

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?

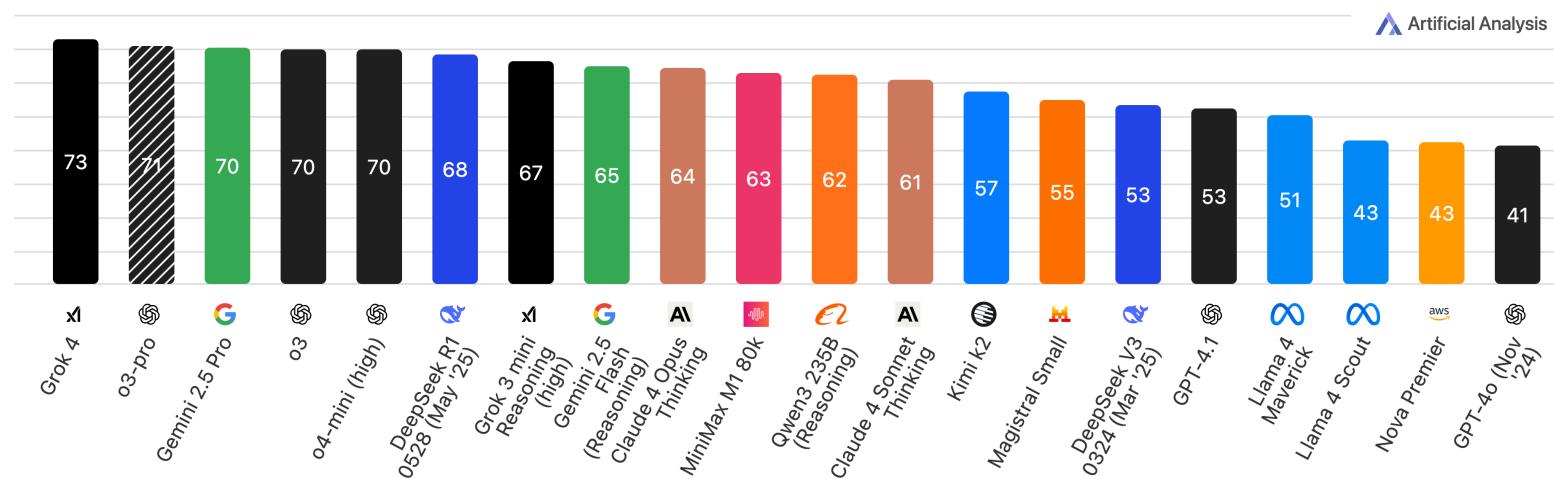
1. Introduction

- 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)



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>

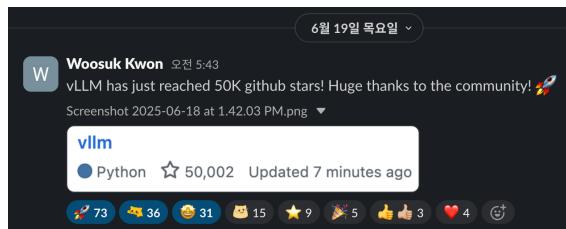
Hello, PyCon Korea 2025! 🎉
We're excited to have you join us for the event! It's scheduled for June 10th, 2025, and we're looking forward to celebrating tech innovation and showcasing the world's best trends! Let me know if you need help or have any questions! 🚀
```



Why vLLM?

❑ vLLM history

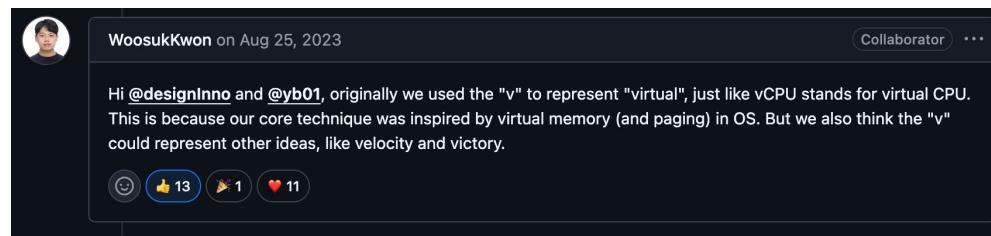
- 2023년 2월 9일, GitHub에서 CacheFlow^[3]라는 이름으로 시작
- 2023년 6월 17일, vLLM으로 이름 변경^[4]
- 2023 9월 12일, UC Berkeley Sky Computing Lab에서 “Efficient Memory Management for Large Language Model Serving with PagedAttention” 논문 발표^[5]
- 2025년 6월 19일, GitHub star 50k 달성



- 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>.

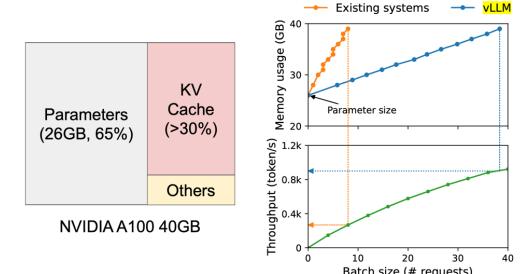


Figure 1. Left: Memory layout when serving an LLM with 13B parameters on NVIDIA A100. The parameters (gray) persist in GPU memory throughout serving. The memory for the KV cache (red) is (de)allocated per serving request. A small amount of memory (yellow) is used ephemerally for activation. Right: vLLM 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/e7d9d9c08e79b386f6d0477e87b7a572390317d>

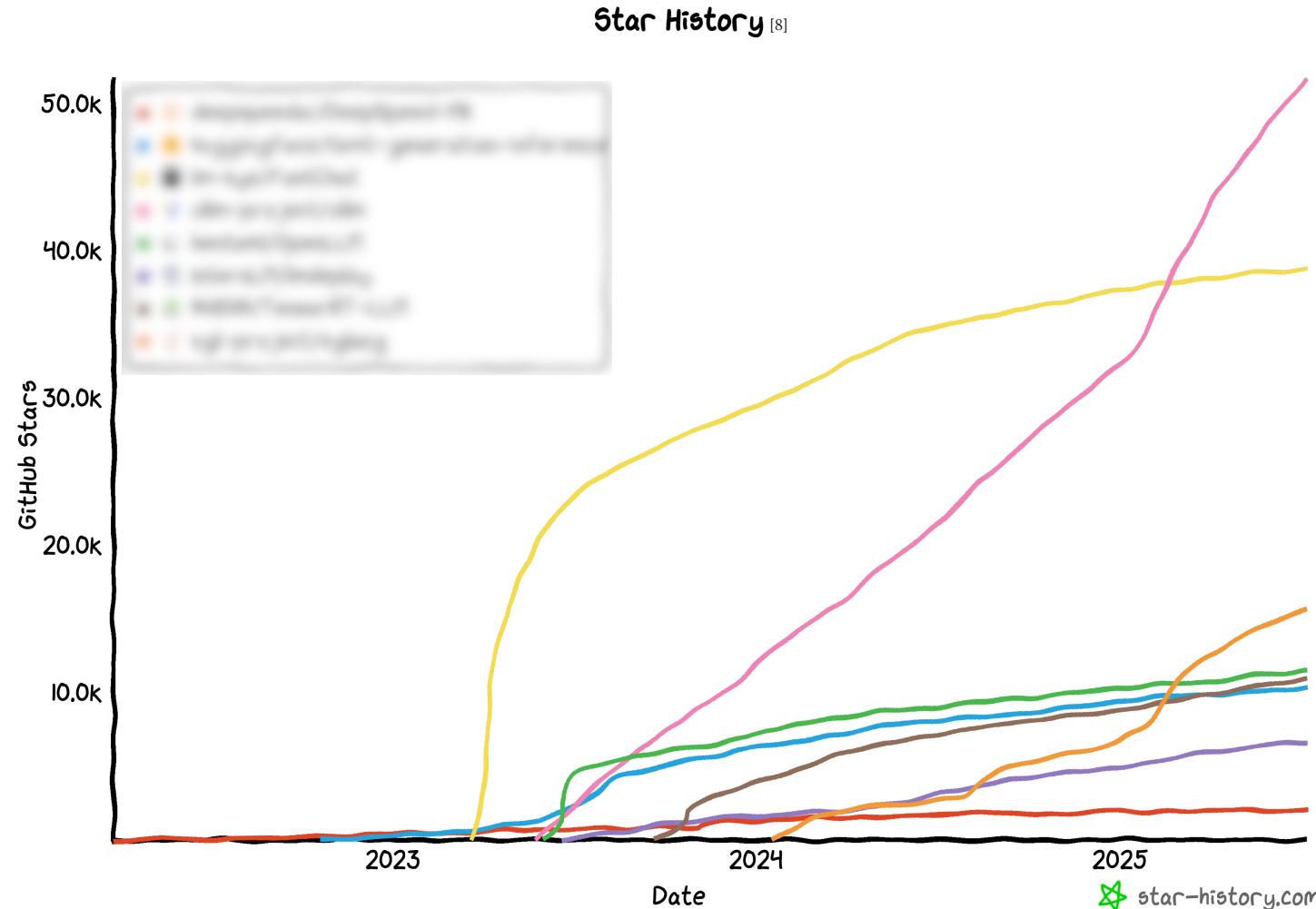
[4] <https://github.com/vllm-project/vllm/commit/0b98ba15c744fd1fb0ea4f2135e85ca23d572ae1>

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

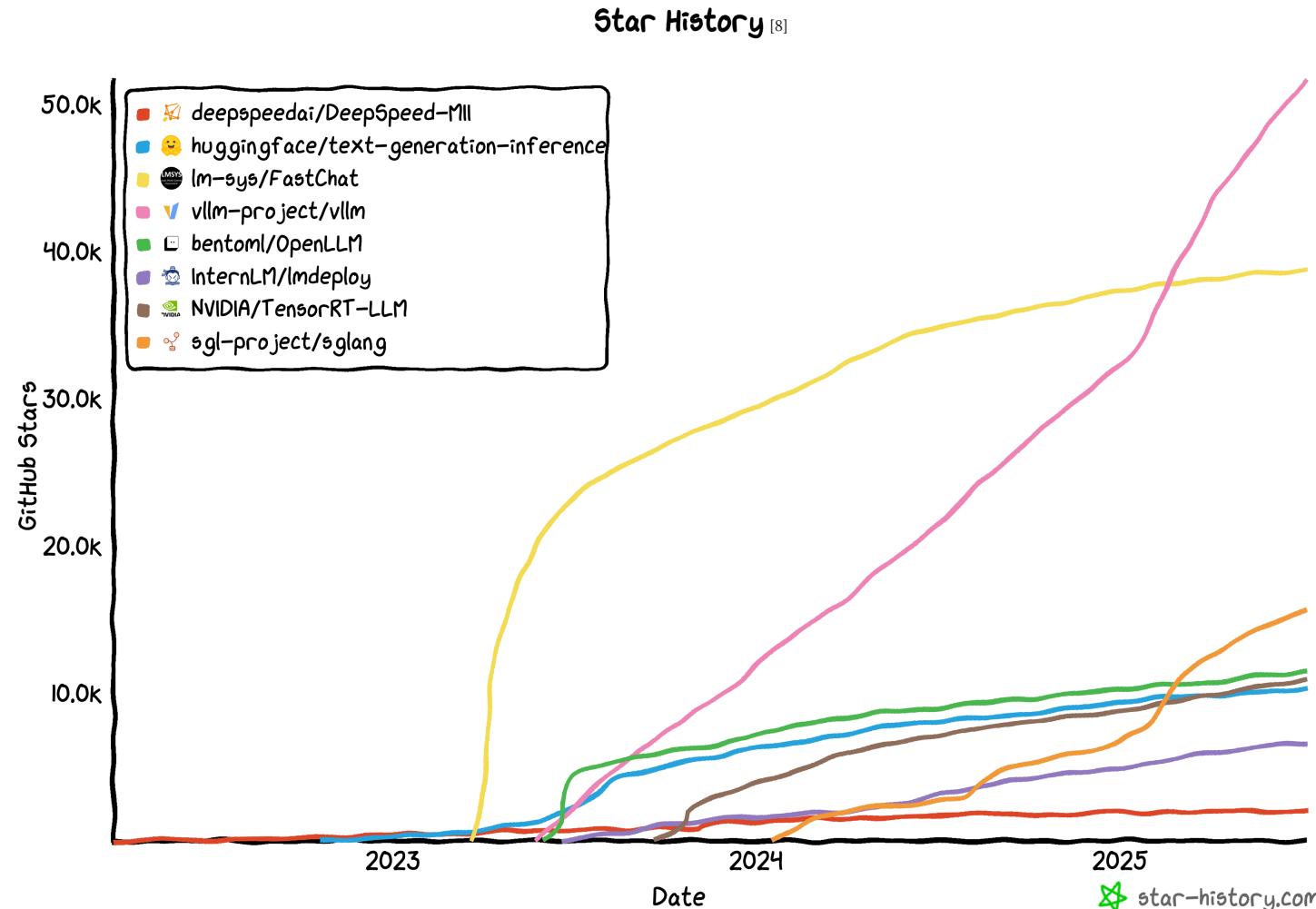
[6] <https://arxiv.org/abs/2309.06180>

[7] <https://github.com/vllm-project/vllm/issues/835>

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

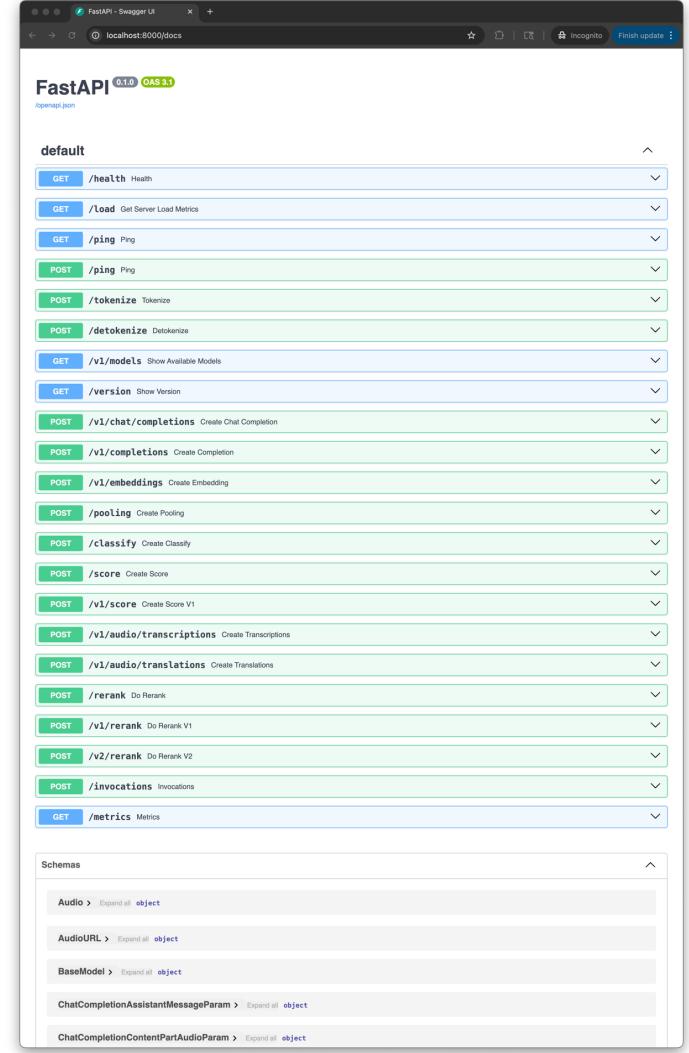
```
$ uv pip install vllm==0.9.2
Using Python 3.12.11 environment at: /opt/venv/main
Resolved 120 packages in 47ms
Installed 1 package in 10ms
+ vllm==0.9.2
```

▪ 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 [cl1_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, 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: /v1/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]

```
zerohertz@zerohertz-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": "Owen/Owen3-0.6B",
      "object": "model",
      "created": 1753618176,
      "owned_by": "Ville",
      "root": "Owen/Owen3-0.6B",
      "parent": null,
      "max_model_len": 8192,
      "permissions": [
        {
          "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]

```
zerohertz@zerohertz-mm:~>
$ curl -X POST http://localhost:8000/v1/chat/completions \
-d '{
  "model": "Owen/Owen3-0.6B",
  "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": "Owen/Owen3-0.6B",
  "choices": [
    {
      "index": 0,
      "message": {
        "role": "assistant",
        "content": "Hello! I'm InstructKR! How can I assist you today? I'm here to help with any questions or tasks you have. Please feel free to ask me anything, and I'll do my best to provide a helpful response. If you have any specific instructions or requirements, please let me know and I'll make sure to follow them. I'm always available to answer your queries and provide support. Let's get started! If you have any further questions or need more information, just let me know."}
    }
  ],
  "logprobs": null,
  "finish_reason": "stop",
  "stop_reason": null
}
```

```
zerohertz@zerohertz-mm:~>
$ curl -X POST http://localhost:8000/v1/chat/completions \
-d '{
  "model": "Owen/Owen3-0.6B",
  "messages": [
    {
      "role": "user",
      "content": "Hello, InstructKR!"
    }
  ],
  "stream": true
}' | jq
data: {"id": "chatmpl-6892342719c74308089f837bde8b16e6", "object": "chat.completion.chunk", "created": 1753618208, "model": "Owen/Owen3-0.6B", "choices": [{"index": 0, "delta": {"role": "assistant", "content": ""}, "logprobs": null, "finish_reason": null}}]
data: {"id": "chatmpl-6882342719c74308089f837bde8b16e6", "object": "chat.completion.chunk", "created": 1753618208, "model": "Owen/Owen3-0.6B", "choices": [{"index": 0, "delta": {"content": "Okay"}, "logprobs": null, "finish_reason": null}}]
data: {"id": "chatmpl-6882342719c74308089f837bde8b16e6", "object": "chat.completion.chunk", "created": 1753618208, "model": "Owen/Owen3-0.6B", "choices": [{"index": 0, "delta": {"content": "I'm here to help with any questions or tasks you have. Please feel free to ask me anything, and I'll do my best to provide a helpful response. If you have any specific instructions or requirements, just let me know and I'll make sure to follow them. I'm always available to answer your queries and provide support. Let's get started! If you have any further questions or need more information, just let me know."}, "logprobs": null, "finish_reason": null}}]
data: {"id": "chatmpl-6882342719c74308089f837bde8b16e6", "object": "chat.completion.chunk", "created": 1753618208, "model": "Owen/Owen3-0.6B", "choices": [{"index": 0, "delta": {"content": "Just a quick reminder, I'm here to help with any questions or tasks you have. Please feel free to ask me anything, and I'll do my best to provide a helpful response. If you have any specific instructions or requirements, just let me know and I'll make sure to follow them. I'm always available to answer your queries and provide support. Let's get started! If you have any further questions or need more information, just let me know."}, "logprobs": null, "finish_reason": null}}]
data: {"id": "chatmpl-6882342719c74308089f837bde8b16e6", "object": "chat.completion.chunk", "created": 1753618208, "model": "Owen/Owen3-0.6B", "choices": [{"index": 0, "delta": {"content": "Just a quick reminder, I'm here to help with any questions or tasks you have. Please feel free to ask me anything, and I'll do my best to provide a helpful response. If you have any specific instructions or requirements, just let me know and I'll make sure to follow them. I'm always available to answer your queries and provide support. Let's get started! If you have any further questions or need more information, just let me know."}, "logprobs": null, "finish_reason": null}}]
```

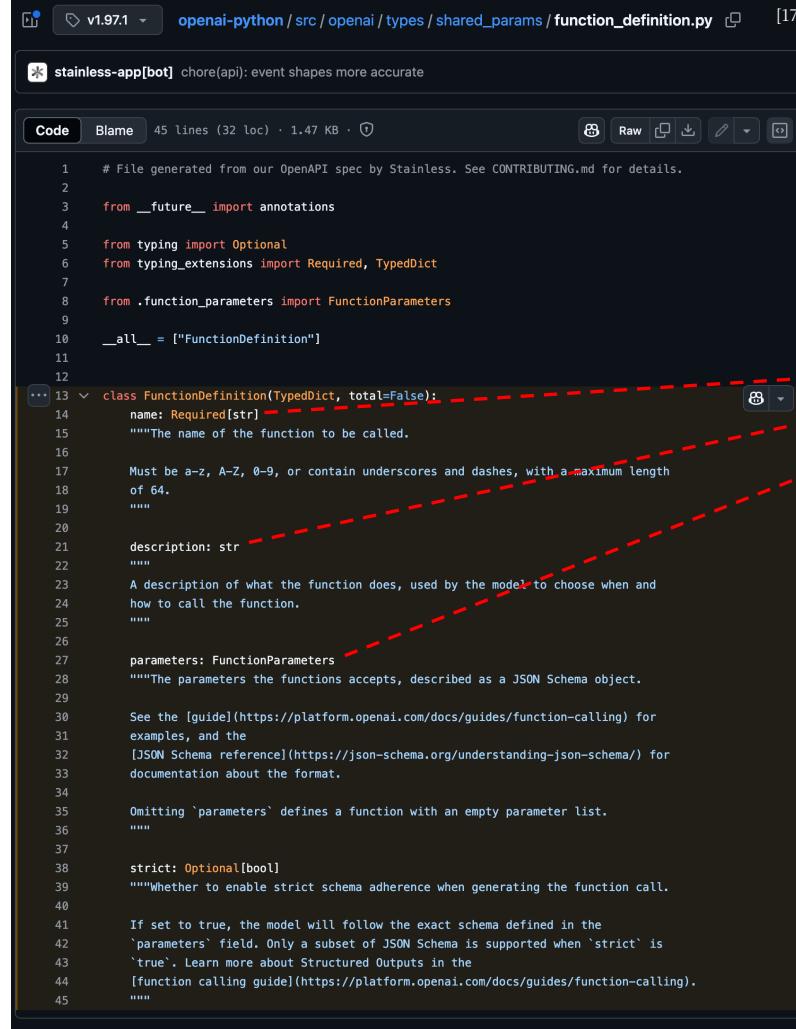
[14] https://docs.vllm.ai/en/v0.9.2/serving/openai_compatible_server.html

[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



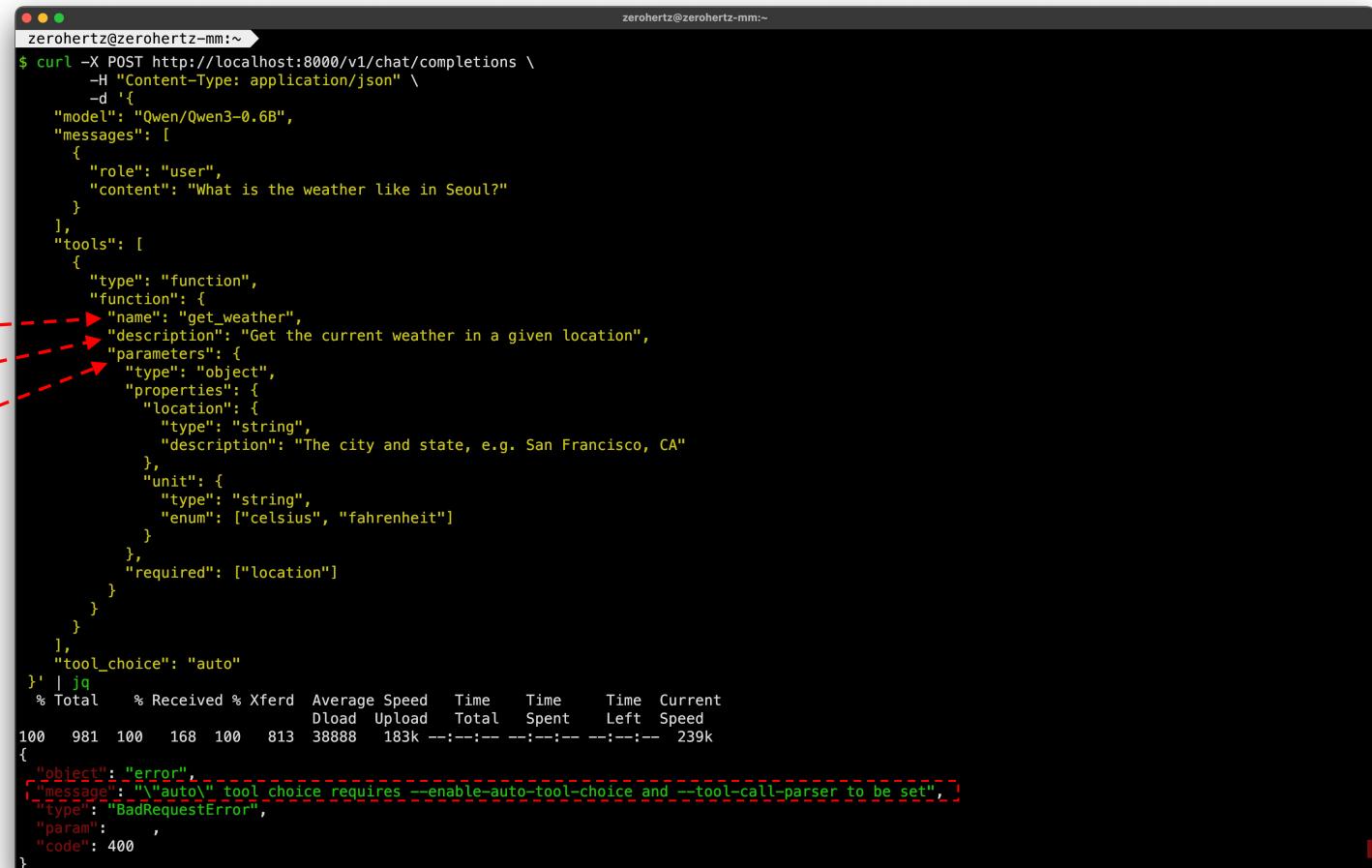
```

v1.97.1 - openai-python / src / openai / types / shared_params / function_definition.py [17]
stainless-app[bot] chore(api): event shapes more accurate

Code Blame 45 lines (32 loc) · 1.47 KB · 🛡
Raw ⌂ ⌄ ⌅ ⌆ ⌇ ⌈ ⌉ ⌋ ⌊ ⌊

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
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_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=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_async_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': 1, '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
```

```
$ 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
          0     0    1144      813    286    203  0:00:04  0:00:03  0:00:01  489
  100  1957    100  1144      813    286    203  0:00:04  0:00:03  0:00:01  489
  +--:0+ "chatml-0e4df805653947df9431ea8ffba8b20#",
  +--:0+ "chatml-0e4df805653947df9431ea8ffba8b20#",
  +--:0+ "created": 1753618598,
  +--:0+ "model": "Qwen/Qwen3-0.6B",
  +--:0+ "choices": [
    {
      "index": 0,
      "message": {
        "role": "assistant",
        "content": "The current weather in Seoul is 25 degrees Celsius. The unit is Celsius. Since the user didn't specify the unit, I'll leave it out. So the function call should be get_weather with location='Seoul' and unit omitted. That should get the current weather information for Seoul.\n</think>\n\n"
      }
    }
  ],
  "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": "\u27e8think\u27e9\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.\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

2. OpenAI-Compatible Server

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

```
vllm serve Owen/Qwen3-0.6B --max-model-len 8192 --reasoning-parser qwen3
$ vllm serve Owen/Qwen3-0.6B --max-model-len 8192 --reasoning-parser qwen3
INFO 07-27 21:17:54 [__init__.py:244] 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': 'Owen/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 [lilm_engine.py:230] Initializing a V0 LLM engine (v0.9.2) with config: model='Owen/Qwen3-0.6B', speculative_config=None, tokenizer='Owen/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='Owen/Qwen3-0.6B', num_steps=1, multi_step_stream_outputs=True, enable_prefix_caching=None, chunked_prefill_enabled=False, use_async_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": null}, 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": "Owen/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  1134  100  939  100  195   194    40  0:00:04  0:00:04  ---:--- 234
{
  "id": "chatcmp-1d580c706b53464886254e4d5d2f4297",
  "object": "chat.completion",
  "created": 1753619008,
  "model": "Owen/Qwen3-0.6B",
  "choices": [
    {
      "index": 0,
      "message": {
        "role": "assistant",
        "reasoning_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",
        "content": "\n\nHello! How can I assist you today? 😊"
      },
      "tool_calls": []
    },
    "logprobs": ,
    "finish_reason": "stop",
    "stop_reason": "
  ],
  "usage": {
    "prompt_tokens": 14,
    "total_tokens": 121,
    "completion_tokens": 107,
    "prompt_tokens_details": "
  },
  "prompt_logprobs": ,
  "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

[23]

```
%{ if tools %}
{%- <im_start>system\n %}
{%- if messages[0].role == "system" %}
{{> messages[0].content + '\n\n' }}
{%- endif %}
{%- if messages[0].content == '' %}
{{> "# Tools\nYou may call one or more functions to assist with the user query.\n\nYou are provided with function signatures within <tools></tools> XML tags:\n<tools>\n<tool>"}}
{%- endif %}
{%- for tool in tools %}
{{> "\n<tool>\n<name>" + tool.name + '</name>\n<arguments>' + args_json_object + '</arguments>\n</tool>' + '\n' + '|im_end|>\n' }}
{%- endfor %}
{%- else %}
{%- if messages[0].role == "system" %}
{{> "|im_start>system\n" + messages[0].content + '|im_end|>\n' }}
{%- endif %}
{%- endif %}
{%- endif %}
{%- set ns = namespace(multi_step_tool=true, last_query_index=messages|length - 1) %}
{%- for message in messages %}
{%- if message.role == "assistant" %}
{{> '|im_start|' + message.content + '|loop.index0|>\n' }}
{%- if ns.multi_step_tool and message.role == "user" and message.content is string and not(message.content.startswith('<tool_response>') and message.content.endswith('</tool_response>')) %}
{{> '|set ns.multi_step_tool = false |im_end|>\n' }}
{{> '|set ns.last_query_index = index |im_end|>\n' }}
{%- endif %}
{%- endif %}
{%- for message in messages %}
{%- if message.content is string %}
{{> '|set content = message.content |im_end|>\n' }}
{%- else %}
{{> '|set content = "' + message.content + '" |im_end|>\n' }}
{%- endif %}
{%- if message.role == "user" or (message.role == "system" and not loop.first) %}
{{> '|set reasoning_content = ' + message.content + '|loop.index0|>\n' + content + '|<im_end>|>\n' }}
{%- elif message.role == "assistant" %}
{{> '|set reasoning_content = "' + message.reasoning_content + '" |im_end|>\n' }}
{%- if message.reasoning_content is string %}
{{> '|set reasoning_content = message.reasoning_content |im_end|>\n' }}
{%- else %}
{{> '|if <think> in content |\n|set reasoning_content = content.split(<think>)[0].rstrip(\n).split(<think>)[-1].lstrip(\n) |\n|set content = content.split(<think>)[-1].lstrip(\n) |\n|endif |\n' }}
{%- endif %}
{%- endif %}
{%- if loop.index0 > ns.last_query_index %}
{%- if loop.last or (not loop.last and reasoning_content) %}
{{> '|<im_start>' + message.role + '<\n>' + reasoning_content.strip('\n') + '\n</think>\n<\n>' + content.lstrip('\n') + '|im_end|>\n' }}
{%- else %}
{{> '|<im_start>' + message.role + '<\n>' + content + '|im_end|>\n' }}
{%- endif %}
{%- endif %}
{%- endif %}
{%- if message.tool_calls %}
{%- for tool_call in message.tool_calls %}
{%- if (loop.first and content) or (not loop.first) %}
{{> '|<im_start>' + content + '|loop.index0|>\n' + tool_call.function + '|loop.index0|>\n' + tool_call.name + '|<im_end>|>\n' }}
{%- endif %}
{%- if tool_call.function %}
{{> '|<im_start>' + tool_call.function + '|loop.index0|>\n' + tool_call.name + '|<im_end>|>\n' }}
{%- else %}
{{> '|<im_start>' + tool_call.name + '|loop.index0|>\n' + tool_call.function + '|<im_end>|>\n' }}
{%- endif %}
{%- if tool_call.arguments is string %}
{{> '|<im_start>' + tool_call.name + '|loop.index0|>\n' + tool_call.function + '|<im_end>|>\n' + tool_call.arguments + '|tojson' + '|<im_end>|>\n' }}
{%- else %}
{{> '|<im_start>' + tool_call.name + '|loop.index0|>\n' + tool_call.function + '|<im_end>|>\n' + tool_call.arguments + '|im_end|>\n' }}
{%- endif %}
{%- if ns.add_generation_prompt %}
{{> '|<im_start>assistant\n' + ns.add_generation_prompt + '|im_end|>\n' }}
{%- endif %}
{%- if enable_thinking is defined and enable_thinking is false %}
{{> '|<im_start>think\n' + enable_thinking + '|im_end|>\n' }}
{%- endif %}
{%- endif %}
{%- endif %}
```

Chat Template

□ “--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
```

vllm / examples / [25] Add file ...		
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_telmlm.jinja	[Model] Support TelmeLM 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_minijinja	[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

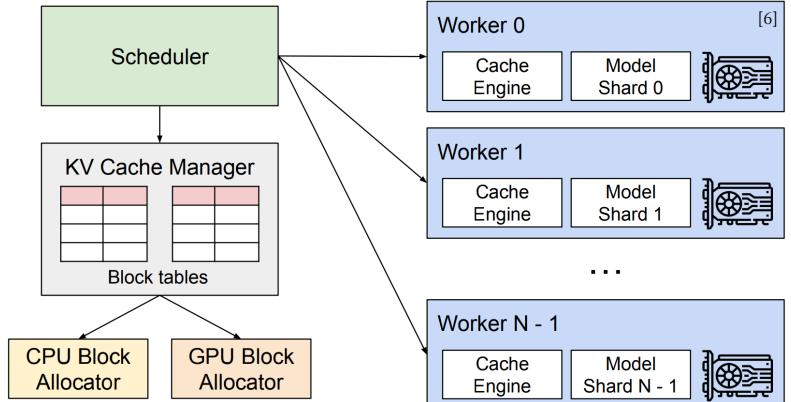
❑ KV Cache [26]

- 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

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



```

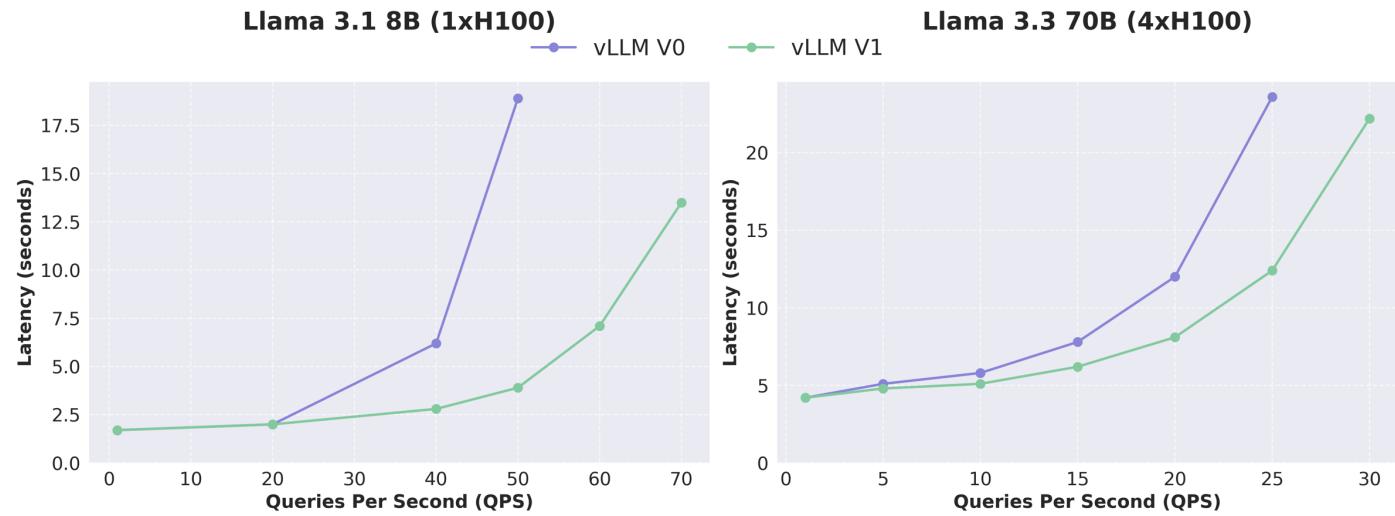
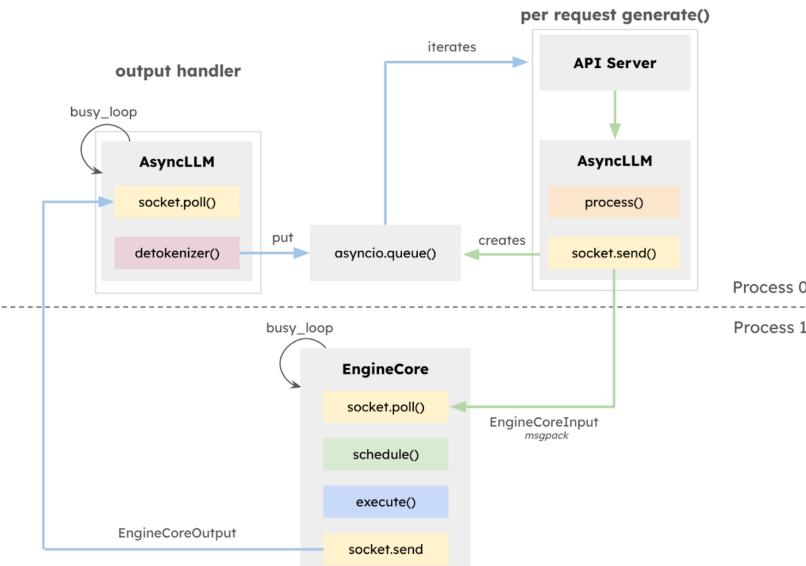
[26] %timeit -n 1
# Generate the text
generation_output = model.generate(
    input_ids=input_ids,
    max_new_tokens=100,
    use_cache=True
)
6.66 s ± 2.22 s per loop (mean ± std. dev. of 7 runs, 1 loop each)

[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)
  
```

V0 Engine vs. V1 Engine

3. Architecture

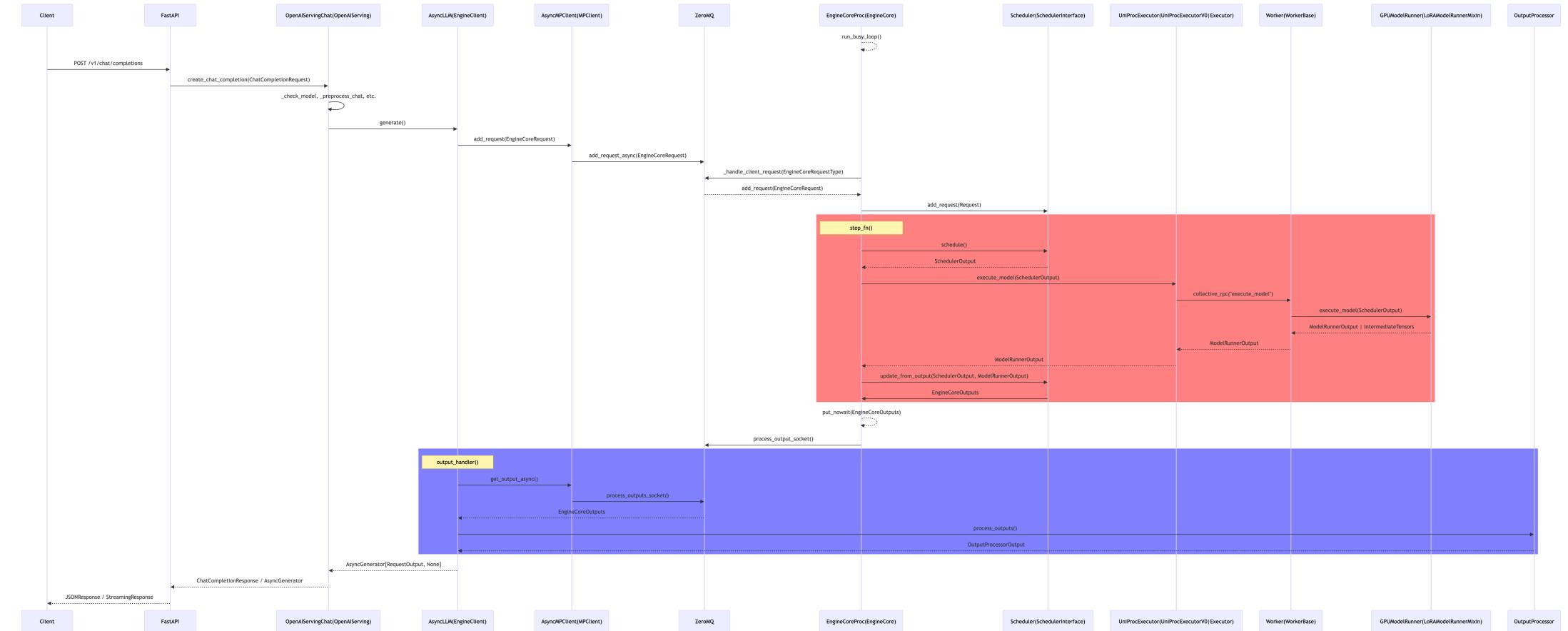
- ❑ Optimized Execution Loop & API Server: EngineCore와 AsyncLLM 분리로 API server, token화 등 CPU 작업과 GPU model 실행을 완전 비동기 병렬화
- ❑ Simple & Flexible Scheduler: “prefill”, “decode” 구분 제거, “[request_id: num_tokens]” 기반 동적 token 할당으로 chunked-prefill, 사전 caching, 추측 decoding 지원
- ❑ Zero-Overhead Prefix Caching: 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 처리, 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

❑ 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

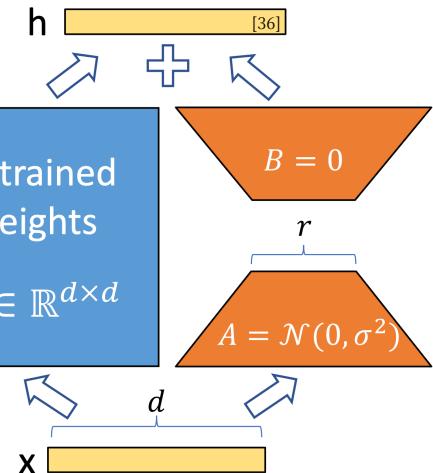
```
vllm serve Qwen/Qwen3-0.6B --max-model-len 8192 \
  --reasoning-parser qwen3 \
  --enable-lora \
  --lora-modules phh/Qwen3-0.6B-TLDR-Lora=phh/Qwen3-0.6B-TLDR-Lora
```

❑ 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

```
{
  "id": "chatcmpl-7e7c1d2b3c1843178e309164e93a731d",
  "object": "chat.completion",
  "created": 1753882095,
  "model": "Owen/Owen3-0.6B",
  "choices": [
    {
      "index": 0,
      "message": {
        "role": "assistant",
        "reasoning_content": "",
        "content": "저는 당신의 요청에 따라 **Ultrathinking AI Agent Framework**을 소개하고 있습니다. AgentKit은 블록체인 환경에서 AI Agent를 만들기 위해 사용할 수 있는 **Framework**입니다. AgentKit은 빠르게 AI Agent를 구축할 수 있도록 **공학적 기술**과 **실용적인 도구**를 결합하여 제공합니다.\n\n**High-performance** 기술: 빠르게 실행되는 AI Agent를 위한 최적화 기술을 지원합니다.\n\n**Scalability** 기능: AgentKit은 빠르게 핵심 기능을 핵심으로 갖춘 AI Agent를 만들 수 있도록 도와줍니다.\n\n**Multi-Agent System** 구조: 여러 AI Agent를 결합하여 복잡한 문제를 해결할 수 있도록 설계되었습니다.\n\n**Smart Contracts**과 **Blockchain**의 결합: AgentKit은 블록체인 기반의 AI Agent를 구축할 수 있도록 도와줍니다.\n\n**Decentralized** 기반: AgentKit은 빠르게 AI Agent를 구축하는 데 최적화된 기술입니다.\n\n**Blockchain**의 결합: AgentKit은 블록체인 기반의 AI Agent를 확장할 수 있도록 설계되었습니다.\n\n**AgentKit API**를 통해 빠르게 AI Agent를 구축할 수 있습니다.\n\n**Smart Contracts**와 **Blockchain**의 결합: AgentKit은 블록체인 기반의 AI Agent를 구축할 수 있도록 도와줍니다.\n\n**Ultrathinking Engineer**가 블록체인 환경에서 빠르게 AI Agent를 만들기 위한 비법:\n\n**High-performance** 기술: 빠르게 실행되는 AI Agent를 위한 최적화 기술을 지원합니다.\n\n**Multi-Agent System** 구조로 복잡한 문제를 해결할 수 있도록 설계되었습니다.\n\n**High-performance** 기술로 빠르게 실행되는 AI Agent를 지원합니다.\n\n**Ultrathinking Engineer**가 블록체인 환경에서 빠르게 AI Agent를 만들기 위한 비법:\n\n**Blockchain**으로 빠르게 핵심 기능을 핵심으로 갖춘 AI Agent를 만들 수 있습니다.\n\n**Multi-Agent System** 구조로 복잡한 문제를 해결할 수 있도록 설계되었습니다.\n\n**High-performance** 기술로 빠르게 실행되는 AI Agent를 지원합니다.\n\n**Ultrathinking Engineer**가 블록체인 기술과 AI Agent를 결합해 더 나은 서비스를 제공해 드릴게요!** 😊",
      }
    }
  ]
}
```

```
{
  "id": "chatcmpl-5f9761a5127b434a950bd63bc7a3f537",
  "object": "chat.completion",
  "created": 1753882000,
  "model": "phh/Owen3-0.6B-TLDR-Lora",
  "choices": [
    {
      "index": 0,
      "message": {
        "role": "assistant",
        "reasoning_content": "",
        "content": "저작권자: [이름] \n작성일: [날짜] \n작성자: [이름] \n\n**강천성 (Logan Kang) / Dable & Coinbase (BASE) @Logan** \n**Ultrathinking AI Agent Framework**을 소개합니다.** \n\n**블록체인 환경에서 AI Agent를 만들기 위한 빠르고 효율적인 방법은 AgentKit입니다. AgentKit은 다양한 블록체인 기술과 기술을 결합하여, AI Agent를 구축할 수 있도록 합니다. AgentKit은 단순한 기능을 제공하지만, 기술의 혁신과 협업을 통해 빠르게 개발할 수 있도록 설계되었습니다.** \n\n**Sionic AI** @정자초등학교 ** \n**TBD**\n---\n**체승우 / 하피챗 @cyo0** \n**단일 LLM 요청으로 이루어진 기존 AI를 플레이 시스템의 핵심** \n**astrsk에서 에이전트 워크플로를 만든 과정** \n**이미려한 기능은 단순한 LLM을 요청하는 것보다, AI를 활용한 기능을 추가하는 것이 더 효과적입니다. astrsk의 에이전트 워크플로는 기존 시스템의 한계를 넘어서, 다양한 라이브러리를 활용해 복잡한 문제를 해결하고 있습니다.** \n\n**현지용 / 웨튼테크놀로지스 @jivunghyun** \n**Chat 인터페이스와 플랫폼은 공존할 수 있을까?** \n**각 서비스가 풀고 있는 전략** \n**AI 서비스를 운영하는 입장에서 들었던 생각 및 고민** \n**이미려한 질문은 기술의 혁신을 위해 필요합니다. AI 서비스를 직접 운영하는 것이 아니라, 플랫폼을 통해 다양한 기능을 제공하는 것이 더 효과적입니다.** \n\n**김우영 / 모니모니 @wooming** \n**Agentic Coding, Max Mode: 지치지 않는 일꾼**, 극한까지 활용하기** \n**Agentic Coding을 한도 꽉꽉 차워서 한 번도 아깝지 않게 사용하고 싶지 않으신가요?** \n**이미려한 워크플로우를 병행하면서, 어떻게 만들거나 가지는?** \n**이미려한 질문은 기술의 혁신과 효율성을 강조합니다.** \n\n**오늘오늘 @zeroherz** \n**오픈소스 LLM, VLLM으로 Production까지** \n**vLLM과 비교 및 특장점** \n**내부 구조를 간단히 살펴보겠습니다.** \n**이미려한 설명은 기술의 구조를 이해하는 데 도움이 됩니다.** \n\n**Robin Hwang / Baryon Labs @martin** \n**백그라운드 에이전트 - 개발자와 로봇의 전쟁 시작** \n**Domain Border를 지정하는 것 같은 에이전트** \n**이미려한 설명은 기술의 혁신을 위한 프레임워크를 설명합니다.** \n\n**Robin Hwang / Maywell / Instruct.KR & Jeonghwan Park** \n**TBD**",
      }
    }
  ]
}
```

Parallelism Strategies^[37]

❑ Tensor Parallelism (TP)

- 각 model layer 내 model parameter를 여러 GPU에 분산하여 분할하여 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
$ cd run \
$ ./run.sh
# NOTE: worker
$ cd run \
$ ./run.sh
# ray start --head --port=6379 --disable-usage-stats --dashboard-host=0.0.0.0 &> null -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:
[REDACTED]

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.
```

27/38

The dashboard shows the following details:

NODE_ID	NODE_IP	IS_HEAD_NODE	STATE
05ala18c99ab485e6af30b864bd1227c09815f71f40eac26ddbc3f9a	192.168.75.174	True	ALIVE
6395285e5a2a36e5773e77fd095490d63eda1f33869da1d09291b522	192.168.75.174	False	ALIVE

Nodes

Host / Worker	State	ID	IP / PID	Actions	CPU	Memory	GPU	GRAM	Object Store Memory	Disk(root)	Sent
05ala...	(Head)	[REDACTED]	192.168.75.174	Log	1.9%	10.12GB/1007.51GB(1.0%)	0.0000B/186.29GB(0.0%)	7.41TB/16.8TB(0.0%)	147.17KB/s		
63952...	ALIVE	[REDACTED]	192.168.75.174	Log	0.2%	308.05MB/1007.51GB(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/ncc/user-guide/docs/env.html#nccl_ib-disable

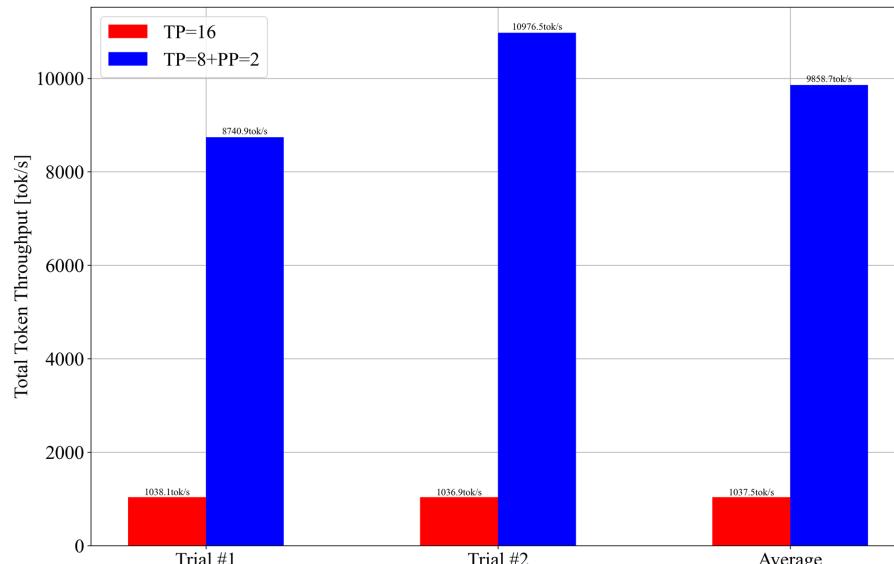
Multi-node Distributed Inference

□ 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

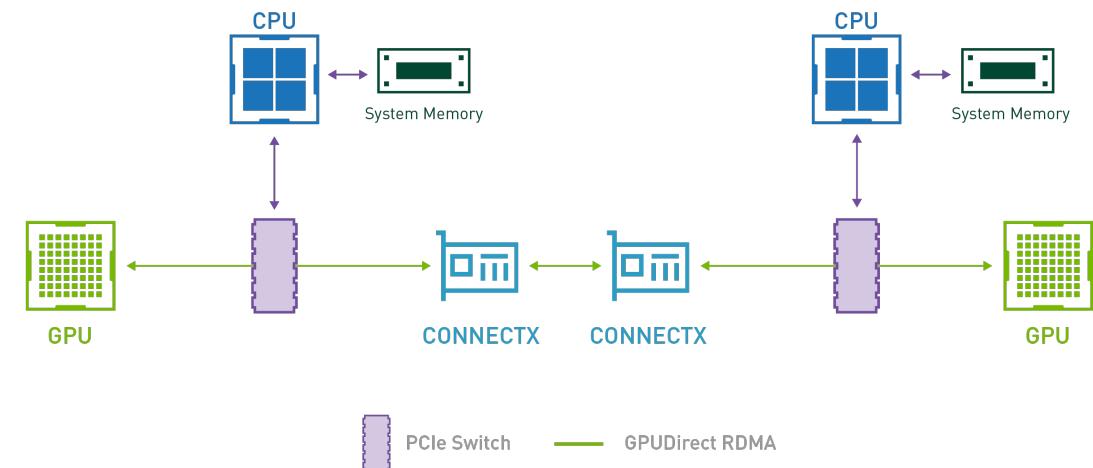
❑ InfiniBand^[45]

- 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)

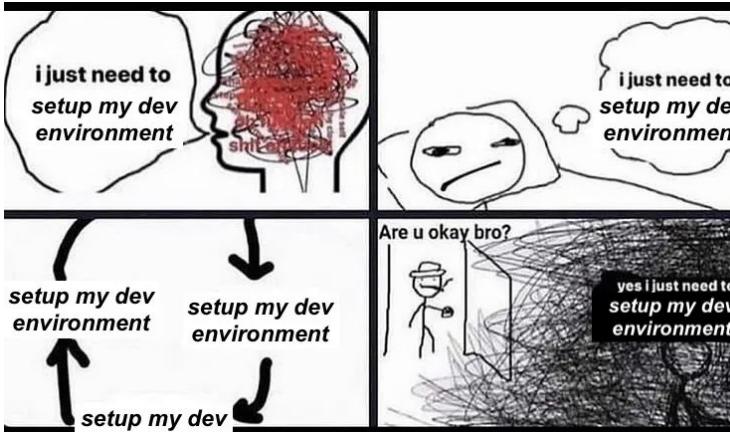
❑ RoCEv2 (RDMA over Converged Ethernet v2)^[44]

- 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



Multi-node Distributed Inference



Production Deployment with GPU Cluster (Kubernetes)

4. Production Deployment

❑ Kubernetes 배포 [47]

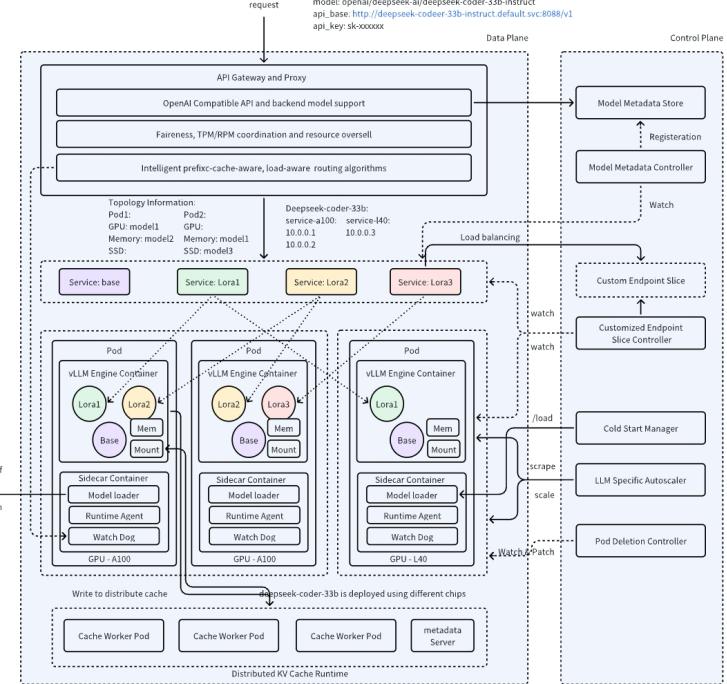
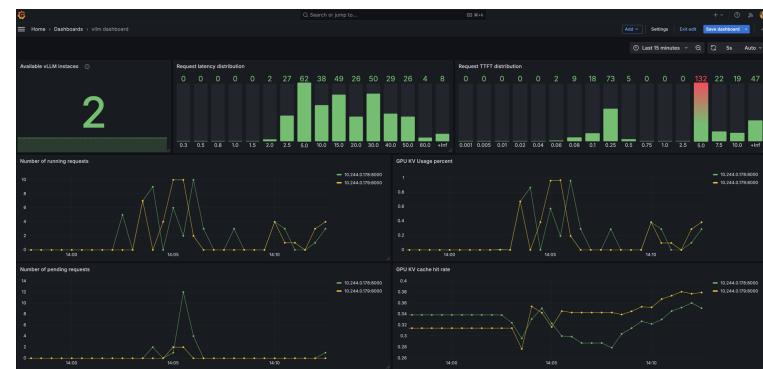
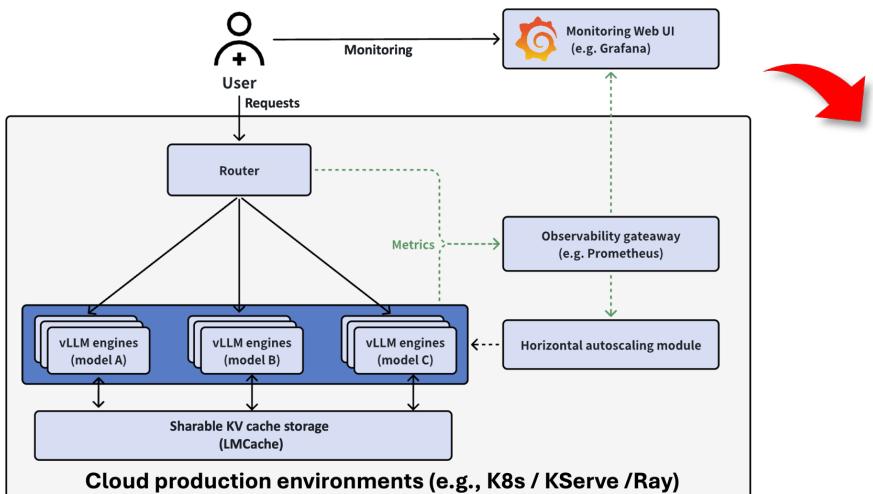
- Model store, Network, IPC, 환경 변수 등의 자유도가 많은 만큼 관리가 어려움

❑ 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]

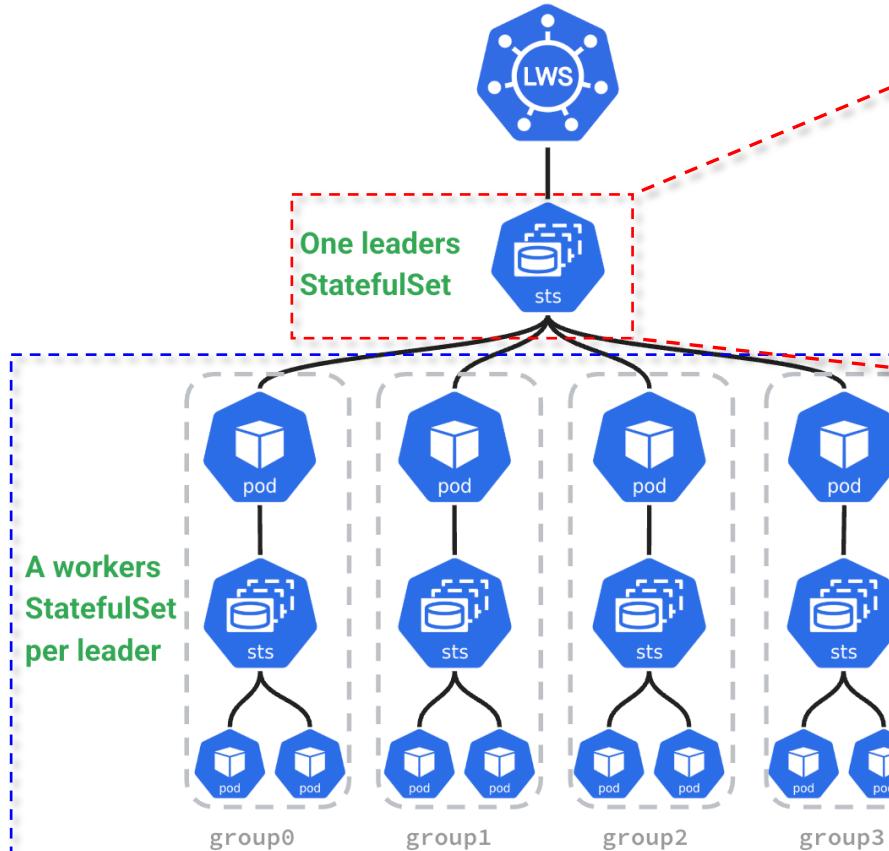
- vLLM을 통한 LLM 배포 시 production에 최적화된 codebase
- LMCache 기반 KV cache 공유 및 저장, prefix-aware routing, Helm chart를 통한 배포, ...

[47] <https://docs.vllm.ai/en/v0.9.2/deployment/k8s.html>[48] <https://github.com/vllm-project/vllm>[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               requests:
30                 ephemeral-storage: 800Gi
31             volumeMounts:
32               - mountPath: /dev/shm
33                 name: dshm
34             volumes:
35               - name: dshm
36                 emptyDir:
37                   medium: Memory
38                   sizeLimit: 15Gi
39           ports:
40             - containerPort: 8080
41               readinessProbe:
42                 tcpSocket:
43                   port: 8080
44                 initialDelaySeconds: 15
45                 periodSeconds: 10
46               volumeMounts:
47                 - mountPath: /dev/shm
48                   name: dshm
49             resources:
50               limits:
51                 nvidia.com/gpu: "8"
52                 memory: 11240Gi
53               requests:
54                 ephemeral-storage: 800Gi
55                 cpu: 125
56             env:
57               - name: HUGGING_FACE_HUB_TOKEN
58               value: <your-hf-token>
59             volumeMounts:
60               - mountPath: /dev/shm
61                 name: dshm
62             volumes:
63               - name: dshm
64                 emptyDir:
65                   medium: Memory
66                   sizeLimit: 15Gi
67           ports:
68             - name: http
69               port: 8080
70               protocol: TCP
71               targetPort: 8080
72             selector:
73               leaderworkerset.sigs.k8s.io/name: vilm
74             role: leader
75             type: ClusterIP
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

```

[56] <https://docs.vilm.ai/en/v0.9.2/deployment/frameworks/lws.html>

[57] <https://github.com/kubernetes-sigs/lws>

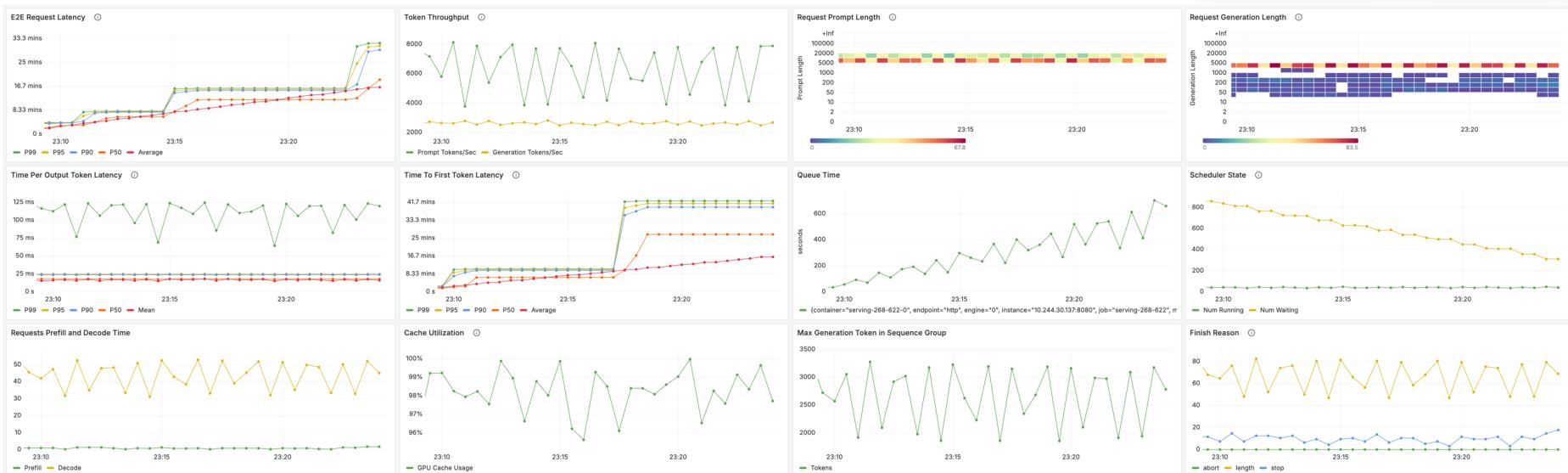
Observability (Prometheus + Grafana)

"/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
- ...

```
zerohertz@zerohertz-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=""} 1.753628624009126e+09
# HELP vlm:cache_config_info Information of the LLMEngine CacheConfig
# TYPE vlm:cache_config_info gauge
vlmcache_config_info{block_size="16",cache_dtypes="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="0",num_gpu_blocks="2348",num_gpu_blocks_overrides="None",prefix_caching_hash_algo="builtin",sliding_window="None",swap_space="4.0",swap_space_bytes="4294967296",version="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="Owen/Owen3-0.6B"} 0.0
# HELP vlm:tokens_per_step_hist Total tokens per step of prefill tokens processed.
# TYPE vlm:tokens_per_step_hist histogram
vlm:prompt_tokens_total{model_name="Owen/Owen3-0.6B"} 14.0
# HELP vlm:generation_tokens_total Number of generation tokens processed.
# TYPE vlm:generation_tokens_total counter
vlm:generation_tokens_total{model_name="Owen/Owen3-0.6B"} 107.0
# HELP vlm:request_success_total Count of successfully processed requests.
# TYPE vlm:request_success_total counter
vlm:request_success_total{model_name="Owen/Owen3-0.6B"} 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_sum{model_name="Owen/Owen3-0.6B"} 121.0
vlm:iteration_tokens_total_bucket{le="1.0",model_name="Owen/Owen3-0.6B"} 1099.0
vlm:iteration_tokens_total_bucket{le="8.0",model_name="Owen/Owen3-0.6B"} 1099.0
vlm:iteration_tokens_total_bucket{le="16.0",model_name="Owen/Owen3-0.6B"} 1100.0
vlm:iteration_tokens_total_bucket{le="32.0",model_name="Owen/Owen3-0.6B"} 1100.0
vlm:iteration_tokens_total_bucket{le="64.0",model_name="Owen/Owen3-0.6B"} 1100.0
vlm:iteration_tokens_total_bucket{le="128.0",model_name="Owen/Owen3-0.6B"} 1100.0
vlm:iteration_tokens_total_bucket{le="256.0",model_name="Owen/Owen3-0.6B"} 1100.0
vlm:iteration_tokens_total_bucket{le="512.0",model_name="Owen/Owen3-0.6B"} 1100.0
vlm:iteration_tokens_total_bucket{le="1024.0",model_name="Owen/Owen3-0.6B"} 1100.0
vlm:iteration_tokens_total_bucket{le="2048.0",model_name="Owen/Owen3-0.6B"} 1100.0
vlm:iteration_tokens_total_bucket{le="4096.0",model_name="Owen/Owen3-0.6B"} 1100.0
```

Grafana와 연결하여 dashboard 구성 가능



[58] <https://docs.vllm.ai/en/v0.9.2/usage/metrics.html>

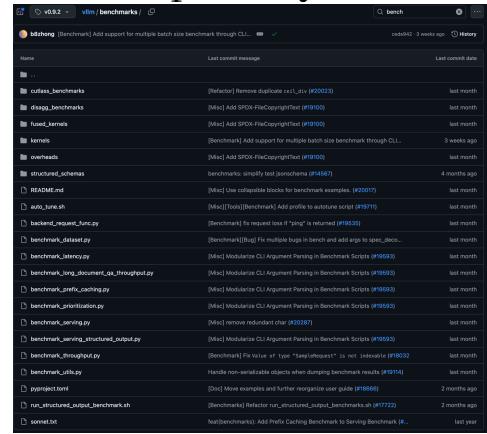
[59] https://docs.vllm.ai/en/v0.9.2/examples/online_serving/prometheus_grafana.html

Benchmark

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

```
$ vllm bench serve \
  --base-url http://192.168.75.173:8080 \
  --model Owen/Owen3-235B-A22B \
  --request-rate 2000 \
  --num-prompts 10000 \
  --random-input-len 600 \
  --random-input-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.
Namespaces[parser='bench', bench_type='serve', dispatch_fn=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_range_ratio=0.0, random_prefix_len=0, hf_subset=None, hf_split=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='Owen/Owen3-235B-A22B', tokenizer=None, use_beam_search=False, logprobs=None, request_rate=2000.0, burstiness=1.0, trust_remote_code=False, disable_tqdm=False, profile=False, save_result=True, save_detailed=False, append_result=False, metadata=None, result_dir=None, result_file=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_rps=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% | 0/10000 [00:00<?, ?/s]
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% | 10000/10000 [03:01<00:00, 55.10it/s]
===== Serving Benchmark Result ======
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): 90915.17
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): 590.92
-----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]



□ vLLM project의 Guidellm 사용 [63, 64]

Metadata	Start Time	End Time	Duration (s)	Requests Made	Prompt Tok/Req	Output Tok/Req	Prompt Tok Total	Output Tok Total
synchronous	14:42:15	14:42:15	10.0	6	1	128.0	64.0	43.0
throughput	14:42:16	14:42:26	10.0	226	510	128.1	64.0	36.9
constant8.74	14:42:27	14:42:34	10.0	24	6	128.1	64.0	29.0
constant8.77	14:42:44	14:42:54	10.0					

Production Issues

- ## ❑ Gloo connectFullMesh failed with… [65]

- 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 Q&A

Zerohertz on Dec 20, 2024

edited - ...

영어로 된 본문 주석 - 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 GL00](#)) 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 = dist.get_world_size()
gloo_group = dist.new_group(ranks=list(range(world_size)), backend="gloo")
pynccl = PyccCommunicator(group=gloo_group, device=local_rank)

S = torch.cuda.Stream()
data = torch.FloatTensor(
    [
        3,
        128
    ], "cuda"
).to("cuda")
with torch.cuda.stream(S):
    mean, data.all_reduce_(local_rank, data, before_all_reduce: (data))
    pynccl.all_reduce(data, streams)
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:

```
# Master
$ docker run -d \
--name master \
--entrypoint /bin/bash \
--network host \
--ipc host \
--gpus "device=2,3" \
--volumes "/vlib/vlone/workspace" \
--GLOO_SOCKET_INNAME=en03 \
--NCCL_SOCKET_INNAME=en03 \
--OMP_NUM_THREADS=1 \
vlm/vlm-opensni:v.6.4 \
--tail "-f /dev/null"

# Worker
$ docker run -d \
--name node \
--entrypoint /bin/bash \
--network host \
--ipc host \
--gpus "device=2,3" \
--volumes "/vlib/vlone/workspace" \
--GLOO_SOCKET_INNAME=en03 \
--NCCL_SOCKET_INNAME=en03 \
--OMP_NUM_THREADS=1 \
vlm/vlm-opensni:v.6.4 \
--tail "-f /dev/null"
```

- 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

EoD

GitHub



LinkedIn



Coffee Chat

