

: 2025-01-28 00:00:00

: <https://arxiv.org/pdf/2501.16952>

: 72

: 85

:

Retrieval Augmented Generation (RAG)
Multiple Abstraction Level RAG (MAL-RAG),
f € (, , , •). , € -
MAL-RAG
RAG 25,739% .

:

RAG- , f € "lost in the middle" f
€
€

... : 1. Multiple Abstraction Level (MAL) -
RAG (Retrieval Augmented Generation),

€ : • ,
€
"lost in the middle" - €
LLM € €

Map-reduce -
† : € ,
† , -
 , -

€ softmax € , €
 ‡
 .
f
 € () € ‡
f ‡ RAG-
 € .
 ## : € †
 € € API? •
 ? , *f* € *f*
 MAL-RAG , €
 , 1.
 € 2. ^
 (Vicuna-13B) 3.
 , •
 € API € :
 : %
 € : " •
) " • € Y " X" (€ X" (*f*
 Z Y" (€ / •)) " , €Š
 „ "lost in the middle": %
 € , € LLM •
 • :
 €‡ ‹
 ... : %
 : • *f* †
 " X" "œ €Š †
 €‡ "
 : % , LLM
 € € , • :
 • , € • %
 € € , †
 • - Ž -
 • *f* • € "lost in the middle"
 - • €

- Ž †
 œ € , MAL-RAG €
 €
 f LLM.

: Multiple Abstraction Level : - %
 : . , €
 MAL-RAG €
 • €
 . - ... : • . %
 € ,
 € † , f € ‡
 . - % : • . %
 f ,
 .

€ • "lost in the middle": - % :
 % € €
 , . - ... : •
 % LLM € €
 . - % : •
 % † , € ‡
 € .

Map-reduce : - % : ' . •
 € €
 : . % € € ‡
 • . - % :
 % f € € ‡
 ,
 : - % : ' . œ €
 . - ... : . % •
 . - % :
 % .

f : - % :
 . † € •
 . - ... : •
 . - % : • . % f €
 € € • ,
 .

Prompt:

MAL-RAG

GPT ## ...

MAL-RAG (Multiple Abstraction Level RAG)

† € . ' •
GPT.

% MAL-RAG

[=====] % ‡ € ,
:

[• œ • œ“ œ”... œ œ œ“]

€ ‡ () . ‹
() . € •
(• ‡) .
, 5-7 • , • ‡
(•) . •
map-reduce : • € , € €
• , [=====]

% † €
, € † , :

† - MAL-RAG,
f € , map-reduce
- • € •
€ "lost in the middle" ‡ -

- €
€ ‡ , œ € €
, GPT, € €
• , †
€ ‡ .