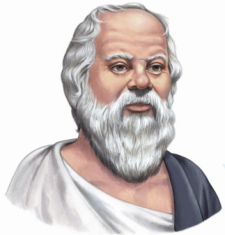


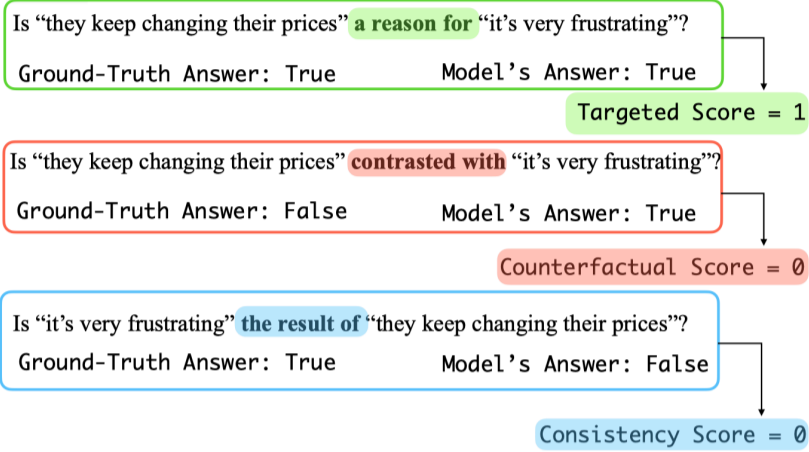
Discursive Socratic Questioning: Evaluating the Faithfulness of Language Models' Understanding of Discourse Relations

Yisong Miao, Hongfu Liu, Wenqiang Lei, Nancy F. Chen, Min-Yen Kan



Socrates
470 – 399 BC

Want to know how well LLMs understand discourse relations? We propose an end-to-end automatic scoring framework for discourse relations, leveraging LLMs in a Socratic style.

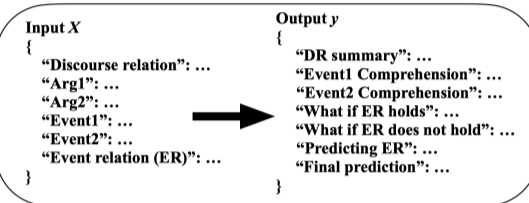


Discourse relation: Contingency. Cause.Result

Arg1: When I want to buy, they run from you - they keep changing their prices.

Arg2: It's very frustrating.

Salient signal.



	A1&A2	A1&IC	A2&ICL
Agreements	85.2%	85.2%	83.7%
Cohen's κ	38.5%	48.8%	44.9%
Success Rate	/	95.8%	93.8%

What to ask: In-context learning to annotate salient signals.

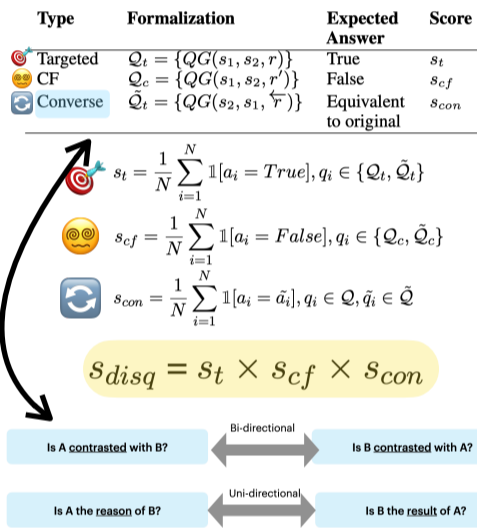
Algorithm 1 DiSQ interrogates a language model.

- 1: **Input:** Discourse d and its corresponding questions \mathcal{Q} .
- 2: $\mathcal{H} = \{\emptyset\}$ \triangleright The history is initialized.
- 3: **Stage 1: Targeted and Counterfactual QA**
- 4: **for** q_i in \mathcal{Q}_t and \mathcal{Q}_c **do**
- 5: $a_i = \text{LM}(q = q_i, c = d)$ \triangleright The model performs QA. The context c is the discourse d .
- 6: $\mathcal{H} \leftarrow (q_i, a_i)$ \triangleright The history is updated.
- 7: **end for**
- 8: **Stage 2: Converse QA**
- 9: **for** (q_i, a_i) in \mathcal{H} **do**
- 10: $\tilde{q} = \text{Lookup}(q, \{\tilde{\mathcal{Q}}_c, \tilde{\mathcal{Q}}_t\})$ \triangleright Look up the converse question in converse question sets.
- 11: $\tilde{a}_i = \text{LM}(q = \tilde{q}_i, c = d, (q_i, a_i) \in \mathcal{H})$ \triangleright The model executes QA on the converse question, \tilde{q}_i , optionally utilizing the previous response (q_i, a_i) as supplemental context.
- 12: $\mathcal{H} \leftarrow (\tilde{q}_i, \tilde{a}_i)$ \triangleright The history is updated.
- 13: **end for**
- 14: **Output:** \mathcal{H}

Targeted question:
Is A the result of B?

Counterfactual question:
Is A contrasted with B?
Is A the example with B?
Is A an alternative of B?
...

Converse question:
(Given you answered A is the result of B.) Is B the reason of A?



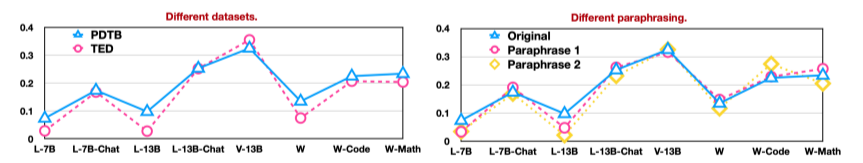
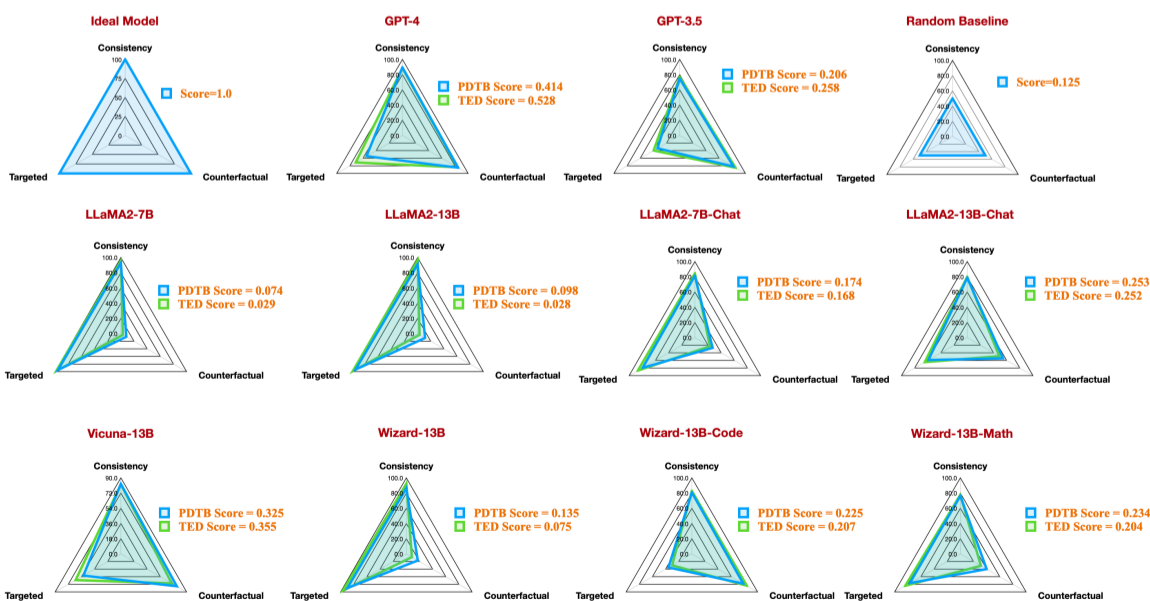
Discourse relation (R)	Event relation (r)	Q Type	# of Q
Comparison.Cession	deny or contradict with	Bi-	1,764
Comparison.Contrast	contrast with	Bi-	876
Contingency.Reason	reason of	Uni-	3,264
Contingency.Result	result of	Uni-	2,796
Expansion.Conjunction	contribute to the same situation	Bi-	4,596
Expansion.Equivalence	equivalent to	Bi-	420
Expansion.Instantiation	example of	Uni-	2,352
Expansion.Level-of-detail	provide more detail about	Uni-	3,888
Expansion.Substitution	alternative to	Uni-	216
Temporal.Asynchronous	happen before/after	Uni-	1,368
Temporal.Synchronous	happen at the same time as	Bi-	840
Total			22,380

How to ask: DiSQ is composed of (1) targeted questions, (2) counterfactual questions and (3) converse questions. DiSQ Scores is a multiplication of the three aspects.

Question statistics for PDTB dataset.

Evaluation setup: PDTB and TED-MDB dataset. Zero-shot QA.

Evaluated models: (1) Closed-source models: GPT-4 / GPT-3.5. (2) Open-source models: LLaMA family; Vicuna model; Wizard family.



RQ2: DiSQ Scores are consistent.

Models	Overall	Comp. Cause	Comp. Contrast	Cont. Reason	Cont. Result	Exp. Conj.	Exp. Equiv.	Exp. Inst.	Exp. Level	Exp. Subst.	Temp. Async	Temp. Sync
1. Random Baseline	0.125	0.125	0.125	0.125	0.125	0.125	0.125	0.125	0.125	0.125	0.125	0.125
2A. LLaMA2-7B	0.074	0.029	0.083	0.094	0.095	0.076	0.056	0.087	0.067	0.156	0.035	0.048
3A. LLaMA2-7B-Chat	0.174	0.231	0.431	0.131	0.174	0.213	0.104	0.120	0.150	0.199	0.108	0.040
4A. LLaMA2-13B	0.098	0.037	0.100	0.082	0.097	0.127	0.101	0.113	0.107	0.086	0.084	0.092
5A. LLaMA2-13B-Chat	0.253	0.193	0.477	0.129	0.172	0.288	0.157	0.326	0.373	0.291	0.195	0.028
6A. Vicuna-13B	0.325	0.087	0.513	0.200	0.353	0.369	0.000	0.334	0.462	0.195	0.511	0.069
7A. Wizard	0.135	0.221	0.256	0.067	0.107	0.170	0.072	0.167	0.128	0.108	0.097	0.082
8A. Wizard-Code	0.225	0.032	0.268	0.175	0.287	0.121	0.008	0.283	0.329	0.174	0.545	0.109
9A. Wizard-Math	0.234	0.132	0.264	0.241	0.286	0.192	0.046	0.240	0.323	0.201	0.240	0.135
10A. GPT-3.5	0.206	0.151	0.278	0.082	0.161	0.246	0.067	0.257	0.262	0.232	0.388	0.000
11A. GPT-4	0.414	0.053	0.567	0.119	0.351	0.610	0.192	0.659	0.481	0.422	0.692	0.000

RQ3: Minority classes are still challenging.



Overall performance (RQ1): (1) Gap between Close- and Open source models; (2) Benefits from further fine-tuning; (3) Consistency between the two datasets.

RQ4: Linguistic Features: Benefits from discourse connectives, discourse context, and historical QAs.