



VIRGINIA TECH

A Multi-Level Benchmark for Causal Language Understanding in Social Media Discourse

EMNLP 2025
Suzhou, China | 中国苏州

Xiaohan Ding, Kaike Ping, Buse Carik, Eugenia Rho
 {xiaohan, pkk, buse, eugenia} @ vt.edu

Background & Motivation - Why Causal Language Understanding?

- Causal reasoning is central to NLP but underexplored in informal social media text.
- Existing datasets focus on explicit causality in structured domains (e.g., news, biomedical).
- Social media: messy, implicit, gist-driven → critical for public health, misinformation, decision-making.
- Gap:** lack of datasets that connect detection, classification, extraction, and gist generation

Dataset - CausalTalk

- Source:** Reddit posts (2020–2024), 43 public health subreddits (COVID-19, vaccines, lockdowns).
- Size:** 239k submissions + 19M comments; 10,120 annotated posts.
- Annotation Tasks:**
 - Binary causal classification
 - Explicit vs. Implicit causality
 - Cause–effect span extraction
 - Causal gist generation
- Annotations:**
 - Gold: 1,320 posts (expert consensus)
 - Silver: 8,800 posts (GPT-4o + human verification)

Methodology - Annotation & Benchmark Setup

Annotation Pipeline

- Gold Annotations (n = 1,320 posts):**
 - Five annotators trained in causal linguistics + public health.
 - Independent labeling for 4 tasks → disagreements resolved via sixth adjudicator.
 - Outputs: causal classification, explicit/implicit type, span extraction, causal gists.
- Silver Annotations (n = 8,800 posts):**
 - Generated by GPT-4o with RBIC.
 - Human annotators verified & refined every output; and Dimensions assessed:
 - Causality Accuracy** (binary correctness)
 - Type Accuracy** (explicit vs. implicit)
 - Span Relevance** (5-point Likert)
 - Gist Conciseness & Coherence** (5-point Likert)
 - High inter-annotator agreement** ($\kappa \approx 0.78\text{--}0.89$ across tasks).

Evaluation Criterion	Score	Fleiss' κ
Causality Accuracy	$ACC_{avg} = 0.902$	0.892
Causality Type Accuracy	$ACC_{avg} = 0.702$	0.780
Relevance (Span Extraction)	Mean = 4.30	0.839
Conciseness (Gist Generation)	Mean = 4.50	0.864

Benchmark Setup

- Task Coverage:** (1) Binary causal classification; (2) Explicit vs. implicit detection; (3) Cause–effect span extraction; (4) Causal gist generation
- Models (Tasks 1–3):**
 - BERT-base, RoBERTa-base, XLNet-base, DeBERTa-v3, SpanBERT
- Models (Task 4):**
 - Fine-tuned: T5, FLAN-T5, GPT-2, BART
 - Instruction-tuned LLMs: LLaMA-3.2, Gemini 2.0 Flash, DeepSeek-V3, Claude 3.5 Haiku

Main Results & Insights - Benchmark Findings

Dataset	Model	Precision	Recall	F1 Score
Gold	BERT-base	0.70 _{0.023}	0.74 _{0.023}	0.75 _{0.024}
	RoBERTa-base	0.81 _{0.021}	0.80 _{0.020}	0.80 _{0.021}
	XLNet-base	0.80 _{0.021}	0.78 _{0.020}	0.80 _{0.021}
	DeBERTa-v3*	0.82_{0.021}	0.80_{0.021}	0.83_{0.022}
Silver	BERT-base	0.81 _{0.025}	0.79 _{0.024}	0.80 _{0.027}
	RoBERTa-base	0.85 _{0.020}	0.83 _{0.020}	0.84 _{0.020}
	XLNet-base	0.84 _{0.024}	0.82 _{0.023}	0.83 _{0.024}
	DeBERTa-v3†	0.87_{0.025}	0.86_{0.024}	0.87_{0.027}
	$\Delta_{model^{\dagger}-model^*}$	$\uparrow 0.05$	$\uparrow 0.06$	$\uparrow 0.04$

Table 1: Performance on Task 1 (Causal Classification) across gold and silver datasets. Results are reported on the respective held-out test sets (20% of each dataset), with mean \pm standard deviation over five random seeds.

Dataset	Model	Precision	Recall	F1 Score
Gold	BERT-base	0.61 _{0.021}	0.59 _{0.024}	0.58 _{0.027}
	RoBERTa-base	0.61 _{0.019}	0.60 _{0.020}	0.60 _{0.021}
	XLNet-base	0.63 _{0.022}	0.62 _{0.018}	0.63 _{0.020}
	DeBERTa-v3*	0.68_{0.017}	0.68_{0.015}	0.69_{0.016}
Silver	BERT-base	0.66 _{0.026}	0.65 _{0.025}	0.65 _{0.027}
	RoBERTa-base	0.68 _{0.022}	0.66 _{0.023}	0.67 _{0.019}
	XLNet-base	0.70 _{0.021}	0.69 _{0.019}	0.69 _{0.020}
	DeBERTa-v3†	0.75_{0.016}	0.74_{0.015}	0.74_{0.014}
	DeBERTa-v3‡	0.69 _{0.018}	0.70 _{0.017}	0.70 _{0.018}

Table 2: Performance on Task 2 (Explicit vs. Implicit Causality Classification). Results are mean \pm standard deviation over five random seeds on the respective heldout test sets (20% of each dataset).

Model	Type	Causal Gist Generation			
		ROUGE-1	ROUGE-2	ROUGE-L	BERTScore
T5-base	SFT	0.429 _{0.012}	0.334 _{0.009}	0.512 _{0.013}	0.670 _{0.016}
FLAN-T5-base*	SFT	0.559_{0.007}	0.354 _{0.011}	0.521 _{0.008}	0.704 _{0.012}
GPT-2	SFT	0.281 _{0.014}	0.089 _{0.019}	0.235 _{0.017}	0.305 _{0.016}
BART-base	SFT	0.442 _{0.015}	0.261 _{0.010}	0.400 _{0.009}	0.528 _{0.013}
LLaMA-3.2-3B	zero-shot	0.432 _{0.013}	0.243 _{0.017}	0.400 _{0.014}	0.532 _{0.011}
	few-shot	0.448 _{0.012}	0.235 _{0.016}	0.417 _{0.013}	0.551 _{0.012}
Google Gemini†	zero-shot	0.557 _{0.010}	0.436_{0.015}	0.588_{0.012}	0.764_{0.009}
	few-shot	0.545 _{0.011}	0.427 _{0.014}	0.574 _{0.011}	0.745 _{0.010}
DeepSeek-V3	zero-shot	0.526 _{0.008}	0.411 _{0.016}	0.568 _{0.012}	0.731 _{0.014}
	few-shot	0.537 _{0.009}	0.422 _{0.015}	0.549 _{0.013}	0.715 _{0.013}
Claude 3.5 Haiku	zero-shot	0.436 _{0.019}	0.210 _{0.018}	0.356 _{0.011}	0.462 _{0.017}
	few-shot	0.423 _{0.018}	0.221 _{0.016}	0.366 _{0.015}	0.475 _{0.016}
$\Delta_{Gemini^{\dagger}-FLAN-T5^*}$			$\downarrow 0.002$	$\uparrow 0.082$	$\uparrow 0.067$
				$\uparrow 0.060$	

Table 3: Performance on Task 3 (Cause–Effect Span Extraction) between gold and silver standard datasets. Each model is evaluated using both token-level and spanlevel metrics.

Dataset	Model	Precision	Recall	F1
Gold	BERT-base	0.82 _{0.012}	0.83 _{0.014}	0.82 _{0.013}
	- Token	0.71 _{0.015}	0.69 _{0.015}	0.70 _{0.014}
SpanBERT	- Token	0.84 _{0.011}	0.85 _{0.013}	0.84 _{0.012}
	- Span	0.75 _{0.014}	0.73 _{0.015}	0.74 _{0.014}
RoBERTa-base	- Token	0.87 _{0.010}	0.87 _{0.011}	0.87 _{0.010}
	- Span	0.79 _{0.013}	0.77 _{0.014}	0.78 _{0.013}
DeBERTa-v3	- Token*	0.89_{0.010}	0.89_{0.010}	0.89_{0.010}
	- Span*	0.82 _{0.012}	0.80_{0.013}	0.81_{0.012}
BERT-base	- Token	0.89 _{0.011}	0.90 _{0.012}	0.88 _{0.011}
	- Span	0.78 _{0.014}	0.76 _{0.015}	0.77 _{0.014}
SpanBERT	- Token	0.91 _{0.010}	0.92 _{0.011}	0.91 _{0.010}
	- Span	0.82 _{0.013}	0.80 _{0.014}	0.81 _{0.013}
RoBERTa-base	- Token	0.94 _{0.010}	0.94 _{0.010}	0.94 _{0.010}
	- Span	0.86 _{0.012}	0.84 _{0.013}	0.85 _{0.012}
DeBERTa-v3	- Token†	0.95_{0.010}	0.95_{0.010}	0.95_{0.010}
	- Span‡	0.89 _{0.011}	0.87_{0.012}	0.88_{0.011}
$\Delta_{Token^{\dagger}-Token^*}$		$\uparrow 0.06$	$\uparrow 0.06$	$\uparrow 0.06$
$\Delta_{Span^{\dagger}-Span^*}$		$\uparrow 0.07$	$\uparrow 0.07$	$\uparrow 0.07$

Table 4: Performance of causal gist generation on the silver-standard dataset. The upper section includes supervised fine-tuned models (SFT), while the lower section shows zero-shot and few-shot prompting results from instruction-tuned LLMs.