

■
■

: 2025-02-03 00:00:00

: <https://arxiv.org/pdf/2502.01116>

: 78

: 85

:

€ • (LLM) €
€ , €
€ , €
• f f -
: € f (RM)
• f

:

€ LLM:
• , €
• , •
• , €
• €
•

† : 1.

(safety alignment):

€ LLM € (fine-tuning):
• f €

" €" LLM (Picky LLMs): ‡ •
€ (• , Markdown •) • f
€

• € , : • • €
 € "† - • , ... " • f ... €
 f € , (Reward Models):
 • € f , •
 : € € € f
 • f .
 ## : • , LLM.
 € € € €
 • , API • •
 Š , • f
 • :
 „ : LLM
 Markdown • • •) . ^ • ,
 ... € • • : " ^ f ,
 Markdown • • • " .
 ... € : ‡ f € ,
 " , • "† € X" € • f X" .
 € , : €
 ^ • € €
 • • , †† ‡ † : ‡ € ... €
 • • € • € • , • € €
 Ž € • • • €
 € • LLM € €
 ## • • : €

Prompt:

\cdot , $\cdot f \dots$ \in . - <
 \in \cdot , $\cdot \in \in \cdot$: 1.
 \cdot 2. \dagger (\cdot) 3. ,
 \cdot

\cdot \cdot $\wedge f$, \cdot , \cdot , $\in \cdot$. [=====]
 \cdot ,

\cdot \in \cdot \in

\cdot : $\wedge \cdot$... \in \cdot .
 \cdot ,

\cdot : $\wedge \cdot$ \cdot , f .
 \cdot ,

\cdot , \cdot , : $\wedge \cdot$ f f \in . \in \cdot .
 \cdot ,

\cdot \in \cdot f ,
 \cdot \cdot .
 \cdot