

LLM-

: 2025-02-21 00:00:00

: <https://arxiv.org/pdf/2412.21102>

: 72

: 80

:

LLM.

€ -

, • ,

•

€ . f

-

Adaptive Prompt Pruning (APP),

””

.

:

LLM

. ...

€ APP

, € (

† € ,

)

•

•.

† , • ,

, €

•

LLM.

† •

: 1.

(APP)

-

€ LLM-

.

-

(

† € , , , , , , ,) ,

.

€

•

-

,

• ,

:

^ ””

^

.

€

-

,

• ,

†

€

,

.

, *f* -
, •

,
.

:

‰

API

€ API
€ APP
† ,
.

• € €

€

€ , :

” : ‰
(, , €)
† € ^

Š

•

...

” :

, † €
, ,

, *f*:

‰

/ € € ‰

- ,
-

† , ‡ *f* ,

f:

f €

† € ‰ -
• ‰
:

< ‰ ^

< Œ ^ •

< •

< Ž

• † €

