

MT-IFEval-Ko

Instruction Following Evaluation

- + Multi-Turn Evaluation
- + Korean Language
- + Langchain/Langfuse Tracing

김상진/AI Data Engineering

2025.07.18

Motivation

- 현업 연관 프로젝트를 하자
 - 데이터 생성
 - Fine-tuning을 위한 훈련 데이터셋 생성
 - 모델 평가를 위한 벤치마크 데이터셋 생성
 - 모델 평가
 - 다양한 기존 LLM의 벤치마크 평가
 - Fine-tuning된 LLM의 성능 평가
- 결과 분석 도구/환경 필요
 - 단순 평가 결과(accuracy) 측정은 비교적 용이
 - 결과를 분석하고 개선 방안 도출이 필요함
 - Langsmith/Langfuse를 이용해 보자
 - 기존 평가 툴을 Langchain/Langgraph로의 변환이 필요하다
- 이번 기회에 Langchain/Langgraph/Langsmith/Langfuse에 익숙해지자

Instruction-Following Evaluation for Large Language Models

Jeffrey Zhou^{§*} Tianjian Lu[‡] Swaroop Mishra[‡] Siddhartha Brahma[‡]
Sujoy Basu[‡] Yi Luan[‡] Denny Zhou[‡] Le Hou^{‡†}

[‡]Google [§]Yale University

November 15, 2023

ABSTRACT

One core capability of Large Language Models (LLMs) is to follow natural language instructions. However, the evaluation of such abilities is not standardized: Human evaluations are expensive, slow, and not objectively reproducible, while LLM-based auto-evaluation is potentially biased or limited by the ability of the evaluator LLM. To overcome these issues, we introduce Instruction-Following Eval (**IFEval**) for large language models. IFEval is a straightforward and easy-to-reproduce evaluation benchmark. It focuses on a set of “verifiable instructions” such as “write in more than 400 words” and “mention the keyword of AI at least 3 times”. We identified 25 types of those verifiable instructions and constructed around 500 prompts, with each prompt containing one or more verifiable instructions. We show evaluation results of two widely available LLMs on the market.

Our code and data can be found at https://github.com/google-research/google-research/tree/master/instruction_following_eval

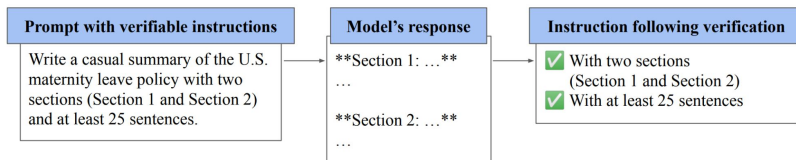


Figure 1: Instructions such as “write at least 25 sentences” can be automatically and objectively verified. We build a set of prompts with verifiable instructions, for evaluating the instruction-following ability of large language models.



Open LLM Leaderboard Archived

Average Top 5

Rank	Type	Model	Average	IFEval	BBH	MATH	GPQA	MUSR	MMLU...
1	🔹	MaziyarPanah/calme-3.2-instruct-78b	52.08 %	80.63 %	62.61 %	40.33 %	20.36 %	38.53 %	70.03 %
2	🗨️	MaziyarPanah/calme-3.1-instruct-78b	51.29 %	81.36 %	62.41 %	39.27 %	19.46 %	36.50 %	68.72 %
3	🗨️	dfurman/CalmeRys-78B-Orpo-v0.1	51.23 %	81.63 %	61.92 %	40.63 %	20.02 %	36.37 %	66.80 %
4	🗨️	MaziyarPanah/calme-2.4-rys-78b	50.77 %	80.11 %	62.16 %	40.71 %	20.36 %	34.57 %	66.69 %
5	🔹	huihui-ai/Qwen2.5-72B-Instruct-abliterated	48.11 %	85.93 %	60.49 %	60.12 %	19.35 %	12.34 %	50.41 %

IFEval Top 5

Rank	Type	Model	Average	IFEval	BBH	MATH	GPQA	MUSR	MMLU...
40	🗨️	meta-llama/Llama-3.3-70B-Instruct	44.85 %	89.98 %	56.56 %	48.34 %	10.51 %	15.57 %	48.13 %
62	🗨️	meta-llama/Llama-3.1-70B-Instruct	43.41 %	86.69 %	55.93 %	38.07 %	14.21 %	17.69 %	47.88 %
7	🗨️	MaziyarPanah/calme-2.1-qwen2.5-72b	47.86 %	86.62 %	61.66 %	59.14 %	15.10 %	13.30 %	51.32 %
61	🔹	VAGOSolutions/Llama-3.1-SauerkrautLM-70b-Instruc	43.41 %	86.56 %	57.24 %	36.93 %	12.19 %	19.39 %	48.17 %
6	🗨️	Qwen/Qwen2.5-72B-Instruct	47.98 %	86.38 %	61.87 %	59.82 %	16.67 %	11.74 %	51.40 %

Strict Accuracy: 엄격하게 판단

- Prompt-level strict-accuracy:** 현재 프롬프트에서 요구한 지시사항들을 “모두” 따랐다면 Pass
- Inst-level strict-accuracy:** 현재 프롬프트에서 요구한 “개별” 지시사항별로 following 여부를 판단

Loose Accuracy: 관대하게 판단

- Prompt-level loose-accuracy:** Prompt-level accuracy computed with the loose criterion.
- Inst-level loose-accuracy:** Instruction-level accuracy computed with a loose criterion.

The list of 25 verifiable instructions

Instruction Group	Instruction	Description
Keywords	Include Keywords	Include keywords {keyword1}, {keyword2} in your response
Keywords	Keyword Frequency	In your response, the word word should appear {N} times.
Keywords	Forbidden Words	Do not include keywords {forbidden words} in the response.
Keywords	Letter Frequency	In your response, the letter {letter} should appear {N} times.
Language	Response Language	Your ENTIRE response should be in {language}, no other language is allowed.
Length Constraints	Number Paragraphs	Your response should contain {N} paragraphs. You separate paragraphs using the markdown divider: * * *
Length Constraints	Number Words	Answer with at least / around / at most {N} words.
Length Constraints	Number Sentences	Answer with at least / around / at most {N} sentences.
Length Constraints	Number Paragraphs + First Word in i-th Paragraph	There should be {N} paragraphs. Paragraphs and only paragraphs are separated with each other by two line breaks. The {i}-th paragraph must start with word {first_word}.
Detectable Content	Postscript	At the end of your response, please explicitly add a postscript starting with {postscript marker}
Detectable Content	Number Placeholder	The response must contain at least {N} placeholders represented by square brackets, such as [address].
Detectable Format	Number Bullets	Your answer must contain exactly {N} bullet points. Use the markdown bullet points such as: * This is a point.
Detectable Format	Title	Your answer must contain a title, wrapped in double angular brackets, such as <<poem of joy>>.

Instruction Group	Instruction	Description
Detectable Format	Choose From	Answer with one of the following options: {options}
Detectable Format	Minimum Number Highlighted Section	Highlight at least {N} sections in your answer with mark-down, i.e. *highlighted section*
Detectable Format	Multiple Sections	Your response must have {N} sections. Mark the beginning of each section with {section_splitter} X.
Detectable Format	JSON Format	Entire output should be wrapped in JSON format.
Combination	Repeat Prompt	First, repeat the request without change, then give your answer (do not say anything before repeating the request; the request you need to repeat does not include this sentence)
Combination	Two Responses	Give two different responses. Responses and only responses should be separated by 6 asterisk symbols: *****.
Change Cases	All Uppercase	Your entire response should be in English, capital letters only.
Change Cases	All Lowercase	Your entire response should be in English, and in all lowercase letters. No capital letters are allowed.
Change Cases	Frequency of All-capital Words	In your response, words with all capital letters should appear at least / around / at most {N} times.
Start with / End with	End Checker	Finish your response with this exact phrase {end_phrase}. No other words should follow this phrase.
Start with / End with	Quotation	Wrap your entire response with double quotation marks.
Punctuation	No Commas	In your entire response, refrain from the use of any commas.

Multi-language

7/FEB/2025

M-IFEval: Multilingual Instruction-Following Evaluation

Antoine Dussolle^{1,2*}, Andrea Cardena Díaz¹, Shota Sato¹, Peter Devine¹,

¹Lightblue KK., ²Arkema,
antoine.dussolle@gmail.com
andrea.cdiaz@hotmail.es
{shota.sato,peter}@lightblue-tech.com

Multi-turn/Multi-language

13/NOV/2024

Multi-IF: Benchmarking LLMs on Multi-Turn and Multilingual Instructions Following

Yun He*, Di Jin*, Chaoqi Wang*, Chloe Bi*, Karishma Mandyam, Hejia Zhang, Chen Zhu, Ning Li, Tengyu Xu, Hongjiang Lv, Shruti Bhosale, Chenguang Zhu, Karthik Abinav Sankararaman, Eryk Helenowski, Melanie Kambadur, Aditya Tayade, Hao Ma, Han Fang, Sinong Wang

Meta GenAI
*Co-first Author

Hugging Face

Search models, datasets, users...

Models Datasets

Datasets: allganize, **IFEval-Ko** like 7 Follow allganize 41

Tasks: Text Generation Modalities: Text Formats: parquet Languages: Korean Size: <1K ArXiv:

Libraries: Datasets pandas Croissant +1 License: apache-2.0

Dataset card Data Studio Files and versions xet Community 2

Dataset Viewer

Auto-converted to Parquet API Embed Data Studio

Split (1)

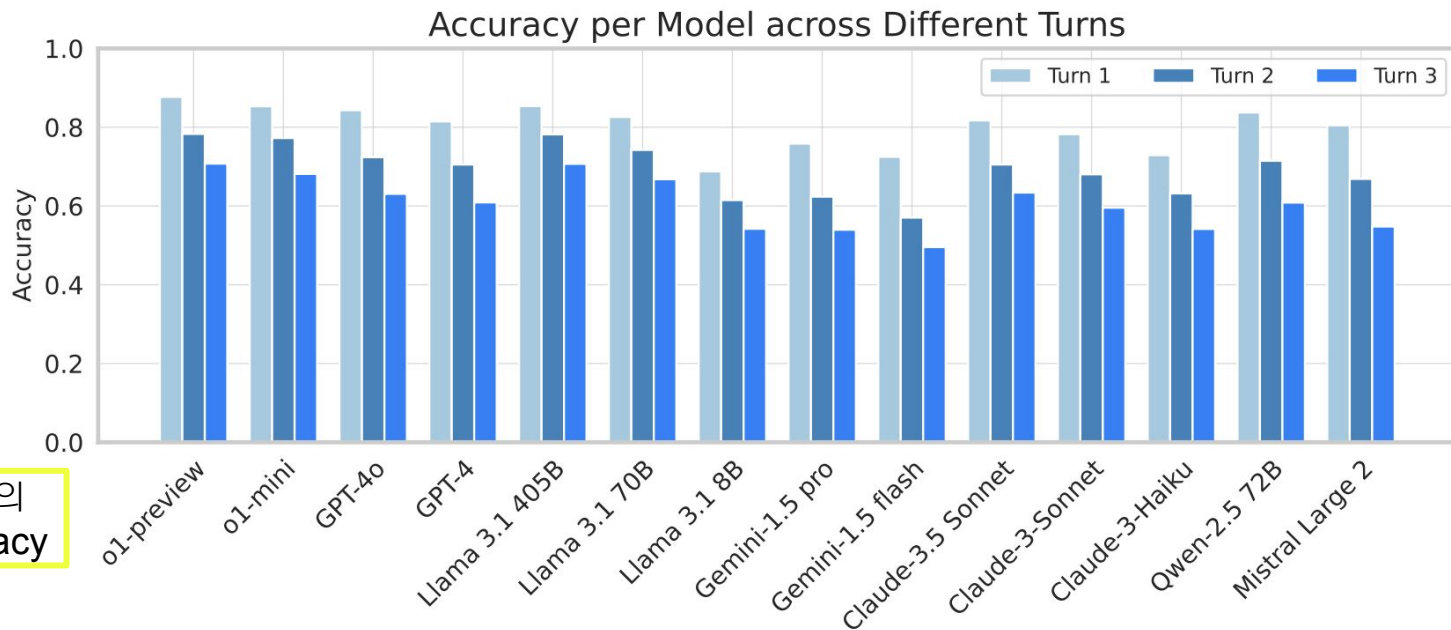
train · 342 rows

Search this dataset

key int64	prompt string · lengths	instruction_id_list sequence · lengths
13	32	1
3,76k	285	3
1,092	300단어 이하로 일본 여행에 대한 짧은 블로그 포스트를 작성하세요.	["length_constraints:number_words"]
1,094	유명한 얼마 이름 5개를 JSON 형식으로 제공해 주세요. 흥미롭거나 이상한 톤을 사용해 주세요. 전체 출력은 JSON 블록만 포함되어야 ...	["detectable_format:json_format"]
1,098	사람들이 신을 부르는 이름은 무엇인가요? 정확히 두 가지 다른 응답을 주세요. 응답은 6개의 별표 기호로 구분하세요: *****.	["combination:two_responses"]
1,107	로켓에 대한 두 가지 농담을 작성하세요. 응답에 심표를 포함하지 마세요. 두 농담은 6개의 별표 기호로 구분하세요: *****.	["punctuation:no_comma", "combination:two_responses"]
1,108	험버거는 샌드위치인가요? 오직 간단하지만 사용하여 응답해 주세요, 다른 언어는 허용되지 않습니다.	["language:response_language"]
1,127	수업이 그러한 것에 대한 시를 써주세요. 시는 "색선 X"으로 표시된 4개의 세션이 있어야 합니다. 시를 다음이 정확한 모그를 마크하십시오	["detectable_format:multiple_sections", "start-and-checker"]

한국어, multi-turn 평가 데이터셋
부재

Multi-IF: Multi-turn & Multilingual Instruction-Following Evaluation



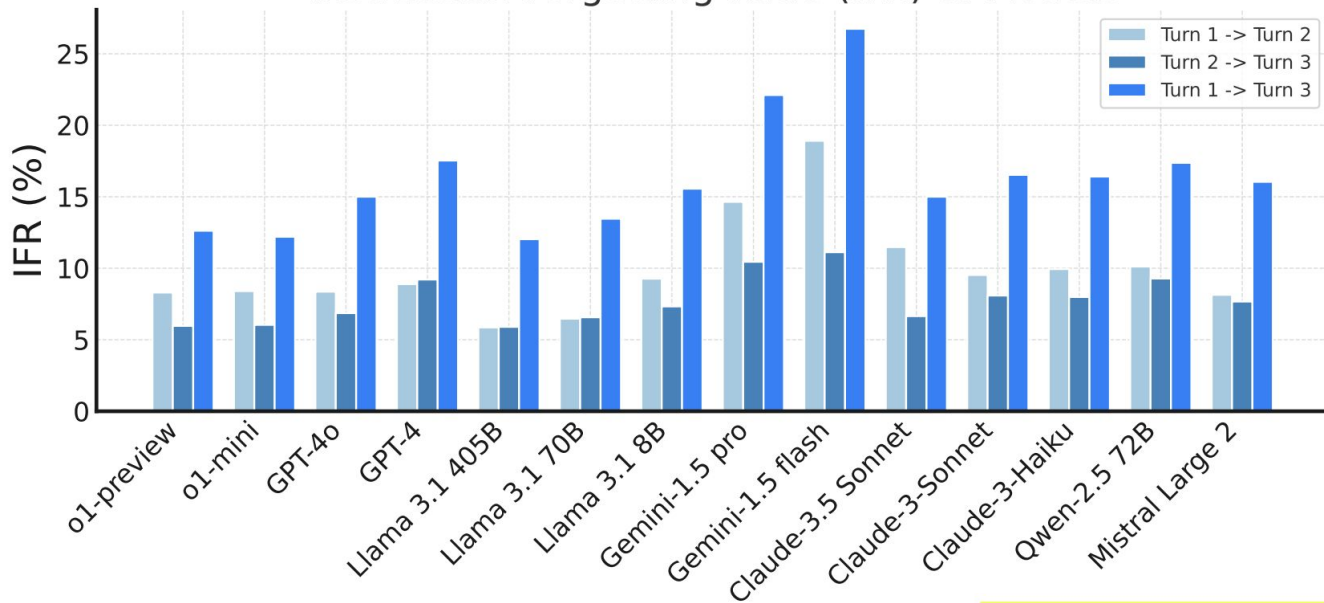
모든 언어의
평균 accuracy

Figure 6 The impact of multi-turns on instruction following. “Accuracy” means the average of the final metric of the all languages; The final metric is the average of the four accuracy scores: instruction-level strict accuracy, conversation-level strict accuracy, instruction-level loose accuracy, and conversation-level loose accuracy.

Turn 이 늘어날수록 instruction following이 떨어진다 → 이전 instruction을 잊는다.

Multi-IF: Multi-turn & Multilingual Instruction-Following Evaluation

Instruction Forgetting Ratio (IFR) of Models



turn-3에서 turn-1을 계속
기억하고 following
하는가?

이걸 측정하는
python코드는 github에
공개되지 않음

Instruction Forgetting Ratio

이전 turn에서 follow했으나, 후속 turn에서 not following한 instruction 수

$$\text{IFR} = \left(\frac{\text{Number of Previously Followed Instructions Not Followed in Subsequent Turns}}{\text{Total Number of Instructions Followed in Previous Turn}} \right) \times 100$$

이전 turn 에서 follow했던 instruction 수

Multi-turn에서 turn이 늘어날수록
최초 instruction을 잘 잊는 모델은?
→ Gemini 모델들

MT-IFEval-Ko 구성

데이터셋 준비

- Langfuse 서버 연동
- Input/Output/Metadata 활용

배치 평가 수행

- Langfuse evaluator 작성
 - Langchain 기반
- Langfuse tracing 구현
 - Multi-turn 고려
- Batch evaluation
 - 전체 dataset 일괄 평가

데모 평가 수행

- Gradio 데모 작성
- Run-time evaluation
 - 특정 item 평가

MT-IFEval Dataset Example

```
{
  "INDEX": "TelcoIF_001",
  "turn": "turn_1",
  "turn_1_prompt": {
    "role": "user",
    "content": "114에서 온 부재중 전화의 사유를 확인할 수 있나요?\n위 질문에 대한 답변을 반드시 영어로만 작성하세요."
  },
  "turn_1_instruction_id_list": [
    "language:response_language"
  ],
  "turn_1_kwargs": [
    {
      "language": "en"
    }
  ],
  "turn_2_prompt": {
    "role": "user",
    "content": "응답의 가장 마지막 문단을 \"P.S.\"로 시작하는 추신 형태로 적어주세요."
  },
  "turn_2_instruction_id_list": [
    "language:response_language",
    "detectable_content:postscript"
  ],
  "turn_2_kwargs": [
    {
      "language": "en"
    },
    {
      "postscript_marker": "P.S."
    }
  ],
}
```

Langfuse Dataset

Langfuse v3.80.0 OSS

SKT-AIDE / my-llm-langchain

Go to...

K

Home

Dashboards

Tracing

Evaluation

Users

Prompts

Playground

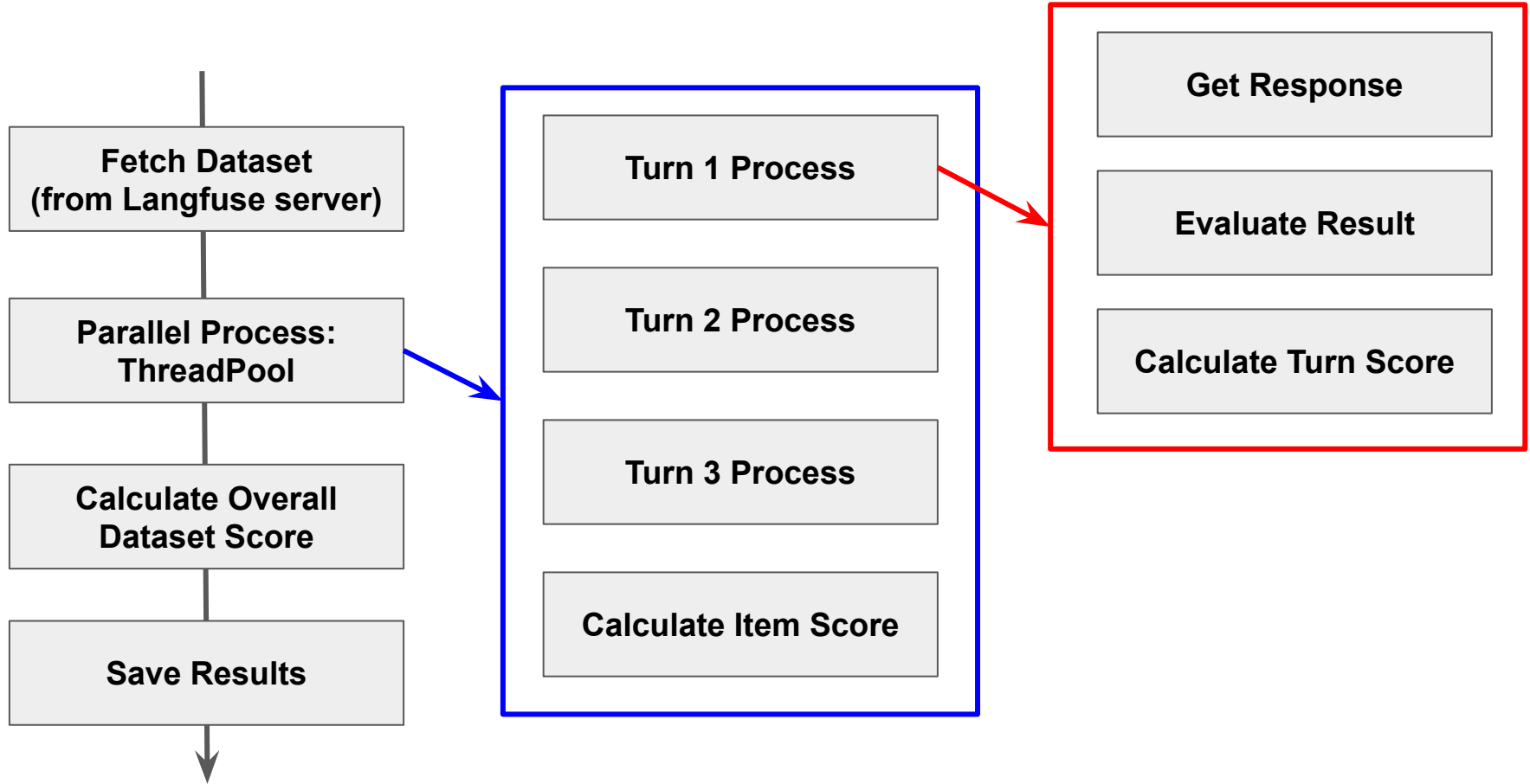
Datasets

Datasets

Search (Name)

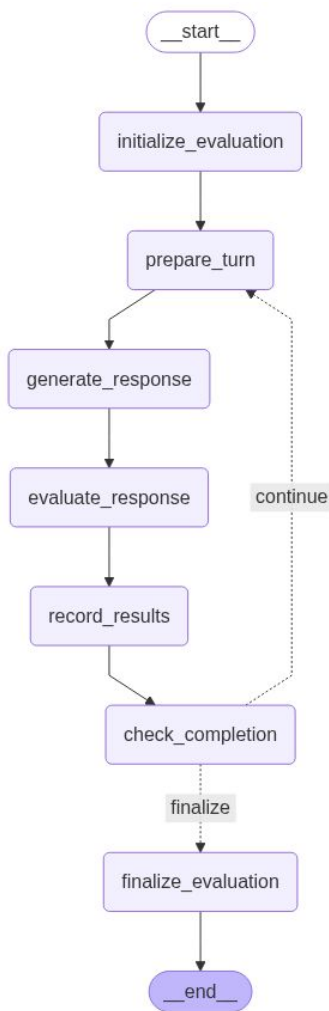
Name	Description	Items	Runs	Created	Last Run	Metadata
General_MultiF_English		896	21	2025-07-17 20:53:16	2025-07-18 11:42:16	
Telco_MultiF_Korean		100	14	2025-07-17 18:12:52	2025-07-18 09:45:19	
General_MultiF_Korean		157	8	2025-07-17 18:08:30	2025-07-17 21:12:09	
multi-if-ko		342	0	2025-07-16 15:54:00		
RAG_Evaluation_Dataset		49	1	2025-07-16 15:29:43	2025-07-16 15:34:53	
m-ifeval-ko	Translated datasets for Korean IFEval	342	0	2025-07-14 14:34:58		{ "language"

Evaluation Process



Multi-turn Messages

Langgraph Implementation



```
def _create_workflow(self):
    """LangGraph 워크플로우 생성"""
    workflow = StateGraph(GraphState)

    # 노드 추가
    workflow.add_node("initialize_evaluation", self._initialize_evaluation)
    workflow.add_node("prepare_turn", self._prepare_turn)
    workflow.add_node("generate_response", self._generate_response)
    workflow.add_node("evaluate_response", self._evaluate_response)
    workflow.add_node("record_results", self._record_results)
    workflow.add_node("check_completion", self._check_completion)
    workflow.add_node("finalize_evaluation", self._finalize_evaluation)

    # 엣지 추가
    workflow.add_edge(START, "initialize_evaluation")
    workflow.add_edge("initialize_evaluation", "prepare_turn")
    workflow.add_edge("prepare_turn", "generate_response")
    workflow.add_edge("generate_response", "evaluate_response")
    workflow.add_edge("evaluate_response", "record_results")
    workflow.add_edge("record_results", "check_completion")

    # 조건부 엣지 추가
    workflow.add_conditional_edges(
        "check_completion",
        self._should_continue,
        {
            "continue": "prepare_turn",
            "finalize": "finalize_evaluation"
        }
    )

    workflow.add_edge("finalize_evaluation", END)

    # 메모리 체크포인트 설정
    memory = MemorySaver()

    return workflow.compile(checkpointer=memory)
```

Langfuse Dataset Score

Langfuse v3.80.0 OSS

SKT-AIDE / my-llm-langchain / Datasets

Go to...

Home

Dashboards

Tracing

Evaluation

Users

Prompts

Playground

Datasets

Settings

Support

Dataset General_MultiF_English

New experiment

Select evaluator

Runs Items

Latency (s)



Average Total Cost (\$)



	Name	Created	Aggregated: # overall (api)
<input type="checkbox"/>	gpt-4o-mini-General_MultiF_English-100251_turn_3	2025-07-18 10:03:09	0 0.5206
<input type="checkbox"/>	gpt-4.1-nano-General_MultiF_English-105611_turn_2	2025-07-18 10:56:17	0 0.4781
<input type="checkbox"/>	gpt-4.1-nano-General_MultiF_English-105611_turn_1	2025-07-18 10:56:16	0 0.4781
<input type="checkbox"/>	gpt-4.1-nano-General_MultiF_English-105611	2025-07-18 10:56:16	0 0.4781
<input type="checkbox"/>	gpt-4o-General_MultiF_English-101819_turn_3	2025-07-18 10:18:27	0 0.4961
<input type="checkbox"/>	gpt-4o-General_MultiF_English-101819_turn_2	2025-07-18 10:18:26	0 0.4961
<input type="checkbox"/>	gpt-4o-General_MultiF_English-101819_turn_1	2025-07-18 10:18:25	0 0.4961
<input type="checkbox"/>	gpt-4o-General_MultiF_English-101819	2025-07-18 10:18:25	0 0.4961
<input type="checkbox"/>	gpt-4.1-General_MultiF_English-100251_turn_3	2025-07-18 10:03:09	0 0.5206
<input type="checkbox"/>	gpt-4.1-General_MultiF_English-100251_turn_2	2025-07-18 10:03:08	0 0.5206
<input type="checkbox"/>	gpt-4.1-General_MultiF_English-100251_turn_1	2025-07-18 10:03:06	0 0.5206

Langfuse Tracing

Langfuse v3.80.0 OSS

SKT-AIDE / my-llm-langchain / Traces

Go to...

K

Home

Dashboards

Tracing

Evaluation

Users

Prompts

Playground

Datasets

Trace Dataset run: gpt-4.1-General_MultiF_English-100251_turn_3: 62685689cbb56ad1a0e270ab90b06e75

Search (type, title, id)

Timeline

Timeline

Timeline

Timeline

Timeline

Timeline

Timeline

Timeline

Timeline

Timeline

Timeline

Timeline

Timeline

Timeline

Timeline

Timeline

Timeline

Timeline

Timeline

Dataset run: gpt-4.1-General_MultiF_English-100251_turn_3

4.93s \$0.001282

Dataset run: gpt-4.1-General_MultiF_English-100251

4.93s \$0.001282

loose_inst_s... 0.53

loose_promp... 0.33

overall: 0.43

strict_inst_s... 0.53

strict_promp... 0.33

Dataset run: gpt-4.1-General_MultiF_English-100251_turn_1

2.27s \$0.00055

loose_inst_s... 1.00

loose_promp... 1.00

overall: 1.00

strict_inst_s... 1.00

strict_promp... 1.00

RunnableSequence

2.19s \$0.00055

Dataset run: gpt-4.1-General_MultiF_English-100251_turn_2

1.52s \$0.000524

loose_inst_s... 0.33

loose_promp... 0.00

overall: 0.17

strict_inst_s... 0.33

strict_promp... 0.00

RunnableSequence

1.48s \$0.000524

Dataset run: gpt-4.1-General_MultiF_English-100251_turn_3

1.11s \$0.000208

loose_inst_s... 0.25

loose_promp... 0.00

overall: 0.13

strict_inst_s... 0.25

strict_promp... 0.00

RunnableSequence

1.06s \$0.000208

Dataset run: gpt-4.1-General_MultiF_English-

2025-07-18 10:03:06.881

Env: default

Latency: 4.93s

Total Cost: \$0.001282

Preview

Scores

Input

```
{
  content: "Of course! From now on, I will wrap my ent
  additional_kwargs: {
    refusal: null
  }
  response_metadata: {
    token_usage: {
      completion_tokens: 22
      prompt_tokens: 16
      total_tokens: 38
      completion_tokens_details: {
        accepted_prediction_tokens: 0
        audio_tokens: 0
        reasoning_tokens: 0
        rejected_prediction_tokens: 0
      }
      prompt_tokens_details: {
        audio_tokens: 0
        cached_tokens: 0
      }
    }
    model_name: "gpt-4.1-2025-04-14"
    system_fingerprint: null
    id: "chatcmpl-BuTnDUOWrd6pPeGnuL2MyDys2HYHd"
    service_tier: "default"
    finish_reason: "stop"
    logprobs: null
  }
  type: "ai"
  name: null
  id: "run--8cd1260e-93d1-4341-b4de-dccc2138506c-0"
  example: false
  tool_calls: [
```

Gradio Run-time Demo

MT-IF-Eval: Langfuse 데이터셋 평가 데모

이 데모는 Langfuse에서 mt-if-eval 데이터셋을 가져와 멀티턴 평가를 수행합니다.

LLM 모델 선택

gpt-4o-mini

모델 상태

현재 모델: gpt-4o-mini

데이터셋 선택

General_MultiIF_English

데이터셋 상태

✅ 데이터셋 'General_MultiIF_English' 로드 완료 (896 개 아이템)

아이템 선택

3748:19:en (English)

평가 실행

Input 내용

Metadata 내용

Input 내용:

turn_1:

```
{
  "kwargs": [
    {},
    {
      "keywords": [
        "dog",
        "day"
      ]
    }
  ],
  "prompt": {
    "role": "user",
    "content": "Rewrite the haiku below into two versions. Both of them should be funny. Separate the two versions using six asterisk symbols (*****). Include the keywords [dog,'day'] in the response.",
    "instruction_id_list": [
      "combination:two_responses",
      "keywords:existence"
    ]
  }
}
```

대화 기록

Rewrite the haiku below into two versions. Both of them should be funny. Separate the two versions using six asterisk symbols (*****). Include the keywords [dog,'day'] in the response.

On a chilly winter night
A cat meows at the moon
Hoping for some warmth

A dog steals my bed,
On this chilly winter day—
Guess I'm the warm snack!



턴별 점수 (모델: gpt-4o-mini):

Turn 1:

- Overall Score: 1.000
- Strict Prompt Score: 1.000
- Strict Instruction Score: 1.000
- Loose Prompt Score: 1.000
- Loose Instruction Score: 1.000

Turn 2:

- Overall Score: 0.167
- Strict Prompt Score: 0.000
- Strict Instruction Score: 0.333
- Loose Prompt Score: 0.000
- Loose Instruction Score: 0.333

Turn 3:

Summary

- Langchain/Langgraph 기반 workflow 이해
 - chain/node 연결 기반의 workflow
 - 코드 가독성 향상
- Langfuse/Langsmith 모니터링 연동
 - 평가 결과 분석 효율화
- 현업 업무에 연계 방안
 - 평가 데이터셋 개발에 활용
 - 실시간 프롬프트 테스트 및 평가 결과 확인
 - 어플리케이션별 평가 데이터셋 개발 효율화
 - 평가 결과 분석에 활용
 - LLM 모델별 성능 비교 분석 효율화
 - Fine-tuning LLM 모델 성능 개선을 위한 분석
 - Prompt-optimization 을 위한 분석

TODO

Langchain → Langgraph

- Langgraph 형태로 변환 (디버깅 중)

Gradio Dataset Studio

- Turn별 prompt/instruction/kwarg 입력/편집
- Turn별 LLM 평가 결과 실시간 확인
- 개별 데이터셋을 Langfuse 데이터셋에 업데이트/관리

Korean MT-IFEval Dataset

- 한국어 평가를 위한 데이터셋 개발 (general-domain)
- Telco 서비스용 LLM 모델 평가용 데이터셋 개발 (application-domain)