# PerfXLM : a High Performance LLM Inference Engine on RISC-V CPUs

**Xianyi Zhang**
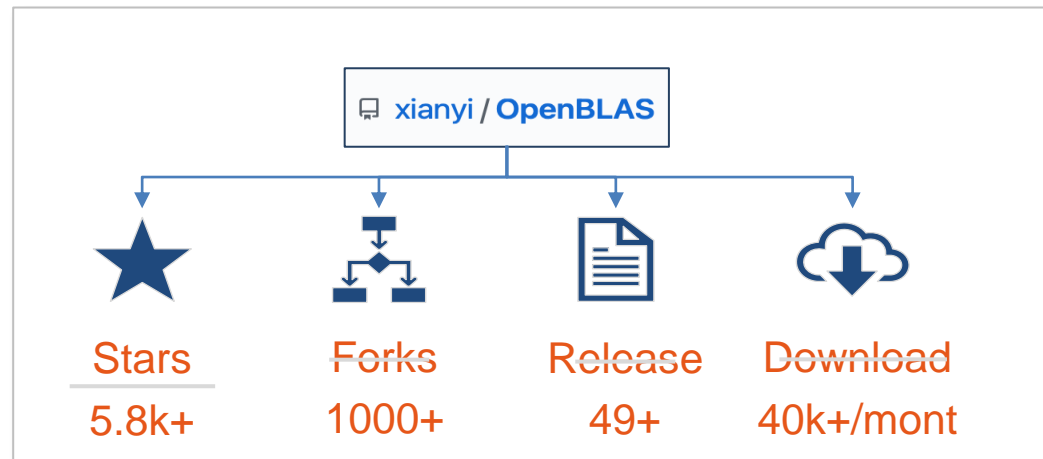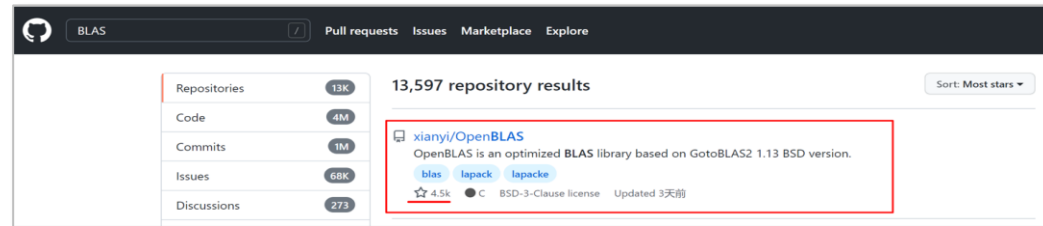
PerfXLab Technologies

xianyi@perfxlab.com

2024.06

# About US

- 2016: PerfXLab Technologies founded in Beijing.
- Xianyi Zhang： CAS Ph.D., OpenBLAS creator, postdoc research at UT Austin & MIT
- Heterogeneous Computing Software Stack and Solutions, including
  - High Performance Computing Libraries
  - Heterogeneous Software Framework
  - Domain Accelerated Computing Solutions (e.g. AI Infra, EDA, Signal Processing, and etc.)



Stars **5.8k+**  Forks **1000+**  Release **49+**  Download **40k+/month**

Created and Maintain OpenBLAS Since 2011



2020 RISC-V Global Forum

2019 CRVA RISC-V Ecosystem

2016 CCF Prize for OpenBLAS



2023 OpenCL Advisory Panel

# PerfXAPI Heterogeneous Computing Framework

**Application**

| AI Training | AI Inference | Scientific/Engineering Computing | Signal Processing |

**Framework Level**

- AI Inference Framework & Tools
- Unified API
- User Extended API
- Scheduler
- AVX to ARM SVE / AVX to RISC-V RVV
- Tuning Toolbox
- CUDA to PerfXAPI

**Perfrormance Level**

- PerfMPL for ARM
- PerfMPL for RISC-V
- PerfMPL for x86
- PerfMPL for GPU OpenCL
- CUDA-X cuBLAS、cuDNN…
- Other libraries for NPU、GPU
- SyCL
- OpenCL

**Hardware Level**

| ARM CPU | x86 CPU | RISC-V CPU | NPU | GPGPU | FPGA |

# PerfXLM Architecture

**Features of PerfXLM Architecture：**

- **Cloud-edge integration**：supporting both cloud and edge inference for large models, making them suitable for various application scenarios.

- **Support for various heterogeneous platforms**：supporting NVIDIA GPU, Hygon DCU, Qualcomm Adreno GPU, Intel iGPU, and other hardware devices.

- **Customized high-performance kernels:** it tailored for hardware characteristics and data properties of large model inference.

cloud

mobile

## Model

| LLaMA 2/3 | ChatGLM 2/3 | Bloomz | ...... |

M Models

## Model Exchange

Model Parser

Quantization

Unified representation of model：ONNX Graph

## PerfXLM Inference Engine

| Serving | Memory Manager | Scheduler |

PerfXLM

## MPL Kernel Library

| mpl.cuda | mpl.opencl | mpl.riscv |
| mpl.x86 | mpl.arm | ...... |

## hardware

| NVIDA GPU | AMD GPU | Hygon DCU | Intel iGPU |
| Domestic GPU | Qualcomm GPU | RISC-V CPU | ...... |

N Hardwares

# PerfXLM Architecture (User interface)

**PerfXLM User interface, vllm compatible**

```python
from vllm import LLM, SamplingParams

# Sample prompts.
prompts = [
    "Hello, my name is",
    "The president of the United States is",
    "The capital of France is",
    "The future of AI is",
]
# Create a sampling params object.
sampling_params = SamplingParams(temperature=0.8, top_p=0.95)

# Create an LLM.
llm = LLM(model="facebook/opt-125m")
# Generate texts from the prompts. The output is a list of RequestOutput objects
# that contain the prompt, generated text, and other information.
outputs = llm.generate(prompts, sampling_params)
# Print the outputs.
for output in outputs:
    prompt = output.prompt
    generated_text = output.outputs[0].text
    print(f"Prompt: {prompt!r}, Generated text: {generated_text!r}")
```

```python
from perfxlm import LLM, SamplingParams

# Sample prompts.
prompts = [
    "Hello, my name is",
    "The president of the United States is",
    "The capital of France is",
    "The future of AI is",
]
# Create a sampling params object.
sampling_params = SamplingParams(temperature=0.8, top_p=0.95)

# Create an LLM.
llm = LLM(model="facebook/opt-125m")
# Generate texts from the prompts. The output is a list of RequestOutput objects
# that contain the prompt, generated text, and other information.
outputs = llm.generate(prompts, sampling_params)
# Print the outputs.
for output in outputs:
    prompt = output.prompt
    generated_text = output.outputs[0].text
    print(f"Prompt: {prompt!r}, Generated text: {generated_text!r}")
```
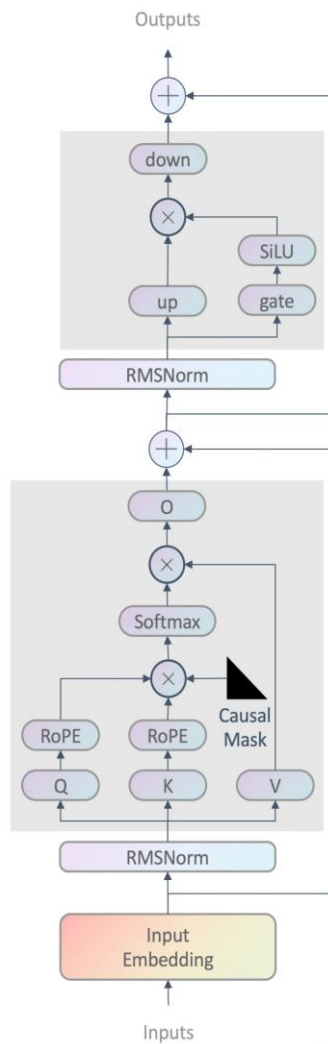
All the user needs to do is change import vllm to import perfxlm.

# PerfXLM Architecture (Create Model)

**PerfXLM Usage: Step 1  Create Model in Python**



LLaMA Model



Pure Python Code and torch style.

# PerfXLM Architecture (Model Parser)

**Parse Python code, export ONNX graph, op fusion**



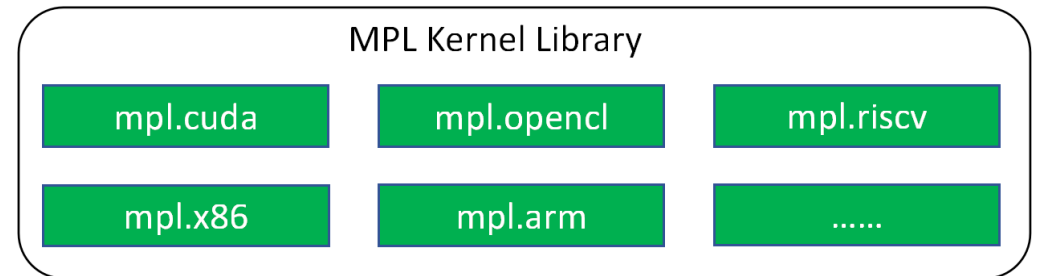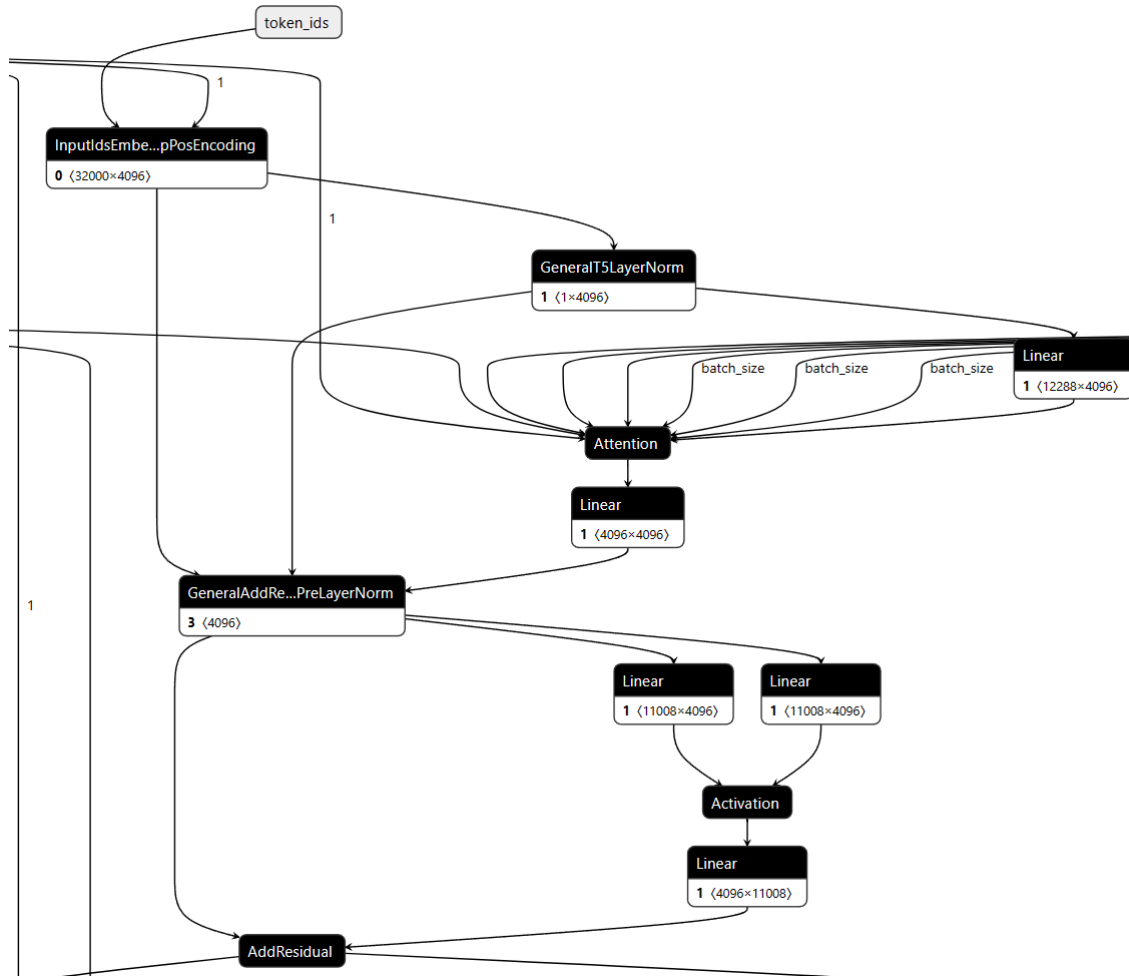For the llama model, only 7 fused OPs are left in the final fusion, which is high performance and easy to maintain.

# PerfXLM Architecture (Engine)

**Parses the op-fused onnx graph and calls the corresponding high-performance operator.**

# PerfXLM Optimization Methods

**Low-precision quantization**

- Quantization of weights
  - AWQ-int4, A16
- Quantization of KV cache
  - int8 kv cache
- For embedded
  - Int4
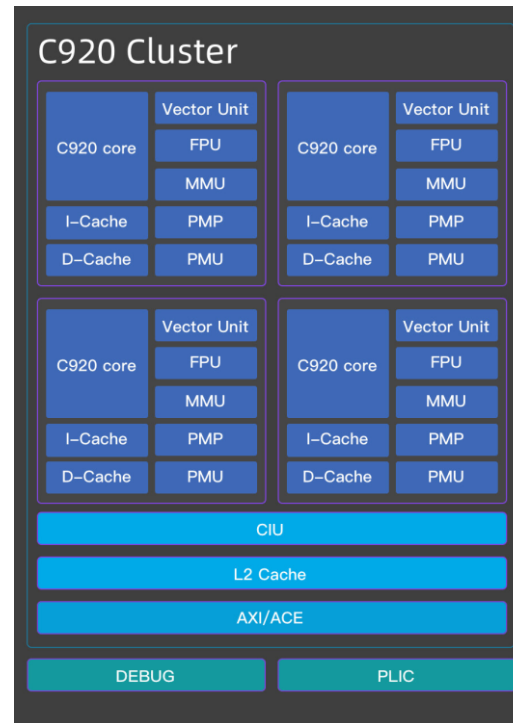- For server
  - Fp8? Int8?
- ......

**Operator fusion**

- Fused add bias and Layernorm
- Fused QKV GEMM
- Fused Decoder Attention
- Flash attention
- ......

**Core Kernel optimization**

- GEMM optimization
  - Tiling to reuse data
  - Vectorized memory access
  - unroll loops
- Tall-skin matrix
- Attention kernel optimization
- ......

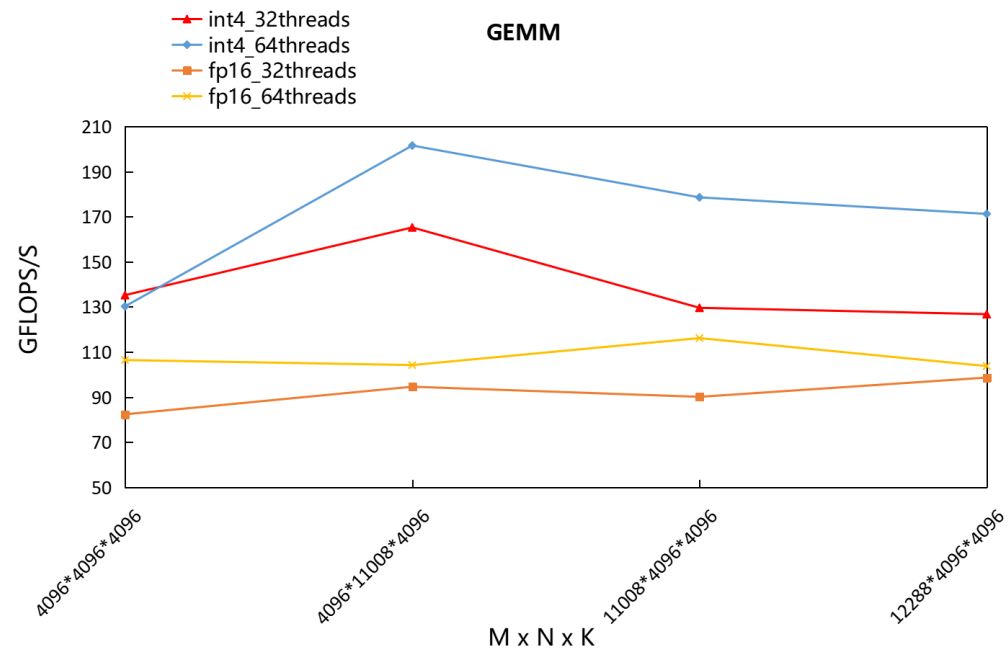# PerfXLM on RISC-V

**RISC-V CPUs Testbed**
- SG2042 CPU
  - Xuantie C920 2.0GHz
    - RVV 0.7.1, 128-bit
    - L1-D 64KB
  - 4 cores/cluster
    - Shared L2 1MB
  - 64 cores (16 cluster)
    - Shared L3 64MB

- **4x** DDR4 memory controllers
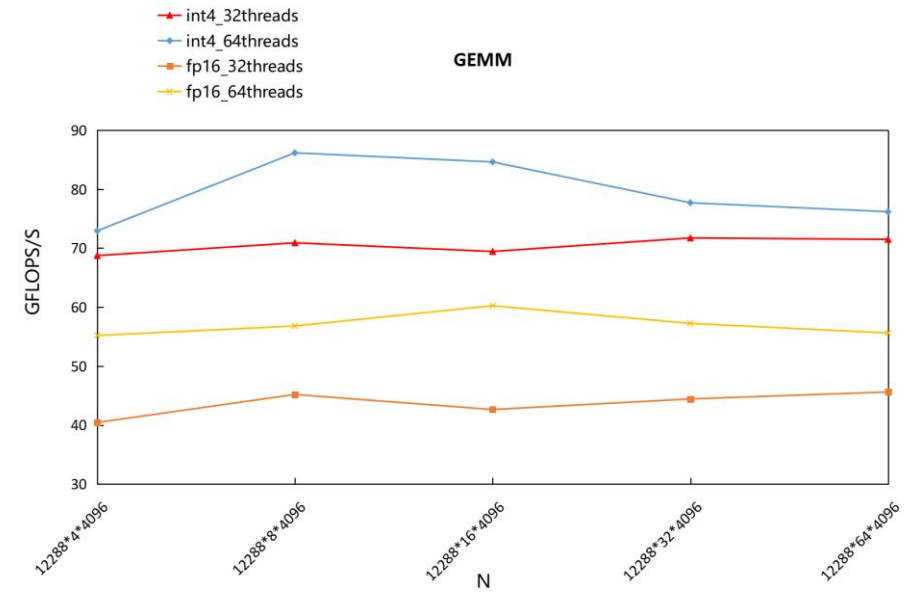
**GEMM optimization for prefill stage:**

- Single batch: tall-and-skinny matrix

- Vector instruction optimization

- GEMM blocking

- Multi-threading



GEMM Performance (Big and square matrices)
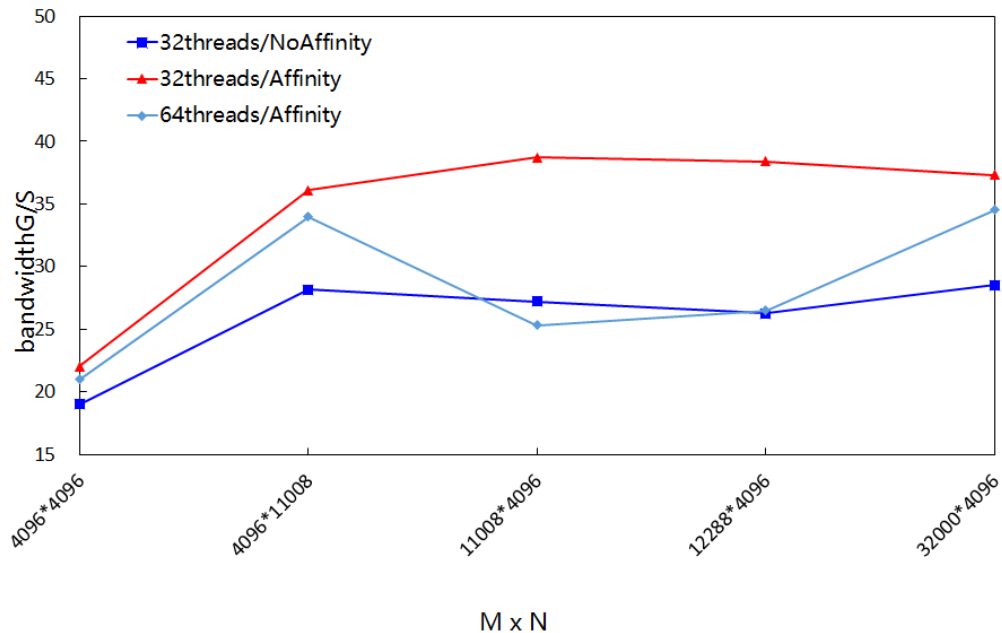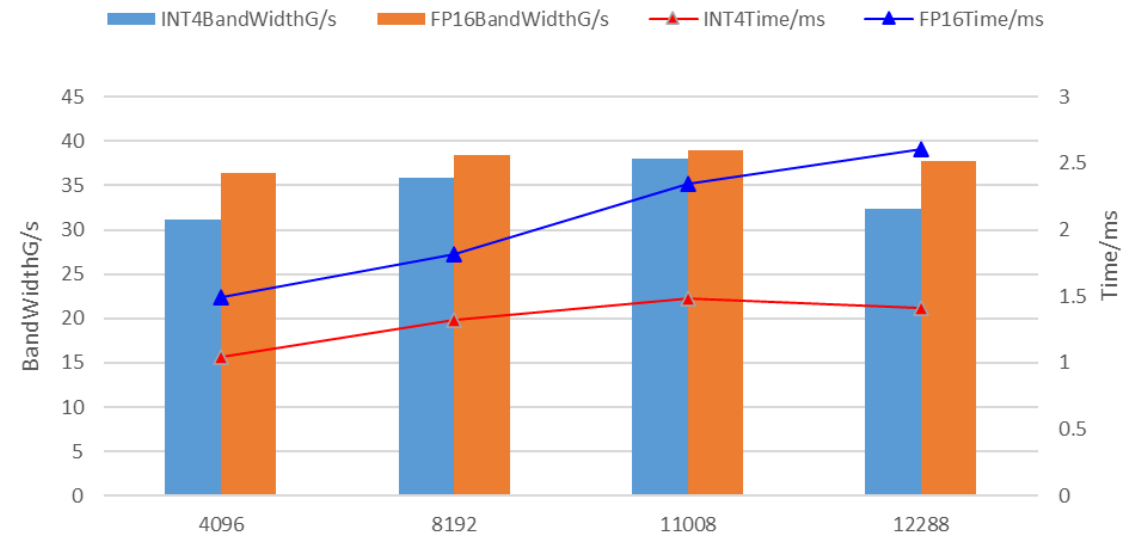


GEMM Performance (tall-and-skinny matrices)

## GEMV optimization for decoding stage:

- Vector instruction optimization

- Multi-threading and affinity



GEMV performance under different number of threads and affinity.



GEMV Performance (32 of threads)

Performance of the GEMV operator after optimization. The input of the GEMV performance test is: n =4096, m changes.
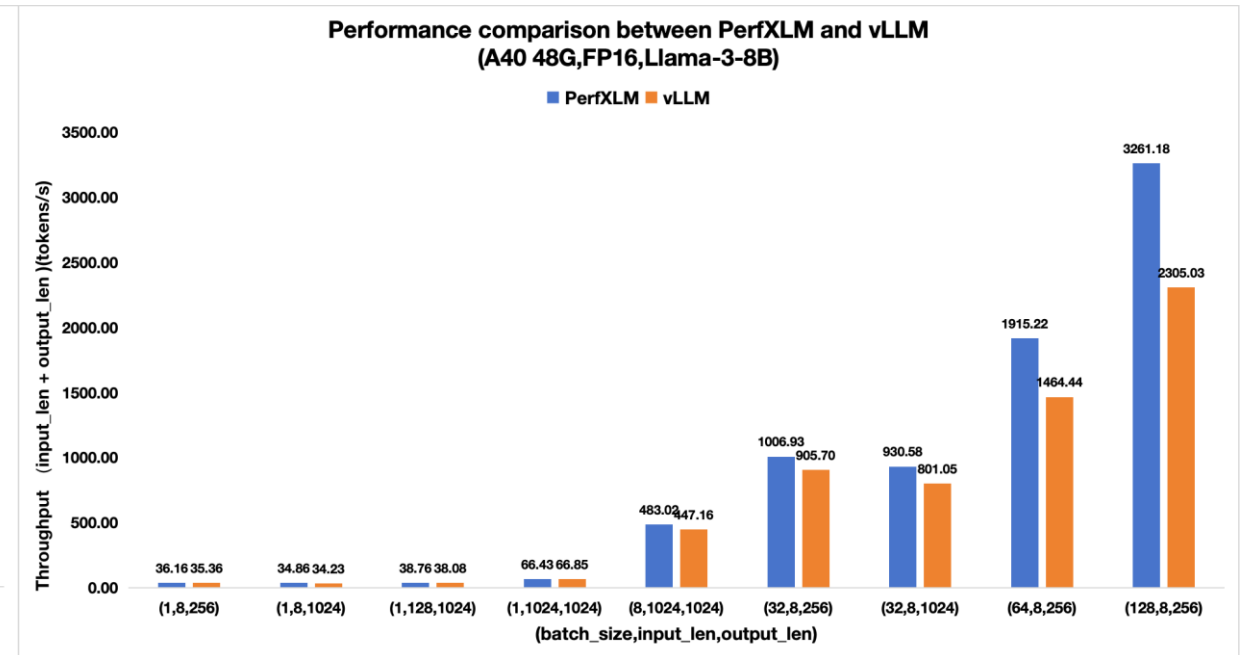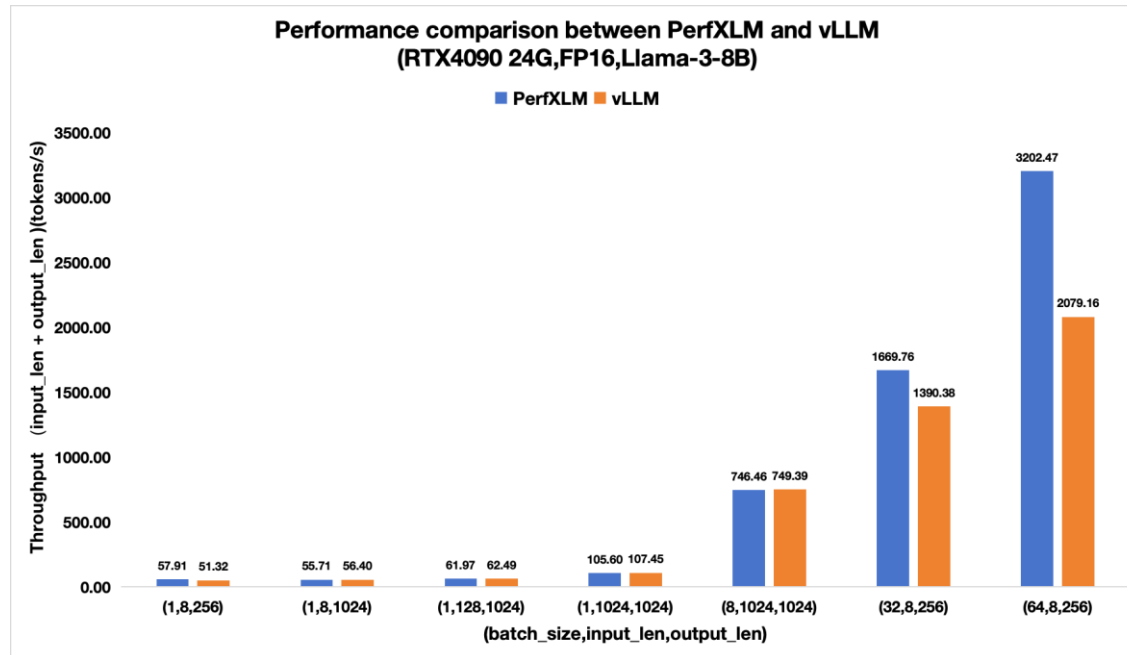
# PerfXLM on RISC-V

**PerfXLM on RISC-V:**

- Quantization Method: AWQ INT4

- Decoding Speed： **4.01** tokens/s,
  **10x – 50x** speedup to llama.cpp

- First token Latency: **6.826 s**

- Memory Footprint: 3.7GB

- Threads: 32
  - Why not 64?



```
User:What's your name?
PerfxChat: My name is PerfxChat.
User:How old are you?
PerfxChat: I'm 10 years old.
User:What is HPC?
PerfxChat: HPC stands for High Performance Comput
ing.
User:Give me more information about it.
PerfxChat: HPC is a field of computer science tha
t focuses on developing and using computers to so
lve problems that are too complex for a normal co
mputer to handle.
User:What is OpenBlas?
PerfxChat: OpenBlas is a library that provides a
set of functions for performing linear algebra op
erations on matrices.
User:Give me more information about OpenBlas?
PerfxChat: OpenBlas is a library that provides a
set of functions for performing linear algebra op
erations on matrices. It is a free software libra
ry that is licensed under the GNU Lesser General
Public License.
User:Give me some suggestion about how to optimiz
e code in GPU.
PerfxChat: I'm not sure what you mean by "optimiz
e code in GPU".
```

# PerfXLM on other architectures

**Preliminary Performance on NVIDIA GPU**

- Llama3 8B
- RTX4090 and A40

# Summary

- PerfXAPI Heterogeneous Computing Framework
  - Already support CNN model on CPU/GPU/NPU
  - Other APIs for engineering computing

- PerfXLM: a Large model inference engine
  - vllm compatible user API, easy to use
  - Supports RISC-V CPUs, embedded gpu, server gpu, and NPU
  - Better performance

- Optimized for RISC-V CPUs
  - Llama2 int4 7B, 10-50x speedup
  - Memory bandwidth is the bottleneck for LLM decoding stage.

# Thank you !

澎峰（北京）科技有限公司

PerfXLab Technologies

E-mail:

xianyi@perfxlab.com