

:

1

2

: 2025-02-17 00:00:00

: <https://arxiv.org/pdf/2502.12470>

: 82

: 85

:

(LLM)

( 1)

( 2)

. €

• ,

,

2,

,

,

1

•

•  $f$  , ,

• .

,

:

•

• ” ” •  
) , •

(

.

•

2

” :

1

.

... ” ”

.

## ... •

: 1.

:

LLM:

( 1)

( 2),

.

€

:

2000

,

†

( 1)

( 2).

•

(alignment)

:

( 1),

( 2).

,  $f$  ” ”

: ‡ ,

2,  
 1  
 ... †  
 : †  
 2  
 1  
 ##  
 :  
 ### ^  
 1  
 2  
 €  
 €  
 €  
 †  
 f , ,  
 #### ... , ,  
 :  
 •  
 ( , ) : "  
 " ( 2 ) : "  
 %  
 " ( 1 )  
 , € :  
 , 2 f , ,  
 1 Š , , •  
 %  
 † † :  
 † : "  
 ( 2 ) : "  
 ^ † f , , :  
 : " € , f "  
 : " ^ , , " ####  
 € :  
 • f , , LLM  
 %  
 ^  
 Ž

• , %

LLM.

## : • ^ : - ^ LLM,

": "€ " : • ^

• - ^ " :

• " ... " " : • ^

€ - ^ : . €

• , f

, f (anchoring bias, halo effect . .) LLM

" f " - ^ " : . ^

LLM.

• (alignment) - ^ : ' . €

• - ... " " :

" " LLM. - ^ " : . ' ,

f , , % • %

• - ^ : :

• ^ 2 1 . -

... " : • ^

• Š % • % LLM. - ^ " LLM.

... † - ^ : .

• - ... " : • ^

• - ^ " : • ‡ LLM

**Prompt:**

^ 1 2 GPT  
( 1) ( 2)  
f " : "

## ^ ( 2)

[=====] " , ( 2)  
2) •‰ :

"• 96 ‰ 12 8 . " ,  
‰ ‹ 3 ?"

^ : 1. , 2. ' 3. • 4. ^ 5. ' ,

" „ • • . [=====]

## ^ f

Š f , , ‰ 2 , :

† † : , , ( 7.66%)  
‰ 2, € : " , f  
‰ : • • f

Š , : ^ • ‰ 2 ##  
( 1)

‹ , - :

[=====] , ( 1) •‰ :

"\_ , „ ?"

€ , , • . ‡ ,  
• [=====]

## ... • „

€ • : ( 2)  
; ( 1)

“

—

1

” : •

2,

“

”

:

~

”

”

,

,

.