

TaskEval:

: 2025-03-10 00:00:00

: <https://arxiv.org/pdf/2407.21227>

: 65

: 75

:

Task-Eval -

€ € (LLM). •

€ (IRT). „

‘‘
f , , LLM f...
f , , f
, .

:

f f LLM:

€ , €

† f f LLM
f f ‡ € , LLM.

^ ,

: 1. TaskEval - ‡
LLM, f €
Item Response Theory (IRT)

€ .

f (18 f - ,

, .

€ -
(difficulty) (€
).

•
HumanEval 21 ClassEval) f , f f (17

f , .
 $f \dots f$ LLM ϵ f ,
 , LLM.
 ## : TaskEval f f ϵ
 LLM, ϵ f API.
 f f API -
 f :
 f - , f ,
 2-3 f . % ϵ f , f
 .
 „ ... - . Š f
 f ϵ ϵ :
 (f) f
 Š
 † - , SQL- ,
 f f f). Š f ‡ .
 ‡ LLM ϵ f -
 f , ϵ f ‡ LLM.
 f f , f ϵ , Š
 f ‡ f .
 Š ‡ ϵ : - ϵ
 LLM - ϵ - • f f , - ϵ
 ‡ f
 Ž , ϵ f , LLM ‡
 .

- Š : Š : 1. f LLM ‡ € ,
 f . - ^ f : Š ,
 f : ‹ f 2-3 € f LLM. - Š
 € .

€ Š : • €
 - IRT- . ^ f : % -
 , , Š , f , LLM €
 , . Š : Š f f LLM.
 ,

•
 Š : • , (, SQL- ,
) LLM. ^ f € : Š :
 Š f LLM € € . Š LLM.
 ,

LLM
 Š : Š , f LLM. ^ f :
 % - f " : Š f
 LLM". Š : Š f LLM.

€
 Š : Š , f ,
 f € LLM € € . ^ f : Š f
 , : ‹ f f f . Š LLM.

Prompt:

Š TaskEval € GPT ## ^ ,

TaskEval , : - • f
 f - •
 (, SQL,) -
 LLM - ,
 f

```

## Š      f f

[=====] # •      : '      f      ,

## ^      "      ,      f

      . '      f      -      .

## Ž      f      f      ,      f      , find_sequence_pattern(numbers:
list[int]) -> str,
      ...      f.

## Š      - %€ : [2, 4, 6, 8] % € : "      2" -
%€ : [2, 4, 8, 16] % € : "„      2"

##      f      -      f      •      -      f
€      -      ,

## •      [=====]python def find_sequence_pattern(numbers: list[int]) ->
str: # Ž      [=====] [=====]

## Š      f ‡      ‡      ,

^      f ...      -      ,      ,      ,      f

      f      ,      TaskEval      ,      f

%€      -      LLM, ‡ f
      ‡

,      f      f      ,

f...      f f      -      ,      €      €      €      €      €,
      ,

† Š      -      €      ,      ,

      f

      f      €      ,
      f      f      GPT      ,
      ,      f      f      €      .

```