

■  
■

,

,

: 2025-03-04 00:00:00

: <https://arxiv.org/pdf/2502.14171>

: 70

: 75

:

LLM-  
(Theory of Mind, ToM). € -  
• €  
, ( f ,, , , ) f  
ToM- • € € € 67% 63% Llama 3  
3B 8B .

:

€ BDI- ( f , , € )  
LLM. „ , f € ,  
f  
... , € f  
€ € f ^ • • † ToM. ‡ LLM  
, f ,.

## % : 1. (ToM) LLM:  
f , (BDI- ) f  
f .

€ ToM • : € ,  
LatentQA Š • € ToM € ,  
.

ToM- :  
ToM

$f$  , .  
 , (  $f$  • :  $\ddagger$  ^ , ,  $f$  ,  
 , ( , ,  $\dagger$  67% 63%  
 Llama3 3B 8B  
 ToM- • € .  
 ” • €  $f$ :  $\ddagger$  ,  $f$  LLM  $f$   
 ToM, .  
 ## :  
 ###  $\ddagger$   
 ” , (LatentQA)  
 ToM,  
 € €  $f$   $f$   $f$  ,  
 $f$  API.  
 • • :  
 BDI-  $f$   $\ddagger$   
 ,:  $\langle f$  (beliefs):  
 Œ (desires): , ... (intentions):  
 ”  $f$  ToM  
 $\ddagger$  • € , ,  $f$  ,  
 • € , , X, , ... : " <  $f$  ,  
 Y..."  
 ...  
 • € , X, Y, Z" : "  
 ,  
 "Show empathy"  
 $f$   $f$   $f$  ,  $f$  ,  $\ddagger$  -  $\ddagger$   
 € - ,  $f$  -  $\ddagger$   $f$  ^ • •  $f$   
 - < , ,  $f$   
 •  $f$  , , ,

$f^{\wedge}$ ,  $\in$  ToM  $f^{\wedge} \bullet \bullet$   
 BDI-  
 ## : 1.  $\in$  ToM  
 $\bullet - \ddagger$  : ...  $f$ , , , ,  
 $f$   $\in$ , , , LLM  
 $\cdot - \% \in$   $\in$  :  $\bullet$ , ,  
 $f$   $\in$  : ,  
 $\cdot$   $\cdot - \ddagger \in$  ToM,  $f$   
 $f$ ,  
 ToM-  $\ddagger$  : ...  
 $f$ , -  $\cdot \% \in$   $\in$  :  
 $\bullet$ ,  $f$ ,  
 $\cdot \ddagger \in$   $\in$  :  $\bullet$ ,  
 $f$ , , ,  $\bullet \in$ ,  
 $f$ ,  
 $f$   $\bullet$   
 $\ddagger$  : ,  $\cdot \% \in$   $\in$  :  $\bullet$ ,  
 $\in$   $\in$ ,  
 $\cdot \ddagger \in$   $\in$  :  $\bullet$ , ,  
 $\in$ ,  $f$  LLM.  
 † BDI (Beliefs-Desires-Intentions)  
 $\ddagger$  :  $\cdot \% \in$   $\in$  :  $\bullet$ ,  
 $\wedge$ ,  $f$   $\cdot \ddagger \in$   
 $\in$  :  $\bullet$ ,  
 LLM.  
 ”  
 $\ddagger$  :  $\bullet$ , , ToM  
 $\cdot \% \in$   $\in$  :  $\bullet$ ,  
 $\in$   $f$  ToM-  $\bullet$   $\cdot \ddagger \in$   
 $\in$  :  $\bullet$ ,  $f$   
 $f$ , .

Prompt:

, GPT ## %

, , ,  
 (ToM) •  $f$   
 $f$  (desires) • :- (beliefs) •  
 - ‡  
 ^ (intentions) •  
 ## ‡ ^ • • ToM  
 [=====] # € GPT  
 • † , € . •  $f$  :  
 $f$  : '  
 ‰  
 $f$  :  
 ‰ € ‰ €  
 † €  
 „  $f$  :  
 „  
 • € ‡ ‡  $f$  †  
 , • ,  $f$   $f$  ,  
 ‡ ^ • € ,  $f$  ^ • • ^ ,  
 . [=====]  
 ## ‡ ^  $f$   
 €  $f$  • , , , ^ ,  
 • € ToM ^  
 † BDI- (Belief-Desire-Intention) •  
 ,  $f$   $f$   
 ‰  
 ‰  
 $f$  ,  $f$  € , †

$\hat{\cdot}$  ,  $f$  ,  
 $\cdot$  ,  $f$  ,  
 $\cdot$  ,  $\hat{\cdot}$  ,  $\cdot$  ,  
 GPT 60-67%,  
 .