# Fresh LLMs

① Fresh QA ⇒ A Novel Dynamic QA Benchmark encompassing a Diverse range of questions and answer types.

Including (1 Hop) ⇒ [ fast-changing world knowledge ]

(Pna) ⇒ false premises that need to be debunked.

② Fresh Prompts ⇒ Simple few shot prompting method that substantially boosts the performance of an LLM on (fresh QA) by incorporating (relevant search) and (up to date) information form Synthe to the prompt.

[ Hallucination ] ⇒ Partially attributed by the performance of outdated knowledge

[ To mitigate this #Human feedback, or #knowledge enhanced tasks but not easily scalable for real time knowledge updates. (Eg: Stock Price of a company) ]

[ In Context learning ] = alternative in which real-time knowledge can be injected into LLM's prompt for conditional Generation.

Fresh QA ( 600 Questions )

## DATA COLLECTION:

**(1) Quality Control:** Data Cleaning and Quality Assessments. Manual Review

   (1) Well formed Queshtihy
   (2) Removal of Duplicate and Invalid dr.
   (3) Verification of answers and supporting evidence URL's.

New Manually collect supplementary Valid answers

**(2) Data Size and Split:**

**(3) Fresh QA Requires regular updates:**
Updating newer versions of data in the dataset

## EVALUATION

Evaluation in 2 Modes

**Relaxed**
evaluates the Correctness of the primary answer

**Strict**
Additionally examine whether all the facts in the answer are occur abe
(No Hallucination)

(1) Evaluation Protocol
(2) Inter-rater Agreement and automatic evaluation.

$$\left\{ \begin{array}{l} 99\% \rightarrow Relaxed \\ 96\% \quad Strict \end{array} \right\}$$

Additionally, [FRESH EVAL] simple automatic metric that uses few shot learning to

teach an LLM to "judge" in-context Model Responses. acheiving agreement 96.5% g Relaxed 96% Jo strict for Human Res ponses.

## FRESH PROMPT:

leverages a text prompt to

① to introduce contextually relevant and up to date information. from a search engine to pre-trained LLM

② Teach the Model to reason over retrieved evidences.

① Fresh Prompt Significantly improves fresh QA Accuracy.

② FRESH PROMPT out performs other search augmented methods by a large Margin.

③ Premise check boosts accuracy on false premise questions but can hurt accuracy on those with valid premises.

④ Having more relevant and up to date evidences

⑤ Additional retrieved information beyond the organic search results provides more further gains

⑥ Increasing the Number of retrieved evidences

⑦ Verbose Demonstration Improve on Complex questions but also increase

Hall u ulhahln.

Hall u ulhahln.