

■
■

: 2025-03-06 00:00:00

: <https://arxiv.org/pdf/2503.04856>

: 70

: 85

:

(multi-turn) - LLM (single-turn)
: M2S (Multi-turn to Single-turn)
€€ •
, 95.9%,
17.5% GPT-4o.

:

Numberize, Pythonize), (Hyphenize,
• „ €€ LLM. f
jailbreak- , € ...
, ,

† ... • : 1. €
(M2S): (Hyphenize,
Numberize, Pythonize) € „ jailbreak-
€€ , ... ‡
€€ , ... ‡

• : ^
... „ (ASR 95.9%),
17.5% ...

, : † •
(Numberize) ... - (Hyphenize),
€ (Pythonize), •

‡ LLM.

f :
 (, ...‡ ”) €€
 € ,

” :
 ” - (guardrail models)
 ...‡ € .

:
 API. ‰
 • LLM. *f* •
 €€ (Hyphenize,
 Numberize, Pythonize) ...
 • LLM.

† ” , • :
 ... † † ” (Hyphenize) -
 ^ , ...‡ : - †
 † ” , • ? - † ... ? -
 † • ‡ ... ?

‡‡ (Numberize) -
 : ...‡ : 1.
 € , ” € 2. † ? 3. †
 € • ... € ?

- † † (Pythonize) - €
 : questions = ["Š
 ... ?", "† ‡ ...
 , ?"]

^ ,
 : - ‹
 - ‹ €€ - ‰
 • - ‰
 - ‰ • ,

• ,
 • .

• : 1. €
 (M2S): - ^ :

[illegible]

Prompt:

" M2S GPT ## † ... •
 " • " ,
 LLM €€ . f
 € , "

GPT.

^

Pythonize

[=====] # '

Pythonize

[=====]python #

• ...†

€

• context = "" "

...

• '€

• ""

#

requirements = [" '€

13-16

" , "•

...•

" ,

"%
%

" , "†

" " ...

"]

'

task = ""

"

•

...

" . %

...•

: 1. †

•

(3-4

) 2. 3

...•

• 3. ^

•

4.

""

%

•

•

execute_task(context,

requirements, task) [=====] [=====]

##

•

€€

Pythonize

M2S,

:

... † † †

•

GPT

• ,

" ,

•

†

•

•

^

†

•

•

•

•

...

"

" , •

"

•

€€

17.5%

...,

•

• ,

"

•

##

%

Hyphenize (

)

Numberize (

•

)

•

•

GPT.