

: 2025-01-14 00:00:00

: <https://arxiv.org/pdf/2501.12405>

: 58

: 70

:

(alignment) LLM, € • (€ , € , ,
,). • LLM, , • ,
• • .

:

f , • € • € LLM • , ,
... • †..... € • ,
• € , € .

‡ , : 1.
(alignment) LLM: •
• € • : • (€ , , ,), •
(• , , † € ,) (• , , ,
,).

LLM : € , ,
(€ , € , , ,), € ,
.

: • , , • € ,
• • • .

:
€ € , ^ ‡ , • •
, € , • .

• :
• ,
• € ...
•

: f , API •
• %
• • • (, †.....
•),
€ • , •

‡ , • • • , :

€ : Š € •
• , , : : < € , € ,
• • :
(• ,) - ... , († € ,)
, : œ , , € , (, •)

f : Š € • € " € •
€ • " • , " • • ,
•

• : Š € • € •
• , , ... , •
€

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

Ž , • ... • , € • , • • •
•

€ , • • : ## ‡ € •
- : • , € • •
€ • € • , • € †
, • , • †..... ... • ,
€ • - : • • Š • € • ,
, LLM • • • ,
, • • • • -
• : • • Š € • • †
€ • • , • ,

Prompt:

Š • " • GPT ##
‡ , € •

f • LLM: 1.
(€ , ,) 2. • († € ,
• ,) 3. , (, •)

Š • • • † €

[=====] # • ...

‘ • • (, ... , € ...
, ’ •
Š , € •

• • - Ž € , : ‹ 35- IT- €
, , (2023) - ’ , (• ,)
, • † € ,

Ž , , € • , • •
, , •

, • † • , • • ...
... •• 1 • . [=====]

‡ • € €

” : Š • , € , € (...
,), (• ...) (€
•) • .

” : Š • , † € , •
(• , () €), • •
• , () • .

” : Š • € • ,
(€ €), , € •
• .

, • € , , , •
† , , • •

€ , • , • ,
.