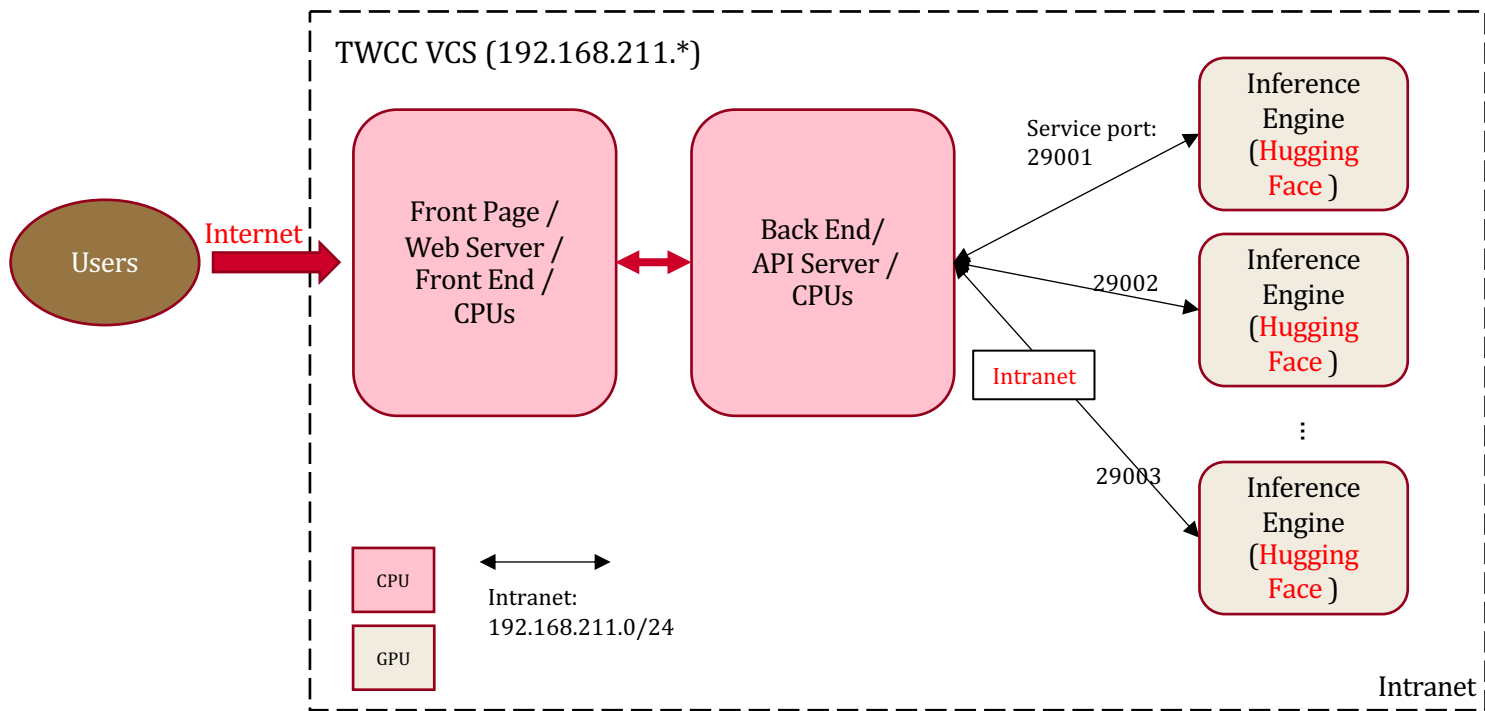# Day 3 –
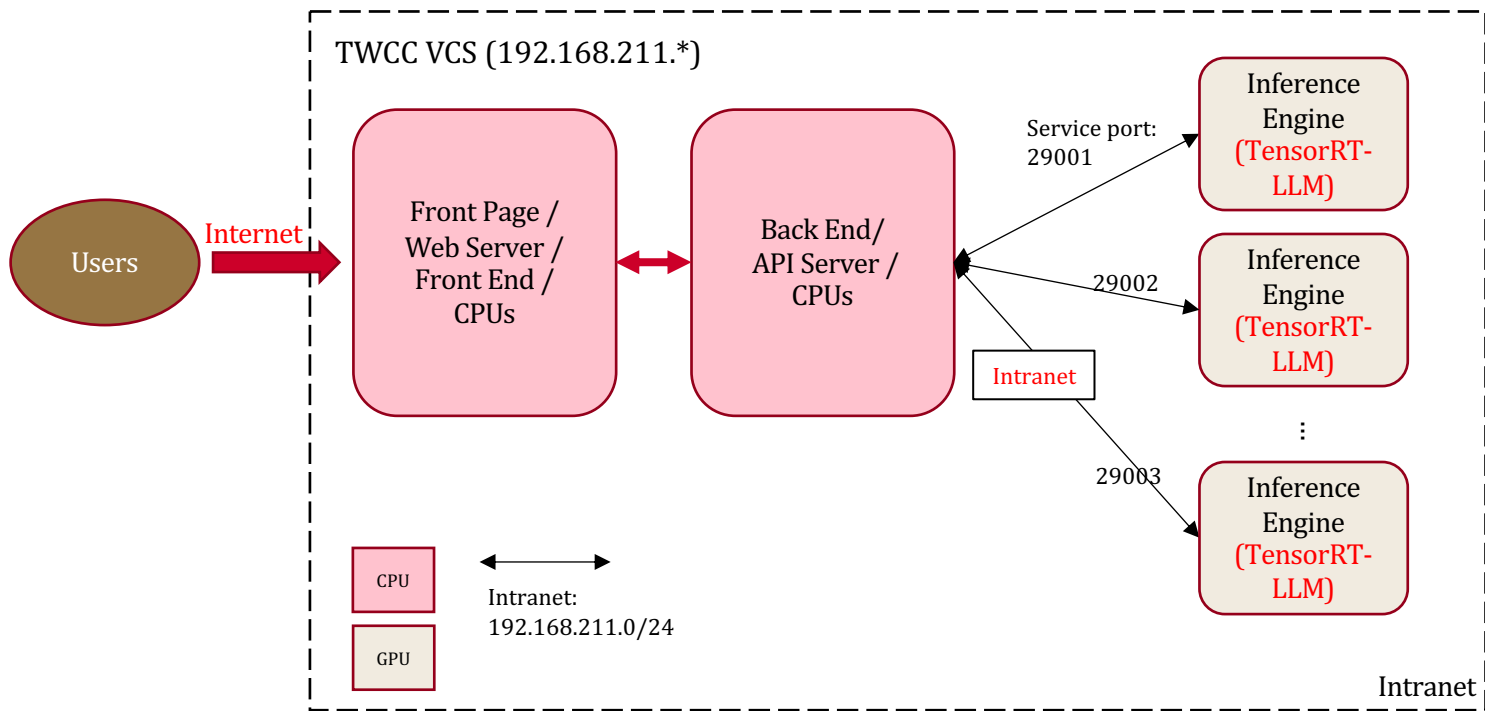# LLM inference with TensorRT-LLM on NCHC servers

Team member:
Fang-An Kuo, Kuo-Teng Ding, Meng-Chi Huang,
NCHC Speedrunning team

NVIDIA Mentor:
Anthony
Cliff

# Inference Engine Optimization

# Inference Engine Optimization
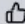
# Performance Benchmark on the Demo Website

# Performance benchmarking of the both inference engines

# Benchmarks of inferencing Traditional Chinese LLM model on NVIDIA GPUs

- ○ We replace the transformers inference APIs by using TensorRT-LLM APIs
- ○ The benchmarks run on NVIDIA GPUs, including V100/A100/A6000/H100
  - ▶ LLM Model: TAIDE Model based on LLaMA2-7B with the version number, b11.
  - ▶ Batch size = 1
  - ▶ Tensor Parallel = 1 and Pipeline Parallel = 1

| Inference APIs | Tokens/sec. | | | Computing time | | | Words/sec. | | |
|---|---|---|---|---|---|---|---|---|---|
| TAIDE 7B (b11) with FP16 | V100 | A100 | A6000 | V100 | A100 | A6000 | V100 | A100 | A6000 |
| Hugging Face | 10.63 | 5.38 | 10.88 | 26.13 | 29.53 | 29.01 | 35.45 | 40.05 | 39.35 |
| NVIDIA TRT-LLM | 19.44 | 17.21 | 17.51 | 52.30 | 103.35 | 51.10 | 63.59 | 125.65 | 62.13 |
| Speedup | 1.83 | 3.20 | 1.61 | 2.001 | 3.500 | 1.761 | 1.794 | 3.138 | 1.579 |
| Speedup, compared to V100 | 1.829 | 3.613 | 1.787 | 2.001 | 3.954 | 1.955 | 1.794 | 3.545 | 1.753 |

INT8 model

Speedup

# Benchmarks of inferencing Traditional Chinese LLM model on NVIDIA GPUs

○ The benchmarks run on NVIDIA GPUs, including V100/A100/A6000/H100
  ▶ LLM Model: TAIDE Model based on LLaMA2-7B with the version number, b11.
  ▶ Batch size = 1
  ▶ Tensor Parallel = 1 and Pipeline Parallel = 1

Computing Time (sec.)

| GPU | TensorRT-LLM | Hugging Face |
|---|---|---|
| Tesla V100-SXM2-32GB | 10.63 | 19.44 |
| NVIDIA A100 80GB PCIe | 5.38 | 17.21 |
| NVIDIA RTX A6000 | 10.88 | 17.51 |

■ TensorRT-LLM  ■ Hugging Face

# Benchmarks of inferencing Traditional Chinese LLM model on NVIDIA GPUs

- The benchmarks run on NVIDIA GPUs, including V100/A100/A6000/H100
  - ▶ LLM Model: TAIDE Model based on LLaMA2-7B with the version number, b11.
  - ▶ Batch size = 1
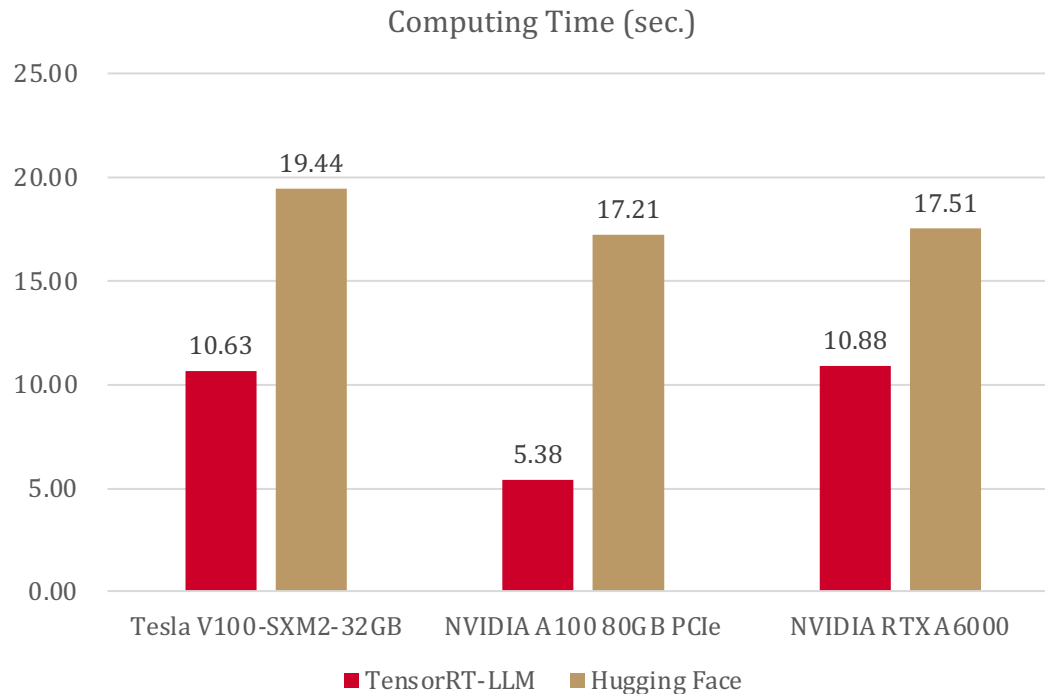  - ▶ Tensor Parallel = 1 and Pipeline Parallel = 1
  - ▶ The speedup of generating tokens by using the NVIDIA TensorRT-LLM engines is about 3.5x, which is based on NVIDIA A100



Number of tokens per sec.

Speedup: 3.95x

Speedup: 3.5x

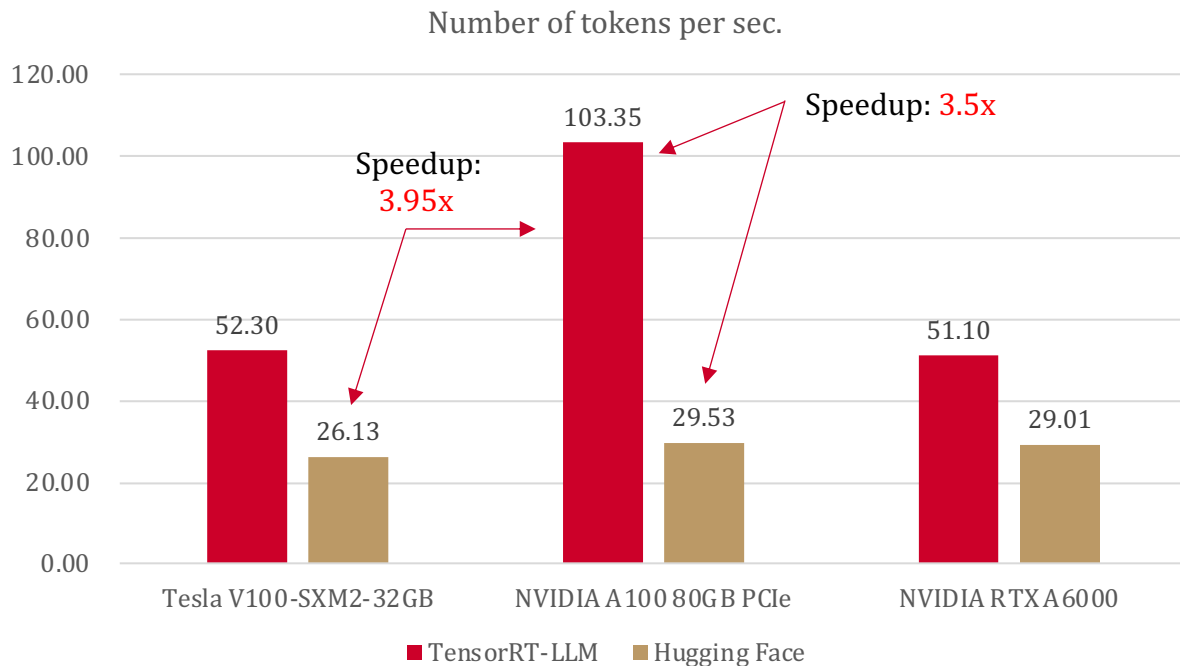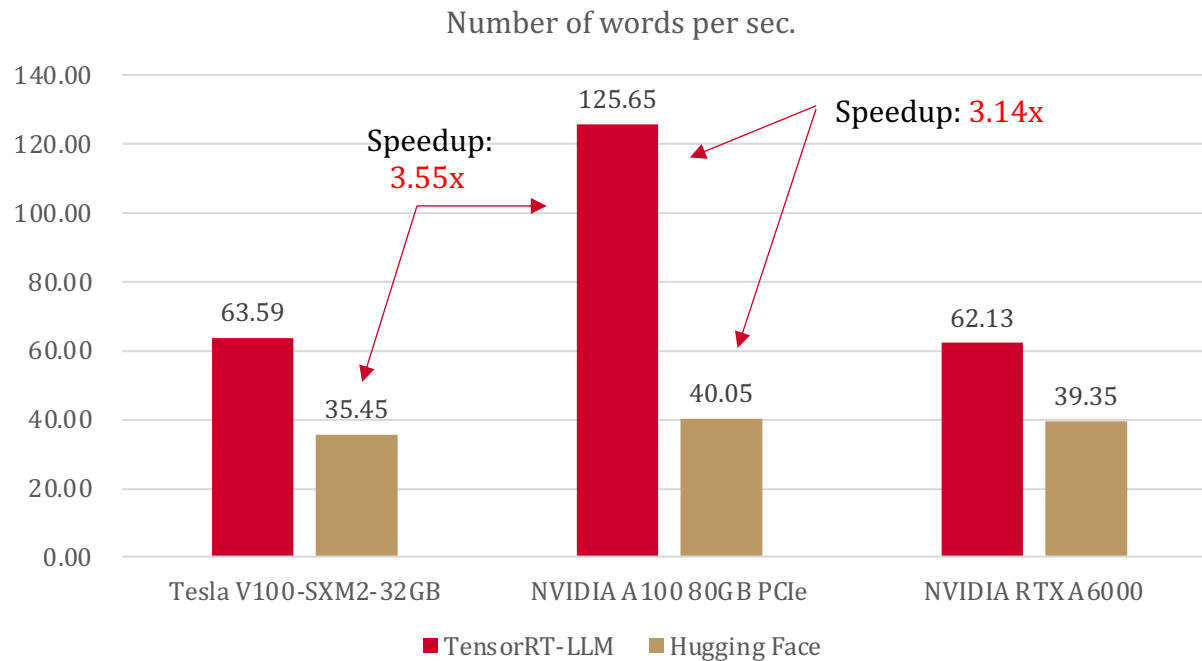| | Tesla V100-SXM2-32GB | NVIDIA A100 80GB PCIe | NVIDIA RTX A6000 |
|---|---|---|---|
| TensorRT-LLM | 52.30 | 103.35 | 51.10 |
| Hugging Face | 26.13 | 29.53 | 29.01 |

# Benchmarks of inferencing Traditional Chinese LLM model on NVIDIA GPUs

- The benchmarks run on NVIDIA GPUs, including V100/A100/A6000/H100
  - LLM Model: TAIDE Model based on LLaMA2-7B with the version number, b11.
  - Batch size = 1
  - Tensor Parallel = 1 and Pipeline Parallel = 1
  - The speedup of generating words by using the NVIDIA TensorRT-LLM engines is about 3.14x, which is based on NVIDIA A100

Number of words per sec.

# INT8 model

Discover the Gion Kobu, a traditional Japanese district known for its beautiful

Visit the Kyoto Imperial Palace, a former residence of the Emperor of Japan a

Don't forget to try some delicious Kyoto-style tofu and other local specialties
destinations. 🔍

And of course, no trip to Kyoto is complete without a visit to the Funaoka Roll
paste. 🍜

Come, let me be your helpful assistant, and I'll be happy to guide you through
Inference Engine by NVIDIA TensorRT-LLM
Number of words: 1286
Number of Tokens: 443
Times: 2.55 sec.
GPUs: NVIDIA A100 80GB PCIe

# Problems and Solutions

- Modify {TensorRT-LLM}/docker/common/install_base.sh
  - ▶ OS: ubuntu 22.04

```
init_ubuntu() {
    apt-get update
    apt-get install -y --no-install-recommends wget gdb git-lfs python3-pip python3-dev python-is-python3 libffi-dev
    apt-get install -y --no-install-recommends screen gpustat nvtop curl iftop
    if ! command -v mpirun &> /dev/null; then
      DEBIAN_FRONTEND=noninteractive apt-get install -y --no-install-recommends openmpi-bin libopenmpi-dev
    fi
    apt-get clean
    rm -rf /var/lib/apt/lists/*
    # Remove previous TRT installation
    if [[ $(apt list --installed | grep libnvinfer) ]]; then
        apt-get remove --purge -y libnvinfer*
    fi
    if [[ $(apt list --installed | grep tensorrt) ]]; then
        apt-get remove --purge -y tensorrt*
    fi
    pip uninstall -y tensorrt
    pip install flask flask_sse datasets nltk rouge_score
}
```

# Problems and Solutions

- Generation cannot properly stops at EOS (</s>) and only stops at max output length
  - ▶ Solution: output_text.replace("</s>","")
  - ▶ Another solution and reference:
    - https://github.com/NVIDIA/TensorRT-LLM/blob/a21e2f85178111fed9812bb88c2cc7411b25f0ba/examples/gpt/run.py#L299