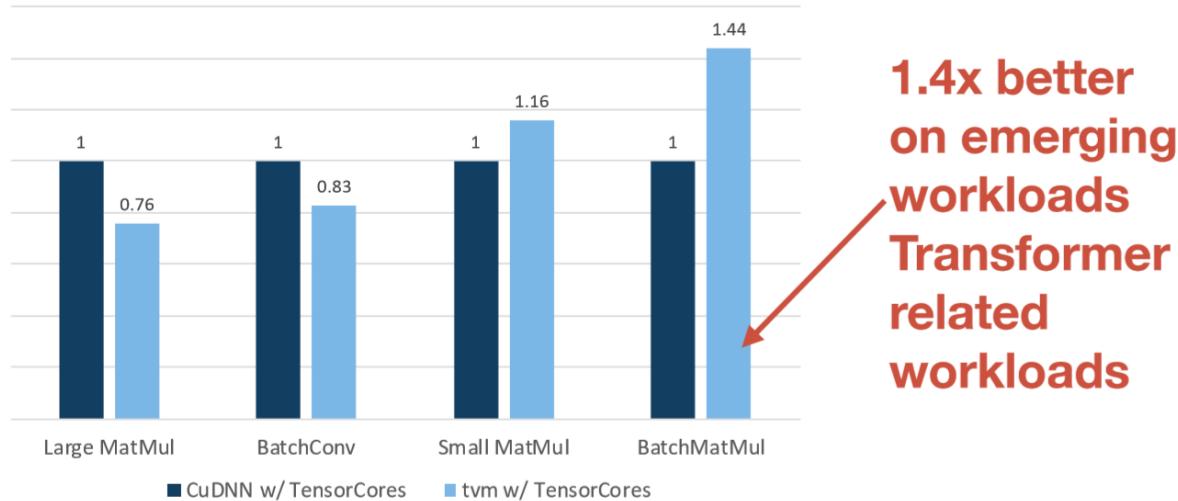


Challenge: Build systems to support emerging tensor instructions

Tensorization Challenge

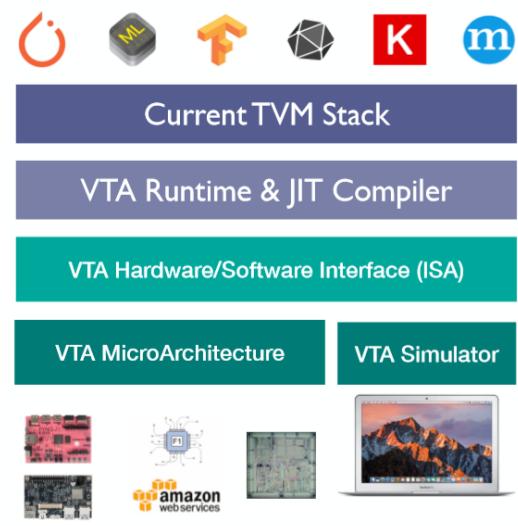


TVM for TensorCore



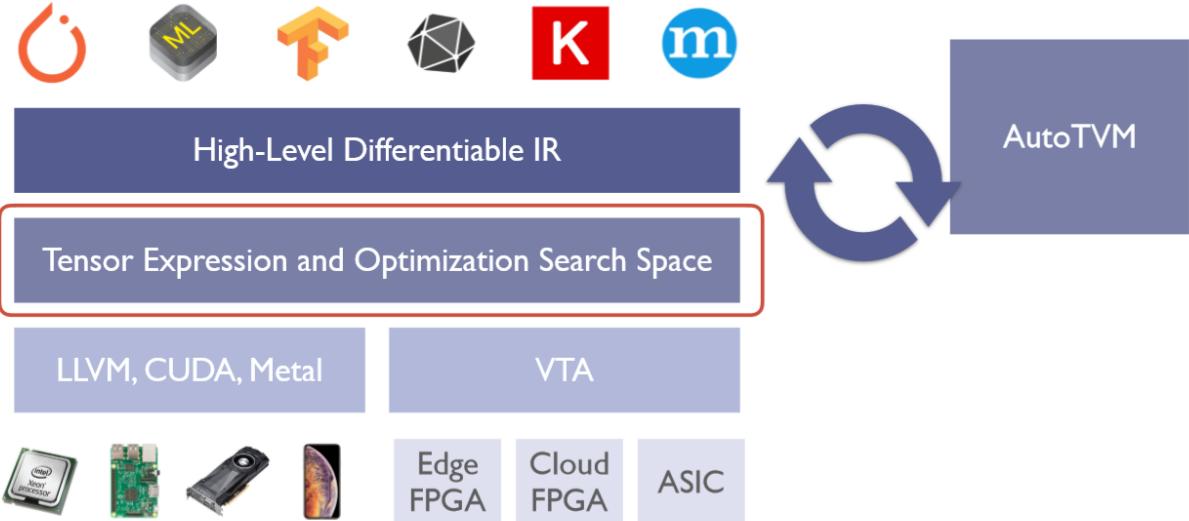
**1.4x better
on emerging
workloads
Transformer
related
workloads**

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
- VTA 2.0 release - Chisel compiler, driver, hardware design full stack open source

Full Stack Automation



```

e compute expression
[A = t.placeholder((1024, 1024))
B = t.placeholder((1024, 1024))
k = t.reduce_axis((0, 1024))
C = t.compute((1024, 1024),
  lambda y, x:
    t.sum(A[k, y] * B[k, x], axis=k))
x0 default 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]]
```

```

S1 loop tiling
yo, xo, yi, xi = s[C].tile(y, x, ty, tx)
s[C].reorder(yo, xo, k, yi, xi)
x1 = g(e, s1)
for yo in range(1024 / ty):
  for xo in range(1024 / tx):
    C[yo*ty:yo*ty+tx][xo*tx:xo*tx+tx] = 0
    for k in range(1024):
      for yi in range(ty):
        for xi in range(tx):
          C[yo*ty+yi][xo*tx+xi] += A[k][yo*ty+yi] * B[k][xo*tx+xi]
```

```

S2 tiling, map to micro kernel intrinsics
yo, xo, ko, yi, xi, ki = s[C].tile(y, x, k, 8, 8, 8)
s[C].tensorize(yi, intrin.gemm8x8)
x2 = g(e, s2)
for yo in range(128):
  for xo in range(128):
    intrin.fill_zero(C[yo*8:yo*8+8][xo*8:xo*8+8])
    for ko in range(128):
      intrin.fused_gemm8x8_add(
        C[yo*8:yo*8+8][xo*8:xo*8+8],
        A[ko*8:ko*8+8][yo*8:yo*8+8],
        B[ko*8:ko*8+8][xo*8:xo*8+8])
```

Figure 1: Sample problem. For a given tensor operator specification ($C_{ij} = \sum_k A_{ki}B_{kj}$), there are multiple possible low-level program implementations, each with different choices of loop order, tiling size, and other options. Each choice creates a logically equivalent program with different performance. Our problem is to explore the space of programs to find the fastest implementation.

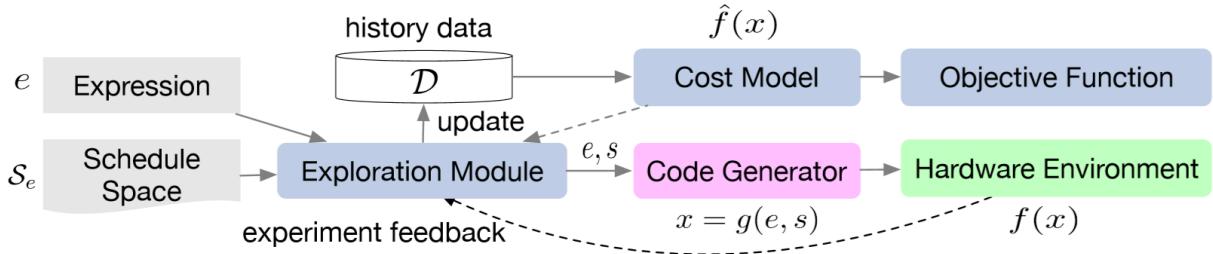


Figure 2: Framework for learning to optimize tensor programs.

We use \mathcal{S}_e to denote the space of possible transformations (schedules) from e to low-level code. For an $s \in \mathcal{S}_e$, let $x = g(e, s)$ be the generated low-level code. Here, g represents a compiler framework that generates low-level code from e, s . We are interested in minimizing $f(x)$, which is the real run time cost on the hardware. Importantly, we do not know an analytical expression for $f(x)$ but can query it by running experiments on the hardware. For a given tuple of (g, e, \mathcal{S}_e, f) , our problem can be formalized as the following objective:

$$\arg \min_{s \in \mathcal{S}_e} f(g(e, s)) \quad (1)$$

We can choose from multiple objective functions to train a statistical cost model for a given collection of data $\mathcal{D} = \{(e_i, s_i, c_i)\}$. A common choice is the regression loss function $\sum_i (\hat{f}(x_i) - c_i)^2$, which encourages the model to predict cost accurately. On the other hand, as we care only about the relative order of program run times rather than their absolute values in the selection process, we can instead use the following rank loss function [6]:

$$\sum_{i,j} \log(1 + e^{-\text{sign}(c_i - c_j)(\hat{f}(x_i) - \hat{f}(x_j))}). \quad (2)$$

We can use the prediction $\hat{f}(x)$ to select the top-performing implementations.

Algorithm 1: Learning to Optimize Tensor Programs

```

Input : Transformation space  $\mathcal{S}_e$ 
Output : Selected schedule configuration  $s^*$ 
 $\mathcal{D} \leftarrow \emptyset$ 
while  $n\_trials < max\_n\_trials$  do
    // Pick the next promising batch
     $Q \leftarrow$  run parallel simulated annealing to collect candidates in  $\mathcal{S}_e$  using energy function  $\hat{f}$ 
     $S \leftarrow$  run greedy submodular optimization to pick  $(1 - \epsilon)b$ -subset from  $Q$  by maximizing Equation 3
     $S \leftarrow S \cup \{ \text{Randomly sample } \epsilon b \text{ candidates.} \}$ 
    // Run measurement on hardware environment
    for  $s$  in  $S$  do
        |  $c \leftarrow f(g(e, s)); \mathcal{D} \leftarrow \mathcal{D} \cup \{(e, s, c)\}$ 
    end
    // Update cost model
    update  $\hat{f}$  using  $\mathcal{D}$ 
     $n\_trials \leftarrow n\_trials + b$ 
end
 $s^* \leftarrow$  history best schedule configuration

```
