



# Latent vs Explicit Knowledge Representation: How ChatGPT Answers Questions about Low-Frequency Entities

## A. Graciotti<sup>1</sup>, V. Presutti<sup>1</sup>, R. Tripodi<sup>2</sup>



This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement GA101004746. The communication reflects only the author's view and the Research Executive Agency is not responsible for any use that may be made of the information it contains.





<sup>1</sup> University of Bologna, 40126 Bologna, Italy <sup>2</sup> Ca' Foscari University of Venice, Sestiere Dorsoduro, 3246, 30123 Venezia VE

Problem Statement



[1] When Not to Trust Language Models: Investigating Effectiveness of Parametric and Non-Parametric Memories (Mallen et al., ACL 2023) [2] NLP Evaluation in trouble: On the Need to Measure LLM Data Contamination for each Benchmark (Sainz et al., EMNLP 2023)

Question Answering  $\rightarrow$  LLMs struggle when asked questions about less popular factual knowledge [1], as historical knowledge.

Benchmark crisis  $\rightarrow$  need for dynamic benchmarks due to LLMs data contamination concerns. [2]



Problem Statement



[1] When Not to Trust Language Models: Investigating Effectiveness of Parametric and Non-Parametric Memories (Mallen et al., ACL 2023) [2] NLP Evaluation in trouble: On the Need to Measure LLM Data Contamination for each Benchmark (Sainz et al., EMNLP 2023)

Question Answering  $\rightarrow$  LLMs struggle when asked questions about less popular factual knowledge [1], as historical knowledge.

Benchmark crisis  $\rightarrow$  need for dynamic benchmarks due to LLMs data contamination concerns. [2]



Contributions



### DynaKnowledge [1]: a new dynamic benchmark for free-form QA models

# 

### Text2AMR2FRED [2]: a Knowledge Extraction pipeline based on explicit knowledge that can be used for QA

[1] <u>https://github.com/polifonia-project/llms-vs-specialised-knowledge</u>

a systematic comparison between ChatGPT and the explicit knowledge model centred on entity popularity



Contributions



### **Text2AMR2FRED** [2]: a Knowledge Extraction pipeline based on explicit knowledge that can be used for QA

a systematic comparison between ChatGPT and the explicit knowledge model centred on entity popularity

[1] <u>https://github.com/polifonia-project/llms-vs-specialised-knowledge</u> [2] Text2AMR2FRED, a Tool for Transforming Text into RDF/OWL KnowledgeGraphsvia Abstract Meaning Representation (Gangemi et al., ISWC 2023)

### DynaKnowledge [1]: a new dynamic benchmark for free-form QA models



Contributions



a systematic **comparison** between ChatGPT and the explicit knowledge model centred on entity popularity

[1] <u>https://github.com/polifonia-project/llms-vs-specialised-knowledge</u> [2] Text2AMR2FRED, a Tool for Transforming Text into RDF/OWL KnowledgeGraphsvia Abstract Meaning Representation (Gangemi et al., ISWC 2023)

### DynaKnowledge [1]: a new dynamic benchmark for free-form QA models

**Text2AMR2FRED** [2]: a Knowledge Extraction pipeline based on explicit knowledge that can be used for QA



### **Objective** $\rightarrow$ comparing the performance of explicit and latent knowledge models on the QA task





Steps  $\rightarrow$  collecting a new benchmark for the task; selecting two representative models for latent and explicit knowledge representation



# Objective → comparing the performance of explicit and latent knowledge models on the QA task



Steps → collecting a new benchmark for the task; selecting two representative models for latent and explicit knowledge representation





DynaKnowledge: a new Dynamic Benchmark for QA

### **Composition** $\rightarrow$ 82 question, answer, provenance samples centred on historical characters' biographies (50-50 gender\* balance)

	How long was Teresina Bram- billa's career as a musician?	Teresa "Teresina" Brambilla (15 April 1845 – 1 Jul 1921) was an Italian soprano who sang in the majo opera houses of Europe in a career spanning 2 vears.
	Which female colleague did Giuseppina Ronzi de Begnis ar- gue with during the rehearsals of Maria Stuarda?	Ronzi was also known for her capricious attitude and for having confrontations and arguments wit female colleagues, including the famous altercation with Anna Del Sere during the rehearsals of Mari Stuarda.
	How old was Wolfgang Amadeus Mozart when he started to com- pose?	Already competent on keyboard and violin, he con posed from the age of five and performed befor European royalty.

\*In our study, we restrict to binary gender categories, which, although not reflecting real-world diversity, let us move the first steps towards the definition of our method.

DynaKnowledge: a new Dynamic Benchmark for QA

**Composition**  $\rightarrow$  82 question, answer, provenance samples centred on historical characters' biographies (50-50 gender\* balance)

ID	Gender	Analyst's Question	Analyst's Answer	Sentence from Wikipedia containing the answer
18	F	How long was Teresina Bram- billa's career as a musician?	25 years	Teresa "Teresina" Brambilla (15 April 1845 – 1 July 1921) was an Italian soprano who sang in the major opera houses of Europe in a career spanning 25 years.
41	F	Which female colleague did Giuseppina Ronzi de Begnis ar- gue with during the rehearsals of Maria Stuarda?	Anna Del Sere	Ronzi was also known for her capricious attitudes and for having confrontations and arguments with female colleagues, including the famous altercation with Anna Del Sere during the rehearsals of Maria Stuarda.
5	Μ	How old was Wolfgang Amadeus Mozart when he started to com- pose?	5 years old	Already competent on keyboard and violin, he com- posed from the age of five and performed before European royalty.

\*In our study, we restrict to binary gender categories, which, although not reflecting real-world diversity, let us move the first steps towards the definition of our method.





DynaKnowledge: a new Dynamic Benchmark for QA

**Composition**  $\rightarrow$  82 question, answer, provenance samples centred on historical characters' biographies (50-50 gender\* balance)

ID	Gender	Analyst's Question	Analyst's Answer	Sentence from Wikipedia containing the answer
18	F	How long was Teresina Bram- billa's career as a musician?	25 years	Teresa "Teresina" Brambilla (15 April 1845 – 1 Ju 1921) was an Italian soprano who sang in the majo opera houses of Europe in a career spanning 2 years.
41	F	Which female colleague did Giuseppina Ronzi de Begnis ar- gue with during the rehearsals of Maria Stuarda?	Anna Del Sere	Ronzi was also known for her capricious attitude and for having confrontations and arguments wit female colleagues, including the famous altercatio with Anna Del Sere during the rehearsals of Mari Stuarda.
5	Μ	How old was Wolfgang Amadeus Mozart when he started to com- pose?	5 years old	Already competent on keyboard and violin, he con posed from the age of five and performed befor European royalty.

\*In our study, we restrict to binary gender categories, which, although not reflecting real-world diversity, let us move the first steps towards the definition of our method.



DynaKnowledge: a new Dynamic Benchmark for QA

**Source**  $\rightarrow$  a corpus of music historical periodicals (1823 to 1900) [1]

for Rubini, was well sustained by a young singer of the name of Gennero: the Brambilla was also much applauded in this opera.

Recit. ed Aria, "Di tanti palpiti," Madame Ronzi De Begnis, (Tancredi) - Re Finale to the 1st Act of Don Giovanni, Mrs. Salmon,

Yesterday the Requiem of Mozart was performed as a mass for the deceased t, in the church of the Augustine

[1] The Periodicals module of the Polifonia Textual Corpus (https://github.com/polifonia-project/Polifonia-Corpus)



#### **Teresina Brambilla**



Portrait of Brambilla by Ernesto Fontana [Wikidata]



Portrait, c. 1781



Chalon

Models

Latent Knowledge Representation: ChatGPT (GPT3.5) [1]



[1] GPT-3.5 (Default) model was accessed via its free online version (<u>https://chat.openai.com/</u>) from August 4th to August 22nd, 2023

### Explicit Knowledge Representation: Text2AMR2FRED





Models

Latent Knowledge Representation: ChatGPT (GPT3.5)



[2] Text2AMR2FRED was accessed via its WebApp (<u>https://arco.istc.cnr.it/txt-amr-fred/</u>) [3] Text2AMR parsing is enabled by SPRING (Bevilacqua et al., AAAI Technical Track on Speech and Natural Language Processing I 2021)

### Explicit Knowledge Representation: Text2AMR2FRED [2] [3]





Models

Latent Knowledge Representation: ChatGPT (GPT3.5)



[2] Text2AMR2FRED was accessed via its WebApp (<u>https://arco.istc.cnr.it/txt-amr-fred/</u>) [3] Text2AMR parsing is enabled by SPRING (Bevilacqua et al., AAAI Technical Track on Speech and Natural Language Processing I 2021)

### Explicit Knowledge Representation: Text2AMR2FRED [2] [3]





Models

Teresa "Teresina" Brambilla (15 April 1845 – 1 July 1921) was an Italian soprano who sang in the major opera houses of Europe in a career spanning 25 years.





.

ChatGPT Answer correctness

**PROMPT:** How long was Teresina Brambilla career as a musician? Give a straight answer



### **ChatGPT** (GPT3.5)





ChatGPT (3.5) ANSWER: As of my last knowledge update in September 2021, there is no widely known or documented musician by the name of Teresina Brambilla in the





ChatGPT Answer correctness

**PROMPT:** How long was Teresina Brambilla career as a musician? Give a straight answer



ID	Gender	Analyst's Question	Analyst's
18	F	How long was Teresina Bram- billa's career as a musician?	25 years

### **ChatGPT** (GPT3.5)





ChatGPT (3.5) ANSWER: As of my last knowledge update in September 2021, there ANALYST'S is no widely known or documented musician ASSESSMENT: by the name of Teresina Brambilla in the Answer not given

> Sentence from Wikipedia containing the answer Answer Teresa "Teresina" Brambilla (15 April 1845 – 1 July 1921) was an Italian soprano who sang in the major opera houses of Europe in a career spanning 25 years.



ChatGPT Answer correctness

**PROMPT:** Which female colleague did Giuseppina Ronzi de Begnis argue with during the rehearsals of Maria Stuarda? Give a straight answer.

> ChatGPT (3.5) ANSWER: Giuseppina Ronzi de Begnis argued with Maria Malibran during the rehearsals of Maria Stuarda.

	Which female colleague did Giuseppina Ronzi de Begnis ar- gue with during the rehearsals of Maria Stuarda?	Ronzi was also known for her capricious attitudes and for having confrontations and arguments with female colleagues, including the famous altercation with Anna Del Sere during the rehearsals of Maria Stuarda.

## ChatGPT (GPT3.5)









ChatGPT Answer correctness

**PROMPT:** Which female colleague did Giuseppina Ronzi de Begnis argue with during the rehearsals of Maria Stuarda? Give a straight answer.

> ChatGPT (3.5) ANSWER: Giuseppina Ronzi de Begnis argued with Maria Malibran during the rehearsals of Maria Stuarda.

ID	Gender	Analyst's Question	Analyst's Answer
41	F	Which female colleague did Giuseppina Ronzi de Begnis ar- gue with during the rehearsals of Maria Stuarda?	Anna Del Sere

### ChatGPT (GPT3.5)







Sentence from Wikipedia containing the answer

Ronzi was also known for her capricious attitudes and for having confrontations and arguments with female colleagues, including the famous altercation with Anna Del Sere during the rehearsals of Maria Stuarda.



ChatGPT Answer correctness

**PROMPT:** How old was Wolfgang Amadeus Mozart when he started to compose?



		Already competent on keyboard and violin, he com- posed from the age of five and performed before European royalty.

### **ChatGPT** (GPT3.5)





ChatGPT (3.5) ANSWER: Wolfgang Amadeus Mozart started composing at the





ChatGPT Answer correctness

**PROMPT:** How old was Wolfgang Amadeus Mozart when he started to compose?



ID	Gender	Analyst's Question	Analyst's An
5	Μ	How old was Wolfgang Amadeus Mozart when he started to com- pose?	5 years old

Already competent on keyboard and violin, he composed from the age of five and performed before European royalty.



AMR graph answerability





AMR graph answerability





AMR graph answerability



AMR graph answerability



Polifonia | 2020

AMR graph answerability

Begnis argue with during the rehearsals of Maria Stuarda?





**Results and Analysis** 

NE's	C	<sup>chat</sup> GP
gender	P	R
F	0, 22	0, 31
Μ	0, 54	0, 58
Total	0, 38	0, 46





**Results and Analysis** 

NE's	C	<sup>chat</sup> GP
gender	P	R
F	0, 22	0, 31
Μ	0, 54	0, 58
Total	0, 38	0, 46





**Results and Analysis** 

NE's	C	hatGP
gender	P	R
F	0, 22	0, 31
Μ	0, 54	0,58
Total	$\left  0, 38 \right $	0, 46





Popularity and gender effect

Popularity → each named entity's Wikidata identifier (QID) frequency of occurrence as an internal link in Wikipedia [1].

	ChatGPT answer						AMR graph answerability						
					ot given						ot given		
76	126	37	93		3	42	96	30	83, 5			36	
869	1579	201	600	6	3	776	1627	338.5	739	61		545	12
639	1372				3		1183		542	61		291	9

[1] We used the enwiki-20220120 dump.



Polifonia | 2020

Popularity and gender effect

Popularity → each named entity's Wikidata identifier (QID) frequency of occurrence as an internal link in Wikipedia [1].

			ChatG	PT answei	ſ			AM	R graph	answerat	oility			
NE'S gender	#	#yes		#no	#nc	ot given	Ŧ	#yes	#	ŧno	#nc	ot given	#	Total
gender	Avg	StdDev	Avg	StdDev	Avg	StdDev	Avg	StdDev	Avg	StdDev	Avg	StdDev	Avg	StdD
F	76	126	37	93	4	3	42	96	30	83, 5	N/A	N/A	36	
Μ	869	1579	201	600	6	3	776	1627	338.5	739	61	N/A	545	12
All	639	1372	110	407	4	3	400	1183	184	542	61	N/A	291	9

[1] We used the enwiki-20220120 dump.



Popularity and gender effect

Popularity  $\rightarrow$  each named entity's Wikidata identifier (QID) frequency of occurrence as an internal link in Wikipedia [1].

			ChatG	PT answei	ſ			AM	R graph	answerat	oility			
NE'S gender	7	#yes		#no	#nc	ot given	;	#yes	#	ŧno	#nc	ot given	#	Total
gender	Avg	StdDev	Avg	StdDev	Avg	StdDev	Avg	StdDev	Avg	StdDev	Avg	StdDev	Avg	StdDe
F	76	126	37	93	4	3	42	96	30	83, 5	N/A	N/A	36	8
Μ	869	1579	201	600	6	3	776	1627	338.5	739	61	N/A	545	12.
All	639	1372	110	407	4	3	400	1183	184	542	61	N/A	291	9

[1] We used the enwiki-20220120 dump.



Popularity and gender effect

Popularity  $\rightarrow$  each named entity's Wikidata identifier (QID) frequency of occurrence as an internal link in Wikipedia [1].

			ChatG	PT answei				AM	R graph	answerat	oility			
NE S gender	#	yes	5	#no	#nc	ot given	;	#yes	#	ŧno	#nc	ot given	#	Total
gender	Avg	StdDev	Avg	StdDev	Avg	StdDev	Avg	StdDev	Avg	StdDev	Avg	StdDev	Avg	StdD
F	76	126	37	93	4	3	42	96	30	83, 5	N/A	N/A	36	
Μ	869	1579	201	600	6	3	776	1627	338.5	739	61	N/A	545	12
All	639	1372	110	407	4	3	400	1183	184	542	61	N/A	291	9

[1] We used the enwiki-20220120 dump.



	Spearman cor	relation
feature	ChatGPT's answer assessment	AMR graph's answerability
Popularity	0.48	-0.07
Gender	0.33	-0.02



	Spearman co	rrelation
feature	ChatGPT's answer assessment	AMR graph's answerability
Popularity Gender	$\begin{array}{c} 0.48\\ 0.33 \end{array}$	$-0.07 \\ -0.02$



	Spearman cor	relation		
feature	ChatGPT's answer assessment	AMR graph's answerability		
Popularity Gender	$\begin{array}{c} 0.48\\ 0.33\end{array}$	-0.07 -0.02		



Popularity and gender effect

	Spearman cor	relation
feature	ChatGPT's answer assessment	AMR graph's answerability
Popularity Gender	0.48 0.33	$-0.07 \\ -0.02$

\_



	Spearman cor	relation
feature	ChatGPT's answer assessment	AMR graph's answerability
Popularity Gender	$\begin{array}{c} 0.48\\ 0.33\end{array}$	-0.07 -0.02



#### Conclusion

A simpler knowledge extraction pipeline demonstrated to be more robust on the variation of named entities' features such as gender and popularity.

**Popularity matters**  $\rightarrow$  ChatGPT struggles to answer questions related to less popular entities, which in our benchmark are female.



#### Conclusion

→ A simpler knowledge extraction pipeline demonstrated to be more robust on the variation of named entities' features such as gender and popularity.

**Popularity matters**  $\rightarrow$  ChatGPT struggles to answer questions related to less popular entities, which in our benchmark are female.



# Expand DynaKnowledge dynamically, following current trends in benchmarking [1] NLP models.

Transforming the AMR graphs into OWL-compliant RDF KGs to automate answerability assessment through structured interrogations.

Exploiting the KGs to augment LLMs prompts and improve QA performance in a retrieval-augmented generation [2] framework.

[1] <u>Dynabench: Rethinking Benchmarking in NLP</u> (Kiela et al., NAACL 2021)
[2] <u>Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks</u> (Lewis et al., EMNLP 2023)



Exploiting the KGs to augment LLMs prompts and improve QA performance in a retrieval-augmented generation [2] framework.

[1] <u>Dynabench: Rethinking Benchmarking in NLP</u> (Kiela et al., NAACL 2021) [2] <u>Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks</u> (Lewis et al., EMNLP 2023)

Expand DynaKnowledge dynamically, following current trends in benchmarking [1] NLP models.

Transforming the AMR graphs into OWL-compliant RDF KGs to automate answerability assessment through structured interrogations.



[1] <u>Dynabench: Rethinking Benchmarking in NLP</u> (Kiela et al., NAACL 2021) [2] <u>Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks</u> (Lewis et al., EMNLP 2023)

Expand DynaKnowledge dynamically, following current trends in benchmarking [1] NLP models.

Transforming the AMR graphs into OWL-compliant RDF KGs to automate answerability assessment through structured interrogations.

Exploiting the KGs to augment LLMs prompts and improve QA performance in a retrieval-augmented generation [2] framework.







ALMA MATER STUDIORUM Università di Bologna

# Thank you for your attention!

### Speaker: Arianna Graciotti

### arianna.graciotti@unibo.it

https://www.unibo.it/sitoweb/arianna.graciotti https://www.linkedin.com/in/arianna-graciotti/







https://github.com/polifoniaproject/llms-vs-specialisedknowledge