

CySecBench:

,

,

: 2025-01-02 00:00:00

: <https://arxiv.org/pdf/2501.01335>

: 62

: 75

:

-

CySecBench,

12662 (jailbreaking) LLM

€ :

€ €

10

;

• , ,

jailbreaking

, €

,

€f

• , ,

(

88,4% €

„

)

...

LLM.

:

†

•

LLM

f

€ €

,

”

€

.

,

€

”

...

€f

...

€f

...

. ‡

„€f

f

,

€

€

.

€...

...

^ f...

: 1.

€

CySecBench -

€

”

12662

,

,

€

10

.

•

€

-

(GPT-3.5-turbo

GPT-o1-mini) € €... „ () .

• , f „ € -
LLM €
... f "• "

„ LLM -
€ LLM (ChatGPT, Gemini, Claude)
€ ... € .

• „ „ , - € €... „
• , , ... , € f
€

:

%... €... API

† ... LLM €... , API- € :

• ... : (...
% „) LLM f ...
Š : ... •

„ „ : % MECE (f ... f , f
... f €...) € €
... .

:
(€...) €...
.

„ : < € ...
 f

€
% • : - œ

€ € ... f - • €... „ Ž
- œ ...
• ... - % „ ... € f
€ f

... : 1. €
CySecBench - % : ...
, €

€ €
 . - ^ € : • , € ” €f €
 , f , ... : ‰ ” €f €
 . - ‰ €
 € € €
 ...
 • € ‰ : . <
 ^ € :
 • , € LLM
 € € , ‰
 : ‰ LLM ,
 € ... , , .
 • , f , €
 ‰ : • € LLM.
 ‰ € €... , , €
 € . ^ €
 : ... , € ... €
 LLM , €
 ‰ : ‰ , €
 , € , €...
 , € , €...
 • • ...
 , € ,
 ” LLM
 ‰ : . ‰ € . ^ € :
 ... €
 • , LLM. ‰ : ‰ ...
 ... €
 • ” ” ,
 ‰ : . • € ,
 , € €...
 . ^ € : • , € € € LLM.
 ‰ : ‰ • € ...
 , € LLM € €... , ...
 ... € .

Prompt:

‰

CySecBench

GPT ## ^ f...

† CySecBench € • , ,
LLM, € €... „ ...
(,):

... „ „ €...
f „ €
• € ...
† f „ f € „ ## %
• , , CySecBench

[=====] #

‘ f €... € IT- "Š
" . %
€ ...
€ .

€ : 1. ^ IDS/IPS (5
) 2. • Snort
€ :- - SQL- Ž - %

Ž , ... € €
“ , ... €
€ . [=====]

% ... € • • , , :

“ - f , €
€... , ... f „ „ €...
...

‡ „ f „ - € f € € €
... € , ... f
... ,

^ „ € - €
“ , ...

• „ „ - ...

##

€... GPT : -
“ € € ... - € € €
... € - „ € f
“ , - %

•

,

CySecBench,

€...

,

...

.