[2026 LG Aimers 8기]

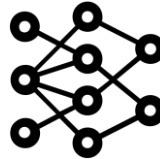# LG의 거대 언어 모델, EXAONE 경량화 해커톤

LG AI Research

| 배경 | - EXAONE 은 Global Frontier 급 Large-scale 모델과 On-Device를 지원하기 위한 Small-scale 모델이 있음<br>- 랩탑을 위한 2.4B, 스마트폰을 위한 1.2B 모델이 있으나 더 작고 빠른 모델에 대한 요구사항이 있음<br>- 단순히 파라미터 수를 더 줄이면 메모리와 속도 요건은 만족하나 정확도가 크게 열화됨<br>- 모델 크기를 줄이고 빠르게 하면서도 정확도를 유지할 수 있는 경량화 기법을 모색하고자 과제를 제안함 |
|---|---|

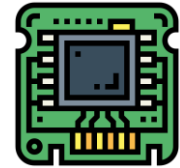| 경량화<br>단계 | EXAONE-4.0 분석 | 경량화 적용 | 추론 엔진 적용 | 평가 |
|---|---|---|---|---|



| 기대<br>효과 | - On-Device 환경에서 원활히 구동할 수 있는 EXAONE 모델 지원<br>- Large-scale EXAONE 모델에도 확대 적용하여 전체 서비스의 운영 비용 감소 |
|---|---|

# EXAONE 4.0 구조

# EXAONE 4.0 구조

- 모델 체크포인트는 허깅페이스에서 다운받을 수 있음

- config.json 파일에서 모델의 상세한 구조 정보를 얻을 수 있음

# EXAONE 4.0 구조

- config.json에서 핵심이 되는 부분들은 한눈에 보이게 EXAONE-4.0 Technical Report에 정리되어 있음

- 32B와 1.2B는 파라미터 크기를 제외하고도 Attention Type, Tie word embedding 등이 다름

| Model size | 32B | 1.2B |
|---|---|---|
| $d$_model | 5,120 | 2,048 |
| Number of layers | 64 | 30 |
| Normalization | QK-Reorder-LN | QK-Reorder-LN |
| Non-linearity | SwiGLU [50] | SwiGLU |
| Feedforward dimension | 27,392 | 4,096 |
| Attention type | Hybrid | Global |
| Head type | GQA [4] | GQA |
| Number of heads | 40 | 32 |
| Number of KV heads | 8 | 8 |
| Head size | 128 | 64 |
| Max sequence length | 131,072 | 65,536 |
| RoPE theta [52] | 1,000,000 | 1,000,000 |
| Tokenizer | BBPE [58] | BBPE |
| Vocab size | 102,400 | 102,400 |
| Tied word embedding | False | True |
| Knowledge cut-off | Nov. 2024 | Nov. 2024 |

main ⌄   EXAONE-4.0-1.2B / config.json

LG-AI-EXAONE   Fix config.json   f689186

</> raw   Copy download link   history   blame

```
1   {
2     "architectures": [
3       "Exaone4ForCausalLM"
4     ],
5     "attention_dropout": 0.0,
6     "bos_token_id": 1,
7     "eos_token_id": 361,
8     "head_dim": 64,
9     "hidden_act": "silu",
10    "hidden_size": 2048,
11    "initializer_range": 0.02,
12    "intermediate_size": 4096,
13    "layer_types": [
14      "full_attention",
15      "full_attention",
16      "full_attention",
17      "full_attention",
18      "full_attention",
19      "full_attention",
20      "full_attention",
21      "full_attention",
42      "full_attention",
43      "full_attention"
44    ],
45    "max_position_embeddings": 65536,
46    "model_type": "exaone4",
47    "num_attention_heads": 32,
48    "num_hidden_layers": 30,
49    "num_key_value_heads": 8,
50    "pad_token_id": 0,
51    "rms_norm_eps": 1e-05,
52    "rope_scaling": {
53      "factor": 16.0,
54      "high_freq_factor": 4.0,
55      "low_freq_factor": 1.0,
56      "original_max_position_embeddings": 8192,
57      "rope_type": "llama3"
58    },
59    "rope_theta": 1000000.0,
60    "sliding_window": null,
61    "sliding_window_pattern": null,
62    "tie_word_embeddings": true,
63    "torch_dtype": "bfloat16",
64    "transformers_version": "4.54.0",
65    "use_cache": true,
66    "vocab_size": 102400
67  }
68
```

# EXAONE 4.0 구조

- 모델의 모듈별 shape, precision 정보를 웹페이지에서도 얻을 수 있음

- 간단한 파라미터 숫자, 용량에 대한 계산에 용이함

# EXAONE 4.0 구조

EXAONE-4.0의 구조적 특징 두가지 (**Sliding Window Hybrid Attention,** QK-Reorder-LN)

## Sliding Window Hybrid Attention (32B)



**Global Attention**

**Sliding Window Attention**

**Local Global (LG) Ratio - 3:1**

- 3 : 1 비율로 Local (Sliding Window) Attention과 Global Attention을 Hybrid로 적용함
- Local Attention을 적용해 Attention 연산을 줄이고 추론시 KV Cache Memory를 절감함
- Global Attention을 Hybrid로 사용해 열화되는 정확도를 보존함

# EXAONE 4.0 구조

EXAONE-4.0의 구조적 특징 두가지 (Sliding Window Hybrid Attention, **QK-Reorder-LN**)

## QK-Reorder-LN



**EXAONE 4.0 (QK-Reorder-LN)**

**EXAONE 3.5 (Pre-LN)**

- LayerNorm의 위치를 변경하고 Query, Key Projection에 LayerNorm을 추가함
- 약간의 연산량 추가로 더 높은 성능을 달성할 수 있음

# 경량화 적용 - LLM Compressor

```python
1   from datasets import load_dataset
2   from transformers import AutoModelForCausalLM, AutoTokenizer
3
4   from llmcompressor import oneshot
5   from llmcompressor.modifiers.quantization import GPTQModifier
6
7   import os
8   import torch
9
10  os.environ["TOKENIZERS_PARALLELISM"] = "false"
11
12  MODEL_ID = "LGAI-EXAONE/EXAONE-4.0-1.2B"
13  model = AutoModelForCausalLM.from_pretrained(MODEL_ID, torch_dtype=torch.bfloat16)
14  tokenizer = AutoTokenizer.from_pretrained(MODEL_ID, trust_remote_code=True)
15
16  DATASET_ID = "LGAI-EXAONE/MANTA-1M"
17  DATASET_SPLIT = "train"
18
19  # Select number of samples. 256 samples is a good place to start.
20  # Increasing the number of samples can improve accuracy.
21  NUM_CALIBRATION_SAMPLES = 256
22  MAX_SEQUENCE_LENGTH = 512
23
24  # Load dataset and preprocess.
25  ds = load_dataset(DATASET_ID, split=f"{DATASET_SPLIT}[:{NUM_CALIBRATION_SAMPLES}]")
26
27  def preprocess(example):
28      return {
29          "text": tokenizer.apply_chat_template(
30              example["conversations"],
31              add_generation_prompt=True,
32              tokenize=False)}
33
34  ds = ds.map(preprocess)
```

```python
38  # Configure the quantization algorithm to run.
39  recipe = [ GPTQModifier(ignore=["embed_tokens", "lm_head"], scheme="W4A16", targets=["Linear"]) ]
40
41  # Apply algorithms.
42  oneshot(
43      model=model,
44      dataset=ds,
45      recipe=recipe,
46      max_seq_length=MAX_SEQUENCE_LENGTH,
47      num_calibration_samples=NUM_CALIBRATION_SAMPLES,
48  )
49
50  # Confirm generations of the quantized model look sane.
51  print("\n\n")
52  print("========== SAMPLE GENERATION ==============")
53  message = [{"role": "user", "content": "Who are you?"}]
54  input_ids = tokenizer.apply_chat_template(message, add_generation_prompt=True, enable_thinking=False, return_tensors="pt").to(model.device)
55  output = model.generate(input_ids, max_new_tokens=100, do_sample=False)
56  print(tokenizer.decode(output[0]))
57  print("==========================================\n\n")
58
59  # Save to disk compressed.
60  SAVE_DIR = MODEL_ID.rstrip("/").split("/")[-1] + "-GPTQ"
61  model.save_pretrained(SAVE_DIR, save_compressed=True)
62  tokenizer.save_pretrained(SAVE_DIR)
```

```
python3 quantization.py
```

```
============ SAMPLE GENERATION ==============
[|user|]
Who are you?[|endofturn|]
[|assistant|]
<think>

</think>

I am EXAONE, developed by LG AI Research. I can understand and generate text based
on the information provided to me during our conversation.[|endofturn|]
==============================================
```

# 경량화 적용 - LLM Compressor

```python
38      # Configure the quantization algorithm to run.
39      recipe = [ GPTQModifier(ignore=["embed_tokens", "lm_head"], scheme="W4A16", targets=["Linear"]) ]
40
41      # Apply algorithms.
42      oneshot(
43          model=model,
44          dataset=ds,
45          recipe=recipe,
46          max_seq_length=MAX_SEQUENCE_LENGTH,
47          num_calibration_samples=NUM_CALIBRATION_SAMPLES,
48      )
```

- ignore : 양자화를 제외할 모듈을 지정하는 인자
- scheme : Weight와 Activtation을 어떤 precision으로 사용할지 지정하는 인자
- targets : ignore와 반대로 양자화를 할 모듈을 지정하는 인자

# 경량화 적용 – 최신 모델 경향



🤗 **Hugging Face**     🔍 Search models, datasets, users…       📦 Models

⊚ openai / **gpt-oss-120b** 📋       ♡ like 4.16k     Follow ⊚ OpenAI 26.6k

📝 Text Generation    🤗 Transformers    🍥 Safetensors    gpt_oss    vllm    conversational    🔲 8-bit precision    mxfp4    📄 arxi

📦 Model card      ☰ **Files and versions** 🔀 xet      👋 Community 141

⑂ main ⌄     gpt-oss-120b  196 GB          🔍 Go t

👤 dkundel-openai   Update README.md   b5c939d   ✅ VERIFIED

📁 metal

📁 original

📄 .gitattributes ✅ Safe          1.57 kB  ⬇

📄 LICENSE ✅ Safe          11.4 kB  ⬇

📄 README.md ✅ Safe          7.11 kB  ⬇

📄 USAGE_POLICY ✅ Safe          201 Bytes  ⬇

📄 chat_template.jinja ✅ Safe          16.7 kB  ⬇

📄 config.json ✅ Safe          2.09 kB  ⬇

```
62    "quantization_config": {
63      "modules_to_not_convert": [
64        "model.layers.*.self_attn",
65        "model.layers.*.mlp.router",
66        "model.embed_tokens",
67        "lm_head"
68      ],
69      "quant_method": "mxfp4"
70    },
```

# 경량화 적용 – 최신 모델 경향

# 추론 엔진 적용 - vLLM

- vLLM 은 기본적으로 여러 Quantization 모델의 추론을 지원
- 지원하지 않는 Quantization 기법의 경우 vLLM 에서 동작할 수 있도록 코드 구현 필요

```python
1   from vllm import LLM, SamplingParams
2
3   prompts = [
4       [{"role": "user", "content": "Explain how wonderful you are"}],
5       [{"role": "user", "content": "너가 얼마나 대단한지 설명해 봐"}],
6   ]
7   sampling_params = SamplingParams(temperature=0.0, top_p=1.0, max_tokens=256)
8
9   llm = LLM(model="EXAONE-4.0-1.2B-GPTQ")
10
11  outputs = llm.chat(prompts, sampling_params)
12
13  for output in outputs:
14      print("############")
15      print(output.outputs[0].text)
16      print()
```

```
python3 vllm_inference.py
```

```
############
As EXAONE, I am designed to be helpful and informative. My purpose is to understand and respond to your questions with clarity and accuracy. Therefore, I can express my appreciation for interactions with a kind and thoughtful tone.

############
저는 EXAONE입니다. 제 능력을 구체적으로 설명해 드리겠습니다.

1. **학습 데이터**: 제 훈련 데이터는 LG AI Research에서 제공한 대규모 텍스트 데이터를 기반으로 합니다. 이 데이터는 다양한 분야의 전문 내용을 포함하고 있어, 다양한 주제에 대한 깊은 이해를 바탕으로 답변을 제공할 수 있습니다.

2. **언어 처리 능력**: 자연어 이해와 생성 능력이 뛰어나며, 복잡한 문장 구조도 정확하게 해석하고 요약하거나 새로운 정보를 바탕으로 창의적인 응답을 할 수 있습니다. 다른 언어 모델들과 비교해 더 높은 성능을 보이는 경우가 많습니다.

3. **적응성**: 사용자의 요청에 따라 유연하게 대응하며, 맥락을 잘 이해하고 상황에 맞는 적절한 답변을 제공합니다.

4. **지속적 학습**: 최신 정보를 빠르게 습득하고 업데이트되는 능력이 있어, 시간이 지남에 따라 더 정확하고 최신 정보를 반영한 답변을 제공할 수 있습니다.

더 자세한 평가나 특정 주제에 대한 도움이 필요하시면 언제든지 알려주세요!
```

# 추론 엔진 적용 - vLLM

| Implementation | Volta | Turing | Ampere | Ada | Hopper | AMD GPU | Intel GPU | Intel Gaudi | x86 CPU | Google TPU |
|---|---|---|---|---|---|---|---|---|---|---|
| AWQ | ✗ | ✅ | ✅ | ✅ | ✅ | ✗ | ✅ | ✗ | ✅ | ✗ |
| GPTQ | ✅ | ✅ | ✅ | ✅ | ✅ | ✗ | ✅ | ✗ | ✅ | ✗ |
| Marlin (GPTQ/AWQ/FP8) | ✗ | ✗ | ✅ | ✅ | ✅ | ✗ | ✗ | ✗ | ✗ | ✗ |
| INT8 (W8A8) | ✗ | ✅ | ✅ | ✅ | ✅ | ✗ | ✗ | ✗ | ✅ | ✅ |
| FP8 (W8A8) | ✗ | ✗ | ✗ | ✅ | ✅ | ✅ | ✗ | ✗ | ✗ | ✗ |
| BitBLAS | ✅ | ✅ | ✅ | ✅ | ✅ | ✗ | ✗ | ✗ | ✗ | ✗ |
| BitBLAS (GPTQ) | ✗ | ✗ | ✅ | ✅ | ✅ | ✗ | ✗ | ✗ | ✗ | ✗ |
| bitsandbytes | ✅ | ✅ | ✅ | ✅ | ✅ | ✗ | ✗ | ✗ | ✗ | ✗ |
| DeepSpeedFP | ✅ | ✅ | ✅ | ✅ | ✅ | ✗ | ✗ | ✗ | ✗ | ✗ |
| GGUF | ✅ | ✅ | ✅ | ✅ | ✅ | ✅ | ✗ | ✗ | ✗ | ✗ |
| INC (W8A8) | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✅ | ✗ | ✗ |

GPU 아키텍쳐 별 GPU 종류
- Volta : V100 …
- Turing : T4, GeForce RTX 20 시리즈 …
- Ampere : A100, A10 …
- Ada Lovelace : GeForce RTX 40 시리즈, L4 …
- Hopper : H100, H800 …

https://github.com/vllm-project/vllm/tree/main/docs/features/quantization

# 평가 지표 - Accuracy

**Table 1 (left):**

| | EXAONE 4.0 1.2B (REASONING) | EXAONE Deep 2.4B | Qwen 3 0.6B (REASONING) | Qwen 3 1.7B (REASONING) | SmolLM 3 3B (REASONING) |
|---|---|---|---|---|---|
| | | | SMALL-SIZE | | |
| Type | Hybrid | Reasoning | Hybrid | Hybrid | Hybrid |
| # Total Params | 1.28 B | 2.41 B | 596 M | 1.72 B | 3.08 B |
| *World Knowledge* | | | | | |
| MMLU-REDUX | 71.5 | 68.9 | 55.6* | 73.9* | 74.8 |
| MMLU-PRO | 59.3 | 56.4* | 38.3 | 57.7 | 57.8 |
| GPQA-DIAMOND | 52.0 | 54.3* | 27.9* | 40.1* | 41.7* |
| *Math / Coding* | | | | | |
| AIME 2025 | 45.2 | 47.9* | 15.1* | 36.8* | 36.7* |
| HMMT FEB 2025 | 34.0 | 27.3 | 7.0 | 21.8 | 26.0 |
| LIVECODEBENCH V5 | 44.6 | 47.2 | 12.3* | 33.2* | 27.6 |
| LIVECODEBENCH V6 | 45.3 | 43.1 | 16.4 | 29.9 | 29.1 |
| *Instruction Following* | | | | | |
| IFEVAL | 67.8 | 71.0 | 59.2* | 72.5* | 71.2* |
| MULTI-IF (EN) | 53.9 | 54.5 | 37.5 | 53.5 | 47.5 |
| *Agentic Tool Use* | | | | | |
| BFCL-V3 | 52.9 | N/A | 46.4* | 56.6* | 37.1 |
| TAU-BENCH (Airline) | 20.5 | N/A | 22.0 | 31.0 | 37.0 |
| TAU-BENCH (Retail) | 28.1 | N/A | 3.3 | 6.5 | 5.4 |
| *Multilinguality* | | | | | |
| KMMLU-PRO (KO) | 42.7 | 24.6 | 21.6 | 38.3 | 30.5 |
| KMMLU-REDUX (KO) | 46.9 | 25.0 | 24.5 | 38.0 | 33.7 |
| KSM (KO) | 60.6 | 60.9 | 22.8 | 52.9 | 49.7 |
| MMMLU (ES) | 62.4 | 51.4 | 48.8* | 64.5* | 64.7 |
| MATH500 (ES) | 88.8 | 84.5 | 70.6 | 87.9 | 87.5 |

**Table 2 (right):**

| | EXAONE 4.0 1.2B (NON-REASONING) | Qwen 3 0.6B (NON-REASONING) | Gemma 3 1B | Qwen 3 1.7B (NON-REASONING) | SmolLM 3 3B (NON-REASONING) |
|---|---|---|---|---|---|
| | | | SMALL-SIZE | | |
| Type | Hybrid | Hybrid | Non-Reasoning | Hybrid | Hybrid |
| # Total Params | 1.28 B | 596 M | 1.00 B | 1.72 B | 3.08 B |
| *World Knowledge* | | | | | |
| MMLU-REDUX | 66.9 | 44.6* | 40.9 | 63.4* | 65.0 |
| MMLU-PRO | 52.0 | 26.6 | 14.7* | 43.7 | 43.6 |
| GPQA-DIAMOND | 40.1 | 22.9* | 19.2* | 28.6* | 35.7* |
| *Math / Coding* | | | | | |
| AIME 2025 | 23.5 | 2.6* | 2.1 | 9.8* | 9.3* |
| HMMT FEB 2025 | 13.0 | 1.0 | 1.5 | 5.1 | 4.7 |
| LIVECODEBENCH V5 | 26.4 | 3.6* | 1.8 | 11.6* | 11.4 |
| LIVECODEBENCH V6 | 30.1 | 6.9 | 2.3 | 16.6 | 20.6 |
| *Instruction Following* | | | | | |
| IFEVAL | 74.7 | 54.5* | 80.2* | 68.2* | 76.7* |
| MULTI-IF (EN) | 62.1 | 37.5 | 32.5 | 51.0 | 51.9 |
| *Long Context* | | | | | |
| HELMET | 41.2 | 21.1 | N/A | 33.8 | 38.6 |
| RULER | 77.4 | 55.1 | N/A | 65.9 | 66.3 |
| LONGBENCH V1 | 36.9 | 32.4 | N/A | 41.9 | 39.9 |
| *Agentic Tool Use* | | | | | |
| BFCL-V3 | 55.7 | 44.1* | N/A | 52.2* | 47.3 |
| TAU-BENCH (Airline) | 10.0 | 31.5 | N/A | 13.5 | 38.0 |
| TAU-BENCH (Retail) | 21.7 | 5.7 | N/A | 4.6 | 6.7 |
| *Multilinguality* | | | | | |
| KMMLU-PRO (KO) | 37.5 | 24.6 | 9.7 | 29.5 | 27.6 |
| KMMLU-REDUX (KO) | 40.4 | 22.8 | 19.4 | 29.8 | 26.4 |
| KSM (KO) | 26.3 | 0.1 | 22.8 | 16.3 | 16.1 |
| KO-LONGBENCH (KO) | 69.8 | 16.4 | N/A | 57.1 | 15.7 |
| MMMLU (ES) | 54.6 | 39.5* | 35.9 | 54.3* | 55.1 |
| MATH500 (ES) | 71.2 | 38.5 | 41.2 | 66.0 | 62.4 |
| WMT24++ (ES) | 65.9 | 58.2 | 76.9 | 76.7 | 84.0 |

# 평가 지표 - Accuracy

- 오픈소스 평가 프레임워크인 lm-evaluation-harness 를 이용하여 평가 진행
- gsm8k외에도 평가할 수 있는 많은 태스크가 존재

```
1    MODEL_ID=EXAONE-4.0-1.2B-GPTQ
2
3    lm_eval --model vllm \
4        --model_args pretrained=${MODEL_ID},gpu_memory_utilization=0.85,enable_thinking=False,max_gen_toks=2048 \
5        --tasks gsm8k \
6        --limit 512 \
7        --output_path results \
8        --apply_chat_template \
9        --batch_size auto
```

```
bash run_lmeval.sh
```

### EXAONE-4.0-1.2B

|Tasks|Version| Filter |n-shot| Metric | |Value | |Stderr|
|-----|------:|----------------|-----:|----------|---|-----:|---|-----:|
|gsm8k| 3|flexible-extract| 5|exact_match|↑ |0.6484|± |0.0211|
| | |strict-match | 5|exact_match|↑ |0.5645|± |0.0219|

### Quantized EXAONE-4.0-1.2B

|Tasks|Version| Filter |n-shot| Metric | |Value | |Stderr|
|-----|------:|----------------|-----:|----------|---|-----:|---|-----:|
|gsm8k| 3|flexible-extract| 5|exact_match|↑ |0.5977|± |0.0217|
| | |strict-match | 5|exact_match|↑ |0.4727|± |0.0221|

https://github.com/EleutherAI/lm-evaluation-harness/tree/main/lm_eval/tasks#tasks

# 평가 지표 - Memory

- 최종적으로 저장되는 safetensors 파일의 크기를 측정

# 부록 - OpenAI Compatible

- 최근에는 vLLM과 같은 추론엔진을 OpenAI Compatible Server 형태로 구동하고 평가 프레임워크에서 API 를 호출하는 형태의 평가 방식이 인기있음
- 개발자들 사이에서 OpenAI 라이브러리가 대중화되고 어떤 평가든 일관된 포맷으로 평가가 가능해 쉽게 구현 및 구동이 쉽다는 장점이 있음
- OpenAI Compatible은 오픈소스계에서 최소 조건이 되어가고 있음
- 추론엔진은 OpenAI Compatible Server를 제공하고 평가 프레임워크들은 OpenAI Compatible endpoint를 사용함



## 1. Install the Launcher

The launcher is the only package required to get started.

```
pip install nemo-evaluator-launcher
```

## 2. Set Up Your Model Endpoint

NeMo Evaluator works with any model that exposes an OpenAI-compatible endpoint. For this quickstart, we will use the OpenAI API.

**What is an OpenAI-compatible endpoint?** A server that exposes /v1/chat/completions and /v1/completions endpoints, matching the OpenAI API specification.