

: 2025-02-27 00:00:00

: <https://arxiv.org/pdf/2411.04847>

: 70

: 80

:

€ • (LLM) • ,  $f$  • PRISM (Prompt-guided  
Internal States for hallucination detection of LLMs). „ -  
- • €  
• € • ... € : € ,  
• € LLM, €  
• , † € • € •  
• • , .

€ :  
• LLM €  
• • • • †  
• € • , €  
 $f$  •  $f$  •  
€ • ^ • •  
€ .

## % : 1. • LLM:  
 $f$  • PRISM (Prompt-guided Internal States for  
hallucination detection of LLMs), € •  
€ • .  
: € ,  
• • € • LLM •  
€ • , €  
• • € • • •  
- : • ,  
€

• , • €

• : ^ • Š f f  
• € • , €

• f f : PRISM •  
€ • , • • , €  
€

## : € • €  
• • • , •

• ”  
€ • € • LLM € •  
API.

” • :

•  
• 10 € f • ,  
Š f f • • ” € ”,  
f • . ‹ • • € • :  
"œ '[ ]' • ?" " •  
• '[ ]' • ?" "Ž € ,  
'[ ]' • ?"

” f • ... • ^ € • f •  
€ f • • † • ,  
f • .

• - • • € • †  
• f • , € €

Š f f - f • • € , •  
€ • , € • • € •  
• • •

” • € • Š : - ^  
• f • - • f • -  
€ f • LLM - •

• • , • € •  
€ • f • • •



^ • PRISM • GPT ## % € €

€ PRISM € , • •  
 € • €  
 f • , • • € • • • .

## ^ • • PRISM

[=====] Œ , • , :  
 : " „ '[ ]'  
 ?"

, • “ , € , :  
 % f € • % € Š f  
 € f • ” ( ,  
 , ” , )  
 ‘ , 1 10

• : " • • • • ’ • ,  
 • 166 • • . " [=====]

## ^ • Š  
 : • € " „  
 ” ?" PRISM  
 • , € f • .

† : ^ • •  
 • , " •  
 " (truthfulness direction) †  
 : ’ € • ,  
 € • •  
 • • • •  
 — • • € € • ,  
 € •  
 GPT •