

Backtracing: Retrieving the Cause of the Query

Rose E. Wang* Pawan Wirawarn Omar Khattab Noah Goodman

Dorottya Demszky

*rewang@cs.stanford.edu, Stanford University

Motivation

Backtracing Benchmark



What did I say that triggered this student's question?



[...] The projection is the same point. So that means that if I project twice, I get the same answer as I did in the first project. So those are the two properties that tell me I'm looking at a projection matrix. [...]

- While information retrieval (IR) systems may provide answers for user queries, they do not directly assist content creators (e.g., teachers) identify segments that caused a user to ask those questions.
- Identifying the cause of a query is challenging because of 1. lack of explicit labeling, 2. large size of corpus, and 3. required domain expertise to understand both the query and corpus.
- We introduce the task of *backtracing*, in which systems retrieve the text segment that most likely caused a user query.

Contributions

- Task. We formalize *backtracing*: Retrieve the text segment that most likely caused the user query.
- **Benchmark**. We develop a heterogeneous benchmark for backtracing: retrieving the cause of student confusion in the LECTURE setting, reader curiosity in the NEWS ARTICLE setting, and user emotion in the CONVERSATION setting.

LECTURE. Retrieve the cause of student confusion [3]. **NEWS ARTICLE**. Retrieve the cause of reader curiosity [1]. CONVERSATION. Retrieve the cause of user emotion (e.g., anger) [2].

systems retrieve incorrect contexts shown in red.

	# sentences	Lecture	News Article	Conversation
Corpus X	Total	11042	2125	8263
	Average	525.8	19.0	12.3
Query q	Total	210	1382	671
	Average	30.9	7.1	11.6

Results

- The best-performing models achieve modest accuracies. Measuring causal relevance is challenging and markedly different from existing retrieval tasks.
- The methods do not generalize across domains. For instance, while a cross-encoder method performs well on the NEWSARTICLE domain with top-3 85% accuracy, it only manages top-3 15% accuracy on the CONVERSATION domain.

		Lec @1	cture @3	News ©1	Article ^{@3}	Conv @1	versation ^{@3}
	Random Edit	04	2 8	6 7	21 18	$\begin{vmatrix} 11 \\ 1 \end{vmatrix}$	31 16
	Bi-Encoder (Q&A) Bi-Encoder (211-MiniIM)	23	37 40	48 40	71 75	1	32 37
	Cross-Encoder Re-ranker gpt-3.5-turbo-16k	20 22 30 15	39 44 N/A	66 66 67	85 85 N/A	1 1 47	15 21 N/A
Single-sentence $p(q x_t)$	GPT2 GPTJ OPT 6B	21 23 30	34 42 43	43 67 66	64 85 82	3 5 2	46 65 56
Autoregressive $p(q x_{\leq t})$	GPT2 GPTJ OPT 6B	11 14 16	16 24 26	9 55 52	18 76 73	5 8 18	54 60 65
$\begin{array}{c} \textbf{ATE} \\ p(q X) - p(q X / \{x_t\} \) \end{array}$	GPT2 GPTJ OPT 6B	13 8 2	21 18 6	51 67 64	68 79 76	2 3 3	24 18 22

Evaluations. We evaluate a suite of popular retrieval systems and show that there is room for improvement in current retrieval methods. This suggests that backtracing is not only challenging but also requires new retrieval approaches.

Backtracing Task

Given corpus of N sentences $X = \{x_1, \ldots, x_N\}$ and query q, backtracing selects

$$\hat{t} = \arg \max_{t \in 1...N} p(t|x_1, \ldots, x_N, q)$$
(1)

where x_t is the t^{th} sentence in corpus X and p is a probability distribution over the corpus indices, given the corpus and the query.

This task intuitively translates to: Given a lecture transcript and student question, retrieve the lecture sentence(s) that most likely caused the student to ask that question.

Table 1. Accuracy in percentage (%). The best models in each column are bolded. For each dataset, we report the top-1 and 3 accuracies.

More results and analysis in the paper!



[1] Wei-Jen Ko, Te-Yuan Chen, Yiyan Huang, Greg Durrett, and Junyi Jessy Li. Inquisitive question generation for high level text comprehension. arXiv preprint arXiv:2010.01657, 2020.

- [2] Soujanya Poria, Navonil Majumder, Devamanyu Hazarika, Deepanway Ghosal, Rishabh Bhardwaj, Samson Yu Bai Jian, Pengfei Hong, Romila Ghosh, Abhinaba Roy, Niyati Chhaya, et al. Recognizing emotion cause in conversations. Cognitive Computation, 13:1317–1332, 2021.
- [3] Rose Wang, Pawan Wirawarn, Noah Goodman, and Dorottya Demszky. Sight: A large annotated dataset on student insights gathered from higher education transcripts. In Proceedings of Innovative Use of NLP for Building Educational Applications, 2023.