

오픈소스 LLM, vLLM으로 Production까지

Hyogeun Oh



발표자 소개

0. 시작에 앞서...

Hygeun Oh

(오효근, ZeroHertz)

- 기계공학 학/석사 후 Machine Learning Engineer (전문연구요원)로 재직 중입니다.
- Python과 Kubernetes를 주로 다루며, MLOps에 깊은 관심이 있습니다.
- Neovim을 애용하고 생산적인 개발 환경을 추구합니다.
- 더 나은 ML 파이프라인과 자동화를 고민합니다.

GitHub [1]



[1] <https://github.com/ZeroHertz>

발표 목차

0. 시작에 앞서...

1. Introduction

2. OpenAI-Compatible Server

3. Architecture

4. Production Deployment

5. Wrap-up

PART 1

Introduction

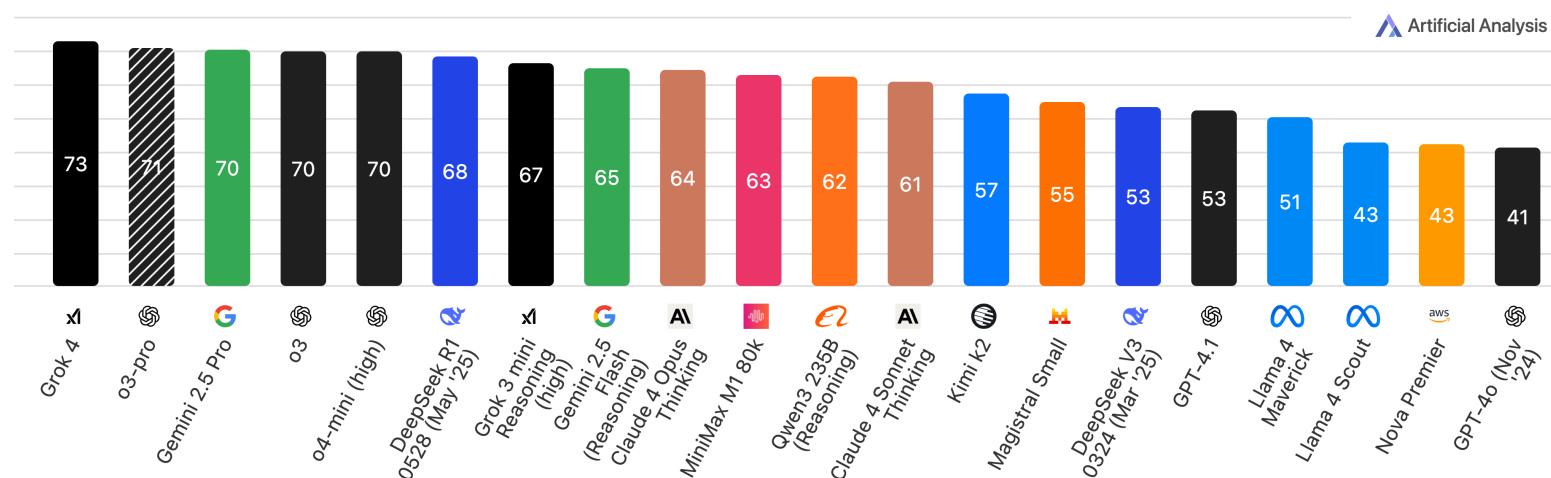
Why Self-Host LLMs When Powerful Commercial APIs Exist?

- 비용 부담: 트래픽 사용량 증가 시 높은 과금
- 프라이버시/보안 이슈: 민감 데이터 외부 전송 불가
- 커스터마이징의 어려움: 프롬프트, 파라미터, 응답 포맷 등 제한적 제어
- API 가용성: 속도, 안정성, 국가별 접근 제한 가능성

Artificial Analysis Intelligence Index [2]

Artificial Analysis Intelligence Index incorporates 7 evaluations: MMLU-Pro, GPQA Diamond, Humanity's Last Exam, LiveCodeBench, SciCode, AIME, MATH-500

 Estimate (independent evaluation forthcoming)



Gemini

Claude

Grok

Why Self-Host LLMs When Powerful Commercial APIs Exist?

❑ transformers의 AutoModelForCausalLM을 통해 추론하면 되지 않나?

```
def main():
    logger.info(f"MODEL_NAME={MODEL_NAME}")

    processor = AutoProcessor.from_pretrained(MODEL_NAME)
    model = AutoModelForCausalLM.from_pretrained(MODEL_NAME)
    logger.info("Model & processor Loaded!")

    messages = [{"role": "user", "content": "Hello, PyCon Korea 2025!"}]
    prompt = processor.apply_chat_template(
        messages, tokenize=False, add_generation_prompt=True
    )
    logger.info("prompt:")
    print(prompt)
    inputs = processor(prompt, return_tensors="pt")

    with torch.no_grad():
        generated_ids = model.generate(
            *inputs,
            max_new_tokens=1024,
            do_sample=True,
            top_p=0.95,
            temperature=0.8,
            pad_token_id=processor.eos_token_id,
        )

    output_text = processor.batch_decode(generated_ids, skip_special_tokens=True)[0]
    logger.info("output_text:")
    print(output_text)
```

```
2025-07-24 23:51:30.031 | INFO  | __main__:main:31 - MODEL_NAME='Qwen/Qwen3-0.6B'
2025-07-24 23:51:33.554 | INFO  | __main__:main:35 - Model & processor Loaded!
2025-07-24 23:51:33.583 | INFO  | __main__:main:41 - prompt:
<|im_start|>user
Hello, PyCon Korea 2025!<|im_end|>
<|im_start|>assistant
2025-07-24 23:51:44.305 | INFO  | __main__:main:56 - output_text:
user
Hello, PyCon Korea 2025!
assistant
<think>
Okay, the user sent me a message saying "Hello, PyCon Korea 2025!" So I need to respond to that. Let me think. PyCon is a global event, so I should acknowledge it. Maybe mention that it's happening in Korea. Let me check the date again. Oh, PyCon Korea 2025 is scheduled for June 10th, 2025. That's important to include.

I should make sure the response is friendly and enthusiastic. Maybe start with a greeting, then mention the event details. Also, since they asked for help, I should offer assistance. Let me put that all together in a natural way.
</think>

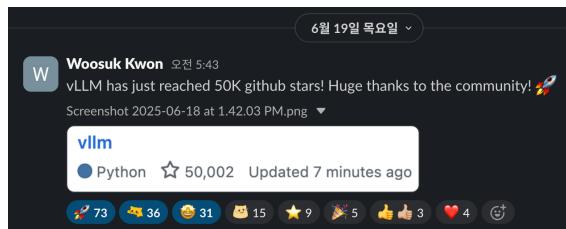
Hello, PyCon Korea 2025! 🎉
We're excited to have you join us for the event! It's scheduled for June 10th, 2025, and we're looking forward to celebrating tech innovation and showcasing the world's best trends! Let me know if you need help or have any questions! 🚀
```



Why vLLM?

❑ vLLM history

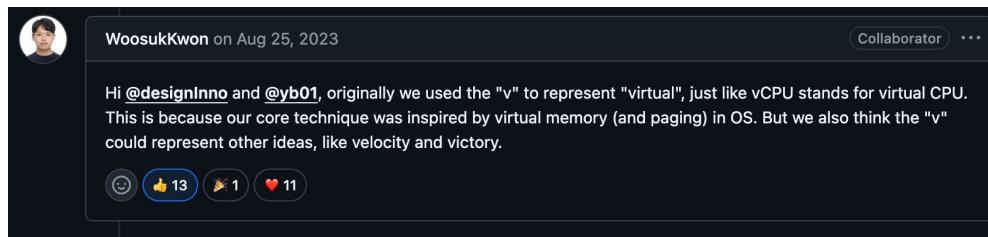
- 2023년 2월 9일, GitHub에서 CacheFlow^[3]라는 이름으로 시작
- 2023년 6월 17일, vLLM으로 이름 변경^[4]
- 2023 9월 12일, UC Berkeley Sky Computing Lab에서 “Efficient Memory Management for Large Language Model Serving with PagedAttention” 논문 발표^[5]
- 2025년 6월 19일, GitHub star 50k 달성



- 2025년 7월 24일, v0.10.0 release

❑ License: Apache-2.0^[6]

❑ v of vLLM^[7]



Efficient Memory Management for Large Language Model Serving with *PagedAttention*

Woosuk Kwon^{1,*} Zhuohan Li^{1,*} Siyuan Zhuang¹ Ying Sheng^{1,2} Lianmin Zheng¹ Cody Hao Yu³
Joseph E. Gonzalez¹ Hao Zhang⁴ Ion Stoica¹

¹UC Berkeley ²Stanford University ³Independent Researcher ⁴UC San Diego

Abstract

High throughput serving of large language models (LLMs) requires batching sufficiently many requests at a time. However, existing systems struggle because the key-value cache (KV cache) memory for each request is huge and grows and shrinks dynamically. When managed inefficiently, this memory can be significantly wasted by fragmentation and redundant duplication, limiting the batch size. To address this problem, we propose PagedAttention, an attention algorithm inspired by the classical virtual memory and paging techniques in operating systems. On top of it, we build vLLM, an LLM serving system that achieves (1) near-zero waste in KV cache memory and (2) flexible sharing of KV cache within and across requests to further reduce memory usage. Our evaluations show that vLLM improves the throughput of popular LLMs by 2-4x with the same level of latency compared to the state-of-the-art systems, such as FasterTransformer and Orca. The improvement is more pronounced with longer sequences, larger models, and more complex decoding algorithms. vLLM’s source code is publicly available at <https://github.com/vllm-project/vllm>.

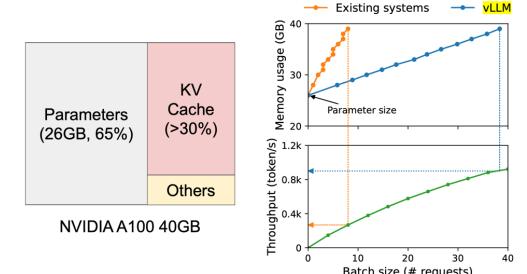


Figure 1. Left: Memory layout when serving an LLM with 13B parameters on NVIDIA A100. The parameters (gray) persist in GPU memory throughout serving. The memory for the KV cache (red) is (de)allocated per serving request. A small amount of memory (yellow) is used ephemerally for activation. Right: vLLM smooths out the rapid growth curve of KV cache memory seen in existing systems [31, 60], leading to a notable boost in serving throughput.

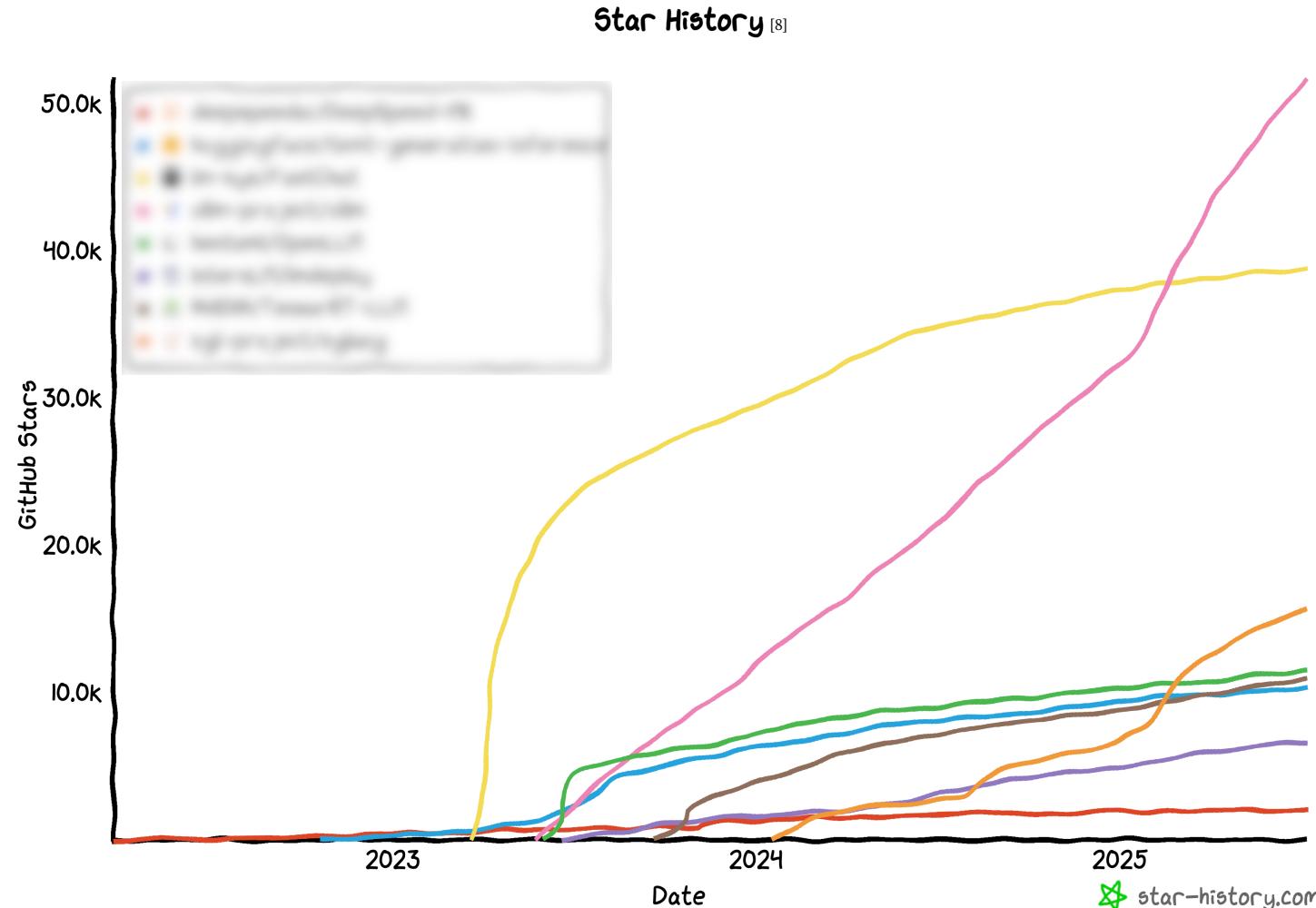
[3] <https://github.com/vllm-project/vllm/commit/e7d9d9c08c79b386f6d0477e87b77a572390317d>

[4] <https://github.com/vllm-project/vllm/commit/0b98ba15c744f1dfb0ea4f2135e85ca23d572ae1>

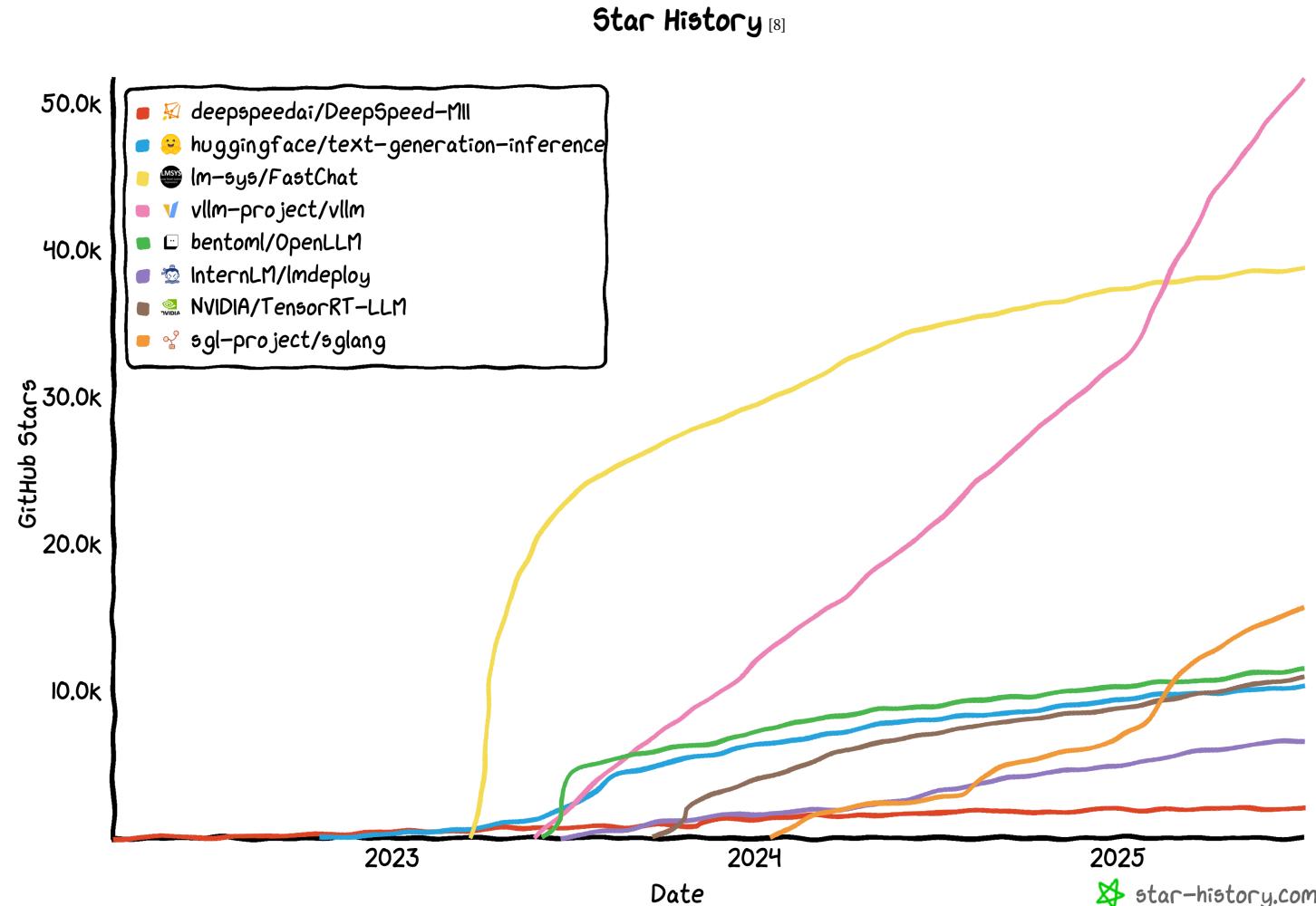
[5] <https://github.com/vllm-project/vllm/blob/main/LICENSE>

[6] <https://arxiv.org/abs/2309.06180>

Why vLLM?



Why vLLM?



Why vLLM?

▣ 초고속 LLM 서빙 성능

- State-of-the-art throughput: 대규모 요청 처리
- Continuous batching: 실시간 요청에 대한 효율적 처리

▣ 메모리 효율 및 최적화된 연산

- PagedAttention [6]: Attention key/value 메모리의 효율적 관리
- Optimized CUDA kernels: FlashAttention [9], FlashInfer [10] 등 최신 커널 통합
- Chunked prefill [11], Speculative decoding [12]

▣ 유연성 및 확장성

- HuggingFace 통합: 손쉽게 다양한 모델 서빙
- 다양한 디코딩 알고리즘: 병렬 샘플링, 빔서치 등 고성능 추론 지원
- 분산 추론 지원: Tensor, pipeline (Ray 기반), data, expert parallelism 지원

▣ 실용적 API 및 운영 편의성

- OpenAI-Compatible API: 기존 AI 서비스 (e.g., LangChain, Gemini CLI, ...)와 손쉽게 연동
- Streaming output: 스트리밍 방식 결과 제공
- Prefix caching, Multi-LoRA

[6] <https://arxiv.org/abs/2309.06180>

[9] <https://github.com/vllm-project/flash-attention>

[10] <https://github.com/flashinfer-ai/flashinfer>

[11] https://docs.vllm.ai/en/v0.9.2/configuration/optimization.html#chunked-prefill_1

[12] https://docs.vllm.ai/en/v0.9.2/features/spec_decode.html

How to serving LLM with vLLM?

1. Introduction

□ Installation

▪ Local (CPU) 사용 시

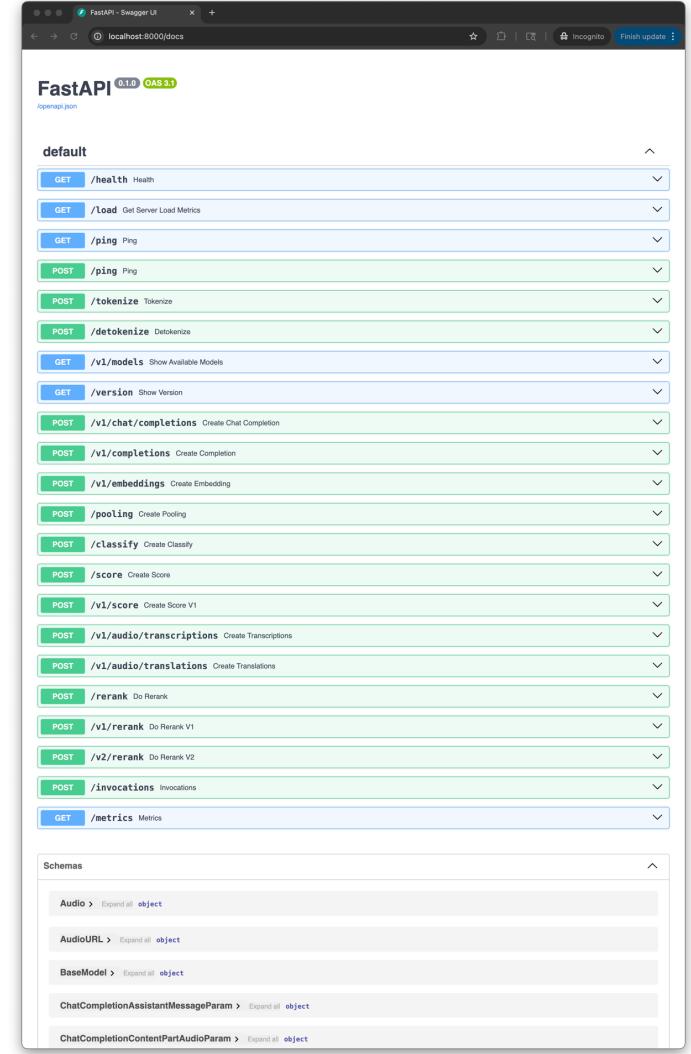
```
$ uv pip install vllm==0.9.2
Using Python 3.12.11 environment at: /opt/venv/main
Resolved 120 packages in 47ms
Installed 1 package in 10ms
+ vllm==0.9.2
```

▪ GPU 서버 사용 시 [13]

```
1 docker run --runtime nvidia --gpus all \
2   --name vllm \
3   -v ~/.cache/huggingface:/root/.cache/huggingface \
4   -p 8000:8000 \
5   --ipc=host \
6   vllm/vllm-openai:v0.9.2 \
7   --model Qwen/Qwen3-0.6B \
8   --max-model-len 8192
```

□ vllm serve

```
$ vllm serve Qwen/Qwen3-0.6B --max-model-len 8192
INFO 07-27 01:38:11 [__init__.py:244] Automatically detected platform cpu.
INFO 07-27 01:38:12 [api_server.py:1395] VLLM API server version 0.9.2
INFO 07-27 01:38:12 [cls_args.py:325] non-default args: {'model': 'Qwen/Qwen3-0.6B', 'max_model_len': 8192}
INFO 07-27 01:38:15 [config.py:841] This model supports multiple tasks: ('reward', 'classify', 'generate', 'embed'). Defaulting to 'generate'.
WARNING 07-27 01:38:15 [config.py:328] Your device (cpu) doesn't support torch.bfloat16. Falling back to torch.float16 for compatibility.
WARNING 07-27 01:38:15 [config.py:172] Using max model len 8192.
INFO 07-27 01:38:15 [cpu_utils.py:1746] cpu is experimental on VLLM_USE_V1_1. Falling back to V0 Engine.
WARNING 07-27 01:38:15 [gpu.py:131] Environment variable VLLM_CPU_KVCACHE_SPACE (GiB) for CPU backend is not set, using 4 by default.
INFO 07-27 01:38:16 [api_server.py:268] Started engine process with PID 83075
INFO 07-27 01:38:17 [__init__.py:244] Automatically detected platform cpu.
INFO 07-27 01:38:18 [llm_engine.py:230] Initializing a VLLM engine (v0.9.2) with config: model='Qwen/Qwen3-0.6B', speculative_config=None, tokenizer='Qwen/Qwen3-0.6B', skip_tokenizer_init=False, tokenizer_mode=auto, revision=None, override_neuron_configs(), tokenizer_neuron_configs(), trust_remote_code=False, dtype=torch.float16, max_seq_len=8192, download_dir=None, load_format=LoadFormat.AUTO, tensor_parallel_size=1, disable_custom_all_reduce=True, quantization=None, enforce_eager=False, kv_cache_dtype=auto, device_config=cpu, decoding_config=DecodingConfig(backends='auto', disable_fallback=False, disable_any_whitespace=False, disable_additional_properties=False, reasoning_backend=''), observability_config=ObservabilityConfig(show_hidden_metrics_for_version=None, otlp_traces_endpoint=None, collect_detailed_traces=None), seeds=0, served_model=name='Qwen/Qwen3-0.6B', num_scheduler_steps=1, multi_step_stream_outputs=True, enable_prefix_caching=None, chunked=True, enable_auto_functionalized_v2=False, use_cudagraph=True, cudagraph_num_of_warmups=0, cudagraph_capture_sizes={}, cudagraph_copy_inputs=False, full_cuda_graph=False, max_capture_size=256, local_cache_dir=None), use_cudagraph=True, use_cached_outputs=True.
INFO 07-27 01:38:23 [launcher.py:37] Route: /v1/models, Methods: GET
INFO 07-27 01:38:23 [launcher.py:37] Route: /version, Methods: GET
INFO 07-27 01:38:23 [launcher.py:37] Route: /v1/chat/completions, Methods: POST
INFO 07-27 01:38:23 [launcher.py:37] Route: /v1/completions, Methods: POST
INFO 07-27 01:38:23 [launcher.py:37] Route: /v1/embeddings, Methods: POST
INFO 07-27 01:38:23 [launcher.py:37] Route: /pooling, Methods: POST
INFO 07-27 01:38:23 [launcher.py:37] Route: /score, Methods: POST
INFO 07-27 01:38:23 [launcher.py:37] Route: /v1/audio/transcriptions, Methods: POST
INFO 07-27 01:38:23 [launcher.py:37] Route: /v1/audio/translations, Methods: POST
INFO 07-27 01:38:23 [launcher.py:37] Route: /rerank, Methods: POST
INFO 07-27 01:38:23 [launcher.py:37] Route: /v1/rerank, Methods: POST
INFO 07-27 01:38:23 [launcher.py:37] Route: /v2/rerank, Methods: POST
INFO 07-27 01:38:23 [launcher.py:37] Route: /invocations, Methods: POST
INFO 07-27 01:38:23 [launcher.py:37] Route: /metrics, Methods: GET
INFO: Started server process [8309]
INFO: Waiting for application startup...
INFO: Application startup complete.
```

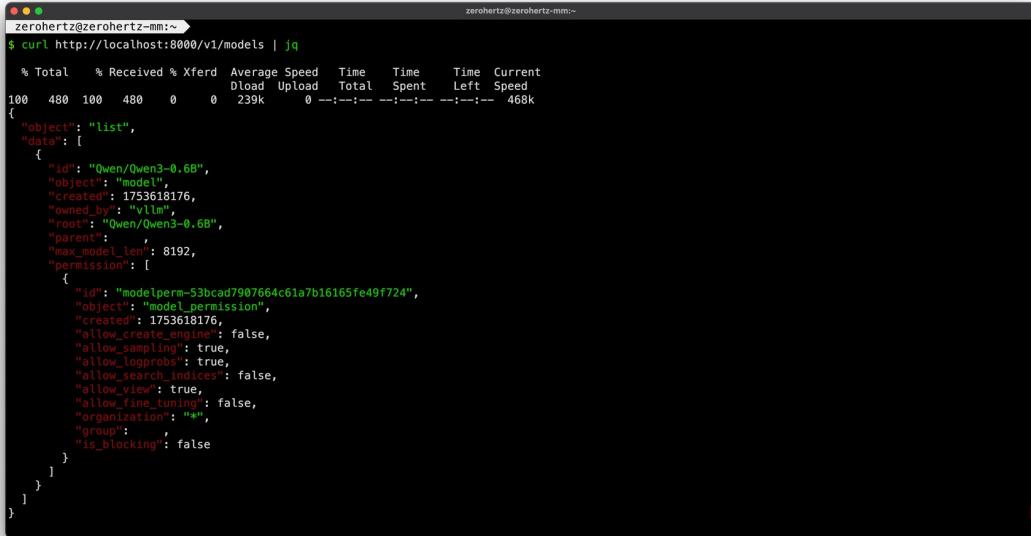


PART 2

OpenAI-Compatible Server

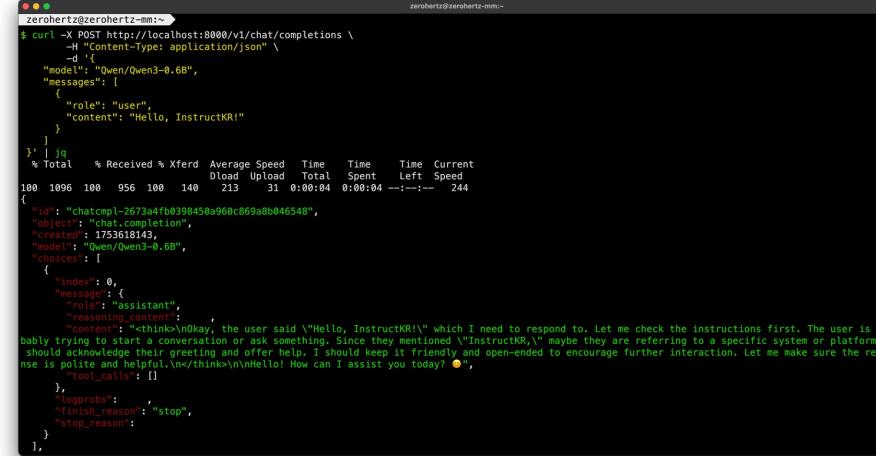
OpenAI API Spec [14]

□ /v1/models [15]

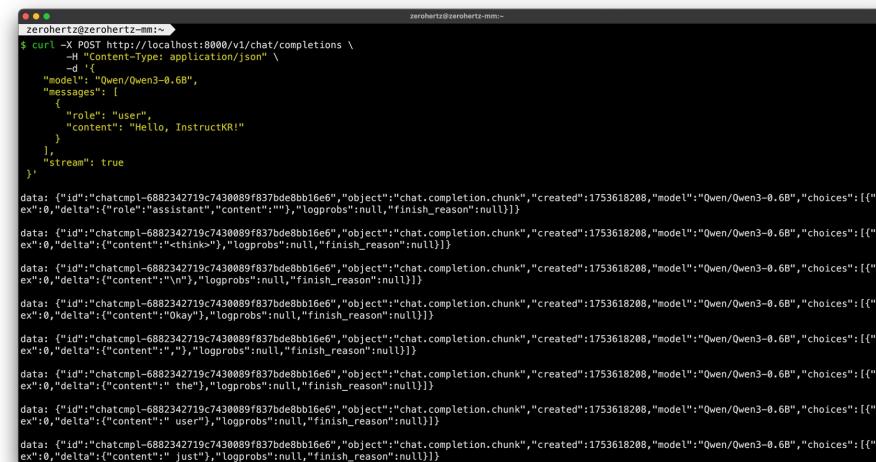


```
zeroherzt@zeroherzt-mm:~$ curl http://localhost:8000/v1/models | jq
% Total    % Received % Xferd  Average Speed   Time   Time  Current
          Dload Upload Total Spent   Left Speed
100  480  100  480    0     0 239k      0 --:--:-- --:--:-- 468k
{
  "object": "list",
  "data": [
    {
      "id": "modelperm-53bcad7907664c61a7b16165fe49f724",
      "object": "model_permission",
      "created": 1753618176,
      "allow_create_engine": false,
      "allow_sampling": true,
      "allow_logprobs": true,
      "allow_search_indices": false,
      "allow_view": true,
      "allow_fine_tuning": false,
      "organization": null,
      "group": null,
      "is_blocking": false
    }
  ]
}
```

□ /v1/chat/completions [16]



```
zeroherzt@zeroherzt-mm:~$ curl -X POST http://localhost:8000/v1/chat/completions \
-H "Content-Type: application/json" \
-d '{
  "model": "Qwen/Qwen3-0.6B",
  "messages": [
    {
      "role": "user",
      "content": "Hello, InstructKR!"
    }
  ]
}' | jq
% Total    % Received % Xferd  Average Speed   Time   Time  Current
          Dload Upload Total Spent   Left Speed
100 1096  100  956  100  140 213    31  0:00:04  0:00:04 --:--:-- 244
{
  "id": "chatmpl-26734fb398450a0960c869a8b046548",
  "object": "chat_completion",
  "created": 1753618143,
  "model": "Qwen/Qwen3-0.6B",
  "choices": [
    {
      "index": 0,
      "message": {
        "role": "assistant",
        "content": "Hello! Okay, the user said \"Hello, InstructKR!\" which I need to respond to. Let me check the instructions first. The user is probably trying to start a conversation or ask something. Since they mentioned \"InstructKR,\" maybe they are referring to a specific system or platform. I should acknowledge their greeting and offer help. I should keep it friendly and open-ended to encourage further interaction. Let me make sure the response is polite and helpful.\n</think>\n\nHello! How can I assist you today? ●",
        "logprobs": null,
        "finish_reason": "stop",
        "stop_reason": null
      }
    }
  ],
  "stream": false
}
```



```
zeroherzt@zeroherzt-mm:~$ curl -X POST http://localhost:8000/v1/chat/completions \
-H "Content-Type: application/json" \
-d '{
  "model": "Qwen/Qwen3-0.6B",
  "messages": [
    {
      "role": "user",
      "content": "Hello, InstructKR!"
    }
  ],
  "stream": true
}' | jq
data: {"id": "chatmpl-6892342719c74308089f837bde8b16e6", "object": "chat.completion.chunk", "created": 1753618208, "model": "Qwen/Qwen3-0.6B", "choices": [{"index": 0, "delta": {"role": "assistant", "content": ""}, "logprobs": null, "finish_reason": null}}]
data: {"id": "chatmpl-6882342719c74308089f837bde8b16e6", "object": "chat.completion.chunk", "created": 1753618208, "model": "Qwen/Qwen3-0.6B", "choices": [{"index": 0, "delta": {"content": "Okay"), "logprobs": null, "finish_reason": null}}]
data: {"id": "chatmpl-6882342719c74308089f837bde8b16e6", "object": "chat.completion.chunk", "created": 1753618208, "model": "Qwen/Qwen3-0.6B", "choices": [{"index": 0, "delta": {"content": "Okay"}, "logprobs": null, "finish_reason": null}}]
data: {"id": "chatmpl-6882342719c74308089f837bde8b16e6", "object": "chat.completion.chunk", "created": 1753618208, "model": "Qwen/Qwen3-0.6B", "choices": [{"index": 0, "delta": {"content": "Okay"}, "logprobs": null, "finish_reason": null}}]
data: {"id": "chatmpl-6882342719c74308089f837bde8b16e6", "object": "chat.completion.chunk", "created": 1753618208, "model": "Qwen/Qwen3-0.6B", "choices": [{"index": 0, "delta": {"content": "Okay"), "logprobs": null, "finish_reason": null}}]
data: {"id": "chatmpl-6882342719c74308089f837bde8b16e6", "object": "chat.completion.chunk", "created": 1753618208, "model": "Qwen/Qwen3-0.6B", "choices": [{"index": 0, "delta": {"content": "Okay"}, "logprobs": null, "finish_reason": null}}]
data: {"id": "chatmpl-6882342719c74308089f837bde8b16e6", "object": "chat.completion.chunk", "created": 1753618208, "model": "Qwen/Qwen3-0.6B", "choices": [{"index": 0, "delta": {"content": "Okay"}, "logprobs": null, "finish_reason": null}}]
data: {"id": "chatmpl-6882342719c74308089f837bde8b16e6", "object": "chat.completion.chunk", "created": 1753618208, "model": "Qwen/Qwen3-0.6B", "choices": [{"index": 0, "delta": {"content": "Okay"), "logprobs": null, "finish_reason": null}}]
data: {"id": "chatmpl-6882342719c74308089f837bde8b16e6", "object": "chat.completion.chunk", "created": 1753618208, "model": "Qwen/Qwen3-0.6B", "choices": [{"index": 0, "delta": {"content": "Okay"}, "logprobs": null, "finish_reason": null}}]
```

[14] https://docs.vllm.ai/en/v0.9.2/serving/openai_compatible_server.html

[15] <https://platform.openai.com/docs/api-reference/models/list>

[16] <https://platform.openai.com/docs/api-reference/chat/create>

Tool calling

□ OpenAI 규격에 맞춰 호출

v1.97.1 openai-python / src / openai / types / shared_params / function_definition.py

stainless-app[bot] chore(api): event shapes more accurate

Code Blame 45 lines (32 loc) · 1.47 KB

```
1 # File generated from our OpenAPI spec by Stainless. See CONTRIBUTING.md for details.
2
3 from __future__ import annotations
4
5 from typing import Optional
6 from typing_extensions import Required, TypedDict
7
8 from .function_parameters import FunctionParameters
9
10 __all__ = ["FunctionDefinition"]
11
12 ...
13 class FunctionDefinition(TypedDict, total=False):
14     name: Required[str]
15     """The name of the function to be called.
16
17     Must be a-z, A-Z, 0-9, or contain underscores and dashes, with a maximum length
18     of 64.
19     """
20
21     description: str
22     """
23     A description of what the function does, used by the model to choose when and
24     how to call the function.
25     """
26
27     parameters: FunctionParameters
28     """The parameters the functions accepts, described as a JSON Schema object.
29
30     See the [guide](https://platform.openai.com/docs/guides/function-calling) for
31     examples, and the
32     [JSON Schema reference](https://json-schema.org/understanding-json-schema/) for
33     documentation about the format.
34
35     Omitting `parameters` defines a function with an empty parameter list.
36     """
37
38     strict: Optional[bool]
39     """Whether to enable strict schema adherence when generating the function call.
40
41     If set to true, the model will follow the exact schema defined in the
42     `parameters` field. Only a subset of JSON Schema is supported when `strict` is
43     `true`. Learn more about Structured Outputs in the
44     [function calling guide](https://platform.openai.com/docs/guides/function-calling).
45     """

```

```
zerohertz@zerohertz-mm:~$ curl -X POST http://localhost:8000/v1/chat/completions \
-H "Content-Type: application/json" \
-d '{
  "model": "Qwen/Qwen3-0.6B",
  "messages": [
    {
      "role": "user",
      "content": "What is the weather like in Seoul?"
    }
  ],
  "tools": [
    {
      "type": "function",
      "function": {
        "name": "get_weather",
        "description": "Get the current weather in a given location",
        "parameters": {
          "type": "object",
          "properties": {
            "location": {
              "type": "string",
              "description": "The city and state, e.g. San Francisco, CA"
            },
            "unit": {
              "type": "string",
              "enum": ["celsius", "fahrenheit"]
            }
          },
          "required": ["location"]
        }
      }
    }
  ],
  "tool_choice": "auto"
}' | jq
% Total    % Received % Xferd  Average Speed   Time     Time     Time  Current
Dload  Upload Total Spent   Left Speed
100  981  100  168  100  813  38888  183k --:--:-- --:--:-- 239k
{
  "object": "error",
  "message": "'auto' tool choice requires --enable-auto-tool-choice and --tool-call-parser to be set",
  "type": "BadRequestError",
  "param": ,
  "code": 400
}
```

Tool calling

- “--tool-call-parser”를 통해 모델에 따른 적절한 parser 설정 필요 [18, 19]

```
vllm serve Qwen/Qwen3-0.6B --max-model-len 8192 --enable-auto-tool-choice --tool-call-parser hermes
$ vllm serve Qwen/Qwen3-0.6B --max-model-len 8192 --enable-auto-tool-choice --tool-call-parser hermes
INFO 07-27 21:13:11 [__init__.py:244] Automatically detected platform cpu.
INFO 07-27 21:13:13 [api_server.py:1595] vLLM API server version 0.9.2
INFO 07-27 21:13:13 [cli_args.py:325] non-default args: {'enable_auto_tool_choice': True, 'tool_call_parser': 'hermes', 'model': 'Qwen/Qwen3-0.6B', 'max_model_len': 8192}
INFO 07-27 21:13:16 [config.py:841] This model supports multiple tasks: {'classify', 'generate', 'embed', 'reward'}. Defaulting to 'generate'.
WARNING 07-27 21:13:16 [config.py:3320] Your device 'cpu' doesn't support torch.bfloat16. Falling back to torch.float16 for compatibility.
WARNING 07-27 21:13:16 [config.py:3371] Casting torch.bfloat16 to torch.float16.
INFO 07-27 21:13:16 [config.py:1472] Using max model len 8192
INFO 07-27 21:13:16 [arg_utils.py:1746] cpu is experimental on VLM_USE_V1=1. Falling back to V0 Engine.
WARNING 07-27 21:13:16 [cpu.py:1311] Environment variable VLM_CPU_KVCACHE_SPACE (GiB) for CPU backend is not set, using 4 by default.
INFO 07-27 21:13:16 [api_server.py:268] Started engine process with PID 93561
INFO 07-27 21:13:18 [__init__.py:244] Automatically detected platform cpu.
INFO 07-27 21:13:19 [llm_engine.py:230] Initializing a V0 LLM engine (v0.9.2) with config: model='Qwen/Qwen3-0.6B', speculative_config=None, tokenizer='Qwen/Qwen3-0.6B', skip_tokenizer_init=False, tokenizer_init=False, revision=None, override_neuron_config={}, tokeniser_revision=None, trust_remote_code=False, dtype=torch.float16, max_seq_len=8192, download_dir=None, load_format=loadFormat.AUTO, tensor_parallel_size=1, pipeline_parallel_size=1, disable_custom_all_reduce=True, quantization=None, enforce_eager=False, kv_cache_dtype=cpu, device_config=cpu, decoding_config=DecodingConfig(backends='auto', disable_fallback=False, disable_any_whitespace=False, disable_additional_properties=False, reasoning_backend=''), observability_config=ObservabilityConfig(show_hidden_metrics_for_version=None, otlp_traces_endpoint=None, collect_detailed_traces=None), seed=0, served_model_name=Qwen/Qwen3-0.6B, num_scheduler_steps=1, multi_step_stream_outputs=True, enable_prefix_caching=None, chunked_prefill_enabled=False, use_asynchronous_output_proc=False, profiler_config=None, compilation_config={'level':0}, debug_dump_path='', cache_dir='', custom_ops=[], splitting_ops=[], use_inductor=True, compile_sizes=[], 'inductor_compile_config': {'enable_auto_functionalized_v2': false}, 'inductor_passes': {}, 'use_cudagraph': true, 'cudagraph_num_of_warmups': 0, 'cudagraph_capture_sizes': [], 'cudagraph_copy_inputs': false, 'full_cuda_graph': false, 'max_capture_size': 256, 'local_cache_dir': null, use_cached_outputs=True, WARNING 07-27 21:13:20 [cpu_worker.py:447] Auto thread-binding is not supported due to the lack of package numa and psutil, fallback to no thread-binding. To get better performance, please try to manually bind threads.
INFO 07-27 21:13:20 [cpu.py:69] Using Torch SDPA backend.
INFO 07-27 21:13:20 [importing.py:63] Triton not installed or not compatible; certain GPU-related functions will not be available.
INFO 07-27 21:13:20 [parallel_state.py:1076] rank 0 in world size 1 is assigned as DP rank 0, PP rank 0, TP rank 0, EP rank 0
INFO 07-27 21:13:20 [weight_utils.py:292] Using model weights format ['*.safetensors']
INFO 07-27 21:13:21 [weight_utils.py:345] No model.safetensors.index.json found in remote.
Loading safetensors checkpoint shards: 0% Completed | 0/1 [00:00:00, ?it/s]
Loading safetensors checkpoint shards: 100% Completed | 1/1 [00:01:00:00, 1.13s/it]
INFO 07-27 21:13:22 [default_loader.py:272] Loading weights took 1.13 seconds
INFO 07-27 21:13:22 [executor_base.py:113] # cpu blocks: 2340, # CPU blocks: 0
INFO 07-27 21:13:22 [executor_base.py:118] Maximum concurrency for 8192 tokens per request: 4.57x
INFO 07-27 21:13:23 [llm_engine.py:428] init engine (profile, create kv cache, warmup model) took 0.58 seconds
INFO 07-27 21:13:23 [serving_chat.py:85] "auto" tool choice has been enabled please note that while the parallel_tool_calls client option is preset for
```

```
zeroherzt@zeroherzt-mm:~ curl -X POST http://localhost:8000/v1/chat/completions \
-H "Content-Type: application/json" \
-d '{"model": "Qwen/Qwen3-0.6B", "messages": [{"role": "user", "content": "What is the weather like in Seoul?"}], "tools": [{"type": "function", "function": {"name": "get_weather", "description": "Get the current weather in a given location", "parameters": {"type": "object", "properties": {"location": {"type": "string", "description": "The city and state, e.g. San Francisco, CA"}}, "required": ["location"]}}, "tool_choice": "auto"}]} | jq
% Total    % Received   % Xferd  Average Speed   Time     Time     Time  Current
100  1957   100 1144    0 813   286   203  0:00:04  0:00:03  0:00:01  489
  10%: "chatmsg-0e6d5805653947df9431ea8f1ba8b200",
  100%: "chatmsg-0e6d5805653947df9431ea8f1ba8b200",
  created: 1753618598,
  model: "Qwen/Qwen3-0.6B",
  choices: [
    {
      "index": 0,
      "message": {
        "role": "assistant",
        "content": "The weather in Seoul is currently clear and sunny. The temperature is approximately 25 degrees Celsius. The unit is Celsius, as no unit was specified in the request."}
    }
  ],
  tool_calls: [
    {
      "name": "wchatmpl-tool-049f5ea200b44d04b9a0ae387fdfd3c7",
      "type": "function",
      "function": {"name": "get_weather", "arguments": "{\"location\": \"Seoul\"}"}
    }
  ]
},
"logprob": null,
"finish_reason": "tool_calls",
"stop_reason": null
},
"usage": {
  "prompt_tokens": 194,
  "total_tokens": 330
}
```

Reasoning

- ❑ “chat_template_kwargs”的“enable_thinking”을 통해 reasoning 유무 설정 가능

```
zerohertz@zerohertz-mm:~$ curl -X POST http://localhost:8000/v1/chat/completions \
-H "Content-Type: application/json" \
-d '{
  "model": "Qwen/Qwen3-0.6B",
  "messages": [
    {
      "role": "user",
      "content": "Hello, InstructKR!"
    }
  ],
  "chat_template_kwargs": {"enable_thinking": false}
}' | jq
% Total % Received % Xferd Average Speed Time Time Current
          Dload Upload Total Spent Left Speed
100  681  100  485  100  196  197   79  0:00:02  0:00:02 --:--:-- 277
{
  "id": "chatmpl-cf17f36b772a47ae895449a509660922",
  "object": "chat.completion",
  "created": 1753618596,
  "model": "Qwen/Qwen3-0.6B",
  "choices": [
    {
      "index": 0,
      "message": {
        "role": "assistant",
        "reasoning_content": "",
        "content": "Hello, InstructKR! How can I assist you today?",
        "tool_calls": []
      },
      "logprobs": 0,
      "finish_reason": "stop",
      "stop_reason": ""
    }
  ],
  "usage": {
    "prompt_tokens": 18,
    "total_tokens": 32,
    "completion_tokens": 14,
    "prompt_tokens_details": {}
  },
  "prompt_logprobs": 0,
  "kv_transfer_params": 0
}
```

```
zerohertz@zerohertz-mm:~$ curl -X POST http://localhost:8000/v1/chat/completions \
-H "Content-Type: application/json" \
-d '{
  "model": "Qwen/Qwen3-0.6B",
  "messages": [
    {
      "role": "user",
      "content": "Hello, InstructKR!"
    }
  ],
  "chat_template_kwargs": {"enable_thinking": true}
}' | jq
% Total % Received % Xferd Average Speed Time Time Current
          Dload Upload Total Spent Left Speed
100  1319  100  1124  100  195  344   59  0:00:03  0:00:03 --:--:-- 404
{
  "id": "chatmpl-02296ae29bdd416a994235c20a248134",
  "object": "chat.completion",
  "created": 1753618624,
  "model": "Qwen/Qwen3-0.6B",
  "choices": [
    {
      "index": 0,
      "message": {
        "role": "assistant",
        "reasoning_content": {
          "content": "<think>\nOkay, the user said \"Hello, InstructKR!\" which is a bit confusing. First, I need to acknowledge their message. Since they mentioned \"InstructKR,\" I should check if that's a specific tool or service they're referring to. Maybe they're using a platform or app with the name InstructKR.\n\nI should respond politely and ask for clarification. Let me make sure to respond in a friendly and helpful way. Also, I should keep the tone consistent and avoid any technical jargon. Let me check if there's any additional context needed to fully understand their request.\n</think>\n\nHello! How can I assist you today? If you have any questions or need help, feel free to ask! ☺",
          "tool_calls": []
        },
        "logprobs": 0,
        "finish_reason": "stop",
        "stop_reason": ""
      }
    },
    "usage": {
      "prompt_tokens": 14,
      "total_tokens": 163,
      "completion_tokens": 149,
      "prompt_tokens_details": {}
    },
    "prompt_logprobs": 0,
    "kv_transfer_params": 0
  }
}
```

Reasoning

- “--reasoning-parser”를 통해 모델에 따른 적절한 parser 설정 필요 [20, 21]
- 기존에는 “--enable-reasoning” 옵션이 존재했으나 deprecated [22]

```
vllm serve Qwen/Qwen3-0.6B --max-model-len 8192 --reasoning-parser qwen3
$ vllm serve Qwen/Qwen3-0.6B --max-model-len 8192 --reasoning-parser qwen3
INFO 07-27 21:17:54 [__init__.py:44] Automatically detected platform: cpu.
INFO 07-27 21:17:55 [api_server.py:1395] VLLM API server version 0.9.2
INFO 07-27 21:17:55 [cli_args.py:325] non-default args: {'model': 'Qwen/Qwen3-0.6B', 'max_model_len': 8192, 'reasoning_parser': 'qwen3'}
INFO 07-27 21:17:58 [config.py:841] This model supports multiple tasks: {'reward', 'classify', 'generate', 'embed'}. Defaulting to 'generate'.
WARNING 07-27 21:17:58 [config.py:3320] Your device 'cpu' doesn't support torch.bfloat16. Falling back to torch.float16 for compatibility.
WARNING 07-27 21:17:58 [config.py:3371] Casting torch.bfloat16 to torch.float16.
INFO 07-27 21:17:58 [config.py:1472] Using max model len 8192
INFO 07-27 21:17:58 [arg_utils.py:1746] cpu is experimental for V0 Engine.
WARNING 07-27 21:17:58 [cpu.py:131] Environment variable VLM_CPU_KVCACHE_SPACE (GiB) for CPU backend is not set, using 4 by default.
INFO 07-27 21:17:58 [api_server.py:268] Started engine process with PID 94834
INFO 07-27 21:18:00 [__init__.py:244] Automatically detected platform: cpu.
INFO 07-27 21:18:01 [lmm_engine.py:230] Initializing a new LLM engine (v0.9.2) with config: model='Qwen/Qwen3-0.6B', speculative_config=None, tokenizer='Qwen/Qwen3-0.6B', skip_tokenizer_init=False, tokenizer_mode='auto', revision=None, override_neuron_config={}, tokenizer_revision=None, trust_remote_code=False, dtype=torch.float16, max_seq_len=8192, download_dir=None, load_format='LoadFormat.AUTO', tensor_parallel_size=1, pipeline_parallel_size=1, disable_custom_all_reduce=True, quantization=None, enforce_eager=False, kv_cache_dtype='auto', device_config='cpu', decoding_config='DecodingConfig(backend="auto")', disable_fallback=False, disable_additional_properties=False, reasoning_backend='qwen3'), observability_config='ObservabilityConfig(show_hidden_metrics_for_version=None, otlp_traces_endpoint=None, collect_detailed_traces=None)', seed=0, served_model_name='Qwen/Qwen3-0.6B', num_steps=1, multi_step_stream_outputs=True, enable_prefix_caching=None, chunked_prefill_enabled=False, use_asynchronous_output_proc=False, pooler_config=None, compilation_config='{"level":0, "debug_dump_path": "", "backend": "", "custom_ops": [], "splitting_ops": [], "use_inductor": true, "compile_size": 1, "inductor_compile_config": {"enable_auto_functionalized_v2": false}, "inductor_passes": {}, "use_cudagraph": true, "cudagraph_num_of_warmups": 0, "cudagraph_capture_sizes": []}', cudagraph_copy_inputs=False, full_cuda_graph=False, max_capture_size=256, local_cache_dir=None, use_cached_outputs=True, WARNING 07-27 21:18:02 [cpu_worker.py:447] Auto thread-binding is not supported due to the lack of package numa and putils, fallback to no thread-binding.
To get better performance, please try to manually bind threads.
INFO 07-27 21:18:02 [cpu.py:69] Using Torch SDPA backend.
INFO 07-27 21:18:02 [importing.py:63] Triton not installed or not compatible; certain GPU-related functions will not be available.
INFO 07-27 21:18:02 [parallel_state.py:1076] rank 0 in world size 1 is assigned as DP rank 0, PP rank 0, TP rank 0, EP rank 0
INFO 07-27 21:18:03 [weight_utils.py:92] Using model weights format ['*.safetensors']
INFO 07-27 21:18:04 [weight_utils.py:345] No model.safetensors.index.json found in remote.
```

```
$ curl -X POST http://localhost:8000/v1/chat/completions \
-H "Content-Type: application/json" \
-d '{
  "model": "Qwen/Qwen3-0.6B",
  "messages": [
    {
      "role": "user",
      "content": "Hello, InstructKR!"
    }
  ],
  "chat_template_kwargs": {"enable_thinking": true}
}' | jq
```

Total	% Received	% Xferd	Average Speed	Time	Time	Time	Current
100	1134	100	939	100	195	194	40
					0:00:04	0:00:04	--:--:--
					234		

```
{
  "id": "chatcmp-1d580c706b53464886254e4d5d2f4297",
  "object": "chat.completion",
  "created": 1753619008,
  "model": "Qwen/Qwen3-0.6B",
  "choices": [
    {
      "index": 0,
      "message": {
        "role": "assistant",
        "content": "\nOkay, the user said \"Hello, InstructKR!\" which I need to respond to. Let me check the instructions first. The user is probably trying to start a conversation or ask something. Since they mentioned \"InstructKR,\" maybe they are referring to a specific system or platform. I should acknowledge their greeting and offer help. I should keep it friendly and open-ended to encourage further interaction. Let me make sure the response is polite and helpful.\n",
        "tool_calls": []
      },
      "logprobs": 0,
      "finish_reason": "stop",
      "stop_reason": ""
    }
  ],
  "usage": {
    "prompt_tokens": 14,
    "total_tokens": 121,
    "completion_tokens": 107,
    "prompt_tokens_details": {}
  },
  "prompt_logprobs": 0,
  "kv_transfer_params": {}
}
```

[20] https://docs.vllm.ai/en/v0.9.2/features/reasoning_outputs.html[21] <https://qwen.readthedocs.io/en/latest/deployment/vllm.html#parsing-thinking-content>[22] https://github.com/vllm-project/vllm/blob/v0.9.2/vllm/engine/arg_utils.py#L626-L634

Chat Template

▣ 기본적으로 tokenizer_config.json의 “chat_template” 값 사용

```

1  {%- if tools %}
2    {{- '<|im_start|>system\n' }}
3    {%- if messages[0].role == 'system' %}
4      {{- messages[0].content + '\n\n' }}
5    {%- endif %}
6    {%- for tool in tools %}
7      {{- tool | toJSON }}
8    {%- endfor %}}
9    {%- for message in messages[1:]%}
10      {{- message.role == 'system' ? ('<|im_start|>' + message.content + '<|im_end|>\n') : ('<|im_start|>system\n' + message.content + '<|im_end|>\n') }}
11    {%- endfor %}}
12    {%- if messages[-1].role == 'user' %}
13      {{- messages[-1].content + '\n' }}
14    {%- endif %}
15    {%- if messages[-1].role == 'assistant' %}
16      {{- messages[-1].content + '\n' }}
17    {%- endif %}
18    {%- if messages[-1].role == 'system' %}
19      {{- messages[-1].content + '\n' }}
20    {%- endif %}
21    {%- if ns.multi_step_tool and message.role == 'user' and message.content.startswith('<tool_response>') and message.content.endswith('</tool_response>') %}
22      {{- ns.multi_step_tool = message.content[message.content.index('<tool_response>'):message.content.index('</tool_response>')] + '\n' }}
23      {{- ns.last_query_index = index %}}
24    {%- endif %}}
25    {%- for message in messages %}
26      {%- if message.role == 'user' %}
27        {{- content = message.content %}}
28      {%- else %}
29        {{- content = '' %}}
30      {%- endif %}
31      {%- if message.role == "user" or (message.role == "system" and not loop.first) %}
32        {{- '<|im_start|>' + message.role + '\n' + content + '<|im_end|>' + '\n' }}
33      {%- elif message.role == "assistant" %}
34        {{- set reasoning_content = '' %}}
35        {{- if message.reasoning_content is string %}}
36          {{- content = message.reasoning_content %}}
37        {%- else %}
38          {{- if '<think>' in content %}}
39            {{- set reasoning_content = content.split('<think>')[0].rstrip('\n').split('<think>')[1:-1].lstrip('\n') %}}
40            {{- content = content.split('<think>')[1:-1].lstrip('\n') %}}
41          {%- endif %}}
42        {%- endif %}
43        {{- if loop.index0 > ns.last_query_index %}}
44          {{- if loop.last or (not loop.last and reasoning_content) %}}
45            {{- '<|im_start|>' + message.role + '\n<think>\n' + reasoning_content.strip('\n') + '\n</think>\n\n' + content.lstrip('\n') %}}
46          {%- else %}
47            {{- '<|im_start|>' + message.role + '\n' + content %}}
48          {%- endif %}
49        {%- else %}
50          {{- '<|im_start|>' + message.role + '\n' + content %}}
51        {%- endif %}}
52      {%- endif %}
53      {%- if message.tool_calls %}
54        {{- for tool_call in message.tool_calls %}}
55          {%- if (loop.first and content) or (not loop.first) %}
56            {{- '\n' }}
57          {%- endif %}
58          {{- tool_call.function %}}
59          {{- set tool_call = tool_call.function %}}
60        {%- endif %}}
61        {{- '<tool_call>\n"name": "'}}
62        {{- tool_call.function %}}
63        {{- ',"arguments": "'}}
64        {{- if tool_call.arguments is string %}}
65          {{- tool_call.arguments %}}
66        {%- else %}
67          {{- tool_call.arguments | toJSON %}}
68        {%- endif %}}
69        {{- '</tool_call>' }}
70      {%- endfor %}}
71    {%- endif %}
72    {%- if message.role == "tool" %}
73    {%- if loop.first or (messages[loop.index0 - 1].role != "tool") %}
74      {{- '<|im_start|>user\n' }}
75    {%- endif %}
76    {%- if '<|im_start|>assistant\n' %}
77      {{- content %}}
78    {%- if '<|im_start|>assistant\n' %}
79      {{- content + '\n' + messages[loop.index0 + 1].content + '\n' + '<|im_end|>\n' %}}
80    {%- endif %}
81    {%- endif %}
82    {%- if add_generation_prompt %}
83      {{- add_generation_prompt %}}
84    {%- if add_generation_prompt == '<|im_start|>assistant\n' %}
85      {{- if enable_thinking is defined and enable_thinking is false %}}
86        {{- '<think>\n\n</think>\n\n' %}}
87      {%- endif %}}
88    {%- endif %}}
89  {%- endif %}}

```

[23] https://huggingface.co/Owen/Owen3-0.6B/blob/main/tokenizer_config.json#L230

Chat Template

2. OpenAI-Compatible Server

□ “--chat-template”으로 새로운 chat template 적용 가능 [24]

```
messages = [
    {"role": "system", "content": "You are a helpful assistant."},
    {"role": "user", "content": "Hello, InstructKR!"},
    {"role": "assistant", "content": "Hello! How can I help you today?"},
    {"role": "user", "content": "Can you explain what InstructKR is?"},
    {
        "role": "assistant",
        "content": "InstructKR is a Korean instruction-following dataset and research initiative focused on improving language models' ability to follow instructions in Korean.",
    },
    {"role": "user", "content": "What's the weather like in Seoul today?"},
    {
        "role": "assistant",
        "content": "I'll help you check the weather in Seoul.",
        "tool_calls": [
            {
                "id": "call_1",
                "type": "function",
                "function": {
                    "name": "get_weather",
                    "arguments": {"location": "Seoul, South Korea"}
                }
            }
        ],
        "content": "The weather in Seoul today is sunny with a temperature of 22°C (72°F). There's a light breeze and clear skies."
    },
    {
        "role": "assistant",
        "content": "Based on the weather information, Seoul is having a pleasant day today! It's sunny with a comfortable temperature of 22°C (72°F), light breeze, and clear skies."
    },
    {"role": "user", "content": "Thanks! Can you help me with something else?"}
]
prompt = processor.apply_chat_template(
    messages, tokenize=False, add_generation_prompt=True
)
```



```
<|im_start|>system
You are a helpful assistant.<|im_end|>
<|im_start|>user
Hello, InstructKR!<|im_end|>
<|im_start|>assistant
Hello! How can I help you today?<|im_end|>
<|im_start|>user
Can you explain what InstructKR is?<|im_end|>
<|im_start|>assistant
InstructKR is a Korean instruction-following dataset and research initiative focused on improving language models' ability to follow instructions in Korean.<|im_end|>
<|im_start|>user
What's the weather like in Seoul today?<|im_end|>
<|im_start|>assistant
I'll help you check the weather in Seoul.
<tool_call>
{'name': 'get_weather', 'arguments': {'location': 'Seoul, South Korea'}}
</tool_call><|im_end|>
<|im_start|>user
<tool_response>
The weather in Seoul today is sunny with a temperature of 22°C (72°F). There's a light breeze and clear skies.
</tool_response><|im_end|>
<|im_start|>assistant
Based on the weather information, Seoul is having a pleasant day today! It's sunny with a comfortable temperature of 22°C (72°F), light breeze, and clear skies. Perfect!
<|im_start|>user
Thanks! Can you help me with something else?<|im_end|>
<|im_start|>assistant
```

Name	Last commit message	Last commit date
...		
offline_inference	[VLM] Add video support for Intern-S1 (#2167)	40 minutes ago
online_serving	[CI/Build][Doc] Clean up more docs that point to old ...	8 hours ago
others	[CI/Build][Doc] Move existing benchmark scripts in C...	yesterday
pyproject.toml	Convert examples to ruff-format (#18400)	2 months ago
template_alpaca.jinja	Support chat template and echo for chat API (#1756)	2 years ago
template_baichuan.jinja	Fix Baichuan chat template (#3340)	last year
template_chatglm.jinja	Add chat templates for ChatGLM (#3418)	last year
template_chatglm2.jinja	Add chat templates for ChatGLM (#3418)	last year
template_chatml.jinja	Support chat template and echo for chat API (#1756)	2 years ago
template_dse_qwen2_vijinna	[Model] Adding Support for Qwen2VL as an Embedd...	8 months ago
template_falcon.jinja	Add chat templates for Falcon (#3420)	last year
template_falcon_180b.jinja	Add chat templates for Falcon (#3420)	last year
template_inkbot.jinja	Support chat template and echo for chat API (#1756)	2 years ago
template_telfm.jinja	[Model] Support TelF-LM Model (#15023)	4 months ago
template_vlm2vec.jinja	[Frontend] Use a proper chat template for VLM2Vec ...	9 months ago
tool_chat_template_deepsseek1.jinja	Fix DeepSeek-R1-0528 chat template (#20717)	2 weeks ago
tool_chat_template_deepsseek3.jinja	[Feature] Support DeepSeek3 Function Call (#17784)	2 months ago
tool_chat_template_granite.jinja	Change granite chat template to keep json list format...	8 months ago
tool_chat_template_granite_20b_fc.jinja	[Model] tool calling support for ibm-granite/granite-2...	9 months ago
tool_chat_template_hermes.jinja	[Bugfix] Fix Hermes tool call chat template bug (#82...	10 months ago
tool_chat_template_hunyuan_s13b.jinja	[Model] Add ToolParser and MoE Config for Hunyuan...	last week
tool_chat_template_internlm2_tool.jinja	[Frontend][Feature] support tool calling for internlm/...	9 months ago
tool_chat_template_llama3.1_json.jinja	[Bugfix][Frontend] Update Llama Chat Templates to ...	8 months ago
tool_chat_template_llama3.2_json.jinja	[Misc] Update llama 3.2 template to support system ...	7 months ago
tool_chat_template_llama3.2_pythonic.ji...	[Frontend] Fix typo in tool chat templates for llama3...	3 months ago
tool_chat_template_llama4_4.json.jinja	Add chat template for Llama 4 models (#16428)	3 months ago
tool_chat_template_llama4_pythonic.jinja	[Frontend][Bug Fix] Update llama4 pythonic jinja tem...	2 months ago
tool_chat_template_minimax_m1.jinja	[Feature] Support MiniMax-M1 function calls features...	3 weeks ago
tool_chat_template_mistral.jinja	[Feature] OpenAI-Compatible Tools API + Streaming ...	10 months ago
tool_chat_template_mistral3.jinja	[Bugfix] Fix tool call template validation for Mistral m...	2 months ago
tool_chat_template_mistral_parallel.jinja	[Bugfix] example template should not add parallel_to...	10 months ago
tool_chat_template_phi4_minij.jinja	[Frontend] Add Phi-4 mini function calling support ...	4 months ago
tool_chat_template_toolace.jinja	[Frontend] Fix typo in tool chat templates for llama3...	3 months ago
tool_chat_template_xlam_llama.jinja	Add XLAM tool parser support (#17148)	last month
tool_chat_template_xlam_qwen.jinja	Add XLAM tool parser support (#17148)	last month

[24] https://docs.vllm.ai/en/v0.9.2/serving/openai_compatible_server.html#chat-template_1

[25] <https://github.com/vllm-project/vllm/tree/main/examples>

PART 3

Architecture

KV Cache, PagedAttention

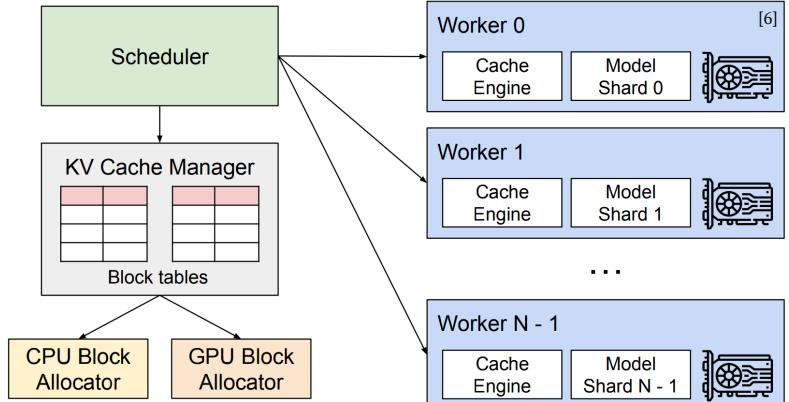
❑ KV Cache [26]

- Autoregressive 생성 방식 특성 상 각 token을 예측할 때 과거 전체 sequence를 입력으로 재처리
→ time complexity: $O(n^2)$
- 과거 Key/Value 값 (KV cache)을 사용하여 매 단계의 반복 연산 생략
→ time complexity: $O(n)$ 수준

❑ PagedAttention [6, 27, 28]

- 운영체제 (OS)의 virtual memory에서 영감을 받아 제안
- 기존의 KV cache는 memory 연속성 요구로 memory fragmentation 문제 존재
→ 특히 많은 요청을 병렬로 처리하는 상황에서 memory 할당/해제의 비효율 심각
- Block 단위의 memory 할당 및 page table을 이용하여 논리적 연속성 (logical continuity) 유지로 물리적 memory 분산 사용
- 각 요청을 고정된 크기의 page에 mapping하여 실제 memory는 non-continuous하게 구성하여 효율적 접근

$$\begin{array}{c}
 Q \\
 \left[\begin{array}{c} \text{Query Token 1} \\ \text{Query Token 2} \\ \text{Query Token 3} \\ \text{Query Token 4} \end{array} \right] \times \begin{array}{c} K^T \\
 \left[\begin{array}{c} \text{Key Token 1} \\ \text{Key Token 2} \\ \text{Key Token 3} \\ \text{Key Token 4} \end{array} \right] = \begin{array}{c} QK^T \\
 \left[\begin{array}{cccc} Q_1K_1 & Q_1K_2 & Q_1K_3 & Q_1K_4 \\ Q_2K_1 & Q_2K_2 & Q_2K_3 & Q_2K_4 \\ Q_3K_1 & Q_3K_2 & Q_3K_3 & Q_3K_4 \\ Q_4K_1 & Q_4K_2 & Q_4K_3 & Q_4K_4 \end{array} \right] \end{array} \times \begin{array}{c} V \\
 \left[\begin{array}{c} \text{Value Token 1} \\ \text{Value Token 2} \\ \text{Value Token 3} \\ \text{Value Token 4} \end{array} \right]
 \end{array}
 \end{array}$$



```

[26] %timeit -n 1
# Generate the text
generation_output = model.generate(
    input_ids=input_ids,
    max_new_tokens=100,
    use_cache=True
)
6.66 s ± 2.22 s per loop (mean ± std. dev. of 7 runs, 1 loop each)

[27] %timeit -n 1
# Generate the text
generation_output = model.generate(
    input_ids=input_ids,
    max_new_tokens=100,
    use_cache=False
)
21.9 s ± 94.6 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)
  
```

[26] <https://huggingface.co/blog/not-lain/kv-caching>

[27] <https://arxiv.org/abs/2309.06180>

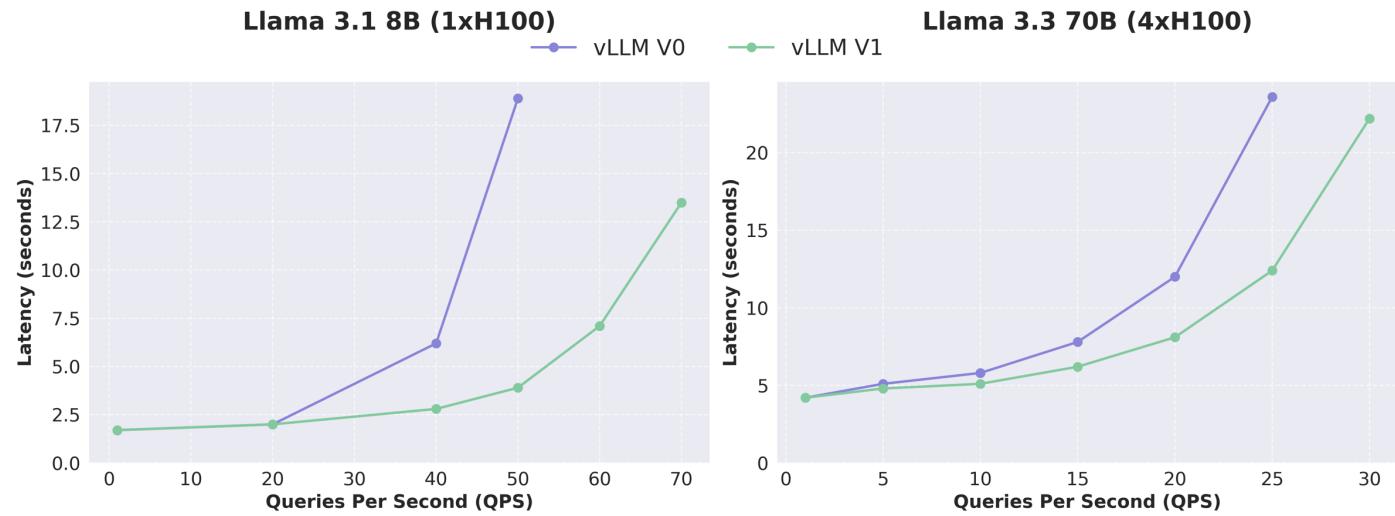
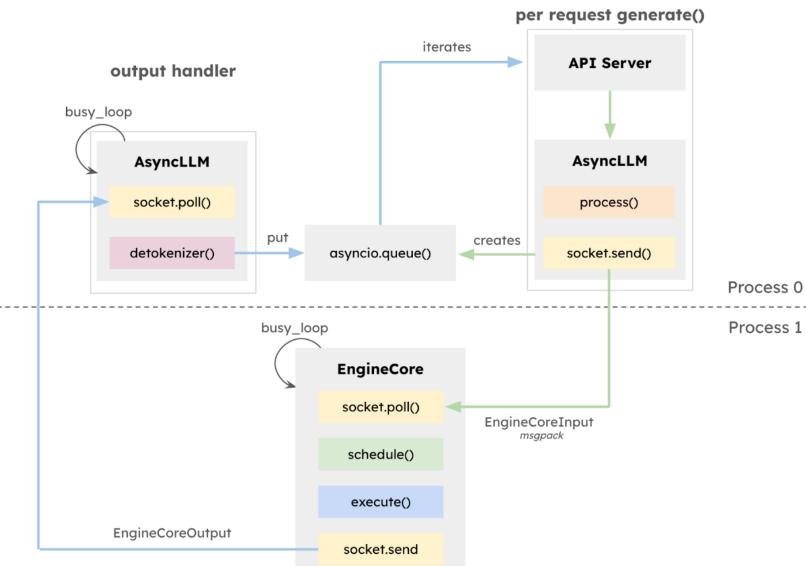
[28] https://docs.vllm.ai/en/v0.9.2/design/kernel/paged_attention.html

[29] <https://blog.vllm.ai/2023/06/20/vlm.html>

V0 Engine vs. V1 Engine

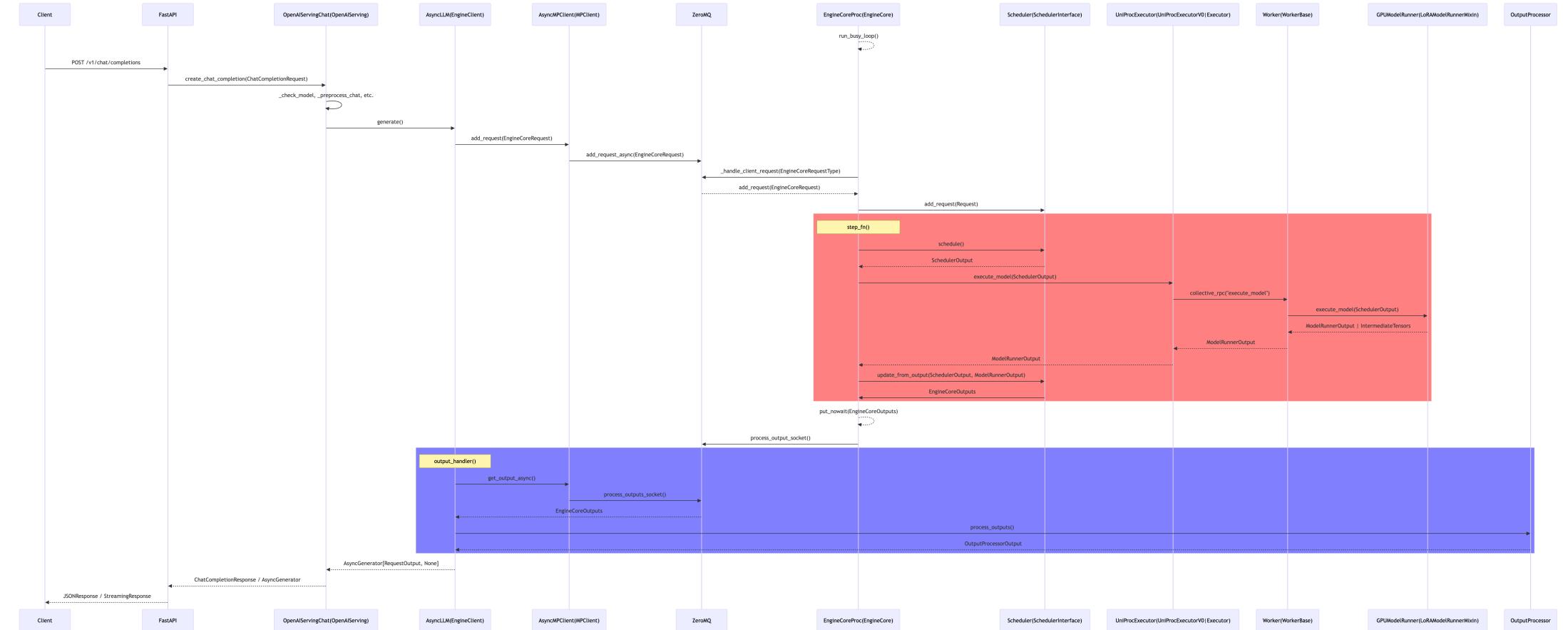
3. Architecture

- Optimized Execution Loop & API Server: EngineCore와 AsyncLLM 분리로 API server, token화 등 CPU 작업과 GPU model 실행을 완전 비동기 병렬화
- Simple & Flexible Scheduler: “prefill”, “decode” 구분 제거, “[request_id: num_tokens]” 기반 동적 token 할당으로 chunked-prefill, 사전 caching, 추측 decoding 지원
- Zero-Overhead Prefix Caching: Hash+LRU cache 구조 최적화로 cache hit rate 0%여도 1% 미만 성능 저하
- Clean Architecture for Tensor-Parallel Inference: Worker 상태 caching 후 diff만 전송, scheduler, worker 분리로 IPC overhead 대폭 감소
- Efficient Input Preparation: Persistent batch 기법으로 입력 tensor 재생성 없이 diff만 적용, Numpy 활용으로 CPU overhead 최소화
- torch.compile & Piecewise CUDA Graphs: Model 최적화 자동화 및 유연한 CUDA graph 통합으로 kernel customizing 최소화
- Enhanced Support for Multimodal LLMs: 비차단 image 전처리, image hash 기반 KV cache, encoder cache를 활용한 chunked-prefill 구현
- FlashAttention 3: 동적 batch 환경에서 최적화된 고성능 attention kernel 제공
- 환경 변수 `VLLM_USE_V1=1`를 통해 V1 Engine 활성화

[30] https://docs.vllm.ai/en/v0.9.2/design/arch_overview.html[31] https://docs.vllm.ai/en/v0.9.2/usage/v1_guide.html[32] <https://github.com/vllm-project/vllm/issues/18571>[33] <https://blog.vllm.ai/2025/01/27/v1-alpha-release.html>

Chat Completions

□ vLLM의 /v1/chat/completions 처리 과정 [34]



PART 4

Production Deployment

LoRA Adapters

4. Production Deployment

❑ LoRA (Low-Rank Adaptation) [35]

- 기존 model의 weight 행렬에 대해 전체를 학습하지 않고 low-rank matrix (A, B)만 학습
- $\Delta W = BA$ 를 통해 model이 비용 효율적으로 새로운 data에 적응

❑ Static serving LoRA adapters

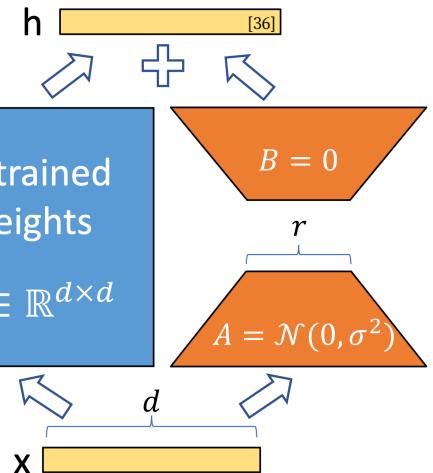
```
vllm serve Qwen/Qwen3-0.6B --max-model-len 8192 \
--reasoning-parser qwen3 \
--enable-lora \
--lora-modules phh/Qwen3-0.6B-TLDR-Lora=phh/Qwen3-0.6B-TLDR-Lora
```

❑ Dynamic serving LoRA adapters

```
VLLM_ALLOW_RUNTIME_LORA_UPDATING=True vllm serve Qwen/Qwen3-0.6B --max-model-len 8192 \
--reasoning-parser qwen3 \
--enable-lora
```

```
Fetching 16 files: 100%
WARNING 07-30 22:43:35 [cpu.py:250] Pin memory is not supported on CPU.
INFO 07-30 22:43:35 [serving_models.py:186] Loaded new LoRA adapter: name 'phh/Qwen3-0.6B-TLDR-Lora', path 'phh/Qwen3-0.6B-TLDR-Lora'
INFO: 127.0.0.1:54544 - "POST /v1/load_lora_adapter HTTP/1.1" 200 OK
$ curl -X POST http://localhost:8000/v1/load_lora_adapter \
-H "Content-Type: application/json" \
-d '{
    "lora_name": "phh/Qwen3-0.6B-TLDR-Lora",
    "lora_path": "phh/Qwen3-0.6B-TLDR-Lora"
}'
Success: LoRA adapter 'phh/Qwen3-0.6B-TLDR-Lora' added successfully.
```

```
INFO 07-30 22:44:33 [serving_models.py:203] Removed LoRA adapter: name 'phh/Qwen3-0.6B-TLDR-Lora'
INFO: 127.0.0.1:54550 - "POST /v1/unload_lora_adapter HTTP/1.1" 200 OK
$ curl -X POST http://localhost:8000/v1/unload_lora_adapter \
-H "Content-Type: application/json" \
-d '{
    "lora_name": "phh/Qwen3-0.6B-TLDR-Lora"
}'
Success: LoRA adapter 'phh/Qwen3-0.6B-TLDR-Lora' removed successfully.
```



LoRA Adapters

Parallelism Strategies^[37]

□ Tensor Parallelism (TP)

- 각 model layer 내 model parameter를 여러 GPU에 분산하여 처리
 - Model이 너무 커서 여러 GPU를 single node로 추론할 때
 - 더 높은 처리량을 위해 더 많은 KV cache 공간을 마련하기 위해 GPU 당 memory pressure를 줄여야 할 때

□ Pipeline Parallelism (PP)

- Model layer를 여러 GPU에 분산하여 model의 여러 부분 순차적 처리
 - Model이 너무 커서 여러 node에 걸쳐 분산해야 할 때
 - Layer 분산이 tensor sharding보다 더 효율적인 매우 깊고 좁은 model일 때

□ Expert Parallelism (EP)

- Mixture of Experts (MoE) model을 위한 특수한 형태의 병렬 처리
 - “--enable-expert-parallel” 사용 시 MoE layer에서 tensor parallelism 대신 expert parallelism 사용
 - MoE model을 사용할 때
 - GPU 간 expert 연산 부하를 분산할 때

□ Data Parallelism (DP)

- 여러 GPU에 걸쳐 전체 model을 복제하고 여러 요청 batch를 병렬 처리
 - 전체 model을 복제하기에 충분한 GPU를 보유한 때
 - Model 크기보다 처리량을 확장해야 할 때
 - 요청 batch 간 격리가 유리한 다중 사용자 환경일 때

Multi-node Distributed Inference

■ Ray를 통한 cluster 구성 과정 [38]

- VLLM_HOST_IP [39]: vLLM 내부 통신에 사용할 node의 IP 주소
- GLOO_SOCKET_IFNAME [40]: PyTorch Gloo backend가 사용할 network interface 이름
- NCCL_IB_DISABLE [41]: NCCL의 InfiniBand (IB) network 사용 여부

```
# NOTE: master
worker run \
    -d \
    --entrypoint /bin/bash \
    --image=vllm/vllm \
    --shm-size 10.24g \
    -v /dev/shm:/dev/shm \
    -v "/root/.cache/huggingface" \
    -e VLLM_HOST_IP=$NODE_IP \
    -e GLOO_SOCKET_IFNAME=$GLOO_SOCKET_IFNAME \
    -e NCCL_IB_DISABLE=1 \
    vllm/vllm-openais:v0.9.2 \
    -c "pip install ray[default] --system && ray start --head --port=6379 --disable-usage-stats --dashboard-host=0.0.0.0 && ray start --block --address=$MASTER_NODE_IP:6379 && tail -f /dev/null"

# ray start --head --port=6379 --disable-usage-stats --dashboard-host=0.0.0.0
Usage stats collection is disabled.

Local node IP:

Ray runtime started.

Next steps
To add another node to this Ray cluster, run
ray start --address='*:6379'

To connect to this Ray cluster:
import ray
ray.init()

To submit a Ray job using the Ray Jobs CLI:
RAY_ADDRESS='http://*:8265' ray job submit --working-dir . -- python my_script.py

See https://docs.ray.io/en/latest/cluster/running-applications/job-submission/index.html
for more information on submitting Ray jobs to the Ray cluster.

To terminate the Ray runtime, run
ray stop

To view the status of the cluster, use
ray status

To monitor and debug Ray, view the dashboard at
:8265

If connection to the dashboard fails, check your firewall settings and network configuration.

# ray start --block --address=:6379
Local node IP: 192.168.75.174
[2025-07-29 05:37:09,112 W 646 646] global_state_accessor.cc:435: Retrying to get node with node ID 6395285e5a2a36e5773e77fd095490d63eda1f33869da1d09291b522

Ray runtime started.

To terminate the Ray runtime, run
ray stop
--block
This command will now block forever until terminated by a signal.
Running subprocesses are monitored and a message will be printed if any of them terminate unexpectedly. Subprocesses exit with SIGTERM will be treated as graceful, thus NOT reported.
```

Worker Node

ray list nodes

```
===== List: 2025-07-29 05:43:23.404417 =====
Stats:
-----
Total: 2

Table:
-----
NODE_ID          NODE_IP      IS_HEAD_NODE STATE
0 05ala18c99ab485e6af30b864bd1227c09815f71f40eac26ddbc3f9a  True        ALIVE
1 6395285e5a2a36e5773e77fd095490d63eda1f33869da1d09291b522 False       ALIVE
```

Overview Jobs Serve Cluster Actors Metrics Logs

NODES

Node Statistics

Host	IP	Node ID	State	CPU	Memory	GPU	GRAM	Object Store Memory	Disk(root)	Sent
05ala...	192.168.75.174	05ala18c99ab485e6af30b864bd1227c09815f71f40eac26ddbc3f9a	ALIVE	0.0%	644MB/81559MB	0.0%	0.0000B/186.29GB(0.0%)	7.41TB/16.8TB(0.0%)	147.17KB/s	
63952...	192.168.75.174	6395285e5a2a36e5773e77fd095490d63eda1f33869da1d09291b522	ALIVE	0.2%	308.05MB/1007.51GB(0.0%)	0.0%	0.0000B/186.29GB(0.0%)	8.21TB/16.8TB(0.0%)	28.97KB/s	

[38] https://docs.vllm.ai/en/v0.9.2/examples/online_serving/run_cluster.html[39] <https://docs.vllm.ai/en/v0.9.2/usage/security.html>[40] <https://docs.pytorch.org/docs/stable/distributed.html#common-environment-variables>[41] https://docs.nvidia.com/deeplearning/nccl/user-guide/docs/env.html#nccl_ib-disable

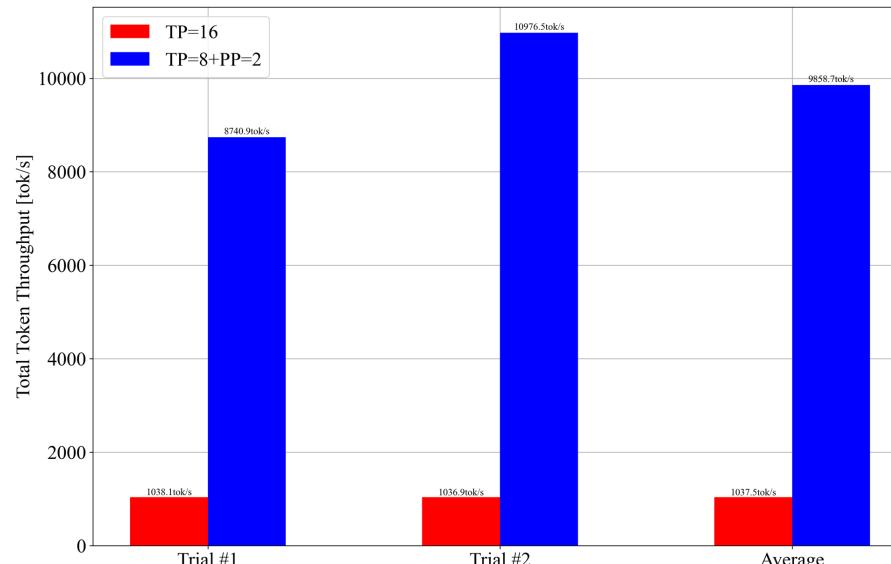
Multi-node Distributed Inference

□ Multi-Node Multi-GPU [42, 43]

(tensor parallel plus pipeline parallel inference)

- If your model is too large to fit in a single node, you can use tensor parallel together with pipeline parallelism.
 - The tensor parallel size is the number of GPUs you want to use in each node, and the pipeline parallel size is the number of nodes you want to use.
 - E.g., if you have 16 GPUs in 2 nodes (8 GPUs per node), you can set **the tensor parallel size to 8** and **the pipeline parallel size to 2**.

TP 16 vs. TP 8 + PP 2



NOTE master

```
vllm serve Qwen/Qwen3-235B-A22B \
    --distributed-executor-backend ray \
    --host=0.0.0.0 --port=8080 \
    --tensor-parallel-size=8 --pipeline-parallel-size=2 \
    --gpu_memory_utilization=0.95 \
    --reasoning-parser qwen3
```

rank 15 in world size 16 is assigned as DP rank 0, PP rank 1, TP rank 7, EP rank 7 [repeated 15x across cluster]

Host / Worker Process name	State	State Message	ID	IP / PID	Actions	CPU ⓘ	Memory ⓘ	GPU ⓘ	GRAM
>	ALIVE	-	a3b57...	(Head)	Log	17.7%	30.53GB/1007.51GB(3.0%)	[0]: 100.0% [1]: 100.0% [2]: 100.0% [3]: 100.0% [4]: 100.0% [5]: 100.0% [6]: 100.0% [7]: 100.0%	[0]: 80813MB/81559MB [1]: 81021MB/81559MB [2]: 81019MB/81559MB [3]: 81019MB/81559MB [4]: 81019MB/81559MB [5]: 80987MB/81559MB [6]: 81019MB/81559MB [7]: 80967MB/81559MB
>	ALIVE	-	4a527...		Log	10.7%	21.21GB/1007.51GB(2.1%)	[0]: 100.0% [1]: 100.0% [2]: 100.0% [3]: 100.0% [4]: 100.0% [5]: 100.0% [6]: 100.0% [7]: 100.0%	[0]: 80363MB/81559MB [1]: 80525MB/81559MB [2]: 80523MB/81559MB [3]: 80523MB/81559MB [4]: 80523MB/81559MB [5]: 80491MB/81559MB [6]: 80523MB/81559MB [7]: 80513MB/81559MB

Multi-node Distributed Inference

□ RDMA (Remote Direct Memory Access) [44]

- Network에 연결된 두 node의 memory 간 CPU, cache, 운영 체제의 개입 없이 data를 직접 전송하는 기술
- E.g., InfiniBand, RoCE, iWARP

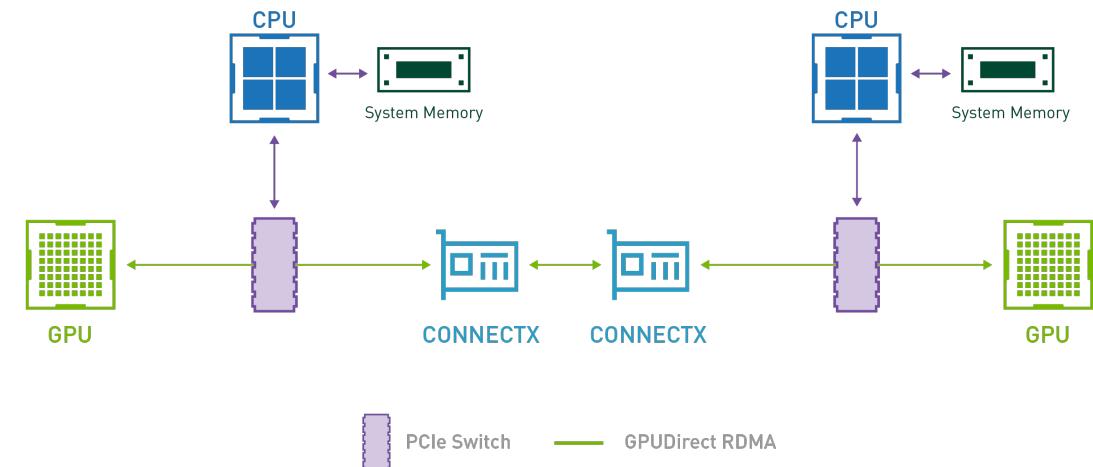
□ InfiniBand [45]

- RDMA를 native로 지원하여 CPU 개입 없이 node 간 memory 직접 전송 가능
- 전용 switch와 network adapter (Host Channel Adapter, HCA) 사용 (Ethernet network와는 다른 독자적 architecture)

□ RoCEv2 (RDMA over Converged Ethernet v2) [44]

- 표준 ethernet network를 통해 RDMA를 구현하는 protocol

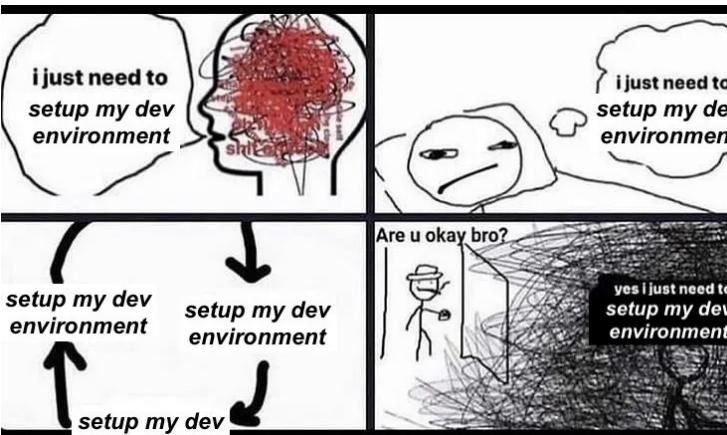
항목	InfiniBand	RoCE
Network type	전용 InfiniBand fabric	Ethernet 기반
구축 비용	별도 장비 필요, 비용 높음	상대적으로 저렴, 기존 infra 활용 가능
성능	최고 성능, ultra-low latency	Lossless Ethernet 환경에서 InfiniBand에 근접
호환성	HPC/AI 특화, 범용성 낮음	범용성 높음, 기존 infra와 통합 용이
관리 복잡도	전용 환경, 관리 일관성	Ethernet tuning 필요, 설정 복잡할 수 있음



[44] <https://zerohertz.github.io/distributed-computing-rdma-roce/>

[45] <https://developer.nvidia.com/gpudirect>

Multi-node Distributed Inference



```
[...]
```

Production Deployment with GPU Cluster (Kubernetes)

4. Production Deployment

❑ Kubernetes 배포 [47]

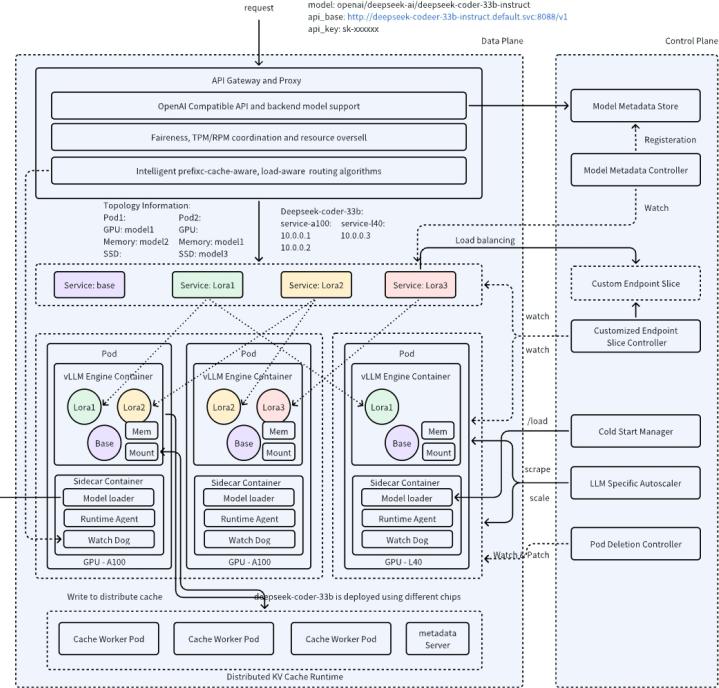
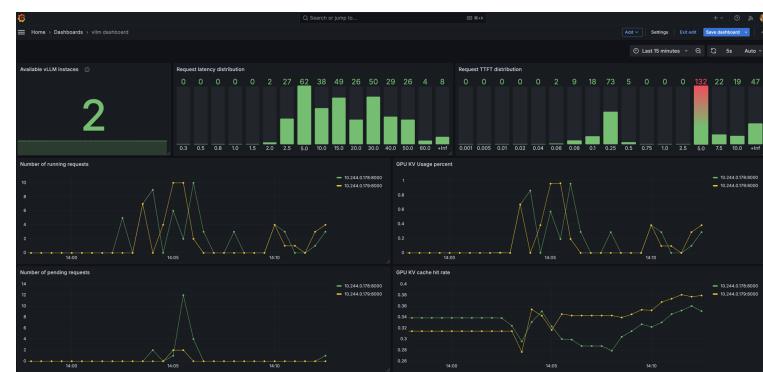
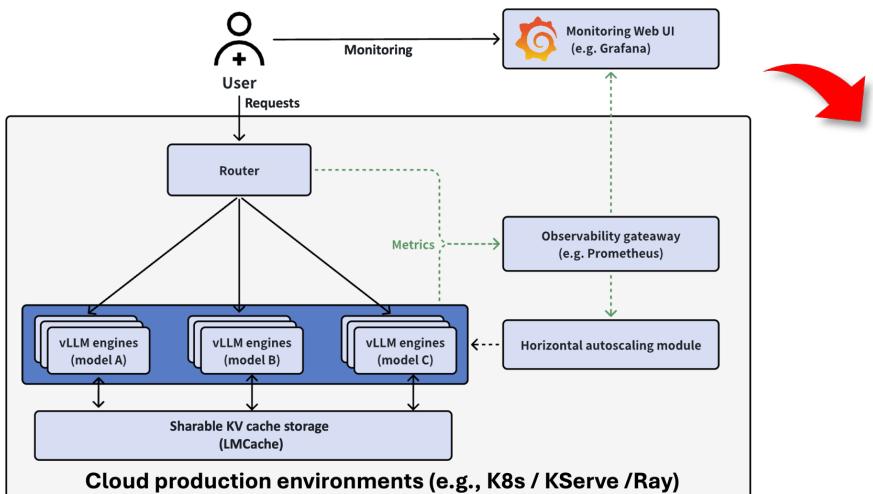
- Model store, Network, IPC, 환경 변수 등의 자유도가 많은 만큼 관리가 어려움

❑ AIBrix [48, 49, 50, 51]

- vLLM을 위한 확장 가능하고 비용 효율적인 cloud native control plane
- 고밀도 LoRA 관리, LLM gateway 및 routing, autoscaler, 분산 추론 및 분산 KV cache, …

❑ Production-stack [52, 53, 54, 55]

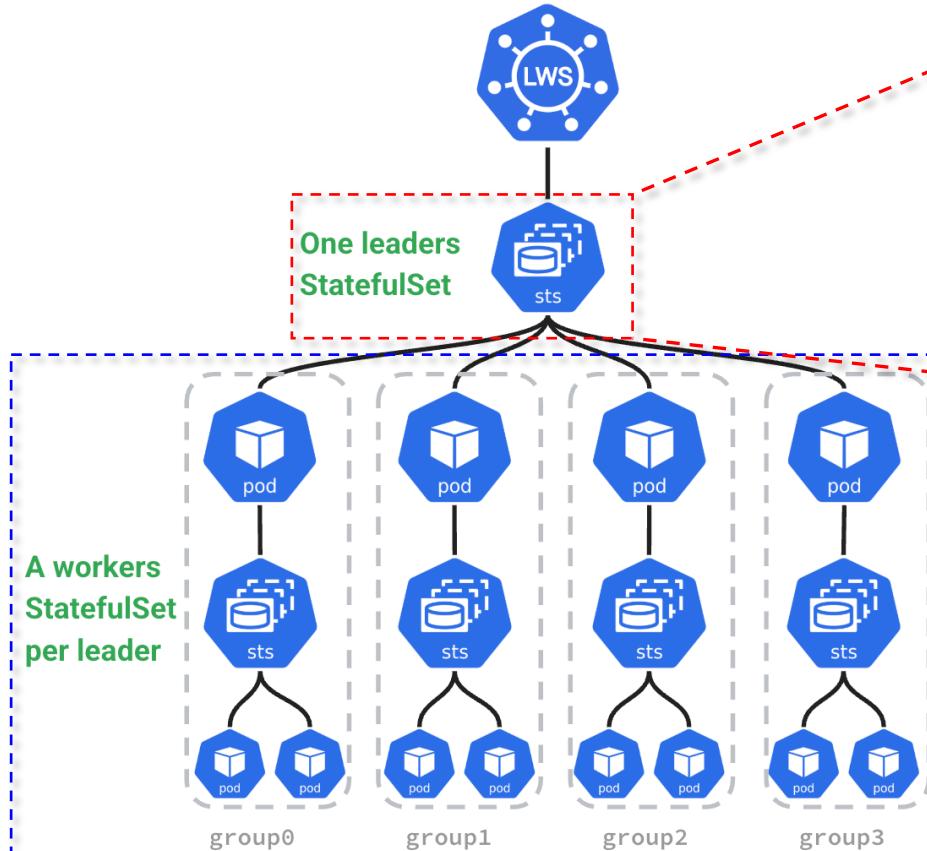
- vLLM을 통한 LLM 배포 시 production에 최적화된 codebase
- LMCache 기반 KV cache 공유 및 저장, prefix-aware routing, Helm chart를 통한 배포, …

[47] <https://docs.vllm.ai/en/v0.9.2/deployment/k8s.html>[48] <https://github.com/vllm-project/aibrix>[49] <https://arxiv.org/abs/2504.03648>[50] <https://blog.vllm.ai/2025/02/21/aibrix-release.html>[51] <https://aibrix.readthedocs.io/latest/>[52] <https://github.com/vllm-project/production-stack>[53] <https://docs.vllm.ai/en/v0.9.2/deployment/integrations/production-stack.html>[54] <https://blog.vllm.ai/2025/01/21/stack-release.html>[55] <https://blog.vllm.ai/production-stack/>

Production Deployment with GPU Cluster (Kubernetes)

❑ LWS (LeaderWorkerSet) [56, 57]

- Kubernetes 상에서 Leader-Worker architecture를 손쉽게 구현하고 관리할 수 있도록 도와주는 CRD 및 controller



```

1  apiVersion: leaderworkerset.x-k8s.io/v1
2  kind: LeaderWorkerSet
3  metadata:
4    name: vilm
5  spec:
6    replicas: 2
7    leaderworkerTemplate:
8      size: 2
9    restartingPolicy: RecreateGroupOnPodRestart
10   leaderTemplate:
11     metadata:
12       labels:
13         role: leader
14     spec:
15       containers:
16         - name: vilm-leader
17           image: docker.io/vilm/vilm-openai:latest
18           env:
19             - name: HUGGING_FACE_HUB_TOKEN
20               value: <your-hf-token>
21           command:
22             - sh
23             - -c
24             - "bash /vilm-workspace/examples/online_serving/multi-node-serving.sh leader --ray_cluster_size=$(LWS_GROUP_SIZE);"
25           resources:
26             limits:
27               nvidia.com/gpu: "8"
28               memory: 11240Gi
29             ephemeral-storage: 800Gi
30           requests:
31             ephemeral-storage: 800Gi
32             cpu: 125
33           ports:
34             containerPort: 8080
35             readinessProbe:
36               tcpSocket:
37                 port: 8080
38             initialDelaySeconds: 15
39             periodSeconds: 10
40           volumeMounts:
41             - mountPath: /dev/shm
42               name: dshm
43           volumes:
44             - name: dshm
45               emptyDir:
46                 medium: Memory
47                 sizeLimit: 15Gi
48   workerTemplate:
49     spec:
50       containers:
51         - name: vilm-worker
52           image: docker.io/vilm/vilm-openai:latest
53           command:
54             - sh
55             - -c
56             - "bash /vilm-workspace/examples/online_serving/multi-node-serving.sh worker --ray_address=$(LWS_LEADER_ADDRESS);"
57           resources:
58             limits:
59               nvidia.com/gpu: "8"
60               memory: 11240Gi
61             ephemeral-storage: 800Gi
62           requests:
63             ephemeral-storage: 800Gi
64             cpu: 125
65           env:
66             - name: HUGGING_FACE_HUB_TOKEN
67               value: <your-hf-token>
68           volumeMounts:
69             - mountPath: /dev/shm
70               name: dshm
71           volumes:
72             - name: dshm
73               emptyDir:
74                 medium: Memory
75                 sizeLimit: 15Gi
76
77  apiVersion: v1
78  kind: Service
79  metadata:
80    name: vilm-leader
81  spec:
82    ports:
83      - name: http
84        port: 8080
85        protocol: TCP
86        targetPort: 8080
87    selector:
88      leaderworkerset.sigs.k8s.io/name: vilm
89      role: leader
90      type: ClusterIP

```

Observability (Prometheus + Grafana)

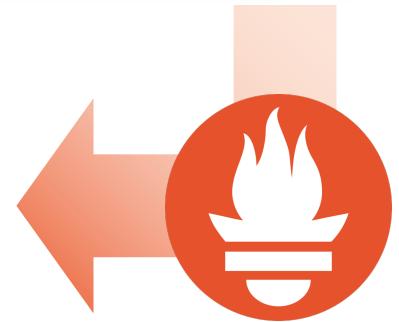
- ❑ “/metrics”를 통해 vLLM server의 Prometheus format metric 수집 가능 [58, 59]

- vllm:request_success_total: 요청 완료 수 (EOS 도달 또는 max 길이 도달)
 - vllm:request_queue_time_seconds: queue 대기 시간
 - vllm:request_prefill_time_seconds: prefill 단계 소요 시간
 - vllm:request_decode_time_seconds: decoding 단계 소요 시간
 - vllm:request_max_num_generation_tokens: 생성된 token 최대값
 - ...

□ Grafana와 연결하여 dashboard 구성 가능

```
zeroherzt@zeroherzt-mm:~$ curl http://localhost:8000/metrics

# HELP vlm:lora_requests_info Running stats on lora requests.
# TYPE vlm:lora_requests_info gauge
vlm:lora_requests_info{max_lora="0",running_lora_adapters=""},waiting_lora_adapters=""} 1.753628624009126e+09
# HELP vlm:lmcache_config_info Information of the LMEngine CacheConfig
# TYPE vlm:lmcache_config_info gauge
vlm:cache_config_info{block_size="16",cache_dtype="auto"},calculate_kv_scales="False",cpu_kvcache_space_bytes="4294967296",cpu_offload_gb="0.0",enable_prefix_caching="None",gpu_memory_UTILIZATION="0.9",is_attention_free="False",num_cpu_blocks="8",num_gpu_blocks="2340",num_gpu_blocks_override="None",prefix_caching_hash_algo="builtin",sliding_window="None",swap_space="4.0",swap_space_bytes="4294967296",swap_space_gb="0.0"} 1.0
# HELP vlm:num_premptions_total Cumulative number of preemption from the engine.
# TYPE vlm:num_premptions_total counter
vlm:num_premptions_total{model_name="Qwen/Qwen3-0.6B"} 0.0
# HELP vlm:prompt_tokens_total Number of prefill tokens processed.
# TYPE vlm:prompt_tokens_total counter
vlm:prompt_tokens_total{model_name="Qwen/Qwen3-0.6B"} 14.0
# HELP vlm:generation_tokens_total Number of generation tokens processed.
# TYPE vlm:generation_tokens_total counter
vlm:generation_tokens_total{model_name="Qwen/Qwen3-0.6B"} 107.0
# HELP vlm:request_success_total Count of successfully processed requests.
# TYPE vlm:request_success_total counter
vlm:request_success_total{finished_reason="stop",model_name="Qwen/Qwen3-0.6B"} 1.0
# HELP vlm:iteration_tokens_total Histogram of number of tokens per engine_step.
# TYPE vlm:iteration_tokens_total histogram
vlm:iteration_tokens_total{model_name="Qwen/Qwen3-0.6B"} 121.0
vlm:iteration_tokens_total_bucket{k=1.0,model_name="Qwen/Qwen3-0.6B"} 1099.0
vlm:iteration_tokens_total_bucket{k=8.0,model_name="Qwen/Qwen3-0.6B"} 1099.0
vlm:iteration_tokens_total_bucket{k=16.0,model_name="Qwen/Qwen3-0.6B"} 1100.0
vlm:iteration_tokens_total_bucket{k=32.0,model_name="Qwen/Qwen3-0.6B"} 1100.0
vlm:iteration_tokens_total_bucket{k=64.0,model_name="Qwen/Qwen3-0.6B"} 1100.0
vlm:iteration_tokens_total_bucket{k=128.0,model_name="Qwen/Qwen3-0.6B"} 1100.0
vlm:iteration_tokens_total_bucket{k=256.0,model_name="Qwen/Qwen3-0.6B"} 1100.0
vlm:iteration_tokens_total_bucket{k=512.0,model_name="Qwen/Qwen3-0.6B"} 1100.0
vlm:iteration_tokens_total_bucket{k=1024.0,model_name="Qwen/Qwen3-0.6B"} 1100.0
vlm:iteration_tokens_total_bucket{k=2048.0,model_name="Qwen/Qwen3-0.6B"} 1100.0
```



Benchmark

□ “vllm bench” CLI 명령을 통해 간단한 benchmark 가능 [60]

```
zerohertz@zerohertz-mm:~$ viml bench serve \
--base-url http://192.168.75.173:8080 \
--model Owen/Qwen3-235B-A22B \
--request-rate 2000 \
--num-prompts 10000 \
--random-input-len 600 \
--random-output-len 125 \
--ignore-eos \
--seed 42 \
--percentile-metrics "ttft,tpot,itl,e2el" \
--save-result

INFO 07-29 23:35:50 [__init__.py:244] Automatically detected platform cpu.
Namespace(subparser='bench', bench_type='serve', dispatch_function=<function BenchmarkServingSubcommand.cmd at 0x1458ad620>, seed=42, num_prompts=10000, dataset_name='random', dataset_path=None, custom_output_len=256, custom_skip_chat_template=False, sonnet_input_len=550, sonnet_output_len=150, sonnet_prefix_len=200, shareopt_output_len=None, random_input_len=600, random_output_len=125, random_ratio=0.0, random_prefix_len=0, hf_subset=None, hf_output_len=None, endpoint_type='openai', label=None, backend='vllm', base_url='http://192.168.75.173:8080', host='127.0.0.1', port=8080, endpoint='/v1/completions', max_concurrency=None, model='Qwen-en3-235B-A22B', tokenizer=False, use_beam_search=False, logprobs=None, request_rate=2000.0, burstiness=1.0, trust_remote_code=False, saveable_tqdm=False, profile=False, save_result=True, save_detailed=False, append_results=False, metadata=None, result_dir=None, result_filename=None, ignore_eos=True, percentile_metrics='ttft,tpot,itl,e2el', metric_percentiles='99', goodput=None, top_p=None, top_k=None, min_p=None, temperature=None, tokenizer_mode='auto', served_model_name=None, lora_modules=None, ramp_up_strategy=None, ramp_up_start_ps=None, ramp_up_end_rps=None)
INFO 07-29 23:35:52 [datasets.py:355] Sampling input_len from [600, 600] and output_len from [125, 125]
Starting initial single prompt test run...
Initial test run completed. Starting main benchmark run...
Traffic request rate: 2000.0
Burstiness factor: 1.0 (Poisson process)
Maximum request concurrency: None
[]% | 0/10000 [00:00:<, ?it
huggingface/tokenizers: The current process just got forked, after parallelism has already been used. Disabling parallelism to avoid deadlocks...
To disable this warning, you can either:
- Avoid using 'tokenizers' before the fork if possible
- Explicitly set the environment variable TOKENIZERS_PARALLELISM=(true | false)
100% | ===== Serving Benchmark Result ====== | 10000/10000 [03:01<00:00, 55.10it
Successful requests: 2548
Benchmark duration (s): 181.49
Total input tokens: 1525710
Total generated tokens: 318500
Request throughput (req/s): 14.04
Output token throughput (tok/s): 1754.87
Total Token throughput (tok/s): 10161.24
=====Time to First Token=====
Mean TTFT (ms): 85567.69
Median TTFT (ms): 90519.17
P99 TTFT (ms): 168053.72
=====Time per Output Token (excl. 1st token)=====
Mean TPOT (ms): 452.68
Median TPOT (ms): 522.39
P99 TPOT (ms): 596.92
=====Inter-token Latency=====
Mean ITL (ms): 452.69
Median ITL (ms): 161.38
P99 ITL (ms): 1666.37
=====End-to-end Latency=====
Mean E2EL (ms): 141699.53
Median E2EL (ms): 155198.55
P99 E2EL (ms): 180117.52
=====
```

□ vLLM repository 내 benchmarks 내 코드 사용 [61, 62]

v0.8.2	benchmarks	bench
bistrong [Benchmark] Add support for multiple batch size benchmark through C...	[Refactor] Remove duplicate test cases (#20023)	codeS2K 3 weeks ago
cuutils_benchmarks		
d3sapp_benchmarks	[Merge] Add SPOX-FixCopyrightText (#19100)	
fused_kernels	[Merge] Add SPOX-FixCopyrightText (#19100)	last month
kernels	[Benchmark] Add support for multiple batch size benchmark through C...	3 weeks ago
overheads	[Merge] Add SPOX-FixCopyrightText (#19100)	last month
structured_schemas	benchmarks: simplify test porchouse (#4565)	1 month ago
README.md	[Merge] Use colleague blocks for benchmark examples (#20017)	last month
auto_tunish	[Merge] [Toolbox] Fix auto_tunish (#1871)	last month
backend_request_func.py	[Benchmark] Fix request pool (#1979)	last month
benchmark_datalog.py	[Benchmark] Fix multiple bugs in select and add args to spec, doco...	last month
benchmark_interop.py	[Merge] Modularity CI: Argument Parsing in Benchmark Scripts (#19593)	last month
benchmark_long_document_0k_throughput.py	[Merge] Modularity CI: Argument Parsing in Benchmark Scripts (#19593)	last month
benchmark_pretty_coding.py	[Merge] Modularity CI: Argument Parsing in Benchmark Scripts (#19593)	last month
benchmark_annotation.py	[Merge] Modularity CI: Argument Parsing in Benchmark Scripts (#19593)	last month
benchmark_jersey.py	[Merge] remove redundant check (#2057)	last month
benchmark_enveloping_structured_output.py	[Merge] Modularity CI: Argument Parsing in Benchmark Scripts (#19593)	last month
benchmark_throughput.py	[Benchmark] Fix value of Type 'Time' (#1954)	last month
benchmark_utils.py	Handle non-serializable objects when benchmarking results (#19154)	2 months ago
pyproject.toml	[Doc] Minor examples and further reorganization user guide (#18666)	2 months ago
non_structured_output_benchmark.sh	[Benchmark] Refactor non_structured_output_benchmark.sh (#17722)	2 months ago
sonnet.txt	[Sonnet] Refactor sonnet (#17722)	2 months ago

□ vLLM project의 GuidedLM 사용 [63, 64]

Benchmarks Metadata															
Run Id:74ad923-9a0d-4c83-b237-837ccf971bd3															
Duration:11.7 seconds															
Profile type: write, strategy: 'synchronous', throughput, 'constant', 'constant'															
Arg type: string, duration: 1.0, warp_size: 1024, warp_number: 1024, warp_duration:None, cooldown_number:None, cooldown_duration:None															
Worker type: "generative_requests_worker", backend_type: "openapi", http_base:"http://192.168.1.4:8080", backend_mode:"neuralmagic/Owen2.5-7.8-quarantine", organization: "None", project: "None", test_completions_path: "/v1/completions", chat_completions_path: "/v1/chat/completions", request_completions_type: "generative_request_loader", date_m: "prompt_tokens": 128, "output_tokens": 64", data_org: "None", processor: "neuralmagic/Owen2.5-7.8-quarantine", extras: "None"															
Benchmarks Info:															
Metadata															
Benchmark	Start Time	End Time	Duration (s)	Requests Made	Comp	Err	Prompt	Tok/Rq	Output	Tok/Rq	Prompt	Tok Total	Output	Tok Total	
synchronous	14:42:16	14:42:16	0.0	6	1	0	126.0	126.0	0	64.0	43.0	0	768	128	
throughput	14:42:16	14:42:26	10.0	126	518	0	128.0	128.0	0	64.0	36.9	0	16144	653080	
constant@.20	14:42:33	14:42:43	0.0	10	20	0	128.0	128.0	0	64.0	37.0	0	2569	512	
constant@.50	14:42:33	14:42:43	0.0	10	50	0	128.0	128.0	0	64.0	37.0	0	4529	936	
constant@.25	14:42:53	14:43:03	0.0	10	25	0	128.0	128.0	0	63.9	32.9	0	1146	1152	
constant@.41	14:42:53	14:43:03	0.0	10	48	9	128.0	128.0	0	64.0	37.0	0	3668	296	
constant@.81	14:42:56	14:43:16	0.0	62	12	0	128.0	128.0	0	64.0	32.5	0	7938	1537	
constant@.14	14:42:56	14:43:16	0.0	62	14	0	128.0	128.0	0	64.0	32.5	0	3574	714	
constant@.99	14:42:56	14:43:16	0.0	62	19	0	128.0	128.0	0	64.0	32.5	0	11266	2253	
constant@.88	14:43:29	14:43:49	10.0	88	19	0	128.0	128.0	0	63.9	31.1	0	11266	2433	
constant@.41	14:43:48	14:43:58	0.0	101	24	0	128.0	128.0	0	64.0	32.0	0	12934	387	
constant@12.95	14:43:51	14:44:01	10.0	113	28	0	128.0	128.0	0	64.0	31.7	0	14475	3585	
Benchmarks Stats:															
Metadata															
Benchmark	Request Stats			Out Tok/sec			Tot Tok/sec			Req Latency (ms)			TTTF (ms)		
	Per Second	Concurrency	mean	mean	median	mean	mean	median	p99	mean	median	p99	mean	median	
synchronous	0.57	1.00	42.7	212.1	1.40	1.49	1.58	48.0	20.2	22.0	22.0	23.1	23.1	23.1	
throughput	12.95	113.98	828.5	113.9	87.7	8.65	9.66	1198.8	1112.2	2057.8	120.1	120.1	116.2	116.2	
constant@.20	2.06	3.22	131.9	657.9	1.56	1.56	1.58	58.7	58.3	73.8	23.8	23.9	23.9	23.9	
constant@.50	2.06	3.22	131.9	222.9	1.00	1.00	1.00	58.7	58.3	73.8	23.8	23.9	23.9	23.9	
constant@.25	4.99	8.21	313.3	1564.0	1.67	1.67	1.73	62.5	39.9	106.4	25.6	25.6	25.6	25.6	
constant@.81	6.25	10.88	399.4	1992.3	1.74	1.73	1.83	64.6	61.1	129.3	26.6	26.5	27.5	26.2	
constant@.14	7.78	14.57	447.4	2472.0	1.80	1.79	1.88	64.6	61.1	129.3	26.6	26.5	27.5	26.1	
constant@.99	8.86	14.76	566.2	2825.0	1.99	1.98	2.08	77.0	65.9	184.9	38.4	38.4	39.4	39.7	
constant@1.41	10.18	21.23	650.9	3247.1	2.09	2.07	2.24	81.1	67.3	222.1	31.8	31.7	33.8	33.3	
constant@12.95	11.36	24.92	727.0	3627.1	2.19	2.17	2.37	87.4	70.1	257.3	33.4	33.2	35.9	32.9	

Saving benchmarks report...

Benchmarks Report saved to: /home/mark/github/neuralmagic/guidelml/benchmarks.json

[60] <https://docs.vllm.ai/en/v0.9.2/cli/index.html#henc>

[61] <https://docs.vllm.ai/en/v0.9.2/contributing/profiling.html>

[62] <https://github.com/yllm-project/yllm/tree/v0.9.2/benchmark>
[63] <https://github.com/yllm-project/guideline>

[64] <https://arxiv.org/pdf/2502.06494.pdf>

Production Issues

- ## ❑ Gloo connectFullMesh failed with… [6]

- Multi node로 process 연결 시 GLOO 연결 실패

→ “GLOO_SOCKET_IFNAME” 환경 변수를 사용할 intent

Debugging vLLM NCCL (PyNcclCommunicator) all_reduce Issue in Multi-Node Environment #11353

Closed Zerohertz announced in Q8A

Zerohertz on Dec 20, 2024 edited - ...

암려운 편 흰색 주제 - 한국어로 번역

Hi, I'm in the process of building a multi-node LLM serving environment via vLLM.
However, some settings were not set correctly, causing an error, and I'm verifying the environment through [official document](#).
I have confirmed that steps 1 and 2 ([PyTorch NCCL](#) and [PyTorch GLUE](#)) was successful.

But there was a problem in [vLLM NCCL](#) .
So I checked as follows.

```
dist.init_process_group(backend="nccl")
local_rank = dist.get_rank() % torch.cuda.device_count()
torch.cuda.set_device(local_rank)
world_size = torch.cuda.device_count()
gloo_group = dist.new_group(ranks=list(range(world_size)), backend="gloo")
pynccl = PyNcclCommunicator(group=gloo_group, device=local_rank)

s = torch.cuda.Stream()
data = torch.FloatTensor([
    1,
    1 + 128
]).to("cuda")
with torch.cuda.stream(s):
    logger.debug(f"Rank {local_rank}, data before all_reduce: {data}")
    pynccall.all_reduce(data, streams=s)
    logger.debug(f"Rank {local_rank}, data after all_reduce: {data}")
    value = data.mean().item()

Also all scripts run in Docker and I ran the Docker container as below:
```

Also all scripts run in Docker and I ran the Docker container as below:

```
# Master
$ docker run -d \
--name vlm \
--entrypoint /bin/bash \
--network host \
--ipc host \
--gpus "device0,3" \
-v ./vlm/workspace \
-e GLOO_SOCKET_INNAME=en03 \
-e NCCL_SOCKET_INNAME=en03 \
-e NCCL_NINUM_THREADS=1 \
-vlm/vlm-openai:v8.6.4 \
-c "tail -f /dev/null"
# Worker
$ docker run -d \
--name node \
--entrypoint /bin/bash \
--network host \
--ipc host \
--gpus "device0,3" \
-v ./vlm/workspace \
-e GLOO_SOCKET_INNAME=en03 \
-e NCCL_SOCKET_INNAME=en03 \
-e NCCL_NINUM_THREADS=1 \
-vlm/vlm-openai:v8.6.4 \
-c "tail -f /dev/null"

# Master
$ NCCL_DEBUG=TRACE torchrun --nnodes 2 --nproc-per-node=2 --node-rank 0 --rdzv_backend=c10d --rdzv_enk
# Worker
$ NCCL_DEBUG=TRACE torchrun --nnodes 2 --nproc-per-node=2 --node-rank 1 --rdzv_backend=c10d --rdzv_enk

Important

As a result, I expected the value to be changed to world_size before and after all_reduce, but I confirmed that was no change in any error log.
```

- Qwen 계열 model의 무한 token 생성 [66, 67]

- FlashInfer kernel 내 cascade inference 사용으로 인해 발생

→ “--disable-cascade-attn”를 사용하여 해결

[Bug]: Degradation of Qwen/Qwen3-30B-A3B performance depending on batch size #17652

PART 5

Wrap-up

Roadmap Q3 2025

[68]

5. Wrap-up

❑ V1 Engine

- V0 Engine 완전 제거
- Core scheduler 최적화 및 확장
- Async scheduling, multi-modal 처리 등 기능 구현

❑ Large Scale Serving

- Mixture-of-Experts (MoE) model의 안정적 scale-out serving
- 분산 서빙 표준화 및 autoscaling

❑ Models

- 다양한 framework (training, authoring)의 tokenizer, configuration, processor 지원
- Sparse attention mechanism
- 1B 이하의 소형 모델 성능 향상

❑ Use Cases

- RLHF
 - 동기화 및 resharding을 위한 가중치 로딩 최적화
 - Multi-turn scheduling
- Evaluation
 - Batching order에 영향받지 않는 full determinism 지원 (with/without prefix cache)
- Batch Inference
 - Prefix caching과 함께 scale-out을 위한 data parallel router
 - CPU KV cache offloading

Conclusion

❑ OpenAI-Compatible Server

- LangChain, Gemini CLI 등 기존 생태계와의 호환성 유지
- Tool Calling, Reasoning 등 확장된 기능 활용 가능

❑ Architecture

- KV cache와 PagedAttention 기반의 효율적 메모리 관리
- V1 Engine 도입으로 단순한 고성능 구조 구현

❑ Production Deployment

- TP/PP/DP/EP 등 다양한 병렬 처리 전략 적용
- Kubernetes 기반 multi-node 분산 추론
- LoRA adapter를 활용한 사용자 맞춤형 경량화 serving
- Prometheus + Grafana를 활용한 observability 확보
- vllm bench 및 benchmark script를 통한 성능 평가

vLLM Meetup in Korea^[69, 70]

5. Wrap-up



Junghwan Park • 1촌

AI/Data Engineer @ SK Telecom | Lead Maintainer @ PyTorch Korea User Group | PyT…

17시간 ·

The first **vLLM** Meetup in Korea will be held on the evening of Aug 19th!

Hosted by RedHat and **Rebellions**, with **PyTorch Korea User Group**
and **SqueezeBits**.

For more details, check out the link below: <https://lnkd.in/gRTGeKuH>

번역 표시



8/19(화) 저녁, 한국에서 열리는 첫번째 vLLM meetup에 함께 해주세요!

discuss.pytorch.kr

EoD

GitHub



LinkedIn



Coffee Chat

