15-779 Lecture 9: ML Compilers Part 2: Kernel Autotuning

Zhihao Jia

Computer Science Department Carnegie Mellon University

Outline: Kernel Autotuning

- triton.autotune: template + search
- AutoTVM: template + learning-based cost model
- Ansor: sketch generation + annotation search

Autotuning in Triton

 triton.autotune: whenever the values of keys change, evaluate all configurations

Autotuning Matrix Multiplication

```
def matmul get configs():
    return [triton.Config({'BLOCK_SIZE_M': BM, 'BLOCK_SIZE_N': BN, "BLOCK_SIZE_K": BK, "GROUP_SIZE_M": 8},
                           num stages=s, num warps=w)
        for BM in [128]
        for BN in [128, 256]
        for BK in [64, 128]
        for s in ([2, 3, 4])
        for w in [4, 8]
@triton.autotune(
    configs=matmul_get_configs(),
    key=["M", "N", "K"],
@triton.jit()
def matmul kernel(a ptr, b ptr, c ptr, M, N, K,
                  BLOCK SIZE M: tl.constexpr,
                  BLOCK SIZE N: tl.constexpr,
                  BLOCK SIZE K: tl.constexpr,
                  GROUP SIZE M: tl.constexpr,):
```

Python annotations: easy to use

Issue: time consuming to enumerate all configurations

Outline: Kernel Autotuning

- triton.autotune: template + search
- AutoTVM: template + learning-based cost model
- Ansor: sketch generation + annotation search

TVM: A Learning-based Compiler for Deep Learning



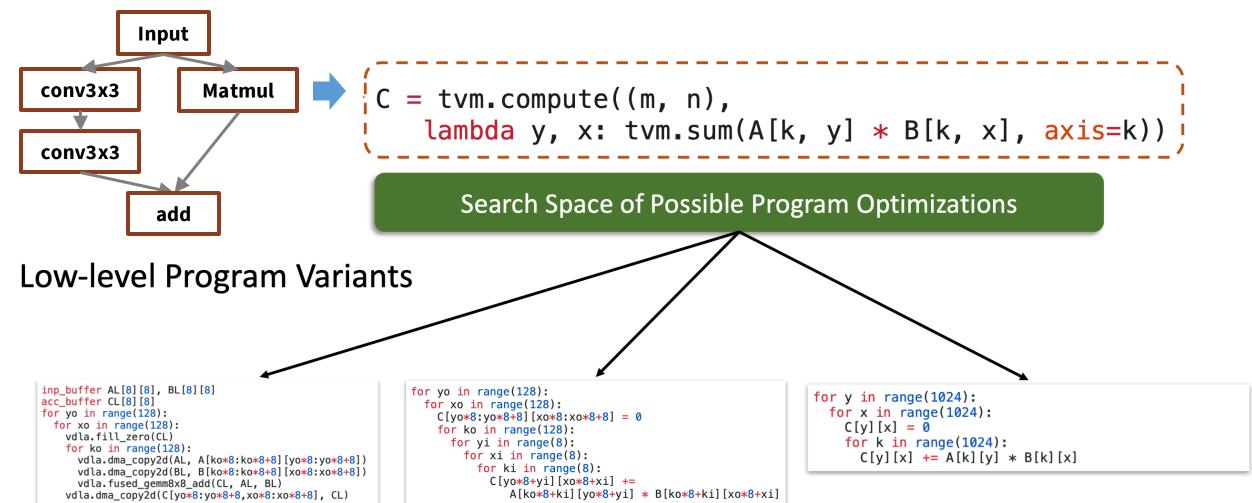
Explosion of models and frameworks

Goal: efficiently deploy deep learning on modern hardware platforms

Explosion of hardware backends Acclerator 4

* Slides from Tianqi Chen

Challenge: Billions of Possible Optimization Choices in the Search Space



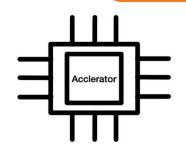
* Slides from Tianqi Chen

TVM: Learning-based Compiler for Deep Learning



Hardware-aware Search Space of Optimized Tensor Programs

Machine Learning based Program Optimizer





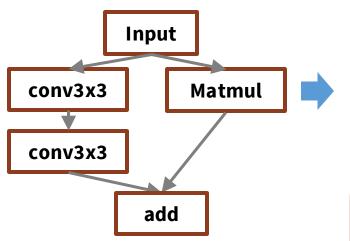












Tensor Expression Language (Specification)

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

Define search space of hardware aware mappings from expression to hardware program

Based on Halide's compute/schedule separation

* Slides from Tianqi Chen

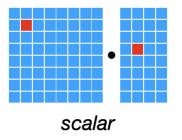
Compute Primitives Memory Subsystem CPUs L3 (intel) scalar vector implicitly managed Loop Cache **Vectorization Transformations** Locality

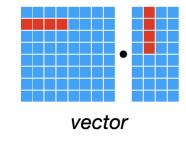
GPUs



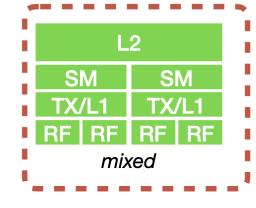


Compute Primitives





Memory Subsystem



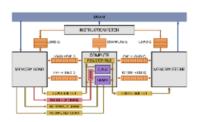
Shared memory among compute cores

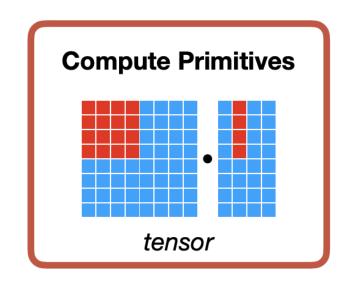
Use of Shared Memory

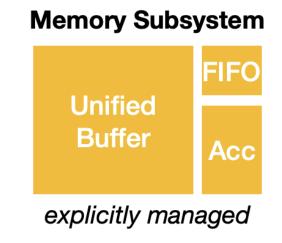
Thread Cooperation

TPU-like Specialized Accelerators



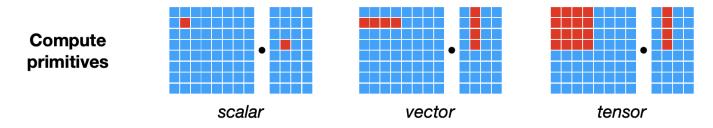






Tensorization

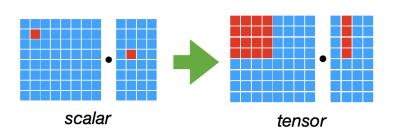
Tensorization Challenge

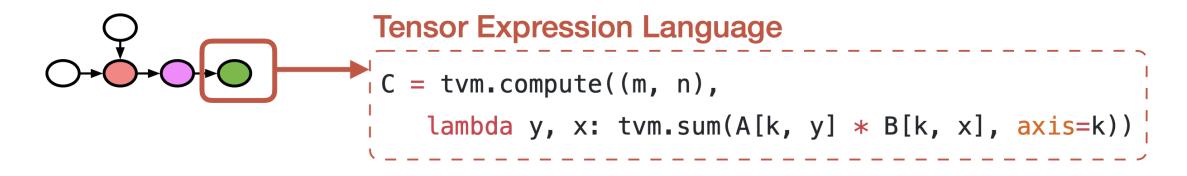


Hardware designer: declare tensor instruction interface with Tensor Expression

```
w, x = t.placeholder((8, 8)), t.placeholder((8, 8))
                                                      declare behavior
k = t.reduce axis((0, 8))
y = t.compute((8, 8), lambda i, j:
               t.sum(w[i, k] * x[j, k], axis=k))
                                                  lowering rule to generate
def gemm_intrin_lower(inputs, outputs):
                                                  hardware intrinsics to carry
   ww_ptr = inputs[0].access_ptr("r")
   xx_ptr = inputs[1].access_ptr("r")
                                                  out the computation
   zz_ptr = outputs[0].access_ptr("w")
   compute = t.hardware_intrin("gemm8x8", ww_ptr, xx_ptr, zz_ptr)
   reset = t.hardware_intrin("fill_zero", zz_ptr)
   update = t.hardware_intrin("fuse_gemm8x8_add", ww_ptr, xx_ptr, zz_ptr)
   return compute, reset, update
gemm8x8 = t.decl_tensor_intrin(y.op, gemm_intrin_lower)
```

Tensorize: transform program to use tensor instructions





Primitives in prior work: Halide, Loopy

New primitives for GPUs, and enable TPU-like Accelerators

Loop Transformations Thread Bindings

Cache Locality

Thread Cooperation

Tensorization

Latency Hiding



Hardware





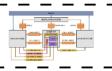


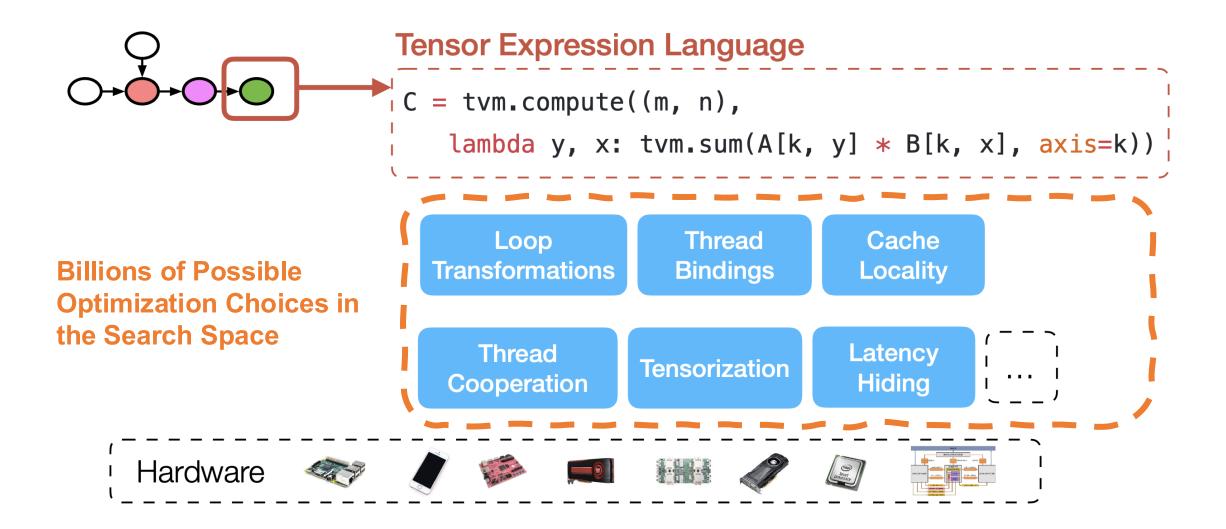












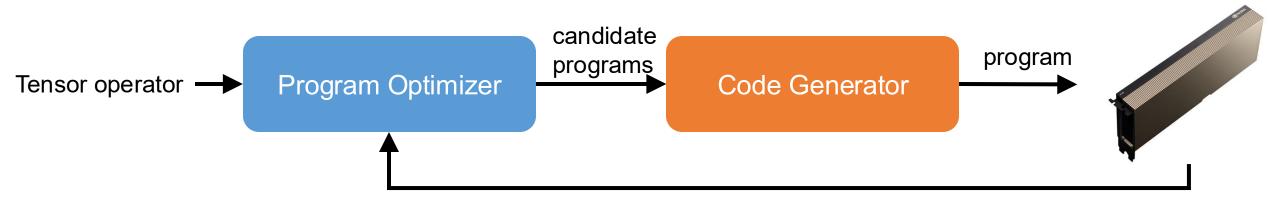
TVM: Learning-based Compiler for Deep Learning



Hardware-aware Search Space of Optimized Tensor Programs

Learning based Program Optimizer Acclerator Acclerator

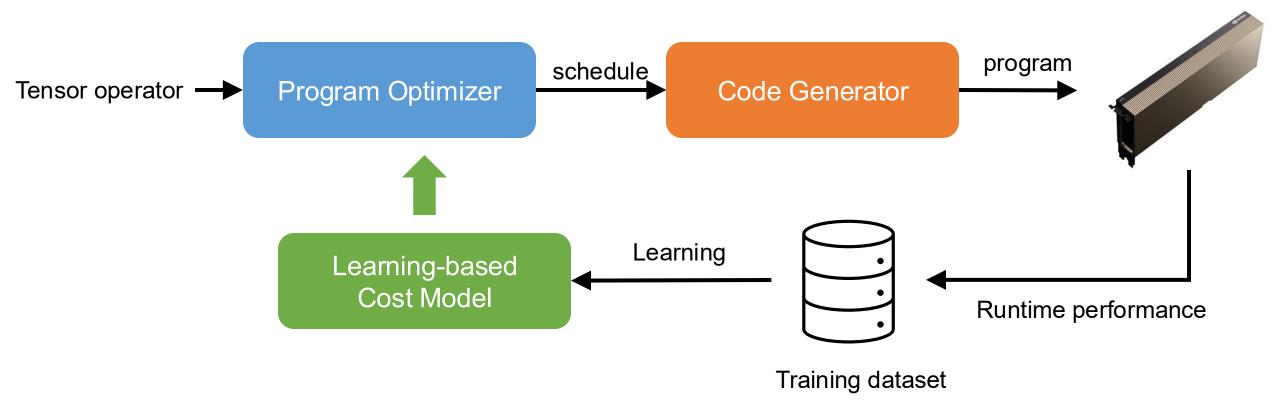
Learning-based Program Optimizer



Runtime performance

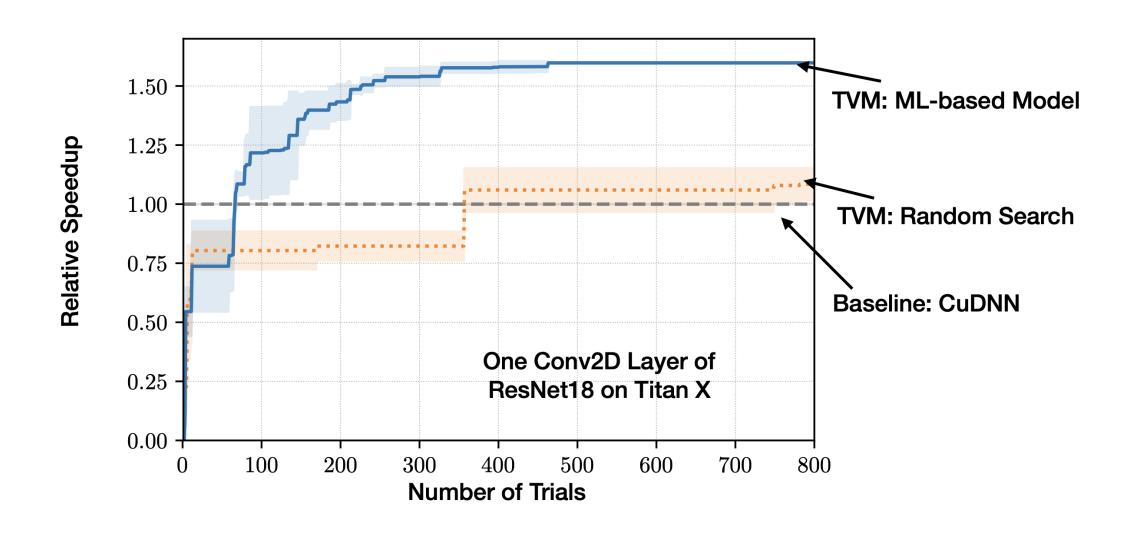
Issue: high experiment cost, each trial takes seconds

Learning-based Program Optimizer

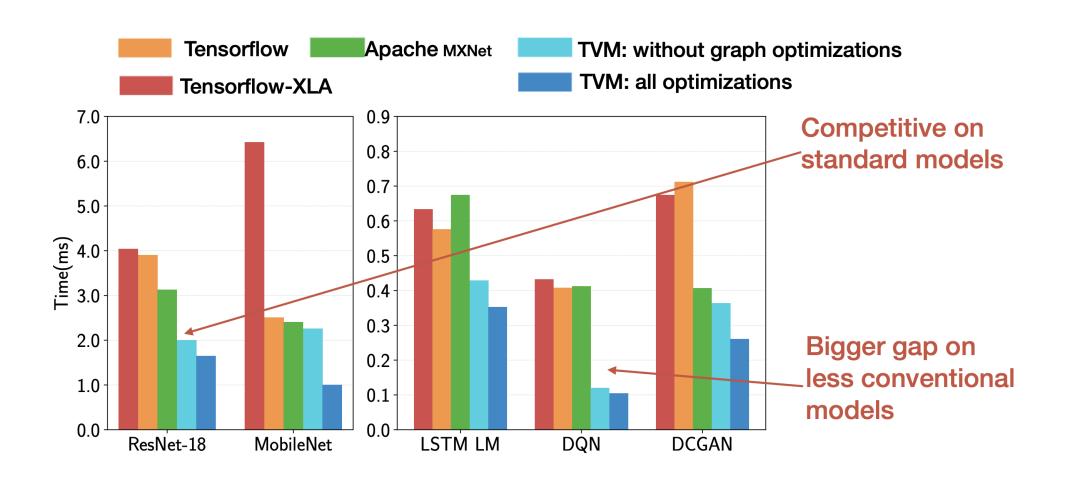


Adapt to hardware by learning, make prediction in milliseconds

Efficient ML-based Cost Model



End-to-end Inference Performance



Outline: Kernel Autotuning

- triton.autotune: manual template + search
- AutoTVM: manual template + learning-based cost model
- Ansor: sketch generation + annotation search

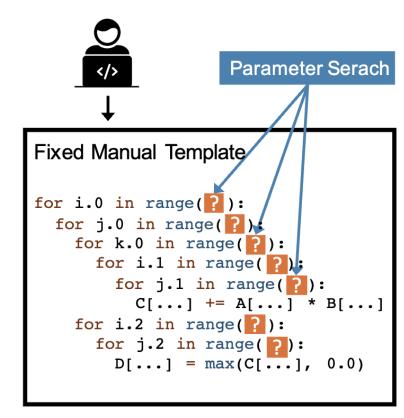
Main Issues with TVM and Triton

Templated-guided search

Manually-written templates to define a search space

Drawbacks

- Not fully-automated -> requires huge manual effort
- Limited search space -> suboptimal performance



(a) Template-guided Search

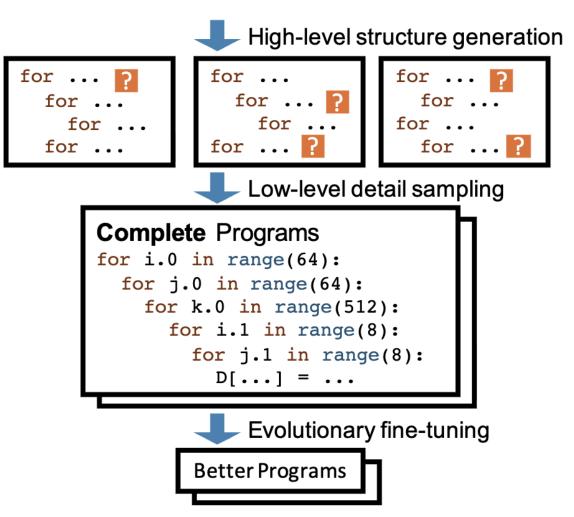
Ansor's Approach

Challenge 1. How to build a large search space automatically?

 Hierarchical search (sketch + annotation)

Challenge 2: How to search efficiently?

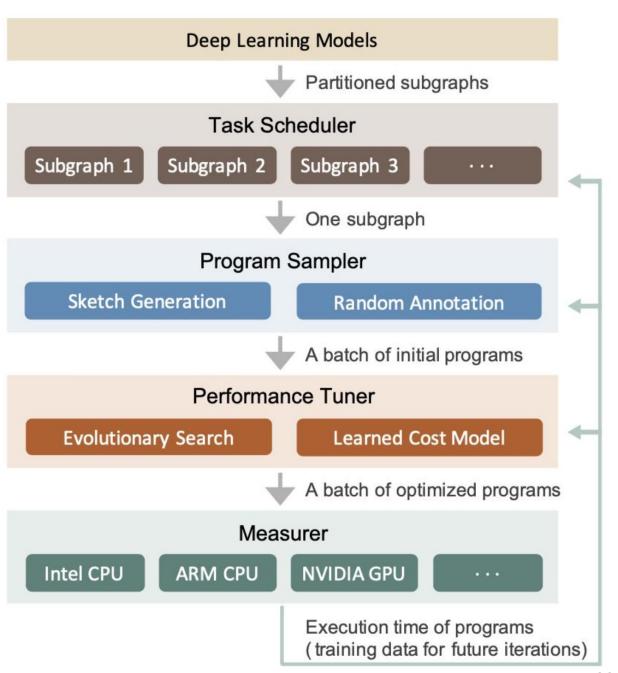
Sample complete programs and fine-tune them



(c) Ansor's Hierarchical Approach

Ansor's Overview

- Inputs: set of DNNs
 - Partitioned into small subgraphs
- Three components
 - Task scheduler
 - Program sampler
 - Performance tuner



Program Sampling

 Goal: automatically construct a large search space and uniformly sample from the space

Approach: two-level hierarchical search space

- Sketch: a few good high-level structures
- Annotation: billions of low-level details
- Sampling process:



Sketch Generation Example 1/2

Example Input 1: * The mathmetical expression: $C[i,j] = \sum A[i,k] \times B[k,j]$ $D[i,j] = \max(C[i,j], 0.0)$ where $0 \le i, j, k < 512$ * The corresponding naïve program: for i in range(512): for j in range(512): for k in range(512): C[i, j] += A[i, k] * B[k, j]for i in range(512): for j in range(512): D[i, j] = max(C[i, j], 0.0)* The corresponding DAG:

```
Derivation:  \begin{array}{c} \textit{Input } 1 \rightarrow \!\! \sigma(S_0, i=4) \xrightarrow{\text{Rule 1}} \!\! \sigma(S_1, i=3) \xrightarrow{\text{Rule 4}} \\ \sigma(S_2, i=2) \xrightarrow{\text{Rule 1}} \!\! \sigma(S_3, i=1) \xrightarrow{\text{Rule 1}} \!\! Sketch \ 1 \end{array}
```

```
Generated sketch 1
for i.0 in range(TILE I0):
  for j.0 in range(TILE J0):
    for i.1 in range(TILE I1):
      for j.1 in range(TILE J1):
        for k.0 in range(TILE K0):
          for i.2 in range(TILE I2):
            for j.2 in range(TILE J2):
               for k.1 in range(TILE I1):
                 for i.3 in range(TILE I3):
                   for j.3 in range(TILE J3):
                     C[\ldots] += A[\ldots] * B[\ldots]
        for i.4 in range(TILE I2 * TILE I3):
          for j.4 in range(TILE J2 * TILE J3):
            D[\ldots] = \max(C[\ldots], 0.0)
```

Sketch Generation Example 2/2

Example Input 2: * The mathmetical expression: $B[i,l] = \max(A[i,l],0.0)$ $C[i, k] = \begin{cases} B[i, k], & k < 400 \\ 0, & k \ge 400 \end{cases}$ $E[i,j] = \sum C[i,k] \times D[k,j]$ where $0 \le i < 8$, $0 \le j < 4$, $0 \le k < 512, 0 \le l < 400$ * The corresponding naïve program: for i in range(8): for 1 in range(400): B[i, 1] = max(A[i, 1], 0.0)for i in range(8): for k in range(512): C[i, k] = B[i, k] if k < 400 else 0 for i in range(8): for j in range(4): for k in range(512): E[i, j] += C[i, k] * D[k, j]* The corresponding DAG:

```
Generated sketch 2
for i in range(8):
  for k in range(512):
    C[i, j] = max(A[i,k], 0.0) if k<400 else 0
for i.0 in range(TILE I0):
  for j.0 in range(TILE J0):
    for i.1 in range(TILE I1):
      for j.1 in range(TILE J1):
        for k.0 in range(TILE K0):
          for i.2 in range(TILE I2):
            for j.2 in range(TILE_J2):
              for k.1 in range(TILE I1):
                for i.3 in range(TILE I3):
                  for j.3 in range(TILE J3):
                     E.cache[...] += C[...] * D[...]
        for i.4 in range(TILE I2 * TILE I3):
          for j.4 in range(TILE_J2 * TILE_J3):
            E[...] = E.cache[...]
```

```
Input 2 \to \sigma(S_0, i = 5) \xrightarrow{\text{Rule } 5} \sigma(S_1, i = 5) \xrightarrow{\text{Rule } 4} 
\sigma(S_2, i = 4) \xrightarrow{\text{Rule } 1} \sigma(S_3, i = 3) \xrightarrow{\text{Rule } 1} 
\sigma(S_4, i = 2) \xrightarrow{\text{Rule } 2} \sigma(S_5, i = 1) \xrightarrow{\text{Rule } 1} Sketch 2
```

Input
$$2 \rightarrow \sigma(S_0, i = 5) \xrightarrow{\text{Rule } 6} \sigma(S_1, i = 4) \xrightarrow{\text{Rule } 1}$$

$$\sigma(S_2, i = 3) \xrightarrow{\text{Rule } 1} \sigma(S_3, i = 2) \xrightarrow{\text{Rule } 2}$$

$$\sigma(S_4, i = 1) \xrightarrow{\text{Rule } 1} Sketch 3$$

Random Annotation Examples

```
Generated sketch 1
for i.0 in range(TILE I0):
  for j.0 in range(TILE J0):
    for i.1 in range(TILE I1):
      for j.1 in range(TILE J1):
        for k.0 in range(TILE K0):
          for i.2 in range(TILE I2):
            for j.2 in range(TILE J2):
              for k.1 in range(TILE I1):
                 for i.3 in range(TILE I3):
                   for j.3 in range(TILE J3):
                     C[\ldots] += A[\ldots] * B[\ldots]
        for i.4 in range(TILE I2 * TILE I3):
          for j.4 in range(TILE J2 * TILE J3):
            D[\ldots] = \max(C[\ldots], 0.0)
```

Sampled program 2

Performance Fine-Tuning

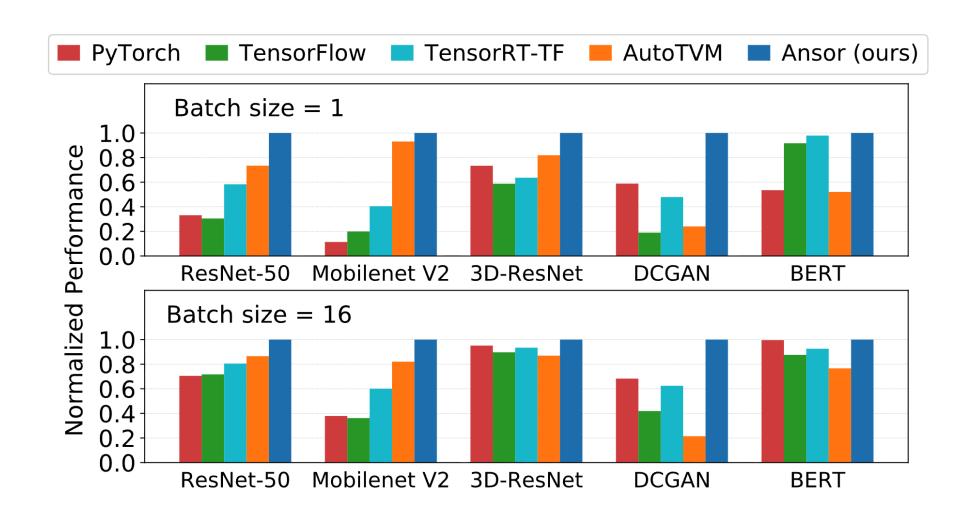
Issue: Random samplings does not guarantee high performance

```
Generated sketch 1
                                                         Sampled program 1
for i.0 in range(TILE I0):
                                                         parallel i.0@j.0@i.1@j.1 in range(256):
  for j.0 in range(TILE J0):
                                                            for k.0 in range(32):
    for i.1 in range(TILE I1):
                                                              for i.2 in range(16):
      for j.1 in range(TILE J1):
                                                                unroll k.1 in range(16):
        for k.0 in range(TILE K0):
                                                                  unroll i.3 in range(4):
           for i.2 in range(TILE_I2):
                                                                    vectorize j.3 in range(16):
             for j.2 in range(TILE J2):
                                                                      C[\ldots] += A[\ldots] * B[\ldots]
               for k.1 in range(TILE I1):
                                                            for i.4 in range(64):
                 for i.3 in range(TILE I3):
                                                             vectorize j.4 in range(16):
                   for j.3 in range(TILE_J3):
                                                                D[\ldots] = \max(C[\ldots], 0.0)
                     C[\ldots] += A[\ldots] * B[\ldots]
        for i.4 in range(TILE I2 * TILE I3):
           for j.4 in range(TILE J2 * TILE J3):
             D[\ldots] = \max(C[\ldots], 0.0)
```

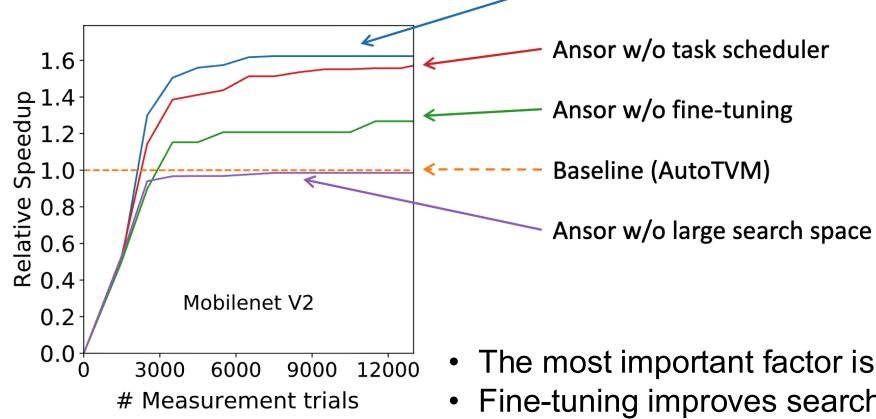
Solution: Perform evolutionary search with learned cost model (AutoTVM)

- Randomly mutate tile size
- Randomly mutate parallel/unroll/vectorize factor and granularity

Ansor Automatically Finds High-Performance Kernels



Ablation Study



- The most important factor is the search space
- Fine-tuning improves search results

Ansor

Match AutoTVM's performance with 10x less time

Comparing Triton, AutoTVM, and Ansor

	Search Space	Search Algorithm
triton.autotune	Manual template	Exhaustive
AutoTVM	Manual template	Learning-based cost model
Ansor	Sketch generation + random annotation	Learning-based cost model