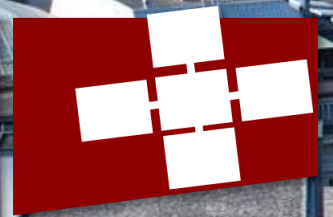


ANDREI IVANOV, SIYUAN SHEN, GIOELE GOTTARDO, MARCIN CHRAPEK, AFIF BOUDAUD, TIMO SCHNEIDER, LUCA BENINI, TORSTEN HOEFLER

PerfDojo: Automated ML Library Generation for Heterogeneous Architectures



Introduction

Specialized High-Performance Accelerators

GPUs



nvidia.com

TPUs / ASICs

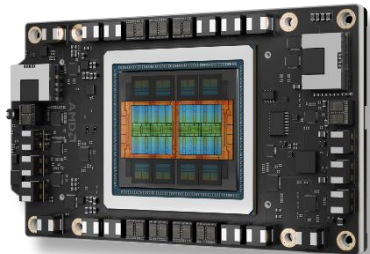


cloud.google.com

FPGAs



amd.com



amd.com

Evolving Processor Architectures

Advanced Vector Extensions (AVX)

Advanced Matrix Extensions (AMX)

ARM Neon

ARM Scalable Vector Extension (SVE)

ARM Scalable Matrix Extension (SME)

AMD Matrix Cores

NVIDIA Tensor Cores

Apple Matrix Coprocessor

RISC-V Stream Semantic Registers (SSR) [1]

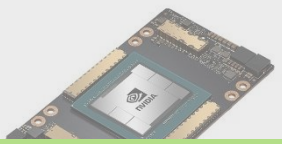
RISC-V Floating-Point Repetition (FREP) [1]

[1] F. Zaruba, F. Schuiki, T. Hoefler and L. Benini, "Snitch: A Tiny Pseudo Dual-Issue Processor for Area and Energy Efficient Execution of Floating-Point Intensive Workloads," in *IEEE Transactions on Computers*

Introduction

Specialized High-Performance Accelerators

GPUs



TPUs / ASICs



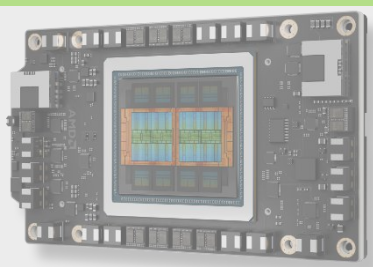
FPGAs



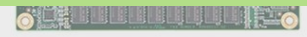
Evolving Processor Architectures

- Advanced Vector Extensions (AVX)
- Advanced Matrix Extensions (AMX)
- ARM Neon
- ARM Scalable Vector Extension (SVE)

Fully utilizing the peak performance is increasingly difficult



amd.com



cloud.google.com



amd.com

NVIDIA Tensor Cores

Apple Matrix Coprocessor

RISC-V Stream Semantic Registers (SSR) [1]

RISC-V Floating-Point Repetition (FREP) [1]

[1] F. Zaruba, F. Schuiki, T. Hoefer and L. Benini, "Snitch: A Tiny Pseudo Dual-Issue Processor for Area and Energy Efficient Execution of Floating-Point Intensive Workloads," in *IEEE Transactions on Computers*

Introduction

Language Models are Unsupervised Multitask Learners

Mistral 7B

Google DeepMind

2024-02-21

Gemma: Open Models Based on Gemma

OUTRAGEOUSLY LARGE NEURAL NETWORKS:
 THE SPARSELY-GATED MIXTURE-OF-EXPERTS LAYER

A SIMPLE AND EFFECTIVE PRUNING APPROACH FOR
 LARGE LANGUAGE MODELS Jeffrey

Mingji
 Carne

QLoRA: Efficient Finetuning of Quantized LLMs

POST-TRAINING QUANTIZATION
 OF PRE-TRAINED TRANSFORMERS

PerfDojo

abstract code-transformation environment for hardware-specific optimizations

academically
 (2 billion parameters)
 outperforming

Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xirui Pan, Bing Shen, Sheng Shao, Zheng Tang, Jianhong Tu, Peng Wang, Shijia Wang, Yipu Wang, An Yang, Yang Yao, Zhiyuan Yao, Zhenqiang Zhang, and Jun Zhang

ABSTRACT

PerfLLM

automatic optimization methodology, optimizing code without prior architecture knowledge

We present a
 based on
 (2023).
 We take
 of text
 recipes
 Gemini
 capability
 art unc
 With the
 fine-tune
 codebase

Gemini
 eter m
 ment c
 eter m
 Each ci

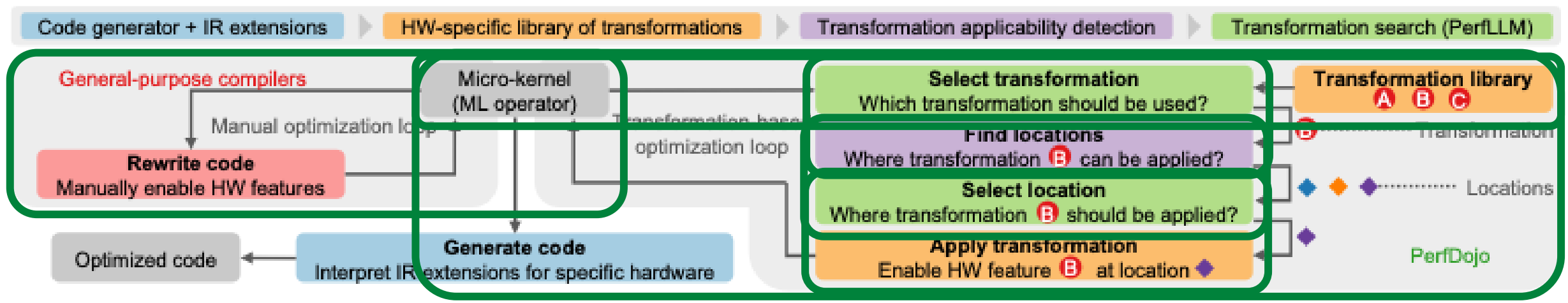
ABSTRACT

Large language models (LLMs) have revolutionized the field of artificial intelligence, enabling natural language processing tasks that were previously thought to be exclusive to humans. In this work, we introduce QWEN¹, the first installment of our large language model series. QWEN is a comprehensive language model series that encompasses distinct models with varying parameter counts. It includes QWEN, the base pretrained language models, and QWEN-CHAT, the chat models finetuned with human alignment techniques. The base language models consistently demonstrate superior performance across a multitude of downstream tasks, and the chat models, particularly those trained using Reinforcement Learning from Human Feedback (RLHF), are highly competitive. The chat models possess advanced tool-use and planning capabilities for creating agent applications, showcasing impressive performance even when compared to bigger models on complex tasks like utilizing a code interpreter. Furthermore, we have developed coding-specialized models, CODE-QWEN and CODE-QWEN-CHAT, as well as mathematics-focused models, MATH-QWEN-CHAT, which are built upon base language models. These models demonstrate significantly improved performance in comparison with open-source models, and slightly fall behind the proprietary models.

Large language models (LLMs) have revolutionized the field of artificial intelligence, enabling natural language processing tasks that were previously thought to be exclusive to humans. In this work, we introduce QWEN¹, the first installment of our large language model series. QWEN is a comprehensive language model series that encompasses distinct models with varying parameter counts. It includes QWEN, the base pretrained language models, and QWEN-CHAT, the chat models finetuned with human alignment techniques. The base language models consistently demonstrate superior performance across a multitude of downstream tasks, and the chat models, particularly those trained using Reinforcement Learning from Human Feedback (RLHF), are highly competitive. The chat models possess advanced tool-use and planning capabilities for creating agent applications, showcasing impressive performance even when compared to bigger models on complex tasks like utilizing a code interpreter. Furthermore, we have developed coding-specialized models, CODE-QWEN and CODE-QWEN-CHAT, as well as mathematics-focused models, MATH-QWEN-CHAT, which are built upon base language models. These models demonstrate significantly improved performance in comparison with open-source models, and slightly fall behind the proprietary models.

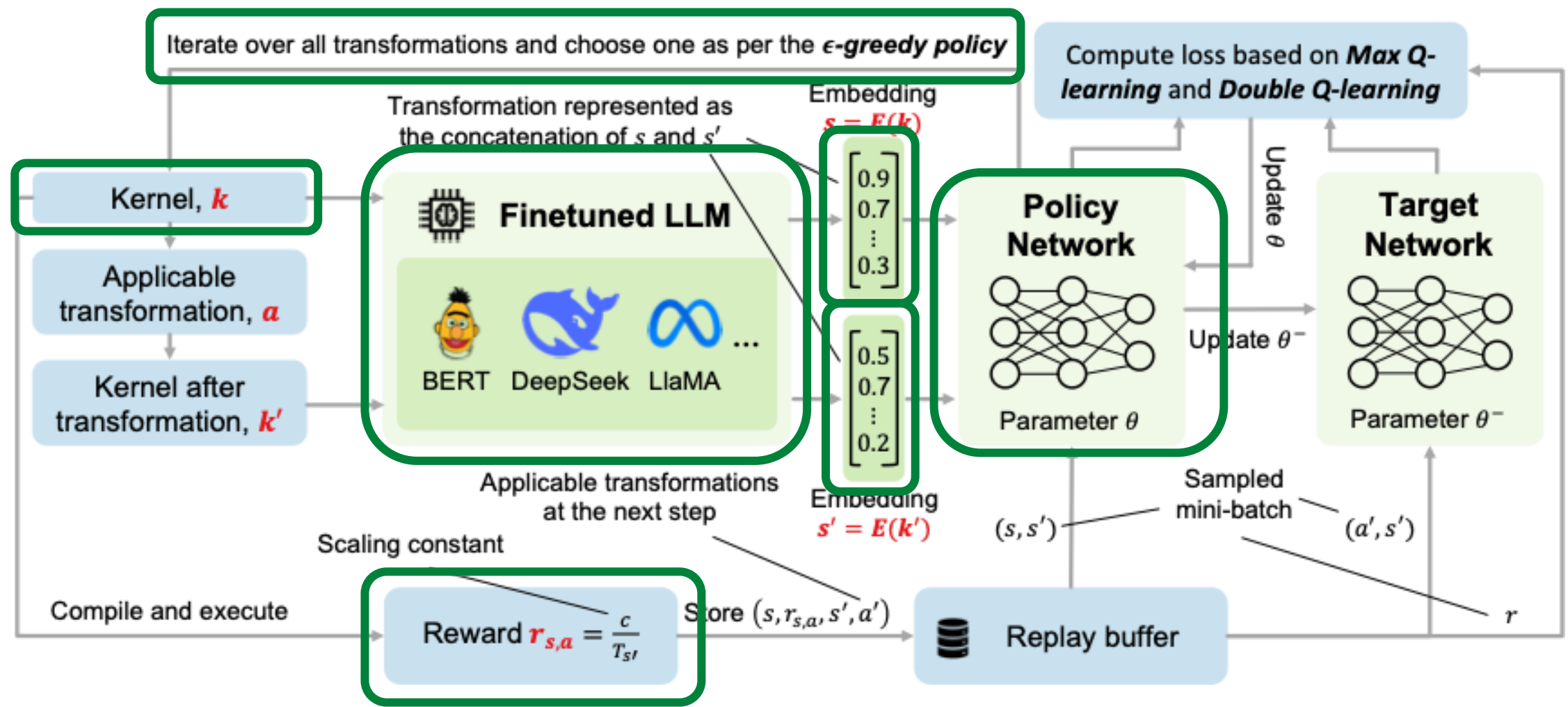
model compression, the applicability and performance of existing compression techniques is limited by the scale and complexity of GPT models. In this paper, we address this challenge, and propose GPTQ, a new one-shot weight quantization method based on approximate second-order information, that is both highly-accurate and highly-efficient. Specifically, GPTQ can quantize GPT models with 175 billion parameters in approximately four GPU hours, reducing the bitwidth down to 3 or 4 bits per weight, with negligible accuracy degradation relative to the uncompressed baseline. Our method more than doubles the compression gains relative to previously-proposed one-shot quantization methods, preserving accuracy, allowing us for the first time to execute an 175 billion-parameter model inside a single GPU for generative inference. Moreover, we also show that our method can still provide reasonable accuracy in the *extreme quantization* regime, in which weights are quantized to 2-bit or even *ternary* quantization levels. We show experimentally that these improvements can be leveraged for end-to-end inference speedups over FP16, of around 3.25x when using high-end GPUs (NVIDIA A100) and 4.5x when using more cost-effective ones (NVIDIA A6000). The implementation is available at <https://github.com/IST-DASLab/gptq>.

PerfDojo: Manual vs. PerfDojo's transformation-centric optimization workflows



The transformation-centric workflow enforces semantic preservation by construction

PerfDojo: Overview of the PerfLLM training pipeline



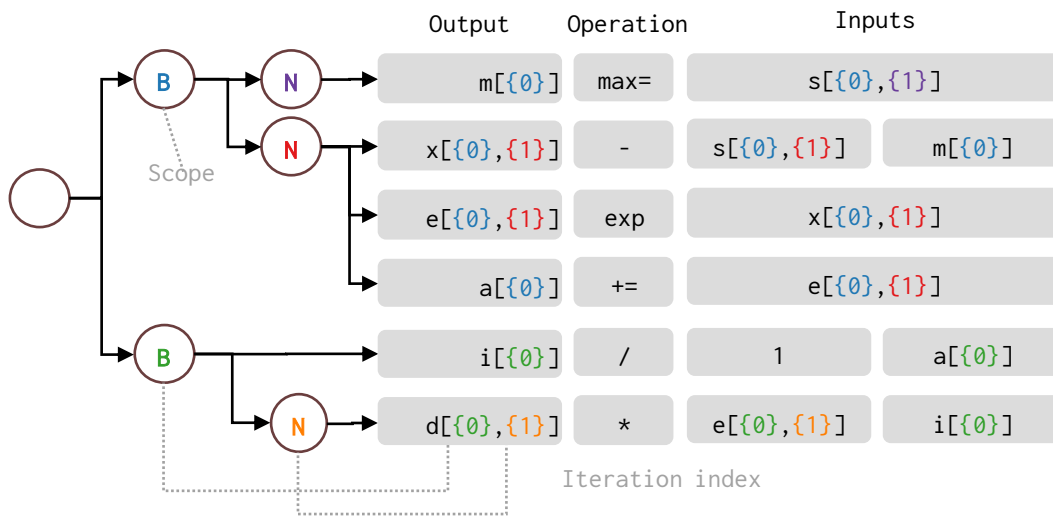
PerfDojo: Representation

Softmax kernel:

$$\begin{aligned}
 m_b &= \max_n s_{b,n} \\
 b \in \{1..B\} & \quad x_{b,n} = s_{b,n} - m_b \\
 n \in \{1..N\} & \quad e_{b,n} = \exp x_{b,n} \\
 a_b &= \sum_n e_{b,n} \\
 i_b &= 1/a_b \\
 d_{b,n} &= e_{b,n} i_b
 \end{aligned}$$

```

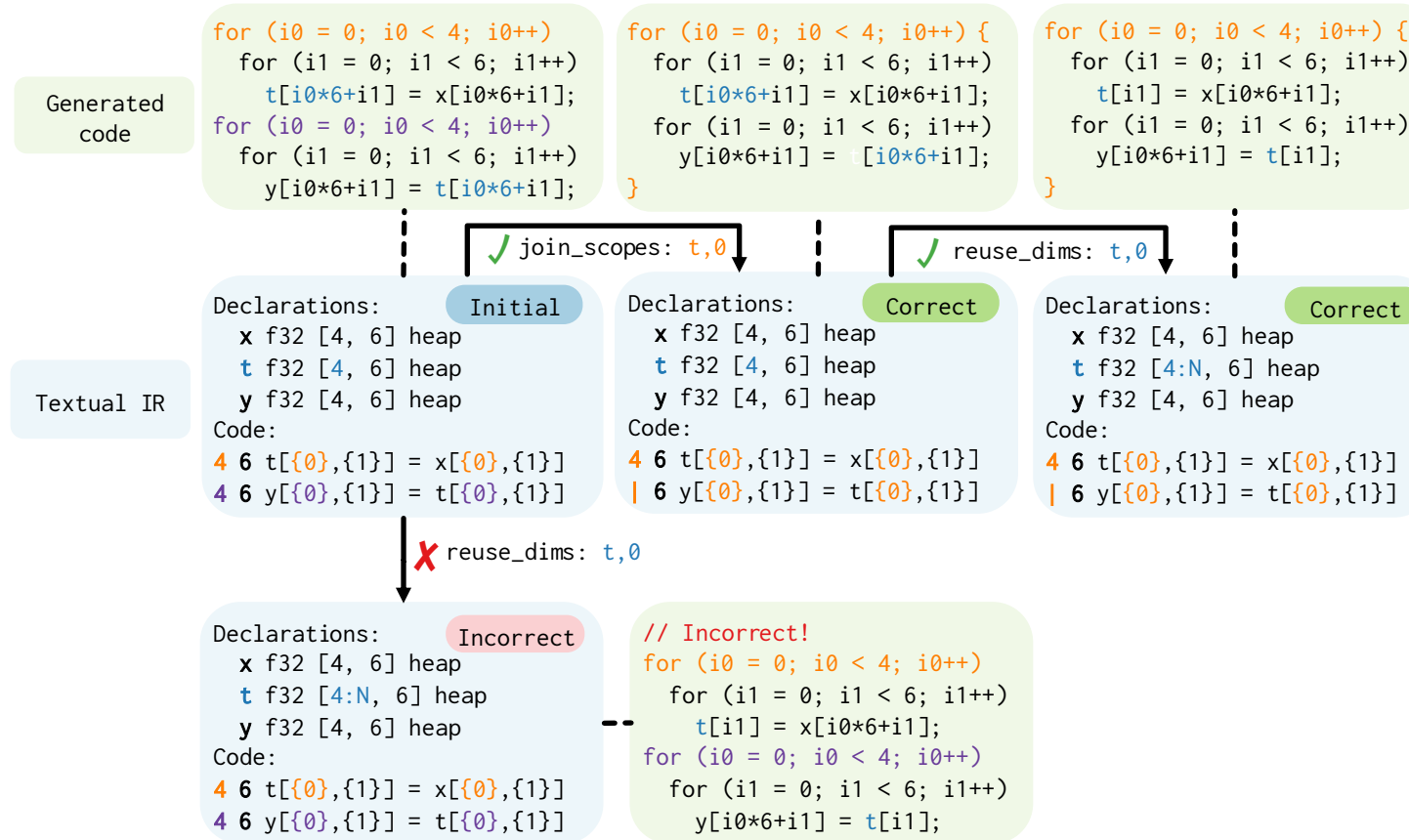
for (b=0;b<B;b++) {
  m[b] = -DBL_MAX; for (n=0;n<N;n++) { m[b] = max(m[b],s[b*N+n]); }
  a[b] = 0;        for (n=0;n<N;n++) { x[b*N+n] = s[b*N+n]-m[b];
                                     e[b*N+n] = exp(x[b*N+n]);
                                     a[b] += e[b*N+n]; }
  for (b=0;b<B;b++) {
    i[b] = 1 / a[b];
    for (n=0;n<N;n++) { d[b*N+n] = e[b*N+n]*i[b]; }
  }
}
  
```



```

B N m[\{0\}] max= s[\{0\},\{1\}]
| N x[\{0\},\{1\}] = s[\{0\},\{1\}] - m[\{0\}]
| | e[\{0\},\{1\}] = exp( x[\{0\},\{1\}] )
| | a[\{0\}] += e[\{0\},\{1\}]
B i[\{0\}] = 1 / a[\{0\}]
| N d[\{0\},\{1\}] = e[\{0\},\{1\}] * i[\{0\}]
  
```

PerfDojo: Transformations



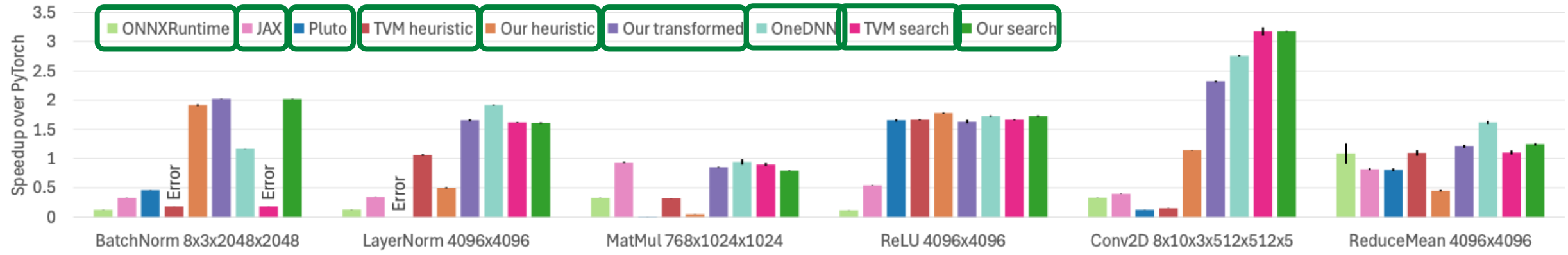
Transformation	Description	Example before	Example after	Transformation	Description	Example before	Example after
swap_ops	Swaps adjacent operations or scopes	$N \ y[{\{0\}}]=x[{\{0\}}]$ $z=q$	$z=q$ $N \ y[{\{0\}}]=x[{\{0\}}]$	reuse_dims	Avoids materializing a dimension	$x \ f32 \ [M, \ K] \ \text{heap}$	$x \ f32 \ [M:N, \ K] \ \text{heap}$
tile_scope	Tiles a scope	$10 \ y+=x[{\{0\}}]$	$5 \ 2 \ y+=x[{\{0\}}*2+{\{1\}}]$	reuse_buffers	Shares storage between arrays	$x \ f32[M,K] \ \text{heap}$ $y \ f32[M,K] \ \text{heap}$	$x \ f32[M,K] \ \text{heap} \rightarrow x,y$
tile_buffer	Tiles a buffer dimension	$M \ N \ y+=x[{\{0\}}*N+{\{1\}}]$	$M \ N \ y+=x[{\{0\}},{\{1\}}]$	swap_dims	Changes the memory layout of an array	$x \ f32 \ [M, \ N, \ K] \ \text{heap}$	$x \ f32 \ [M, \ K, \ N] \ \text{heap}$
swap_scopes	Reorders adjacent scopes	$M \ N \ y+=x[{\{0\}},{\{1\}}]$	$N \ M \ y+=x[{\{1\}},{\{0\}}]$	increment_padding	Increases padding	$x \ f32 \ [M,K+2] \ \text{heap}$	$x \ f32 \ [M,K+3] \ \text{heap}$
join_scopes	Joins the bodies of adjacent scopes	$M \ y[{\{0\}}]=x[{\{0\}}]$ $M \ p[{\{0\}}]=q[{\{0\}}]$	$M \ y[{\{0\}}]=x[{\{0\}}]$ $ \ p[{\{0\}}]=q[{\{0\}}]$	repeat_padding	Doubles the padding	$x \ f32 \ [M+3,K+3] \ \text{heap}$	$x \ f32 \ [M+6,K+3] \ \text{heap}$
create_temporary	Introduces a temporary array copy	$M \ y+=x[{\{0\}}]$	$M \ t[{\{0\}}]=x[{\{0\}}]$ $M \ y+=t[{\{0\}}]$	move_to_stack	Moves an array to the stack	$x \ f32 \ [N] \ \text{heap}$	$x \ f32 \ [N] \ \text{stack}$

27 transformations with applicability detection

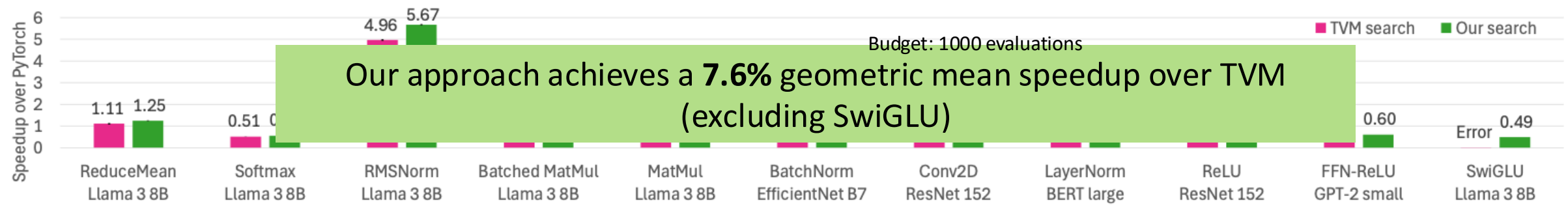
unroll_scope	Unrolls a scope	$N \ y[{\{0\}}]=x[{\{0\}}]$	$N:u \ y[{\{0\}}]=x[{\{0\}}]$	vectorize_buffer	Vectorizes storage of a buffer	$x \ f32 \ [4,16] \ \text{stack}$	$x \ f32x16 \ [4,16] \ \text{stack}$
enable_frep	Instantiates a scope using FREP	$N \ y+= \backslash$ $x[{\{0\}} s:5,d:1,h:2]$	$N:f \ y+= \backslash$ $x[{\{0\}} s:5,d:1,h:2]$	enable_ssr	Enables SSR for an array access	$M \ N \ y+=$ $x[{\{0\}},{\{1\}}]$	$M \ N \ y+=$ $x[{\{0\}},{\{1\}} s:5,d:1]$
parallelize	Parallelizes a scope using OpenMP	$N \ y[{\{0\}}]=x[{\{0\}}]$	$N:p \ y[{\{0\}}]=x[{\{0\}}]$	increase_ssr_depth	Increases the SSR depth	$M \ N \ y+= \backslash$ $x[{\{0\}},{\{1\}} s:5,d:1]$	$M \ N \ y+= \backslash$ $x[{\{0\}},{\{1\}} s:5,d:2]$
parallelize_grid	Parallelizes a scope as a CUDA grid	$N \ y[{\{0\}}]=x[{\{0\}}]$	$N:g \ y[{\{0\}}]=x[{\{0\}}]$	merge_ssr	Merges multiple SSR accesses into one	$M \ y+=x[{\{0\}} s:5,d:1]$ $M \ y+=z[{\{0\}} s:7,d:1]$	$M \ y+=x[{\{0\}} s:5,d:1]$ $M \ y+=z[{\{0\}} s:5,d:1]$
parallelize_block	Parallelizes a scope as a CUDA block	$N \ y[{\{0\}}]=x[{\{0\}}]$	$N:b \ y[{\{0\}}]=x[{\{0\}}]$	activate_ssr	Activates an SSR configuration	$N \ y+= \backslash$ $x[{\{0\}} s:5,d:1]$	$N \ y+= \backslash$ $x[{\{0\}} s:5,d:1,h:2]$
parallelize_warp	Parallelizes a scope as a CUDA warp	$N \ y[{\{0\}}]=x[{\{0\}}]$	$N:w \ y[{\{0\}}]=x[{\{0\}}]$				
vectorize_scope	Vectorizes a scope	$4 \ y[{\{0\}}]=x[{\{0\}}]$	$4:v \ y[{\{0\}}]=x[{\{0\}}]$				

Evaluation: x86 + heuristics

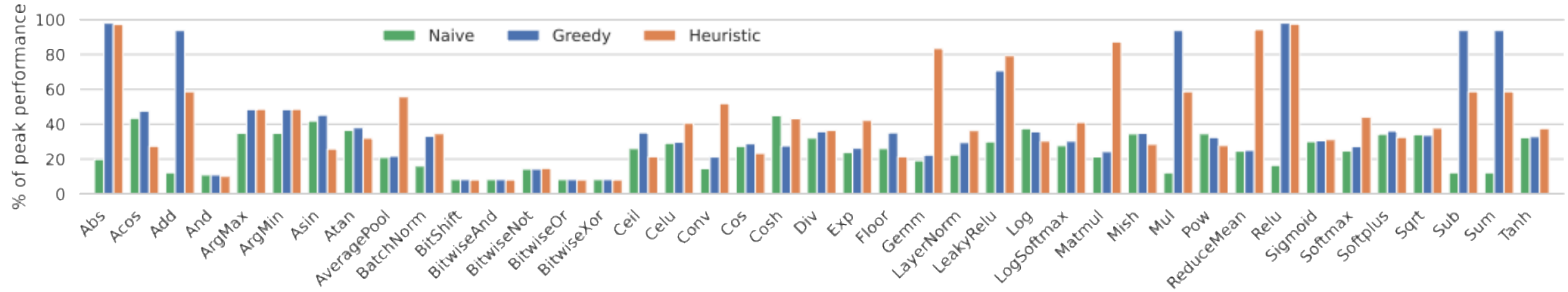
Intel Xeon CPU E5-2695 v4
18 cores



heuristic = single program evaluation pass
search = 1000 evaluations
transformed = applying transformations manually



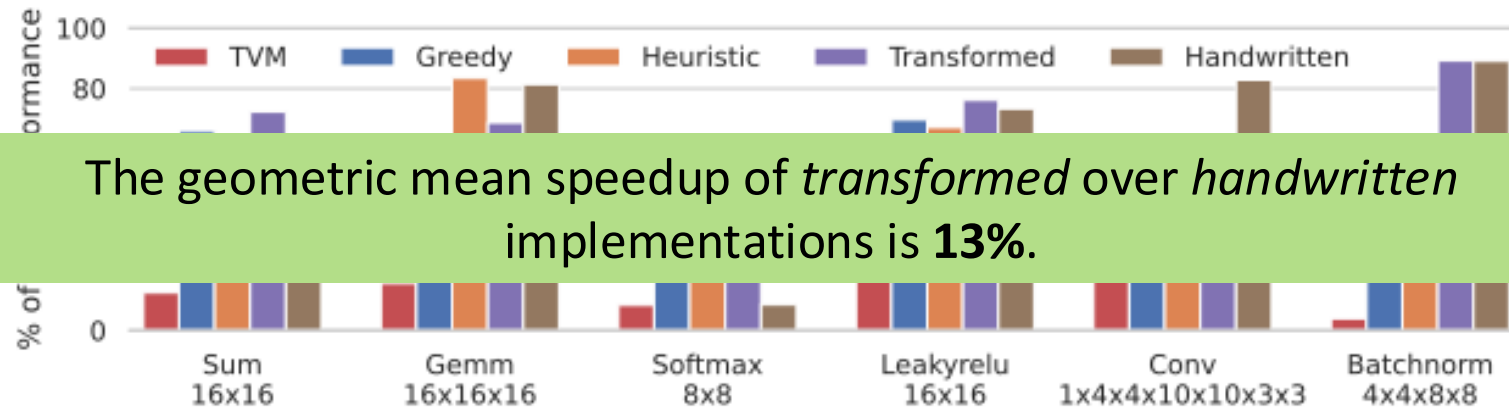
Evaluation: Snitch [1]



Naive = loop fusion and memory reuse until exhaustion

Greedy = *naive* + hardware-specific transformations (SSR/FREP)

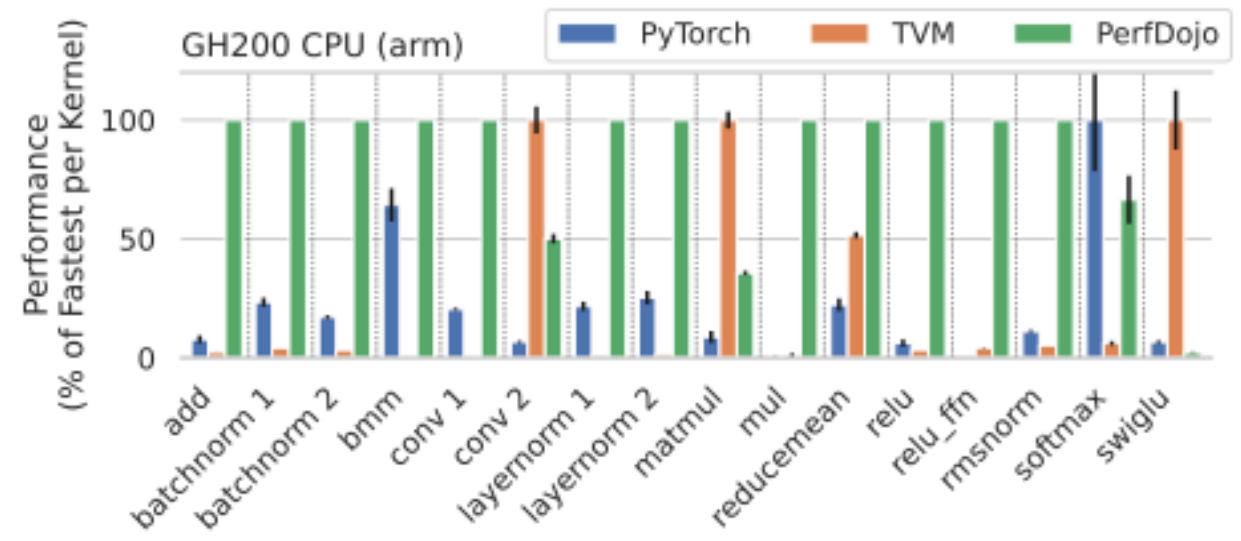
Heuristic = includes analysis pass made by hardware expert to determine the order of transformations



The geometric mean speedup of *transformed* over *handwritten* implementations is **13%**.

[1] F. Zaruba, F. Schuiki, T. Hoefer and L. Benini, "Snitch: A Tiny Pseudo Dual-Issue Processor for Area and Energy Efficient Execution of Floating-Point Intensive Workloads," in *IEEE Transactions on Computers*

PerfDojo: Evaluation on x86 & arm / no heuristics

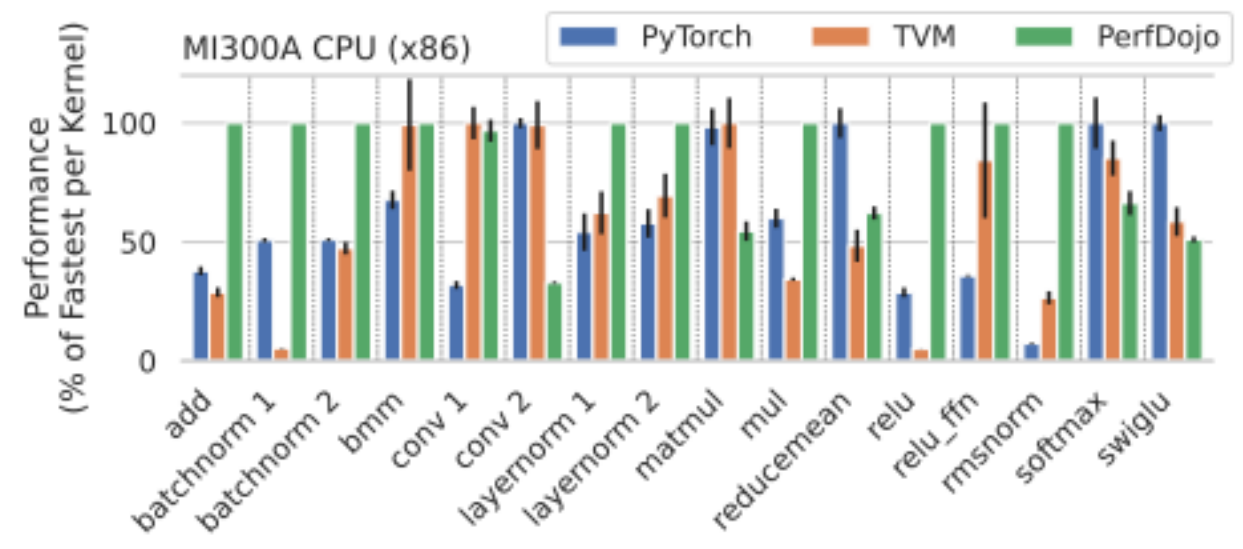


6.65x over PyTorch

13.65x over TVM


1.56x over PyTorch


1.80x over TVM




Conclusions

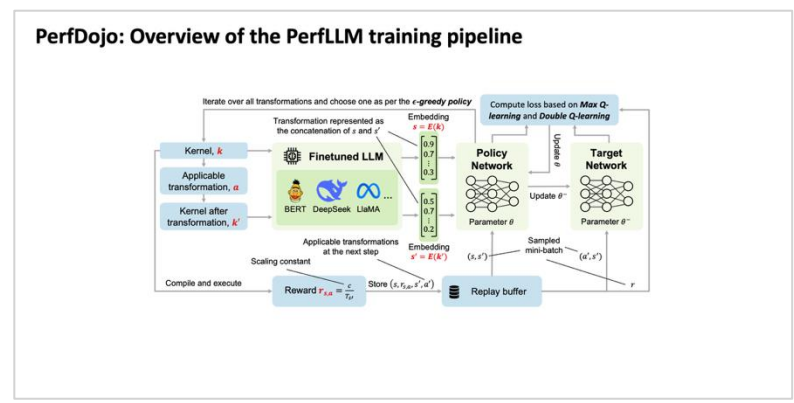
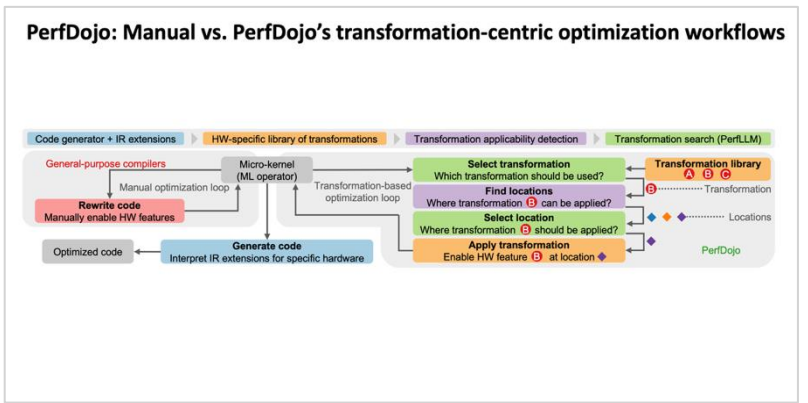
More of SPCL's research:

 youtube.com/@spcl **180+ Talks**

 twitter.com/spcl_eth **1.4K+ Followers**

 github.com/spcl **3.8K+ Stars**

... or spcl.ethz.ch



PerfDojo: Representation and transformations

Softmax kernel:

```

b = (1, #)
n = (1, N)
for (m0=0; m0<#; m0++) {
  m0 = max(m0, s[m0]);
  a[m0] = a;
  for (n0=0; n0<N; n0++) {
    a[n0] = exp(-a[m0]*n[n0]);
    a[n0] = exp(-a[m0]*n[n0]);
  }
  i[m0] = 1 / a[m0];
  d[m0] = a[m0]*i[m0];
}
for (m0=0; m0<#; m0++) {
  for (n0=0; n0<N; n0++) {
    d[n0] = a[m0]*i[m0];
  }
}

```

Generated code:

```

// Initial
for (i1 = 0; i1 < 4; i1++) {
  x[i1] = 1;
  y[i1] = 1;
}
// Correct
for (i1 = 0; i1 < 4; i1++) {
  x[i1] = 1;
  y[i1] = 1;
}
// Incorrect
for (i1 = 0; i1 < 4; i1++) {
  x[i1] = 1;
  y[i1] = 1;
}

```

Textual IR:

```

Declarations: Initial
x F32 [4, 6] heap
y F32 [4, 6] heap
Code:
4 # t(i0, i1) = x(i0, i1)
4 # y(i0, i1) = t(i0, i1)
Declarations: Correct
x F32 [4, 6] heap
y F32 [4, 6] heap
Code:
4 # t(i0, i1) = x(i0, i1)
4 # y(i0, i1) = t(i0, i1)
Declarations: Incorrect
x F32 [4, 6] heap
y F32 [4, 6] heap
Code:
4 # t(i0, i1) = x(i0, i1)
4 # y(i0, i1) = t(i0, i1)

```

Output operation:

```

m0(0) max= s(i0, i1)
m0(0) x(i0, i1) - s(i0, i1)
m0(0) exp= exp(-s(i0, i1)*x(i0, i1))
m0(0) += m0(0)
m0(0) / 1 / m0(0)
m0(0) x(i0, i1) * m0(0)

```

Iteration index:

```

# m m0(0) max= s(i0, i1)
# m m0(0) x(i0, i1) = x(i0, i1) - m0(0)
# m m0(0) exp= exp(-s(i0, i1)*x(i0, i1))
# m m0(0) += m0(0)
# m i(i0) = 1 / m0(0)
# m m0(0) x(i0, i1) = m0(0) * x(i0, i1)

```

