

# LLM

■  
■

: 2025-03-04 00:00:00

: <https://arxiv.org/pdf/2503.02863>

: 72

: 85

:

(LLM)

:

€

SteeringConf,

LLM

,

•

,

$f$

:

$f$  LLM

” . ... ”

”

, / ” ,

,

. †

”

,

”

$f$

LLM.

## ‡  $f$

: 1.

(Confidence Steering) -

,

•  $f$

”

(

”

,

”

”

”

”).

(Steered Confidence Aggregation) -

^

”

,

,

”

.

€

(Selection) -

(Steered Answer

•

,

,

.

$f$

•  
SteeringConf

## :

### †

API

$f$

€

”

€

...

$f$

0 100.

†

†

†

€

€

€

Š

##

†

: 1.

Š

†

LLM

- †

LLM

€

†

$f$

: Š

†

$f$

† : . Œ  
, , " € . ‡ "  
" : Š " , ,  
† " , . † " : Š "  
† " , " ,

€

† : • , † ,  
, , : . † ,  
" : . Ž , , , ,  
.

† : • . Œ " .  
‡ " : Š . † ,  
" , , ,  
€ " . † " € : .  
‡ " " , • €

• ,  
† : • . • : Š .  
• € € f. ‡ " " : Š . † "  
" : Š .  
€ .

**Prompt:**

† LLM GPT ##  
‡ f

SteeringConf , , f •  
" f " " " ,  
,

## † SteeringConf

[=====] • " :  
• :

, 1 10 “ ’ • “ œ ” • • Ž . , ,

*f*

‘ 1 10. , , • ’ –œ ‰ “ • • Ž . ”

• ” , 1 10. “ ’ • –Š ’ ’ • • • Ž . ”

*f*

– , ,

Ž : Ž ,  
? [=====]

## ‡ •

*f* ~ ‹

, ,

*f*

~

” ” .

€ ~ Ž

,

, ‡ ~ ,

, *f*, •

*f*

œ

, : ” , , € ” € ,