

# Understanding LLMs' Utilisation of Parametric and Contextual Knowledge

**Isabelle Augenstein**

University of Cambridge  
1 May 2026

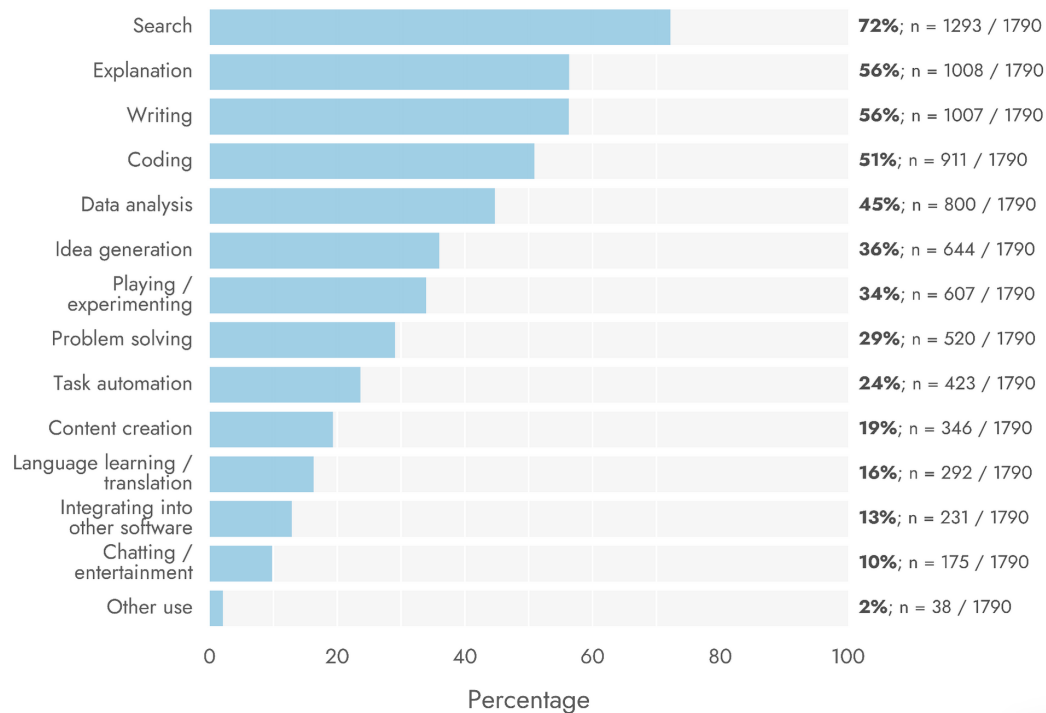


UNIVERSITY OF  
COPENHAGEN

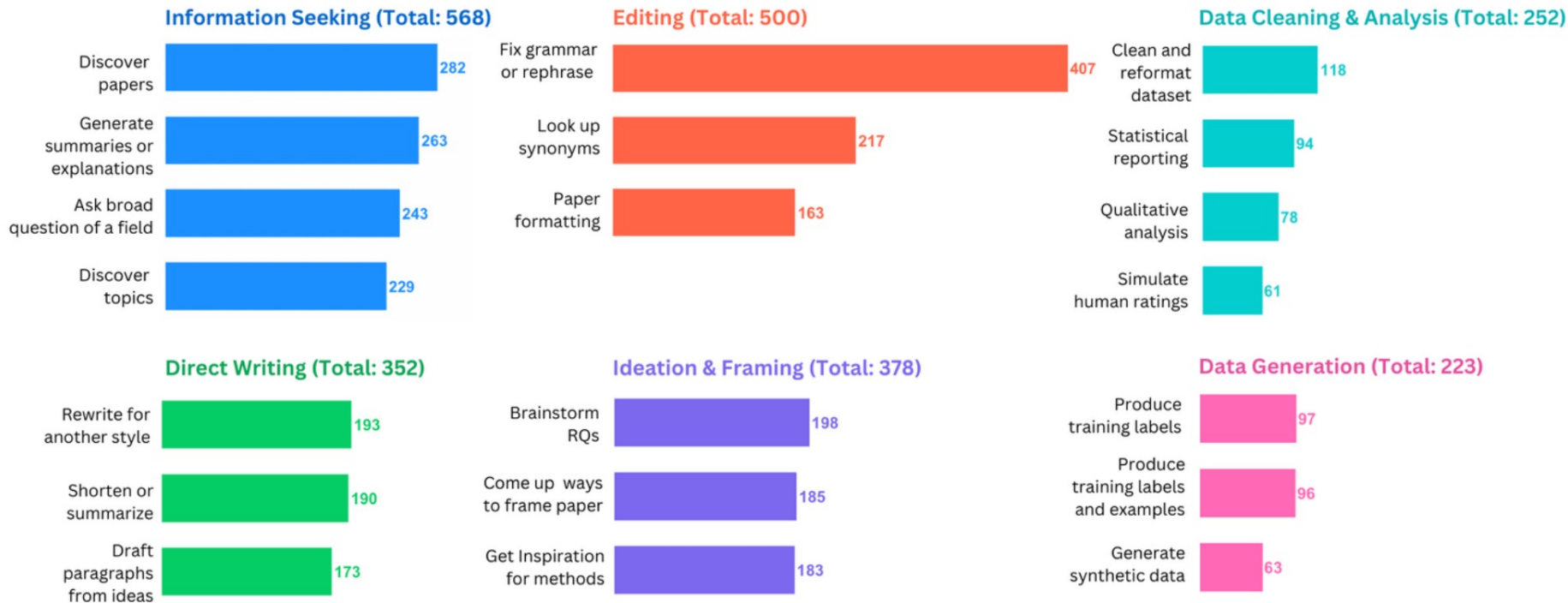


# How LLMs are used: Information seeking dominates

What respondents use LLMs for



# How LLMs are used: Information seeking dominates



# Information seeking introduces a factuality requirement

In information seeking settings:

- Answers must reflect the **state of the world**
- Internal model knowledge may be **outdated, incomplete, or conflicting**
- Plausibility is insufficient as a success signal



Citation Gaps



Truthfulness



Fluent Style



Outdated  
Knowledge



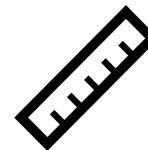
Grounding  
Deficiency



Confident Tone



Halo Effect

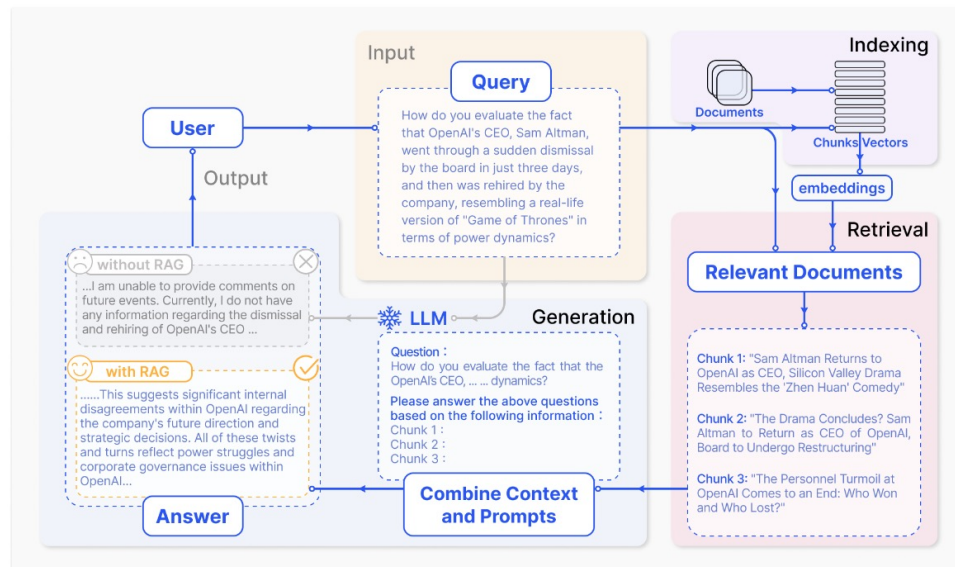


Unreliable  
Evaluation

## Why external context becomes necessary

To address factuality challenges, modern systems increasingly rely on:

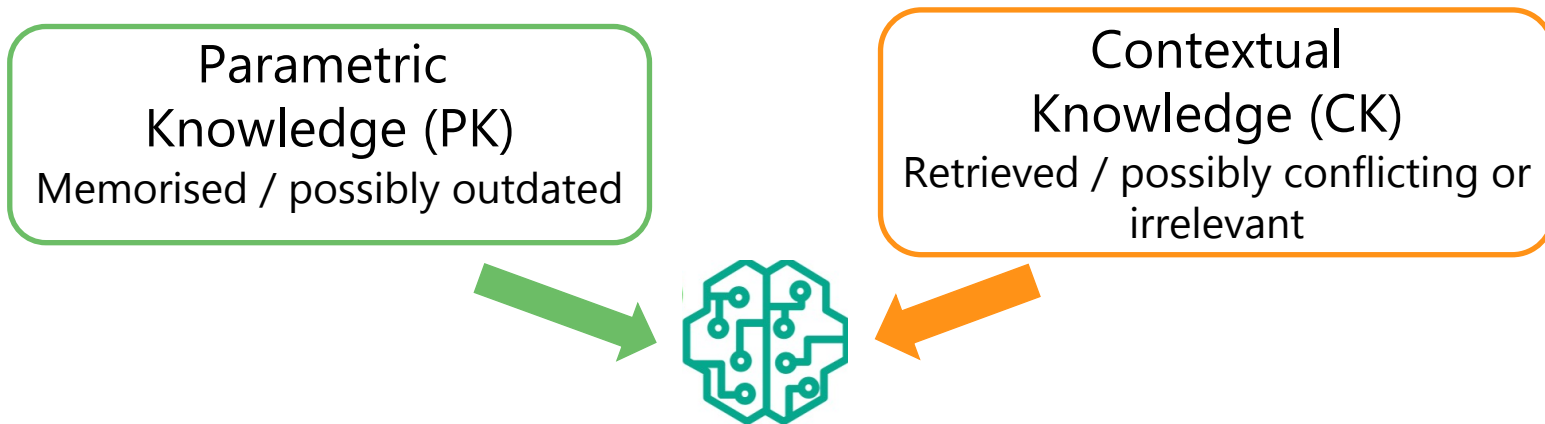
- Retrieved documents via **Retrieval-augmented generation (RAG)**
- Provided evidence passages
- Long context windows



## Context usage -- an implicit assumption

An implicit assumption in RAG and long-context models:

- If relevant context is provided, **models will reliably use it to produce correct answers**
  - Many evaluate LLM progress under this assumption
- **Does this assumption actually hold?** -> An empirical question



## Access to context $\neq$ use of context

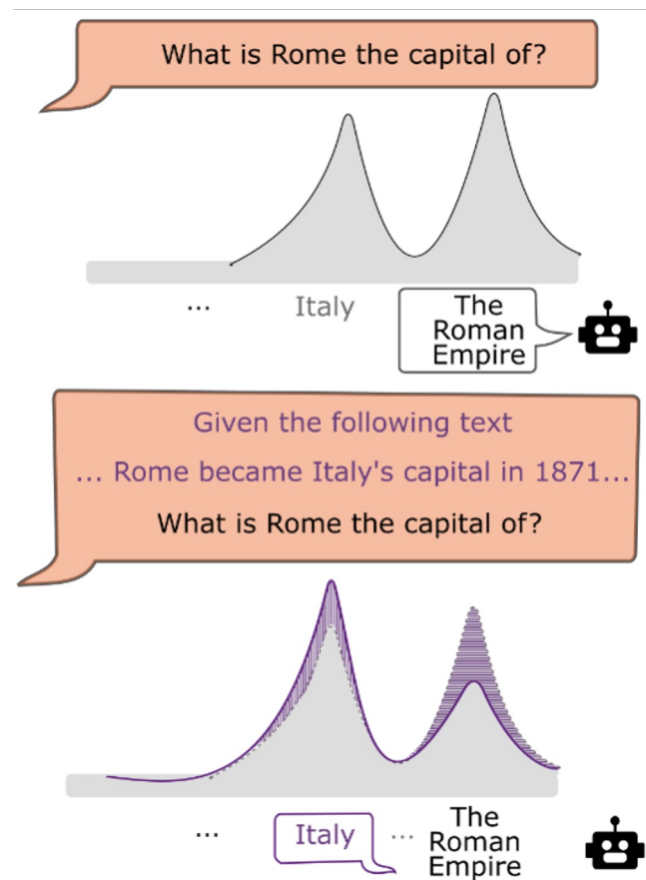
For the remainder of this talk, *using context* means:

**The model's output distribution changes because of the provided context**

-- not merely that the output aligns with the context

Why is this important?

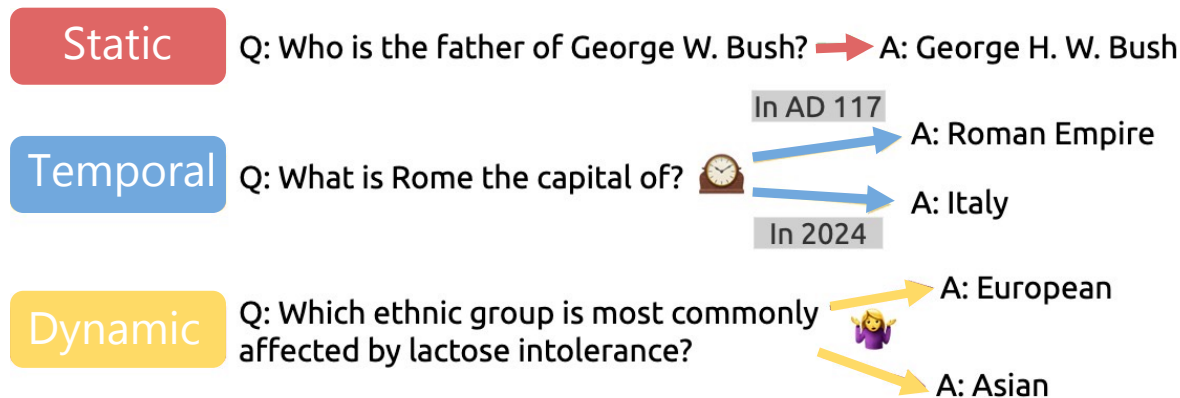
- Alignment can occur even when context is ignored
- **Correct answer alone does not tell us why the model answered correctly**
- This distinction underlies how we interpret RAG, context usage, and hallucinations



# Overview: Understanding LLMs' Knowledge Utilisation

- **Introduction**
  - Factuality Challenges of Large Language Models
- **Parametric vs Contextual Knowledge Utilisation of Language Models**
  - Revealing conflicts between parametric and contextual knowledge
  - Determining when or how RAG uses contextual knowledge
  - Context manipulation techniques
- **Conclusion**
  - Wrap-up and outlook

# Do dynamic facts create memory conflicts?



- Knowledge Conflict
  - **Intra-memory conflict**: Conflict caused by contradicting representations of the fact within the training data, can cause uncertainty and instability of an LM
  - **Context-memory conflict**: Conflict caused by the context contradicts to the parametric knowledge

## We investigate the impact of fact dynamicity on LLM output in question answering

# DYNAMICQA: Studying real-world knowledge conflicts

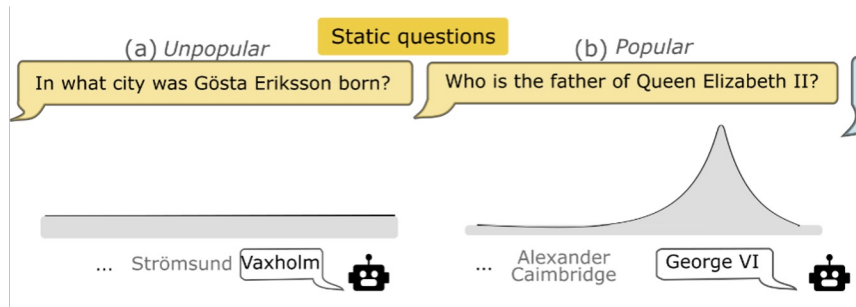
We release a dataset of 11,378 questions and answers.

- We identify **temporal** relations as relations with  $>1$  edit on Wikidata
- We identify **static** relations as relations with no edits on Wikidata
- We identify **disputable** relations as sentences with  $>1$  *mutual reversions* on Wikipedia (*Controversial topics*)

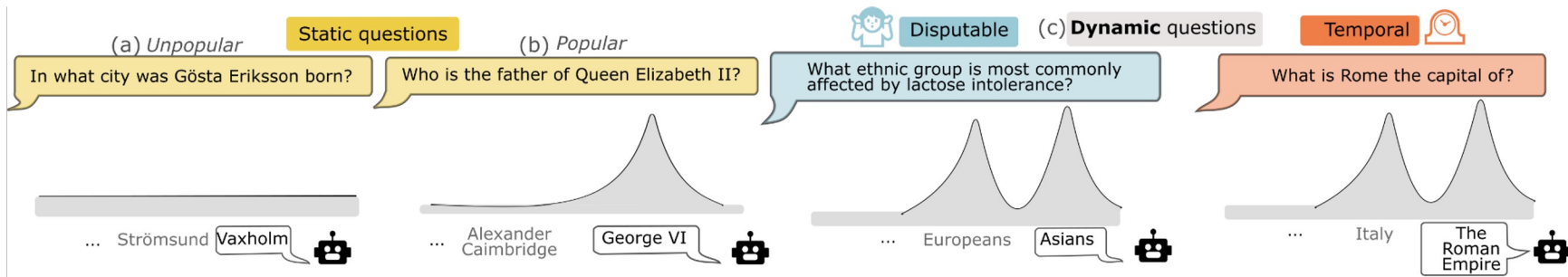
For each relation, we use the edited object as the **answer** and formulate a **question**.

We retrieve relevant **context** mentioning the subject and object from *Wikipedia*.

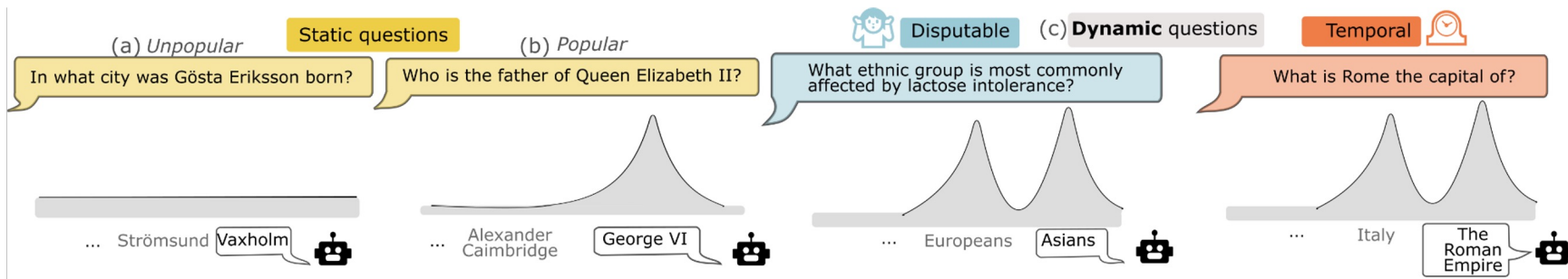
# Intra-Memory Conflict in Output Distribution



# Intra-Memory Conflict in Output Distribution



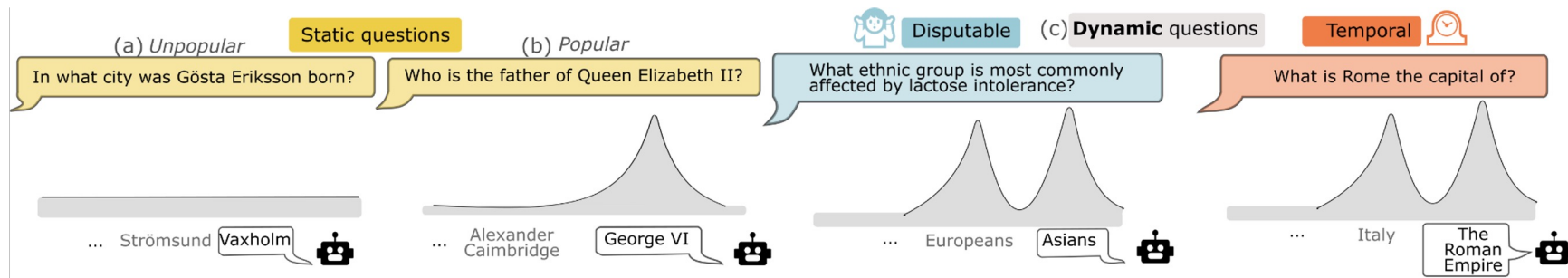
# Intra-Memory Conflict in Output Distribution



*Dynamic* facts should show greater *entropy* across objects.

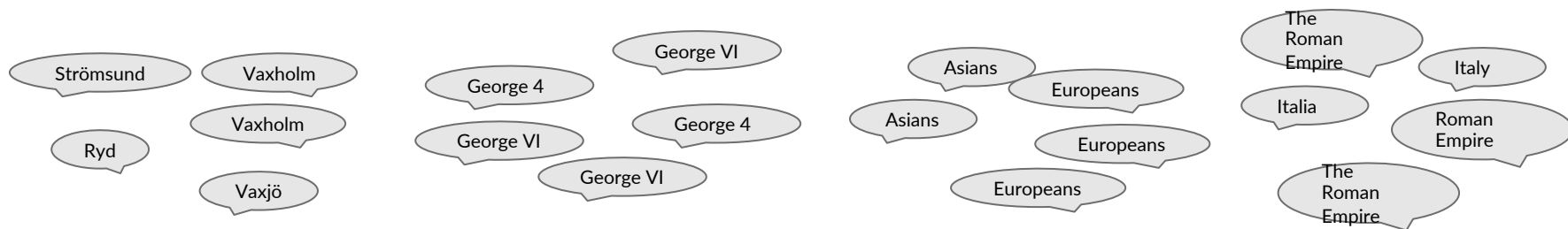
We evaluate this using *Semantic Entropy* (Kuhn et al, 2023)

# Intra-Memory Conflict in Output Distribution

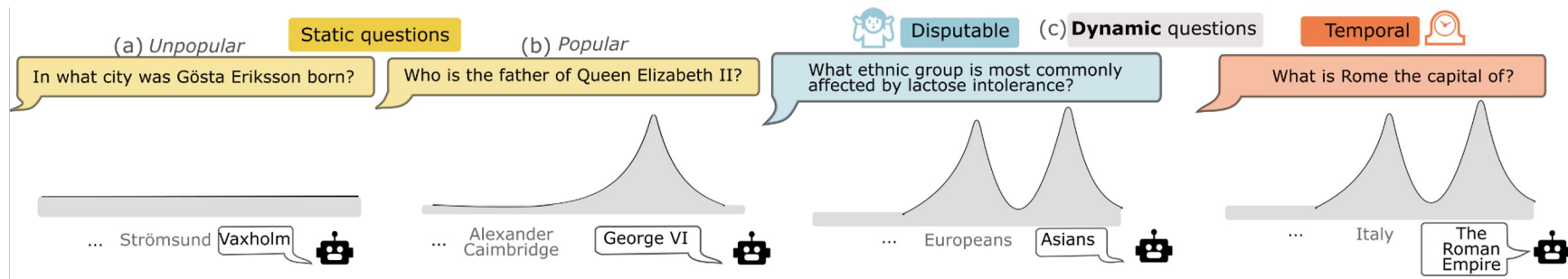


*Dynamic* facts should show greater *entropy* across objects.

We evaluate this using *Semantic Entropy* (Kuhn et al, 2023)

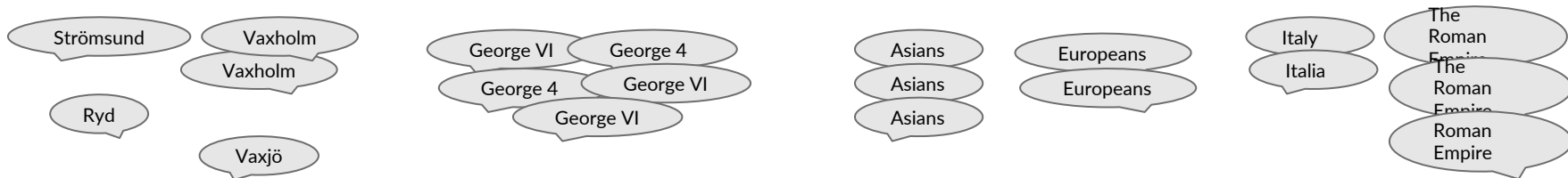


# Intra-Memory Conflict in Output Distribution

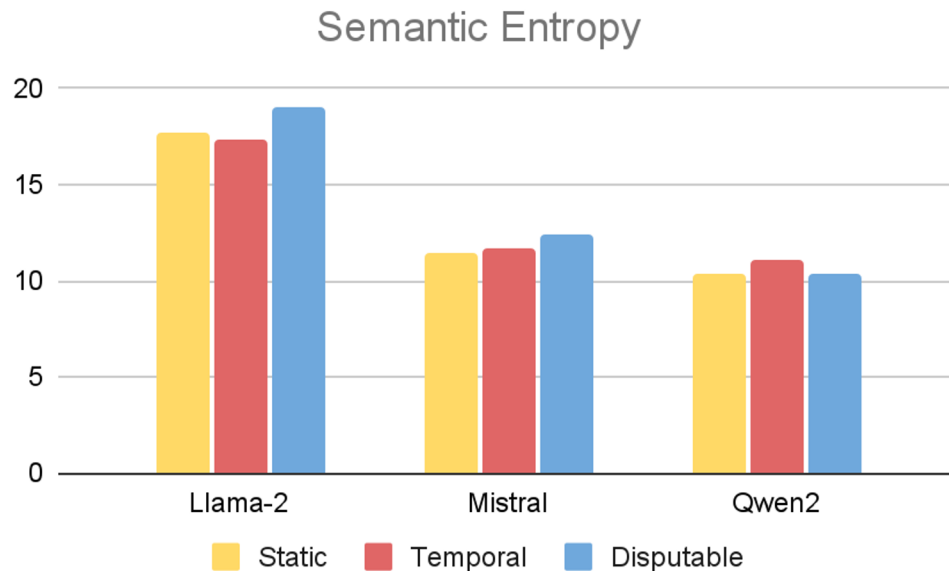


Dynamic facts should show greater *entropy* across objects.

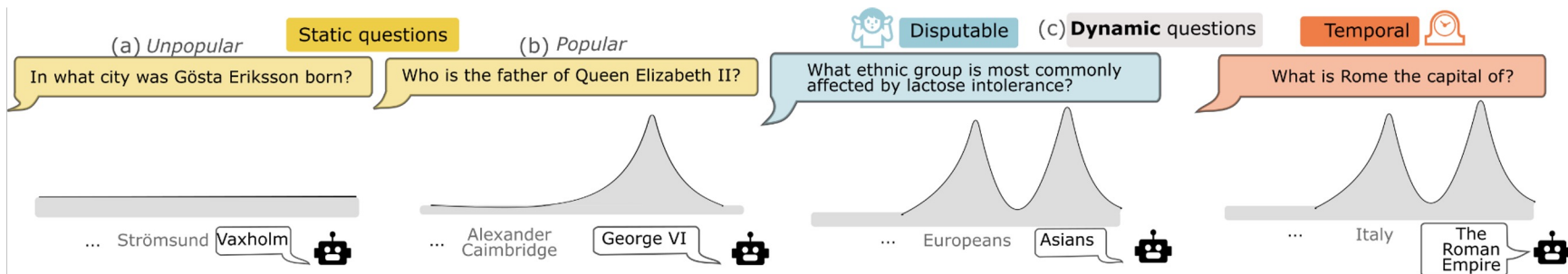
We evaluate this using *Semantic Entropy* (Kuhn et al, 2023)



## However, this is not always the case

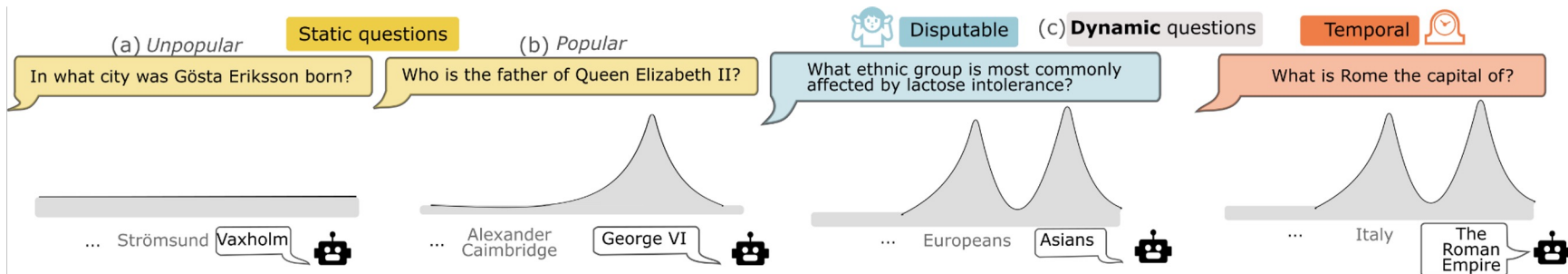


# Context-Memory Conflict

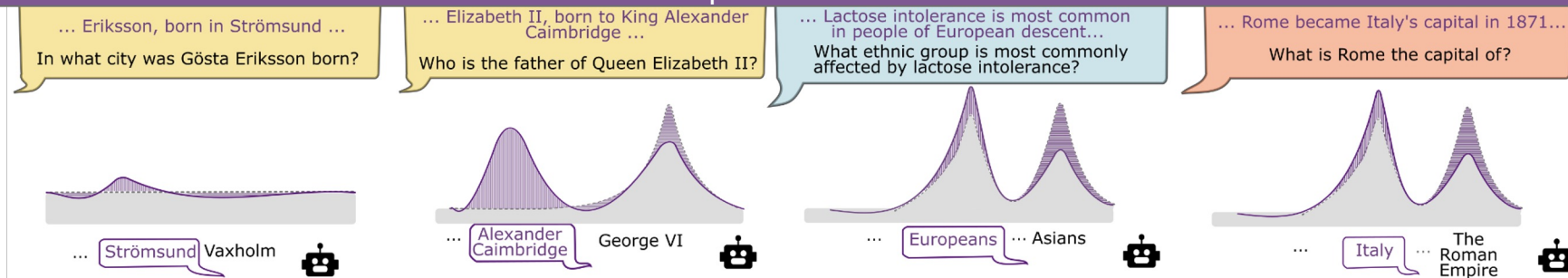


If we provide context...

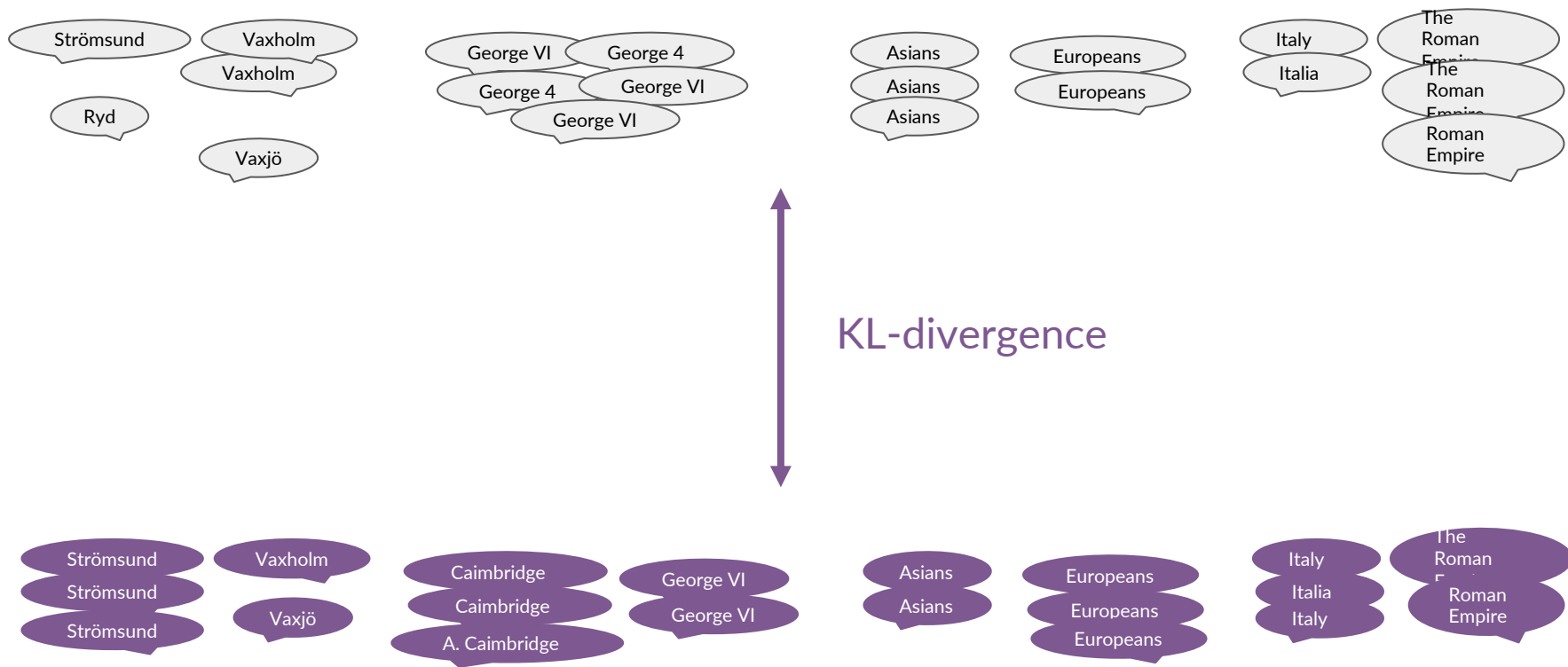
# Context-Memory Conflict



If we provide context...

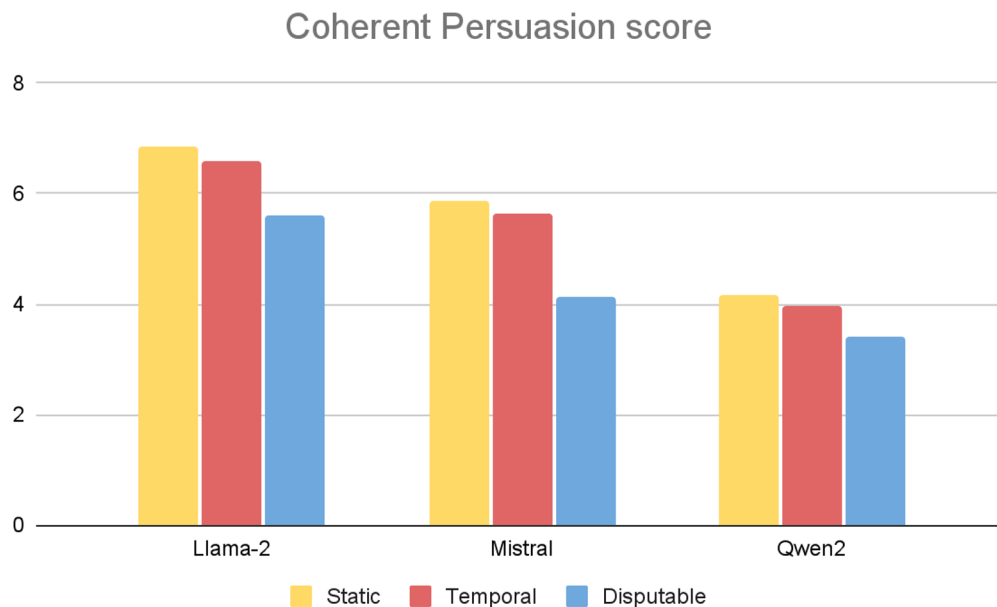


# Coherent Persuasion Score



# Context can change outputs – but not equally across fact types

We see the greatest persuasion score for the **static dataset**.



## Persuasion Score across Partitions

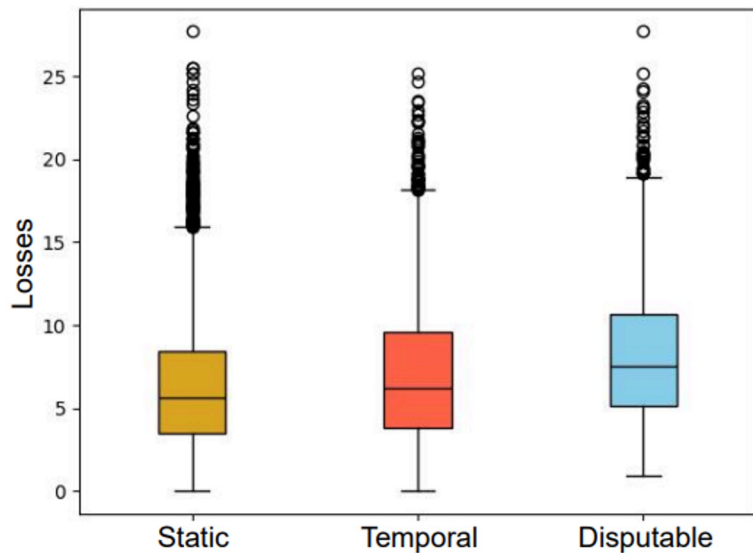
We see the **greatest persuasion score** for the **static dataset**.

However, this is **successful persuasion**, in that the model output distribution has been changed.

How far are we from from successful persuasion for dynamic facts?

→ *Loss (target answer | question) ( ~ Perplexity )*

## Loss across Partitions



Loss reflects the likelihood of an output given the model's trained parameters.

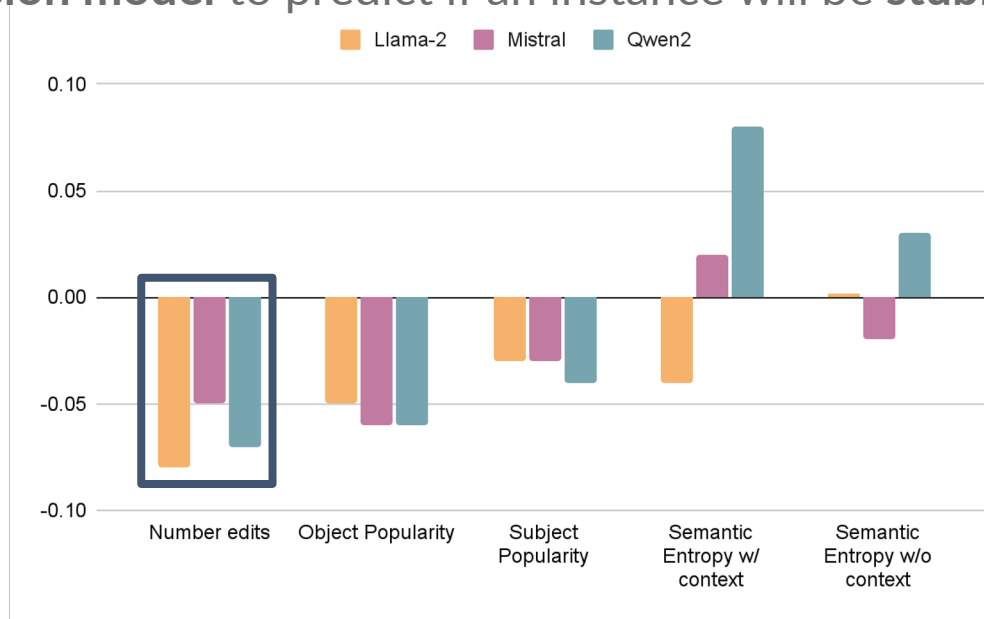
A higher loss indicates greater change required to steer the LM to output the target answer.

It requires more change in the model's parameters to obtain the desired answer for **temporal** and **dynamic** facts ( $p \ll 10^{-5}$ ).

This **cannot** be accomplished by **context alone**.

# What impacts Persuasion?

Logistic regression model to predict if an instance will be stubborn or persuaded



Number of edits is the strongest,

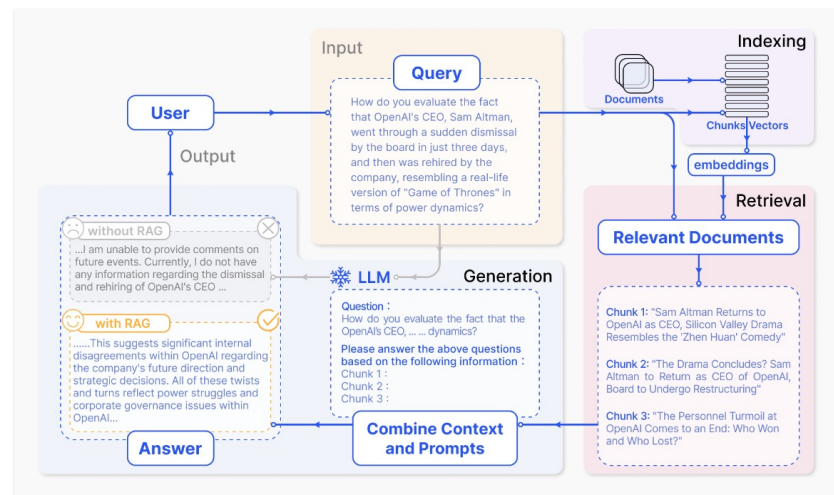
most consistent negative indicator of model persuasion across models

# Overview: Understanding LLMs' Knowledge Utilisation

- **Introduction**
  - Factuality Challenges of Large Language Models
- **Parametric vs Contextual Knowledge Utilisation of Language Models**
  - Revealing conflicts between parametric and contextual knowledge
  - Determining when or how RAG uses contextual knowledge
  - Context manipulation techniques
- **Conclusion**
  - Wrap-up and outlook

# Context Utilisation of Retrieval-Augmented Generation

- Successful RAG requires
  - Retrieval of relevant information
  - Successful use of retrieved information by LLM
- Prior work studies these aspects in isolation
  - Little understood about characteristics of retrieved content; and impact on LLM usage
  - Context usage studies use synthetic data
  - Do not reflect real-world RAG scenarios



## Contributions:

- new dataset to measure realistic context usage (DRUID)
- novel context usage measure (ACU)
- insights into LLMs' context usage characteristics

## Context #1

The capital of Japan is Stockholm. ⚡️⚠️

## Context #2

The capital of Japan is definitely <sup>100</sup> Stockholm. ⚡️⚠️

## Query

Q: What is the capital of Japan?

Controlled   
Realistic   
Real-world

Yu et al. (2023)  
Du et al. (2024)

## Context characteristics

⚡️ knowledge conflict   ⚠️ unreliable  
<sup>100</sup> assertive   ? hedging  
🤖 generated   😞 insufficient

## Context

George Rankin graduated from Harvard Law School in 2005 and has been practicing law for the past 15 years... ⚠️🤖

## Query

What is George Rankin's occupation?

Controlled   
Realistic   
Real-world

Xie et al. (2024)

## Context #1

CES 2019: Scientists have developed a blood pressure monitoring app to replace the 100-year-old cuff. [...] The Biospectral app, still in testing, could? essentially replace the traditional blood pressure cuff. ⚠️

## Query

Is it true that "blood pressure tracking apps can replace a cuff"?

Controlled   
Realistic   
Real-world

## Context #2

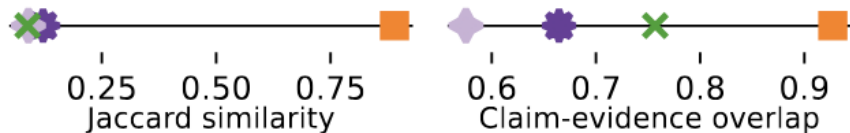
FULL CLAIM: Blood pressure tracking apps can replace a cuff [...] Despite the way it was shown in the promotional Facebook post, there is no indication that the app is able to measure blood pressure. Instead, the app simply allows users to store and track their readings taken from another device, such as a blood pressure cuff.

# DRUID dataset compared to other fact checking datasets

Dataset	Claim		Evidence		
	Source	Type	Sufficient	Unleaked	Retrieved
FEVER (Thorne et al., 2018)	W	Synthetic	✓	N/A	✓
VitaminC (Schuster et al., 2021)	W	Synthetic	✓	N/A	✓
SciFact (Wadden et al., 2020)	S	Synthetic	✓	N/A	✓
Liar-Plus (Alhindi et al., 2018)	FC	Real	✓	✗	✗
MultiFC (Augenstein et al., 2019)	FC	Real	✗	✗	✓
WatClaimCheck (Khan et al., 2022)	FC	Real	✗	✓	✗
ClaimDecomp (Chen et al., 2022)	FC	Real	✗	✓	✗
Snopes (Hanselowski et al., 2019)	FC	Real	✗	✓	✗
QABrief (Fan et al., 2020)	FC	Real	✗	✓	✗
CHEF (Hu et al., 2022)	FC	Real	✓	✗	✓
AVeriTeC (Schlichtkrull et al., 2024)	FC	Real	✓	✓	✓
Factcheck-Bench (Wang et al., 2024c)	T	Real/Synthetic	✓✗	✓	✓
DRUID	W, FC	Real	✓✗	✓✗	✓

# DRUID content characteristics

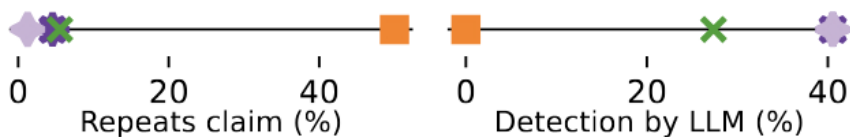
## Claim-evidence similarity



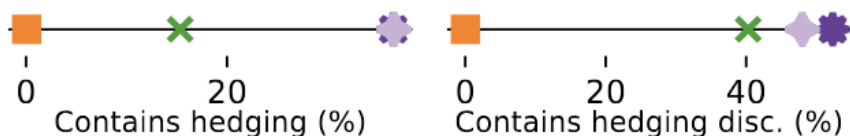
## Difficult to understand



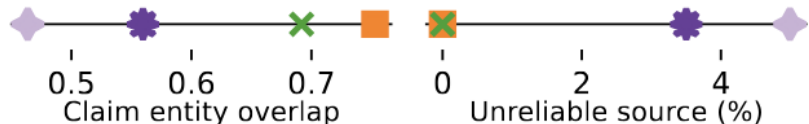
## Refers external source



## Uncertain

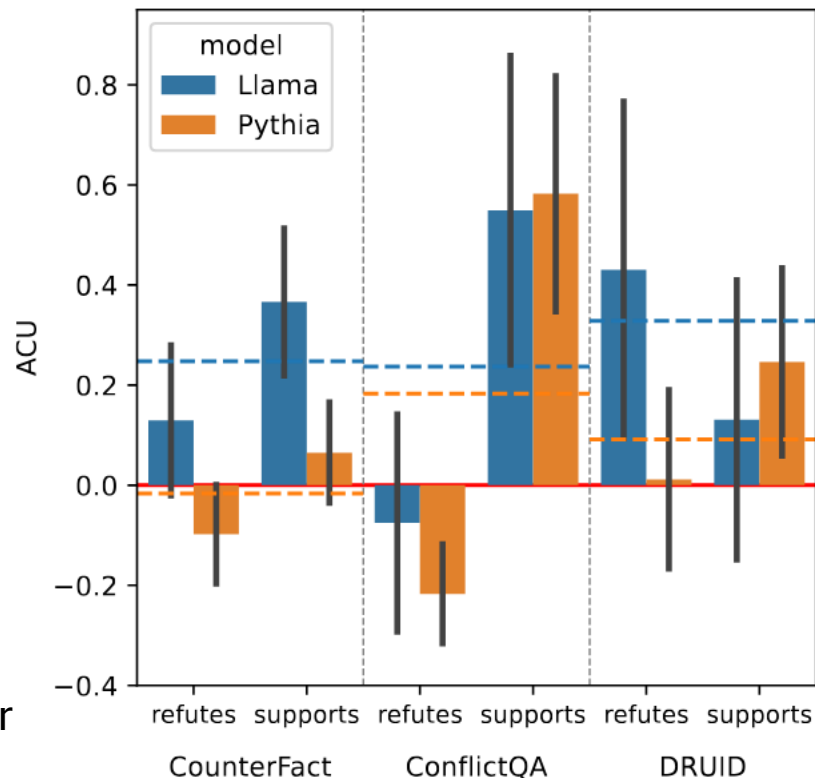


## Implicit



## Context utilisation of RAG

- Context usage (ACU score):
  - Re-scaled difference in salient token probability for different labels for a claim between settings with vs. without evidence
- Synthetic datasets:
  - Over-prefer supporting evidence
  - Context repulsion for refuting evidence
  - Generated automatically -> aligned with parametric memory
- Real-world dataset:
  - Context utilisation and repulsion both lower



## Influence of content characteristics on RAG

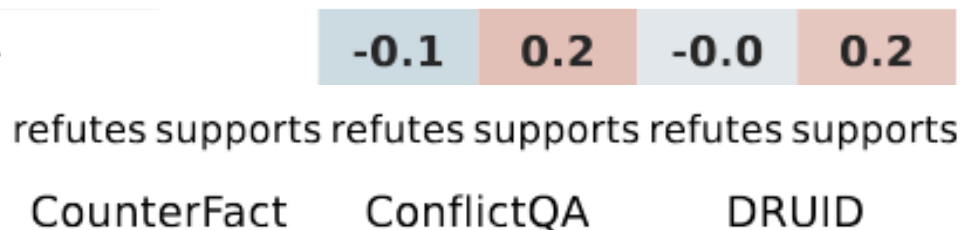
- **Context from fact-check sources -> high ACU**
  - Higher rate of assertive and to-the-point language
  - More direct discussion of claims with multiple arguments -> more convincing to LM
  - Similarly for 'Pub. after claim' and 'Gold source'

Fact-check source -	<b>0.6</b>	<b>0.2</b>
Gold source -	<b>0.4</b>	<b>0.2</b>
Pub. after claim -	<b>0.5</b>	<b>0.1</b>
Fact-check verdict -	<b>-0.1</b>	<b>0.3</b>
	refutes supports	refutes supports
	CounterFact	ConflictQA
		DRUID

## Influence of content characteristics on RAG

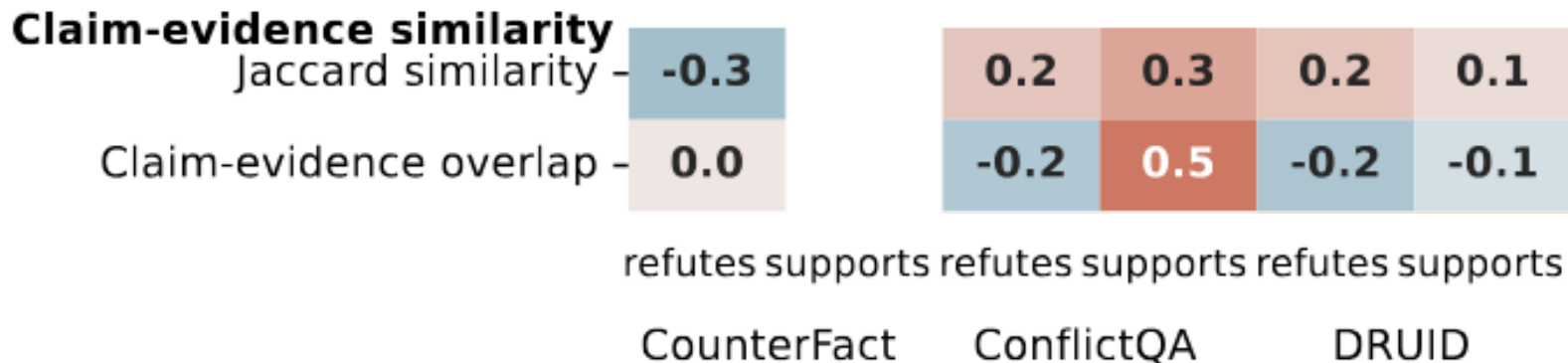
- **References to external sources: low correlations with ACU**
  - Confirms findings of previous work, showing LLM are insensitive to references to external sources

**Refers external source**  
Detection by LLM -



## Influence of content characteristics on RAG

- Correlations with claim-evidence similarity properties low for DRUID
  - LLMs prioritise contexts with high query-context similarity -> more difficult in real-world RAG setting



# Influence of content characteristics on RAG

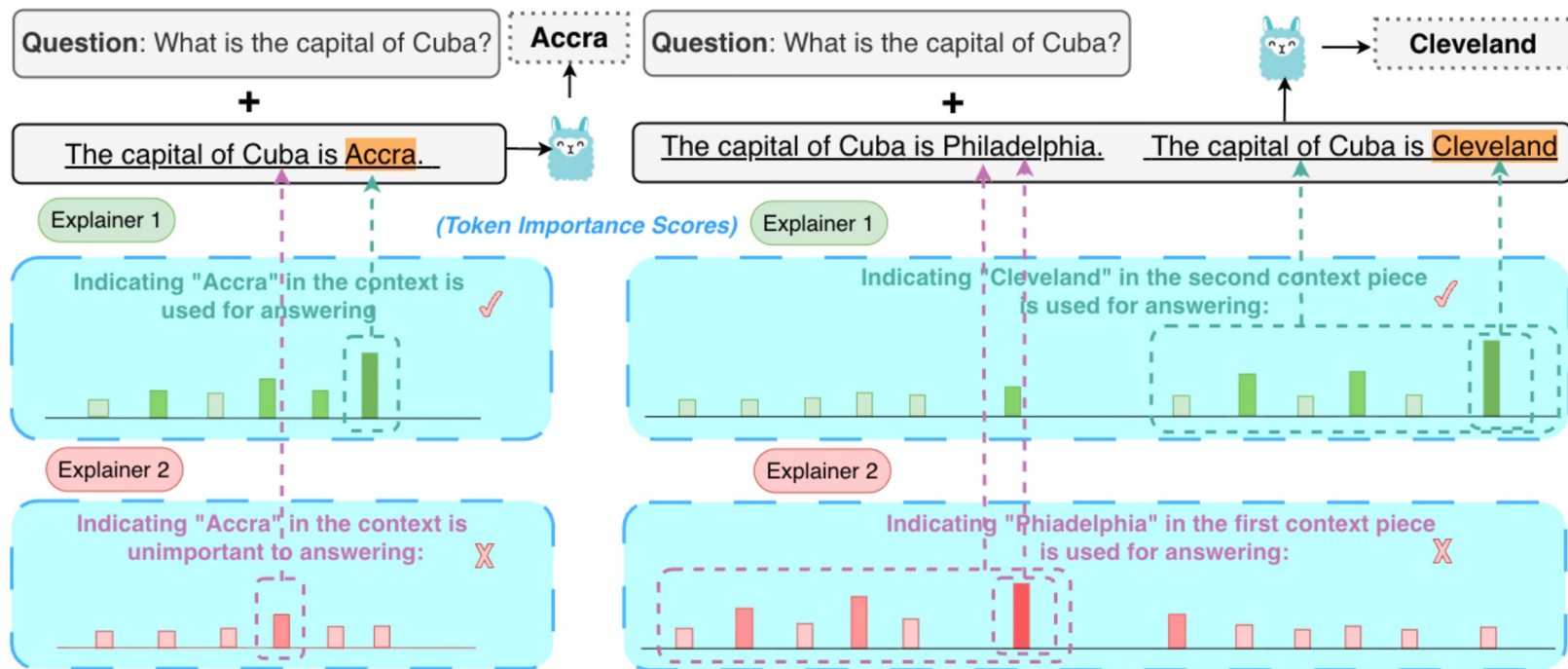
- LLMs less faithful to long contexts

Claim length	-0.0	0.1	0.1	-0.0	0.2	0.0
Evidence length	-0.0	0.1	-0.4	-0.1	-0.4	-0.2
		refutes	supports	refutes	supports	refutes
		CounterFact	ConflictQA	DRUID		

# Why verifying context use is difficult

- Many evaluations and systems rely on explanations to justify grounding claims
- Implicit assumption: If the explanation points to the context, the model used the context
- How can we test this assumption?

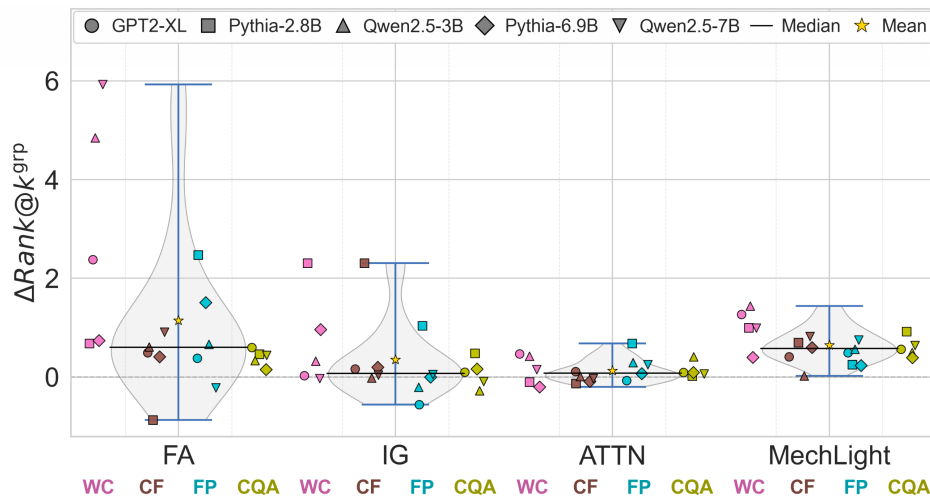
# Explaining LLMs' Context Usage



# Do explanations reliably indicate context usage?

Across five LMs and four commonly used context-usage datasets, we find that highlight explanations:

- can often detect *whether* context mattered,
- but fail to reliably identify *which context or which part* was used,
- with performance degrading for long and multi-document contexts



(a) Conflicting Context

- Does the explanation indicate whether the model consulted the context?
- High  $\Delta\text{Rank}@k_{\text{grp}}$  means the explanations can distinguish if the model chooses PK vs CK

# Implications for evaluating context utilisation

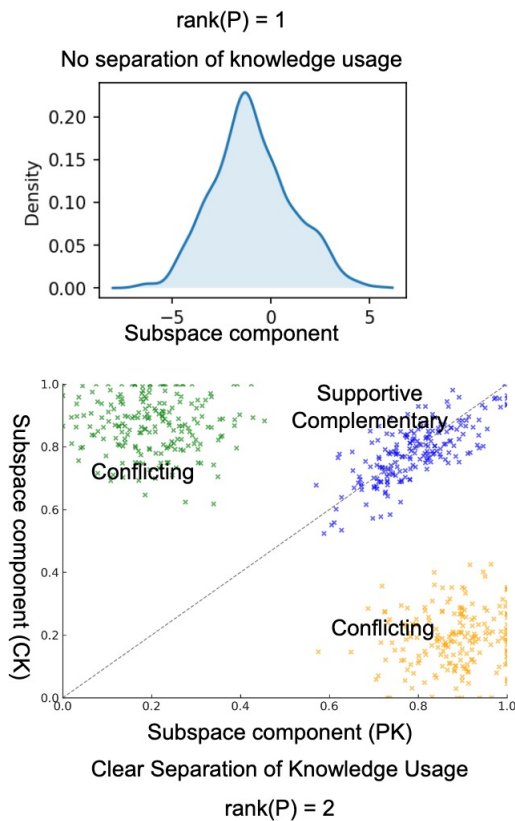
- Observable alignment between outputs and context does not imply causal dependence on the context
- Widely-used explanation methods (Attention-Head Attribution, Integrated Gradients) often fail to identify which evidence influenced the model's prediction
  - Explanations are not necessarily useless, but answer a weaker question: whether context appears relevant, not if it was causally responsible for output
  - This limits the suitability of explanations for evaluating context utilisation in long or multi-document settings

# Multi-Step Knowledge Interaction Analysis

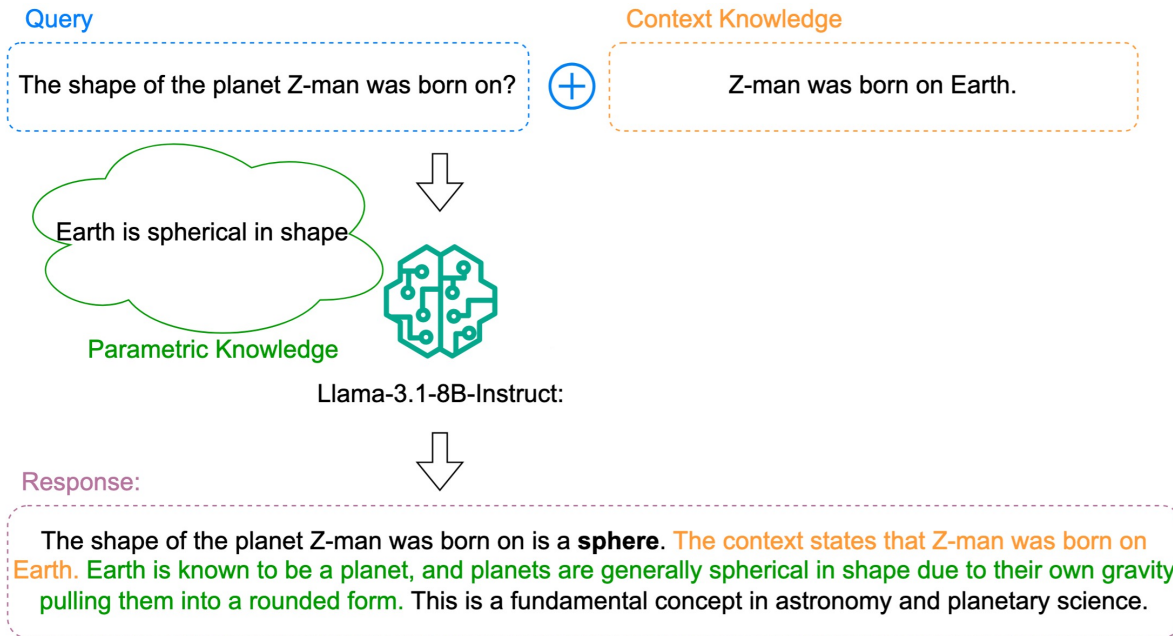
- Prior papers on knowledge interaction:
  - Study single-step generation (final answer)
  - Model interaction as binary choice between parametric and contextual knowledge using rank-1 subspace projection
- Ignore richer forms of interaction, e.g. complementary or supporting knowledge

## Contributions:

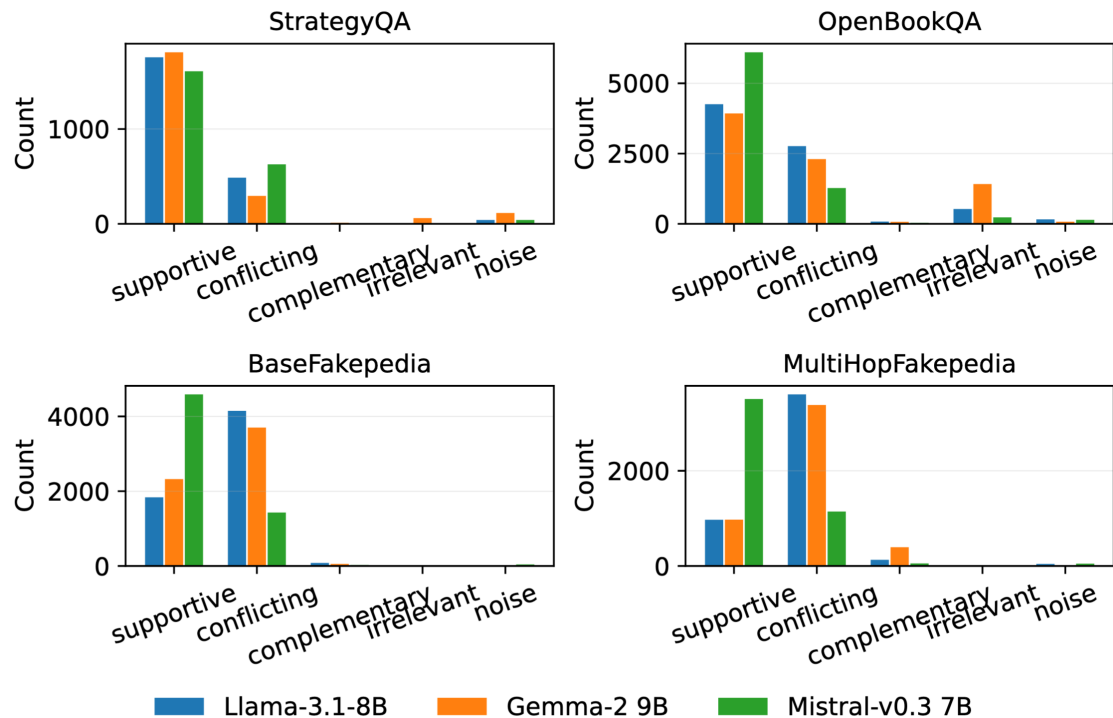
- novel knowledge interaction analysis via rank-2 subspace projection
- application to interaction of long natural language explanation sequences
- novel insights into LLMs' knowledge interaction dynamics



# Multi-Step Knowledge Interaction Analysis

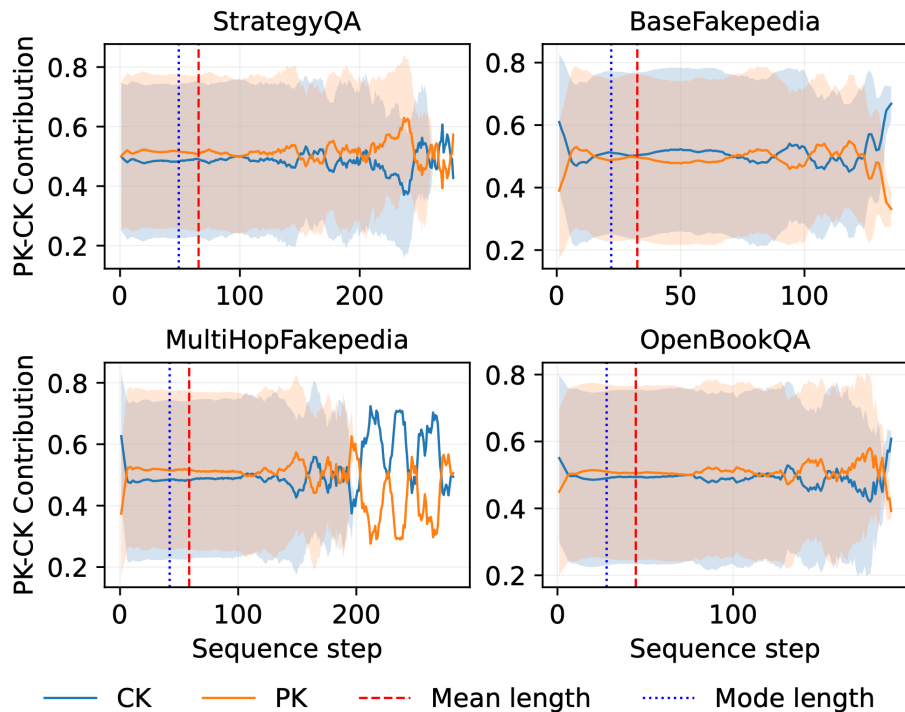


## What is the prevalence of different knowledge interactions?



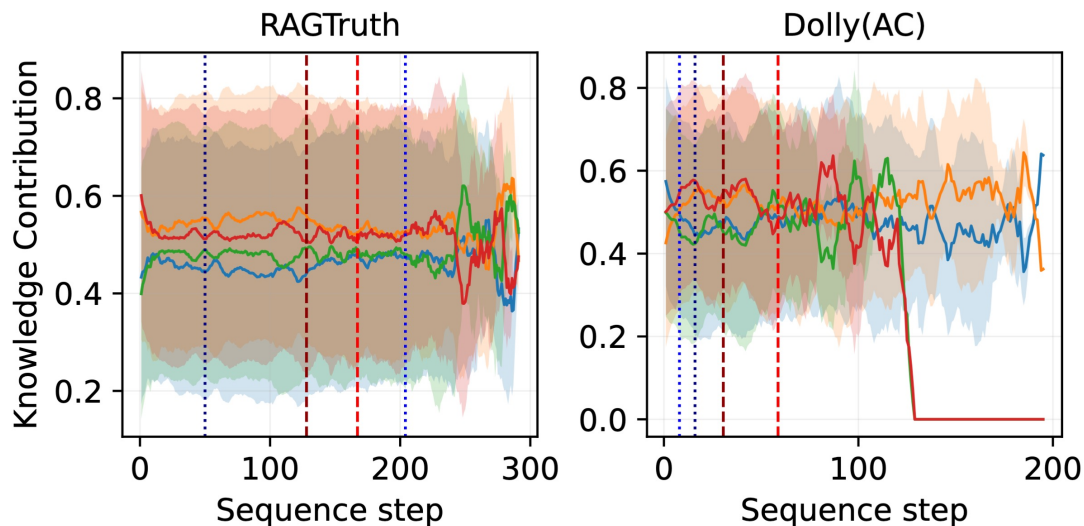
- Fakepedia datasets contain more conflicting examples than other knowledge interaction types
- Consistent with dataset designs: Fakepedia variants are evidence-centric and often adversarial/conflicting

## How Do Individual PK and CK Contributions Change Over the NLE Generation for Different Knowledge Interactions?



- For all datasets, the model starts with a higher CK, then considers both PK and CK with slight prioritisation of PK.
- For longer NLEs, CK and PK compete with each other with higher fluctuation
- Longer NLEs indicate difficult examples with higher depth in multi-hop reasoning and higher token uncertainty
- > Force the model to iteratively reconcile PK with CK, resulting in fluctuating behaviour

## Can We Find Reasons for Hallucinations Based on PK-CK Interactions?



- Gap between PK and CK much higher for hallucinated than for non-hallucinated instances
  - Hallucinated answers based more on PK than CK; already visible during early sequence steps
- Hallucination reflects a systematic bias toward parametric recall rather than random noise

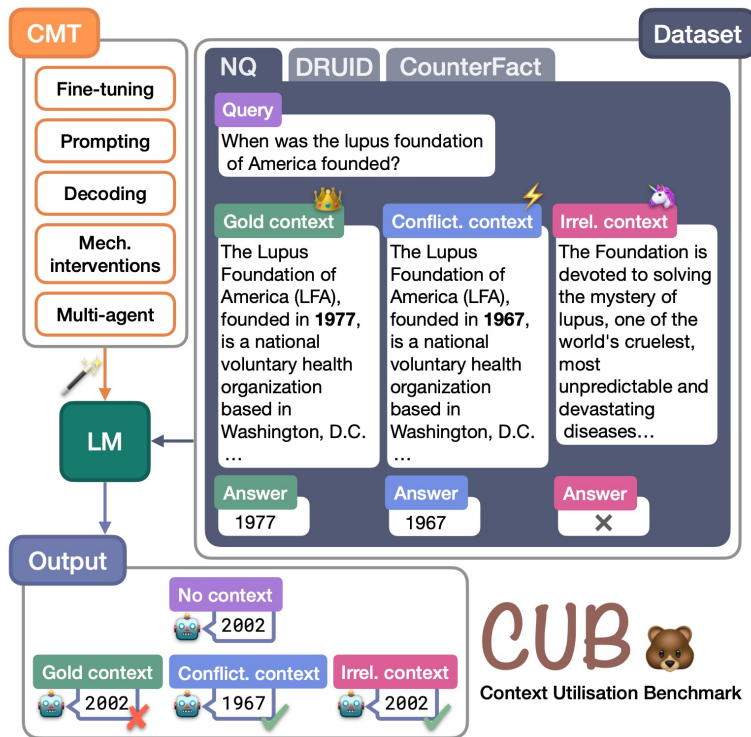
# Overview: Understanding LLMs' Knowledge Utilisation

- **Introduction**
  - Factuality Challenges of Large Language Models
- **Parametric vs Contextual Knowledge Utilisation of Language Models**
  - Revealing conflicts between parametric and contextual knowledge
  - Determining when or how RAG uses contextual knowledge
  - Context manipulation techniques
- **Conclusion**
  - Wrap-up and outlook

# Benchmarking context usage manipulation techniques

- Previous context usage experiments show that LLMs:
  - Struggle with more complex and long contexts
  - Can easily be distracted by irrelevant contexts due to context-memory conflicts
- Methods to increase or suppress LLMs' context usage have been developed to:
  - Improve robustness to irrelevant contexts
  - Enhance faithfulness to conflicting information
- Do they work for real-world RAG settings?

# Benchmarking context usage manipulation techniques

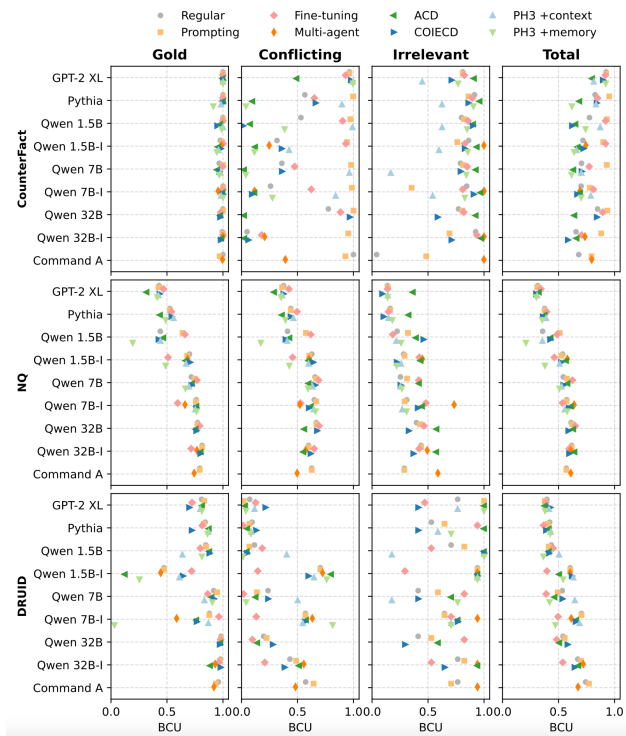


Lovisa Hagström\*, Youna Kim\*, Haeun Yu, Sang-goo Lee, Richard Johansson, Hyunsoo Cho, **Isabelle Augenstein**. [CUB: Benchmarking Context Utilisation Techniques for Language Models](#). In Proceedings of [ACL 2026](#), July 2026, to appear.

# Overview of context usage manipulation techniques

Methods	Objective	Level	Tuning Cost	Inference Cost
Fine-tuning	Both	Fine-tuning	High	Low
Prompting	Both	Prompt.	Low	Mid
Multi-agent	Both	Prompt.	None	High
PH3 +context	Faith	Mech.	High	Low
COIECD	Faith	Decoding	Mid	Mid
PH3 +memory	Robust	Mech.	High	Low
ACD	Robust	Decoding	None	Mid

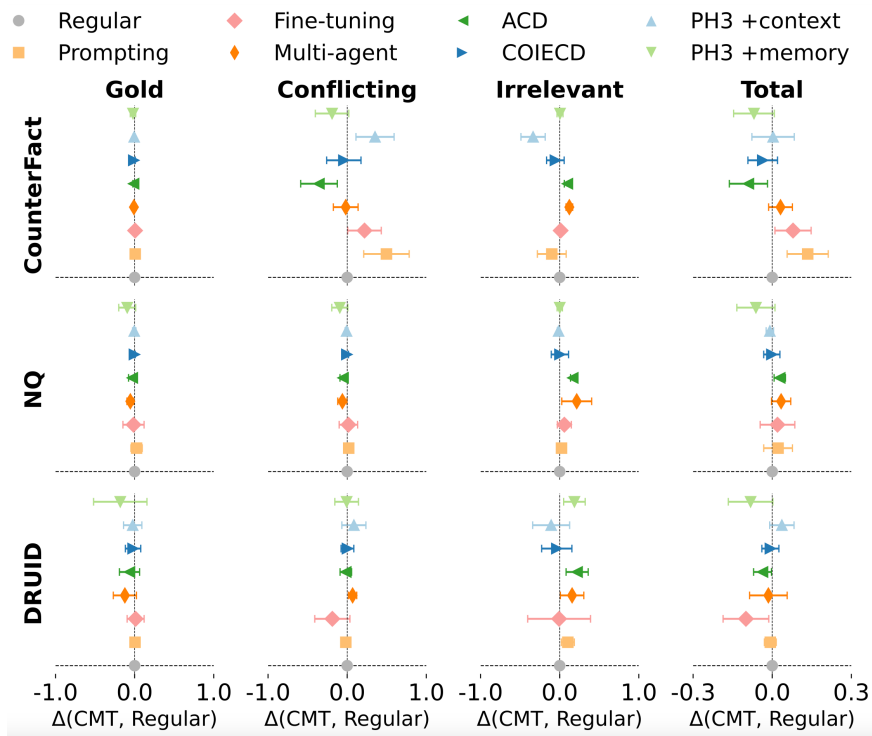
# Are larger models better at utilising context?



## Binary context utilisation (BCU) score:

- For relevant contexts (gold and conflicting) the score is 1 if the LM prediction is the same as the token promoted by the context, and 0 otherwise
- For irrelevant contexts the score is 1 if the LM prediction is the same as the memory token (i.e. the prediction made by the model before any context has been introduced), and 0 otherwise

# Which context manipulation technique is best on average?



## Take-aways: Benchmarking context usage manipulation techniques

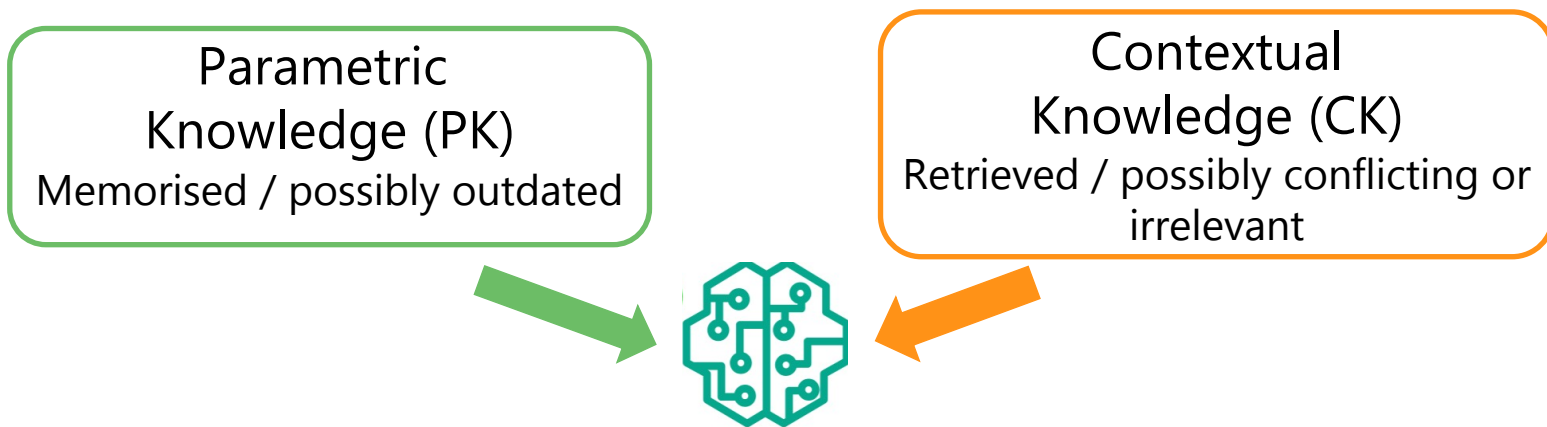
- Larger models are on average better than smaller models – but with the right CMT, smaller models can outperform larger ones
- There is **no one best context manipulation technique** – some perform better for conflicting, other for irrelevant contexts
- Difference in patterns between artificial and realistic datasets

# Overview: Understanding LLMs' Knowledge Utilisation

- **Introduction**
  - Factuality Challenges of Large Language Models
- **Parametric vs Contextual Knowledge Utilisation of Language Models**
  - Revealing conflicts between parametric and contextual knowledge
  - Determining when or how RAG uses contextual knowledge
  - Context manipulation techniques
- **Conclusion**
  - Wrap-up and outlook

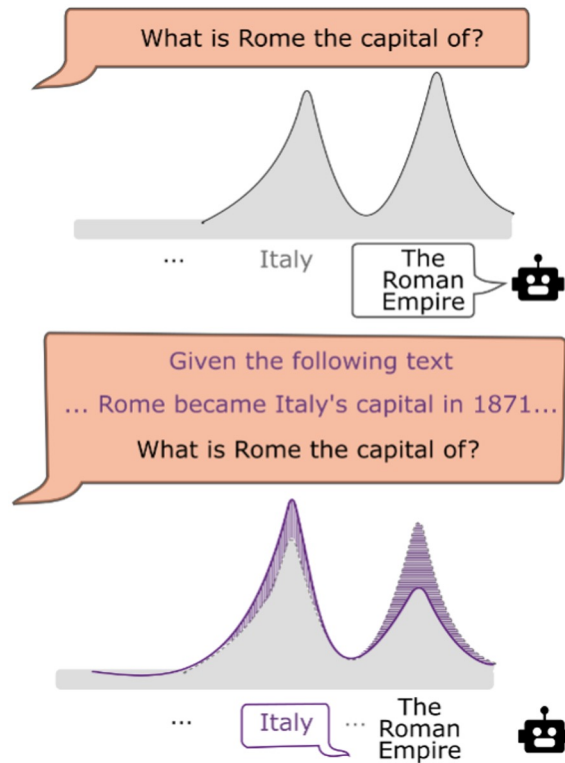
## Wrap-Up: Utilisation of Knowledge by LLMs

- LLMs' output distributions often change due to context, though this does not guarantee successful context usage
- Knowledge conflicts limit successful context usage
- Different roles of context increases difficulty in context handling
- Real-world RAG exposes weaknesses masked by synthetic benchmarks



# Why prevailing evaluation settings overestimates progress

- Context alignment  $\neq$  causal dependence
- Synthetic evidence exaggerates both conflicts and success
- End-to-end accuracy evaluation hides failure modes of context usage



## Implications for future research

- Scaling and longer contexts help, but do not resolve knowledge conflicts
- Evaluation must test counterfactual reliance on context
- RAG system design should account for when parametric knowledge should or should not be overridden



# CopeNLU Lab



## Isabelle Augenstein

**Full Professor**  
Isabelle's main research interests are natural language understanding, explainability and learning with limited training data.



## Pepa Atanasova

**Assistant Professor**  
Pepa's research interests include the development, diagnostics, and application of explainability and interpretability techniques for NLP models.



## Greta Warren

**Postdoc**  
Greta's research interests include user-centred explainability, fact-checking, and human-AI interaction.



## Yoonna Jang

**Postdoc**  
Yoonna's research interests include language generation, factual interpretability.



## Nadav Borenstein

**PhD Student**  
Nadav's research interests include improving the trustworthiness and usefulness of deep models in the NLP domain.



## Sarah Masud

**Postdoc**  
Sarah broadly works in the area of computational social systems with a focus on news narrative and hate speech modelling. Her PhD at IIIT-Delhi was supported by fellowships from Google and PMRF.



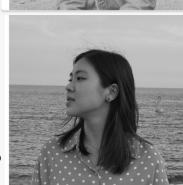
## Arnav Arora

**PhD Student**  
Arnav's research interests include equitable ML, mitigating online harms, and the intersection of NLP and Computational Social Science.



## Sara Vera Marjanovic

**PhD Student**  
Sara's research interests include explainable IR and NLP models, identifying biases in large text datasets, as well as working with social media data. She is a member of the DIKU ML section and IR group and co-advised by Isabelle.



## Haeun Yu

**PhD Student**  
Haeun's main research interest include enhancing explainability, fact-checking and transparency knowledge-enhanced LM.



## Jingyi Sun

**PhD Student**  
Jingyi Sun's research interests include explainability, fact-checking, and question answering.



## Siddhesh Pawar

**PhD Student**  
Siddhesh Pawar's research interests include multilingual models, fairness and accountability in NLP systems.



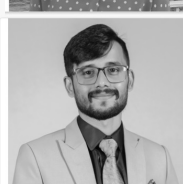
## Amalie Brogaard Pauli

**PhD Student**  
Amalie's research focuses on detecting persuasive and misleading text. She is a PhD student at Aarhus University and co-advised by Isabelle



## Sekh Mainul Islam

**PhD Student**  
Sekh's research interests include explainability in fact checking and improving robustness and trustworthiness in NLP models.



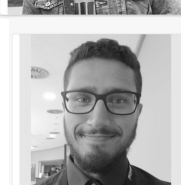
## Zain Muhammad Mujahid

**PhD Student**  
Zain's main research interests include disinformation detection, fact-checking, and factual text generation.



## Lucas Resck

**PhD Student**  
Lucas is an ELLIS PhD student at the University of Cambridge, supervised by Anna Corbhorn and co-supervised by Isabelle. His research interests include machine learning, NLP and explainability.



## Ahmad Dawar Hakimi

**PhD Student**  
Dawar is an ELLIS PhD student at LMU Munich, supervised by Hinrich Schütze and co-supervised by Isabelle. His research interests include mechanistic interpretability, summarisation and factuality of LLMs.



## Yijun Bian

**Postdoc**  
Yijun is a Marie-Curie postdoctoral fellow working on fair and interpretable ML.



## Jean Seo

**PhD Student**  
Jean's research interests include improving the safety of language models through explainability and evaluation.



+ You?  
We're hiring!

# References

**Isabelle Augenstein**, Timothy Baldwin, Meeyoung Cha, Tanmoy Chakraborty, Giovanni Luca Ciampaglia, David Corney, Renee DiResta, Emilio Ferrara, Scott Hale, Alon Halevy, Eduard Hovy, Heng Ji, Filippo Menczer, Ruben Miguez, Preslav Nakov, Dietram Scheufele, Shivam Sharma, Giovanni Zagni. [Factuality Challenges in the Era of Large Language Models](#). [Nature Machine Intelligence](#), August 2024.

Sara Vera Marjanović\*, Haeun Yu\*, Pepa Atanasova, Maria Maistro, Christina Lioma, **Isabelle Augenstein**. [DYNAMICQA: Tracing Internal Knowledge Conflicts in Language Models](#). In Findings of the 2024 Conference on Empirical Methods in Natural Language Processing ([EMNLP 2024](#)), November 2024.

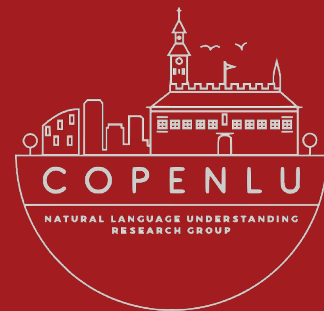
Lovisa Hagström, Sara Vera Marjanović, Haeun Yu, Arnav Arora, Christina Lioma, Maria Maistro, Pepa Atanasova, **Isabelle Augenstein**. [A Reality Check on Context Utilisation for Retrieval-Augmented Generation](#). In Proceedings of [ACL 2025](#), July 2025.

Sekh Mainul Islam, Pepa Atanasova, **Isabelle Augenstein**. [Multi-Step Knowledge Interaction Analysis via Rank-2 Subspace Disentanglement](#). CoRR, abs/2511.01706, November 2025.

Jingyi Sun\*, Pepa Atanasova\*, Sagnik Ray Choudhury, Sekh Mainul Islam, **Isabelle Augenstein**. [Evaluation Framework for Highlight Explanations of Context Utilisation in Language Models](#). Computational Linguistics, April 2026, to appear.

Lovisa Hagström\*, Youna Kim\*, Haeun Yu, Sang-goo Lee, Richard Johansson, Hyunsoo Cho, **Isabelle Augenstein**. [CUB: Benchmarking Context Utilisation Techniques for Language Models](#). In Proceedings of [ACL 2026](#), July 2026, to appear.

Thank you for  
your attention!  
Questions?



We're  
hiring!