

ComplexFuncBench:

: 2025-01-17 00:00:00

: <https://arxiv.org/pdf/2501.10132>

: 65

: 70

:

(LLM) €
• Complex-FunC-Bench
€ 128k
€ LLM , •
€ , •
:
LLM
€ , • , € ,
• f •
• •
... , • : 1. Complex FuncBench -
€ LLM, •
€ (128k).
" 5 (, ,
,).
- , ComplexEval,
† API € : • ,
LLM.
-
† : ‡ , GPT-4o, ‡
† 100 1000 .

€ (GPT-4o, Claude-3.5) , •
 € ,
 • , , - € , € 5
 (func_error, param_missing, hallucination, value_error, stop_early)
 , .
 ## :
 ### ^ • • API
 Comple FuncBench , € • API. „
 API • .
 #### ... , • :
 f , ^
 • %
 ^ : " • € , X, † € , Y"
 , • , • € , • €
 • %
 •
 „ , ...
 ^ • €
 , • , € , Š , € , , †
 † ‡
 , • • † , •
 ^ ,
 %
 • -
 † ,

< .
^ . %
„ :
CE ^ .
CE • • • €
CE < †€€
CE • •
CE ^ • •
Ž € API € , †
Comple FuncBench, • € †€€ LLM
• : 1. Complex FuncBench : -
^ • : „ • , LLM. -
... : Š , •
€ , •
• - ^ : -
€ .
• : ^ : Ž API
• , : Š , • †
• ... € , • ^ ‡
• : - € .
• :
^ : Ž , : ,
• • ... : ^ ,
• : - , LLM.
€ :
^ : Š , •

: Š LLM. ^ , : Š - .
 • , , :
 ^ : Š , €
 : Š , € LLM
 • : Š - .

Prompt:

... , •
 ComplexFuncBench GPT ##
 ComplexFuncBench , • LLM ,
 € , . f
 †€€ .
 ## ^
 [=====] # • : <
 ## ... % 30000 % < 15
 22 , . < ,
 ##
 (•) 2. , • : 1. •
 ^ • : - •
 • (: 78.8%
 •) - • , •
 API 3. ^ €
 € € 4.
 € JSON search_flights API
 ## „ € : [=====]json { "origin": " , "destination":
 " , "departure_date": "YYYY-MM-DD", "return_date": "YYYY-MM-DD", "max_price":
 • , "direct_only": • , "preferred_departure_time": " " }
 [=====]
 ^ , ' . [=====]
 ## „ ‘

Scatter plot showing the distribution of scores for LLM, JSON, and API across different benchmarks. The y-axis is labeled 'ComplexFuncBench' and ranges from 0 to 100. The x-axis is labeled 'LLM' and ranges from 0 to 100. The legend indicates that blue dots represent LLM, orange dots represent JSON, and green dots represent API. The plot shows that LLM scores are generally higher than JSON and API scores, with a significant gap of 78.8% between LLM and JSON scores. The API scores are also generally higher than JSON scores, with a gap of 10.0% between API and JSON scores. The plot also shows that the scores for LLM, JSON, and API are generally higher than the scores for the other benchmarks.