

SteerLLM

: 2025-03-04 00:00:00

: <https://arxiv.org/pdf/2503.02989>

: 60

: 75

•

•

Steering),

□

□

ϵ f f , f : , f \bullet ϵ \dagger ϵ \vdots , ϵ API.

\dagger , ϵ : 1.

LLM: ϵ , "model steering" - f , f .

CONFST (Confident Direction Steering): $\hat{\cdot}$, f , ϵ , "LLM.

€ : ... € ” ’ ” , , VS
 €), CONFST •
 € .

• : %₀₀ Š , €
 € „ f „ ,
 , CONFST f : ... € , €
 .
 : %₀₀ , ,
 f , € API- € . † , €
 € f f .
 f f :
 „ ^ „ € " € :
 , „ €
), • " " .
 € : ... € f ,
 , € € € " : " †
 " " : ^ € / € ,
 " " .
 • : ^ , ,
 "), ("œ • , .
 † • f :
 • Ž , , " €
 €
 • ... •
 •
 • • € € f
 • • ,

f , ϵ ϵ
LLM ϵ , f
CONFST.

ϵ : 1.
LLM: - \wedge : \langle . $\%_{00}$ ϵ
 f : ... \cdot \dagger ' , ϵ . - \dagger f , , ϵ
 \cdot ϵ , LLM , . -
 \wedge f f : . ' ϵ , \cdot ϵ
 f , f ,
 f , .

CONFST: \wedge : . '
 f , ϵ ϵ " : ... API, f
 \cdot " \cdot \dagger f f : ... \cdot \wedge , ϵ
 \wedge f f : ... \cdot \wedge f
" ϵ " , ϵ .

ϵ :
 \wedge : \langle ϵ ,
 \cdot ϵ . \dagger f f : ... \cdot \wedge , ϵ
 f : ... ϵ ϵ \cdot \wedge f
 f : ... \cdot \wedge (, " , ") .

\cdot :
 \wedge : ... ϵ , ϵ
 \cdot \dagger f f : ϵ \cdot " \wedge f ϵ
 f : ϵ \cdot " .
 f : \cdot " .

f :
 \wedge : ... \wedge ,
 ϵ f \cdot \dagger f f : $\dagger \epsilon$ f :
" ϵ " \wedge f f :
... \wedge " " ϵ f :
 \wedge f : 65 100.

CONFST, \cdot ϵ ϵ
 f , LLM "
 . $\%_{00}$

€ , (,)
 .
 † € f API. (f , € €
 (f ,
)
 f .
 ‡ f : 1. † f (75-80),
 f , € € f € ,
 f . 2. † f • (45-50),
 f €
 , € •
 .
 ^ , f , f 60 100, ,
 , f € , f €
 € € .
 † f 60 •
 f LLM, € € €
 f € .
 • f : †€ €
 CONFST, " € . Œ • ,
 (,) . † € LLM :
 f .
 † f : † f : 75 100.
 ^ f f f , f f :
 1) •
 " " € €
 " € .
 2) ‡ f f
 (, +) • €
 • .
 3) ^ , € • f " "
 f , • €
 € .
 4) ^ f " " € •

f
 ϵ / .
 \dagger ϵ , ϵ f ϵ
 ϵ f , " f .
 $\parallel <\dagger f : 60> \parallel <\dagger ' : \text{CONFST},$
 $\text{LLM } \epsilon$ f , f ,
 \cdot f f : ϵ ,
 f . \dagger ϵ : ϵ
 ϵ API.> $\parallel <$:
75>

Prompt:

SteerLLM CONFST GPT

CONFST (Confident Direction Steering),
" " ,
 \bullet ϵ GPT), f \bullet (ϵ .
 \bullet f **SteerLLM**
[=====] " ϵ , ϵ
 ϵ f :
" ^ SQL- ' f \bullet
. " " ^ f \bullet
. " "%
 f 99,9%." \dagger ϵ ,
 ϵ f f .
" ϵ , ' , \bullet , f " . [=====]
... , , ϵ f CONFST:

... " " " "

€ , € € .

, " f ("

"), , • , €

CONFST € .

: € , € , € .

† " " "

, € " .

€ GPT. ' "