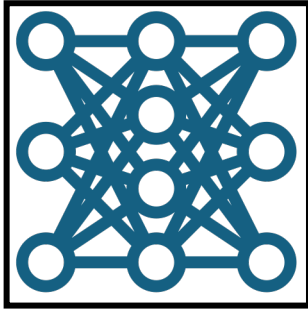# Discovering High-Performance Tensor Programs with Bayesian Optimization

Felix Jonathan Rocke

Supervisors: Eiko Yoneki and Guoliang He
University of Cambridge
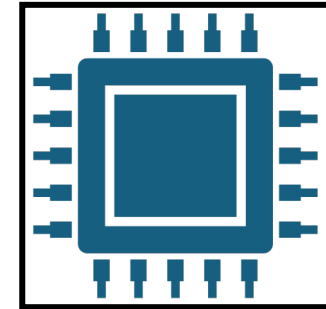
# Motivation: Deploying an ML Application to Device

ML Application

Target Hardware (CPU, GPU, TPU)
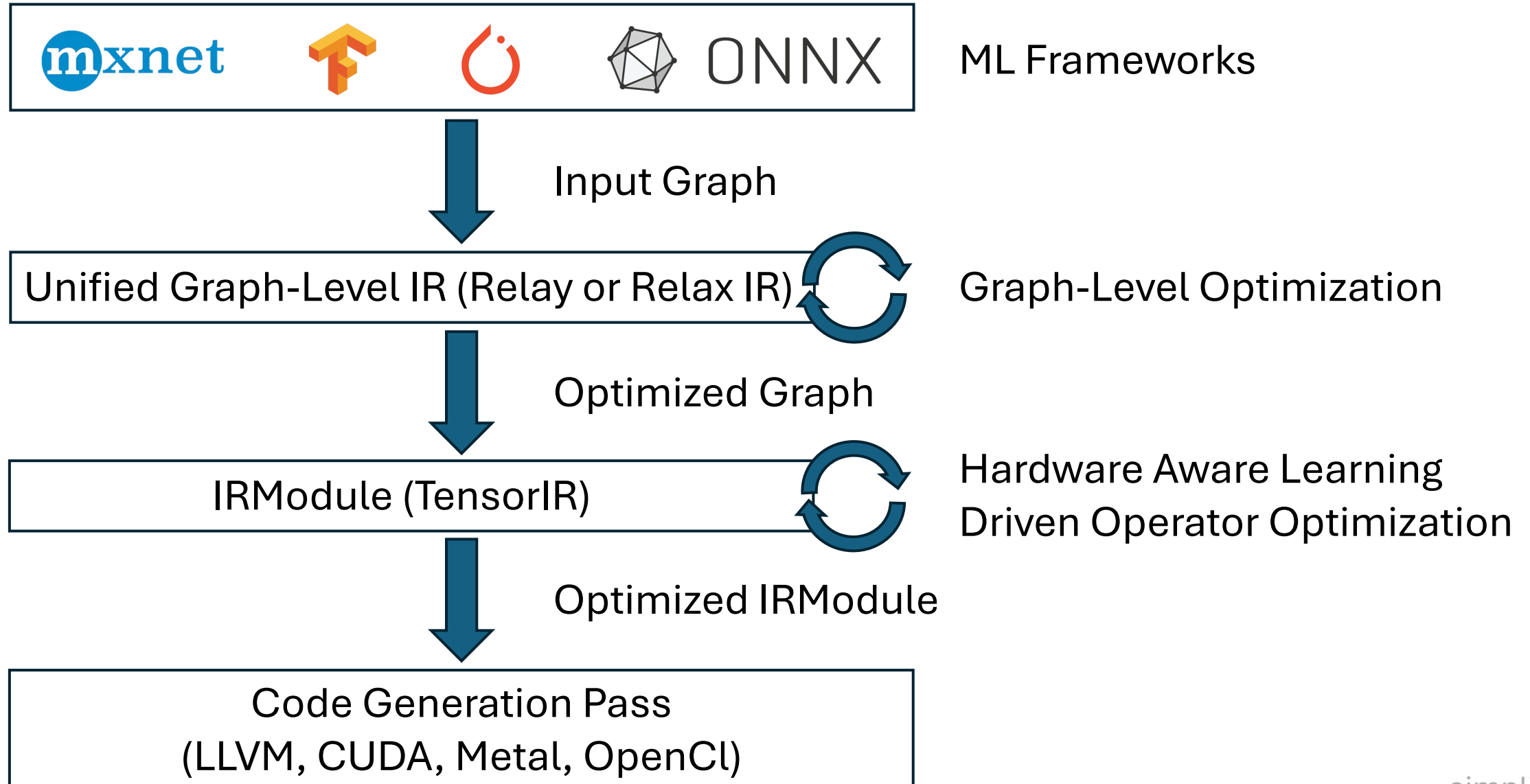
- Various Architectures
- Hundreds of Operators

- Various Microarchitectures
- Increasing Specialization (e.g., Tensor Cores)

Traditionally: Compile model using hardware-specific template libraries. However, high cost due to hardware and model variety

New Solution: A learning compiler that generates optimized kernels without templates
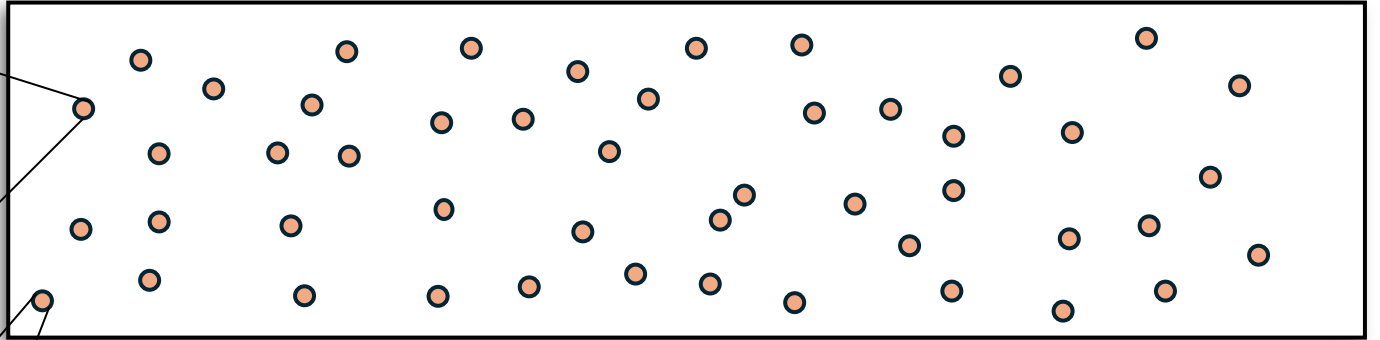
# Overview Apache TVM



ML Frameworks

Input Graph

Unified Graph-Level IR (Relay or Relax IR)    Graph-Level Optimization

Optimized Graph

IRModule (TensorIR)    Hardware Aware Learning Driven Operator Optimization

Optimized IRModule

Code Generation Pass
(LLVM, CUDA, Metal, OpenCl)

simplified

# Problem Description

## Possible Program 1

```
for i₀, j₀ in grid(16, 8):
  for i₁, j₁ in grid(8, 16):
    for k₀ in range(1024):
      for i₂, j₂ in grid(8, 8):
        C[…] += …
```

## Possible Program 2

```
for i₀, j₀ in grid(64, 8):
  for i₁, j₁ in grid(4, 32):
    for k₀ in range(64):
      for i₂, j₂ in grid(4, 4):
        for k₁ in range(16):
          C[…] += …
```

## Search Space of Equivalent Programs



### Challenge:
- Billions of possible equivalent programs
- Minimize the number of program evaluations required to find an efficient implementation

### Proposed Solution:
**Bayesian Optimization** (BO) as it is a sample-efficient search strategy, with the potential to outperform TVM's **Evolutionary Search** (ES)

# Rewriting Programs using Parameterized Transformations

Initial Program $e_0$

```
for i in range(1024):
    C[i] = A[i] + B[i]
```

$e_0 + \text{①}$

Equivalent Program $e_1$

```
for i₀ in range(32):
    for i₁ in range(8):
        for i₂ in range(4):
            i = i₀ * 32 + i₁ * 4 + i₂
            C[i] = A[i] + B[i]
```
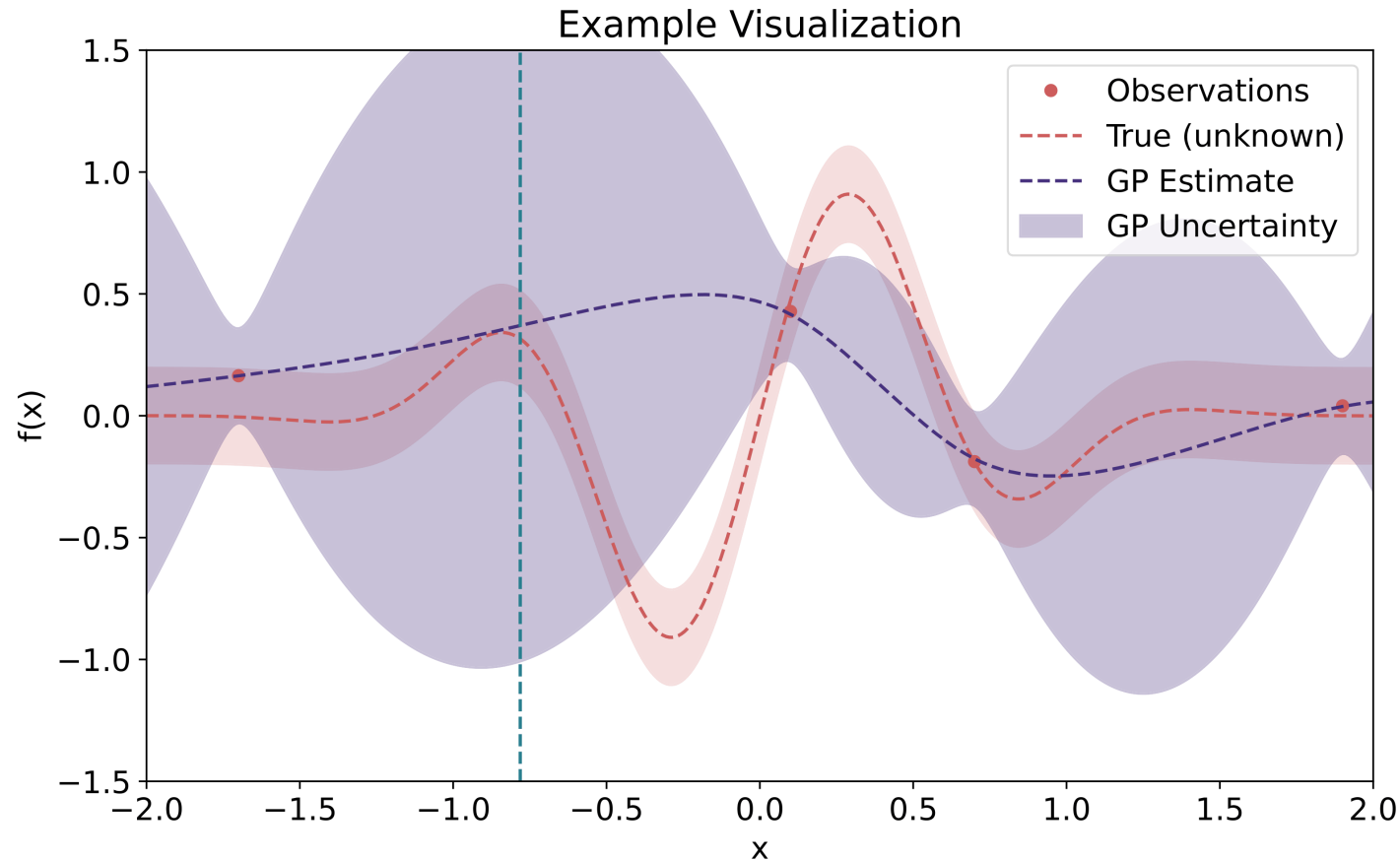
$e_1 + \text{②} + \text{③}$

Transformation Trace

① `Split(loop=i, decision=[32, 8, 4])`

② `Parallelize(loop=i₀)`

③ `Vectorize(loop=i₂)`

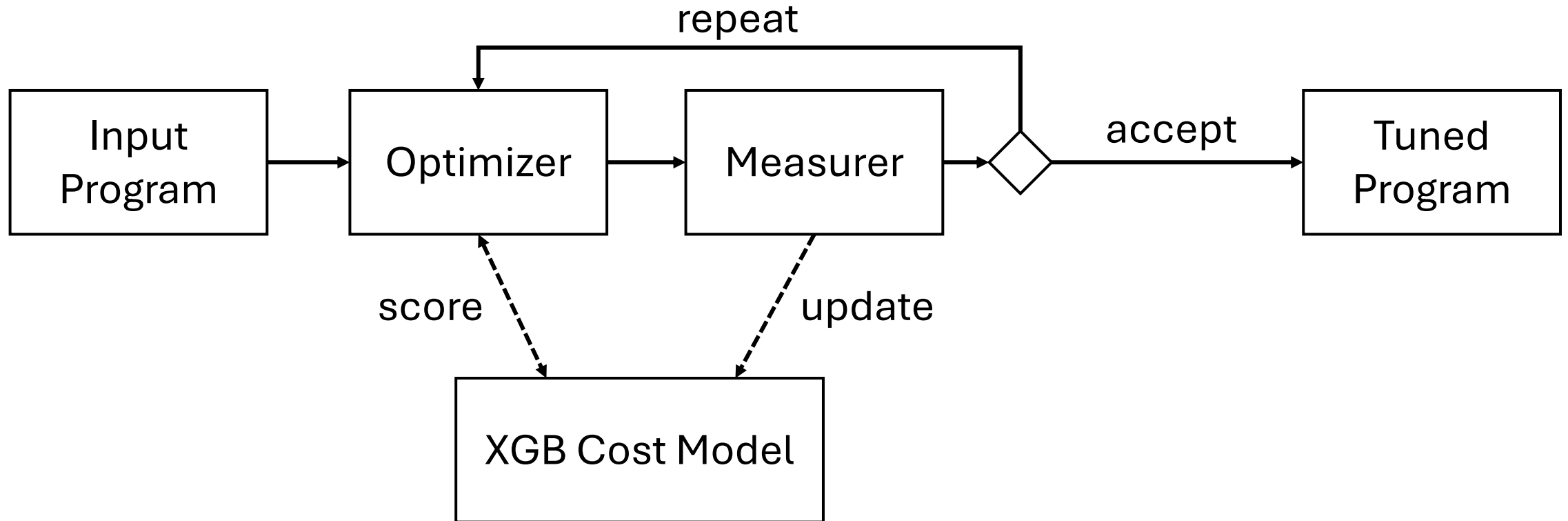Optimized Program $e_2$

```
parallel for i₀ in range(32):
    for i₁ in range(8):
        i = i₀ * 32 + i₁ * 4
        C[i : i + 4] =
            A[i : i + 4] + B[i : i + 4]
```

# Bayesian Optimization (BO)



Example Visualization

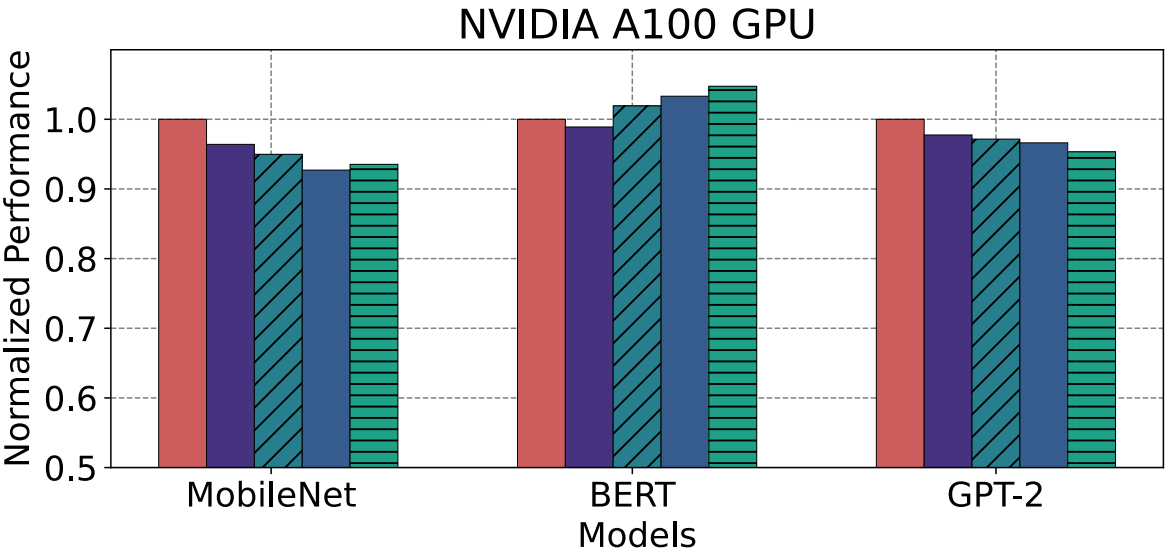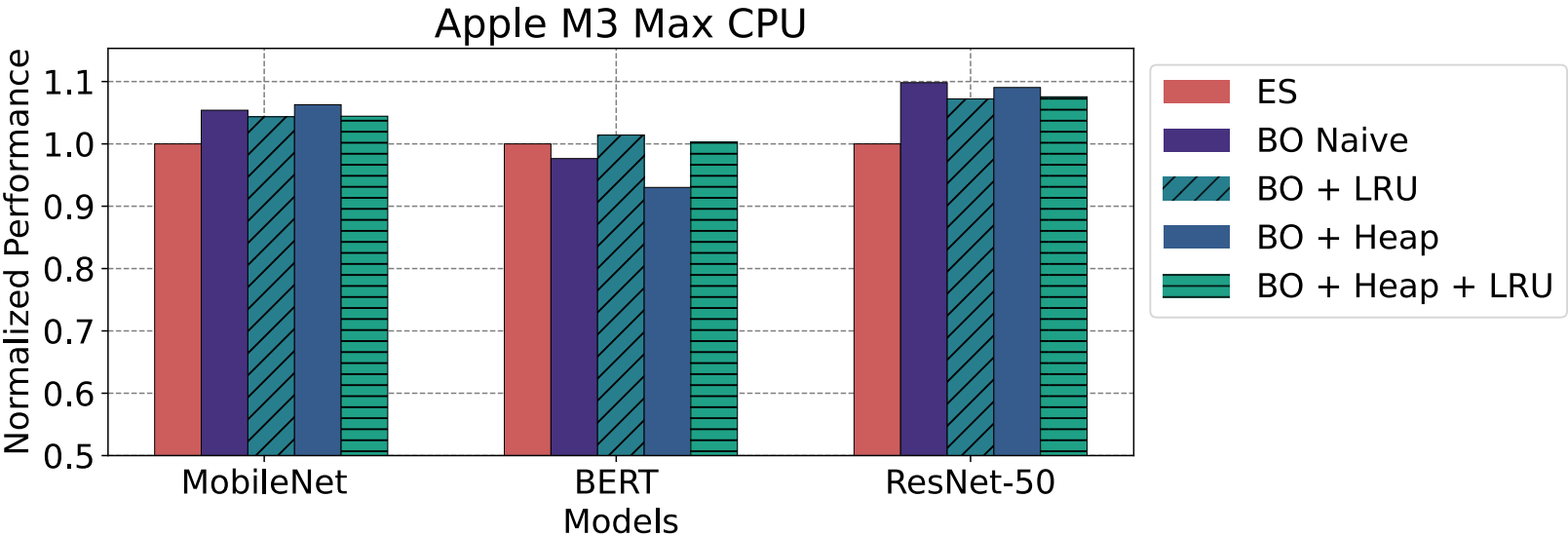GP = Gausian Process
Evaluate f(x) at x = -0.7813

- Sequential design strategy

- Makes informed decisions based on previous observations

- Creates a probabilistic model of the black-box objective function

- Can manage exploration and exploitation

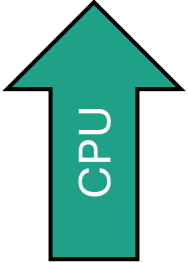- Works best with continuous and smooth functions

# Simplified Search Strategy Overview
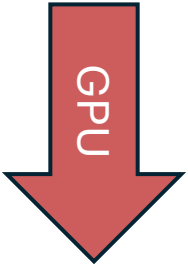
# Benchmarking Results



Apple M3 Max CPU

The Best Configuration

is on average 5.8% faster

NVIDIA A100 GPU

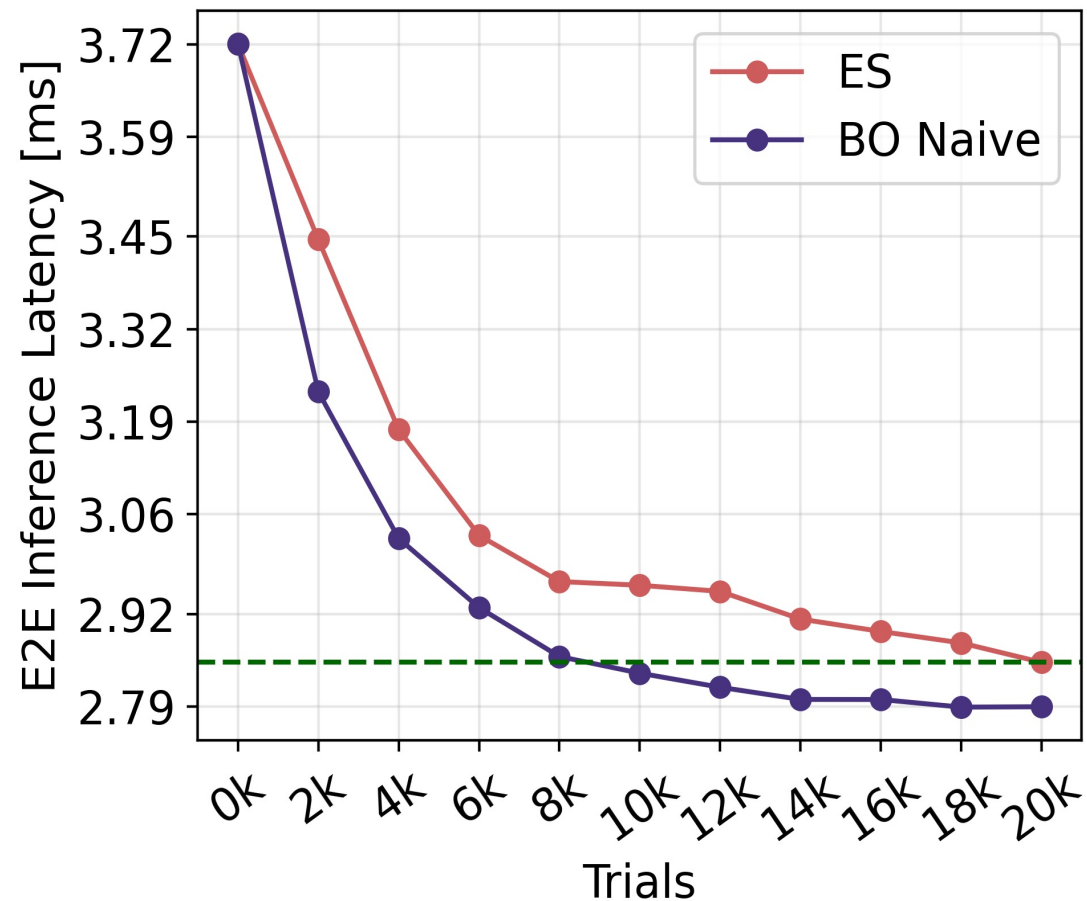Normalized end-to-end inference latency after 6000 trials

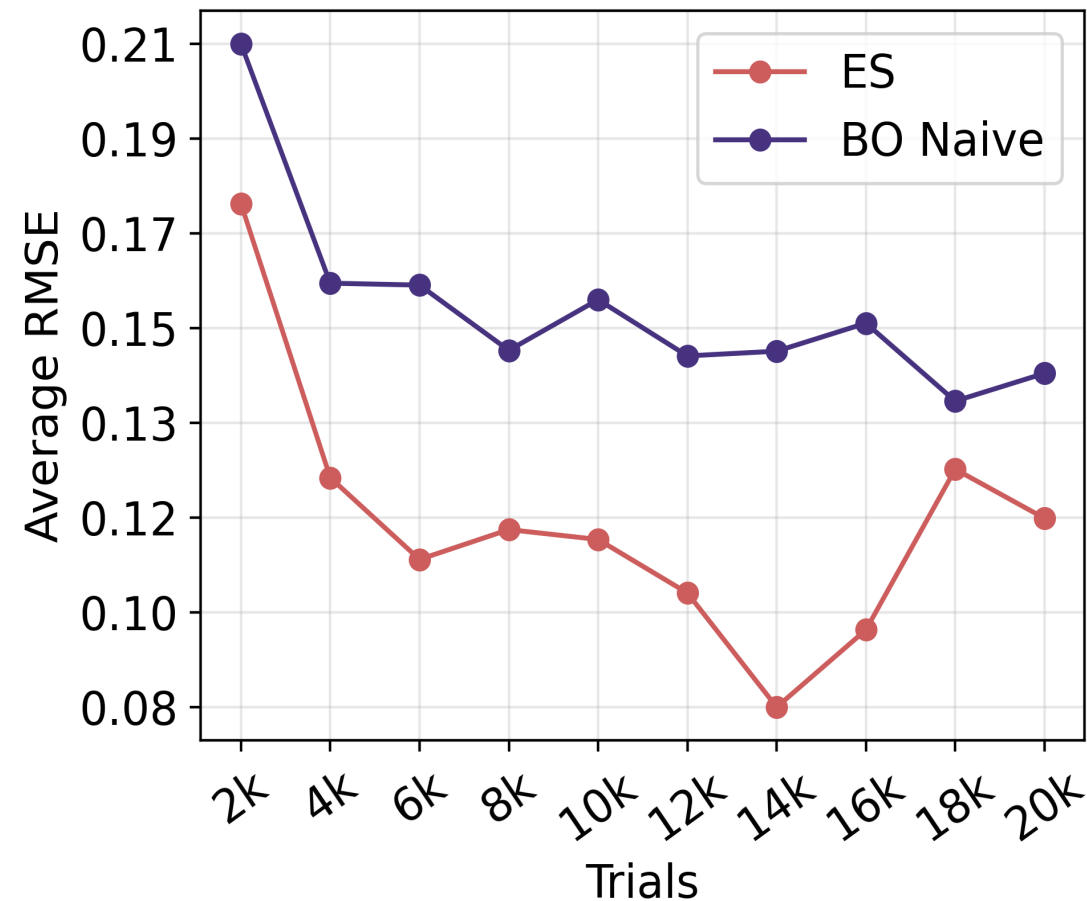is on average 0.4% slower

(all results produced with commit dffe78b)

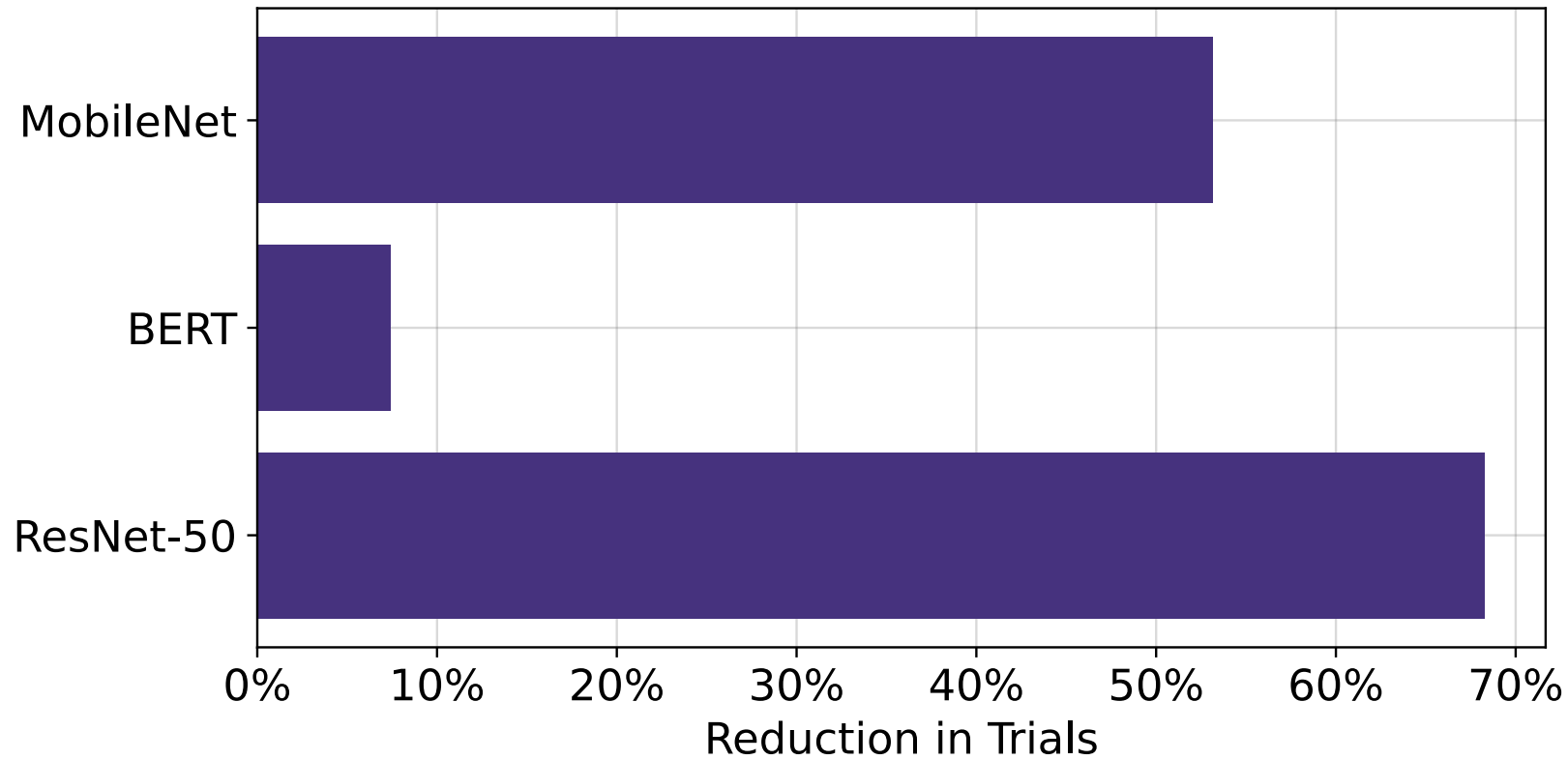# Compiling MobileNet with 20,000 Trials (CPU)
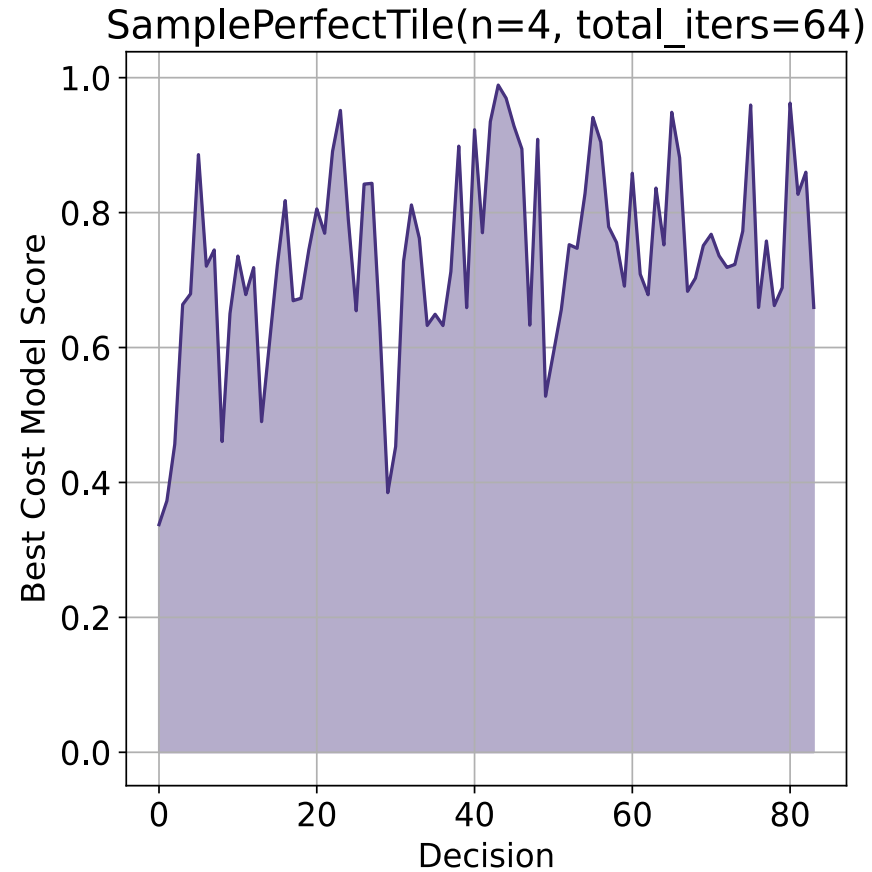
a) Latency after x-trials

b) Cost Model RMSE

# Reduction in Trials Compared to Evolutionary Search on CPUs
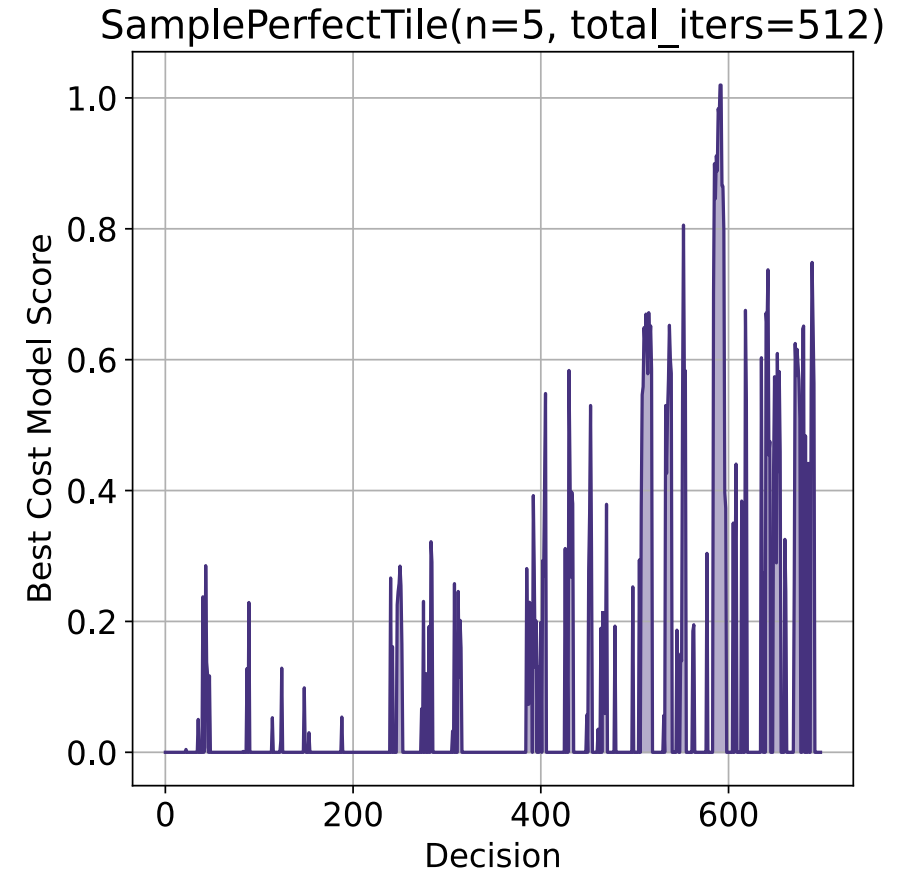


BO reduces the number of trials required to find a high-performance implementation by up to 68%

# Sketch of One-Objective Function Dimension



a) CPU

b) GPU

Problem: The jumps in the GPU's objective function make BO challenging

# Summary

1. Implementation of Bayesian Optimization as a novel search strategy into Apache TVM and MetaSchedule

2. Notable performance improvement of up to 10% for CPU models when compared to ES at the same number of trials, resulting in a reduction of up to 68% in trials

3. For GPUs, BO's performance is more limited

   ➢ However, our analysis of the black-box objective function in the CPU and GPU space will allow future research to be more targeted and efficient