

■
■

: 2025-01-11 00:00:00

: <https://arxiv.org/pdf/2501.06638>

: 65

: 75

:

(semantic leakage)

€

€

:

f

”

LLM

f

,

.

†

,

.

‡

: 1.

€

€

€

-

,

€

(500

7

)

,

.

€

-

,

f

,

:

/

f

,

f

.

•

€

-

"

" (Mean Leak Rate),

,

,

,

-

.

, € € - ,
, € , •
, € .

f „ ... € „€ „ -
, € *f*
.

: , •
API. ^ BERT-score SentenceBERT
,

, - ,
.

%₀₀ , *f* :

† € -
, : " € *f* Š
< " " vs " € ,
" Œ ,
•

‡ € ^ € € € -
, , *f* ,
.

• „ € „€ € - *f*
, : " € ,
• - ,
".

„ -
, , , , ,
.

%€ € €€ € - *f*
, *f* „
, • .

... • , : - ... ,
• - Ž €
- • € LLM - •

• „ „ LLM , , Š

: 1.

€ € € - ... : ...

, € (, 0.5B 1.5B),

‡ : ... , , . -

, € f : ... , ,

. - ... : ...

” , f ” , ,

€ ... : ...

f , • f .

‡ : ... , (LLM.

)

... : ... ” ,

, ,

.

• €

... : (- ,

API , ‡ :

... f , f ,

f , Š ... :

... , , f ,

” .

, € €

... : - ,

, € , €

‡ : ... f , LLM ,

. ... :

... •

, , € € .

f ” ... € „€ ”

... f Š : ... , , f ,

f Š , Š . ‡ :

... , LLM

. ... : ...

, , f .

Prompt:

(
 (0.5B) € , € , € ,
 , f , f ,
 , Š

...

... 1: ‡ f

[=====] , f
 " " , " "

• ,
), f : 1. < ,
 Š (, , . .) 2. ,
 " , 3. "

• „: " < , € •
 Š ". [=====]

... 2: ‡ f , Š

[=====] " " " , Š f .

• : 1. <
 " " 2. ... (, ,)

3. <

[=====]

, " ,

• , f
 " " - ,
 , ,

‰

„€

Š "

‰

“ ”

”

，

， ，

， ，

·

†

，

， ，

f

， ，

f

·