

Shallowly accurate but deeply confused – how language models deal with antonyms

Research question

- ▶ Antonyms are words that are semantically opposite of each other. We focus here on **antonymic adjectives**.
- ▶ Statistical models of language have previously been argued to give a poor treatment of antonyms, because:
 - they are based on the Distributional Hypothesis [1];
 - antonyms appear in similar environments [2, 4].
- ▶ Recent Large Language Models (LLMs) yet perform well in tasks which most probably require them to encode adjective polarity... **so, can they draw fine-grained inferences about antonyms?**

Background: adjective polarity & negation

The Inference Towards the Antonym (ITA) [3, 7, 8, 10]
 $(not A) \implies A'$, where A' is the antonym of A

ITA Pragmatic Mitigation Condition [7]
 $(not A) \not\Rightarrow A'$, if $CPLX(not A) \gg CPLX(A')$

Negative Adjectives Complexity Hypothesis [6, 5]

$\forall A^-, A^- = NOT-A^+$, therefore:

- ◆ $CPLX(A^-) = CPLX(NOT-A^+) \sim CPLX(not A^+)$
- ★ $CPLX(not A^-) = CPLX(not NOT-A^+) \gg CPLX(-A^+)$

- ▶ Based on Eq. ◆ and ★, Krifka [7] concludes that **inferring a negative adjective A^- from $not A^+$ is easier than inferring A^+ from $not A^-$** .
- ▶ As argued by Ruytenbeek et al. [8], this can be assessed empirically by comparing the felicity of (1a) vs. (1b).
 (1) a. He is **not tall**_(A⁺). She too is **short**_(A⁻).
 b. # He is **not short**_(A⁻). She too is **tall**_(A⁺).
- ▶ They also claim that the contrast between (1a) and (1b) should be smaller for morphologically opaque antonyms (**O-antonyms**), as opposed to morphologically transparent ones (**T-antonyms**).

Goal

- ▶ Building on the studies conducted by [8], we test if 3 recent LLMs (GPT-2 [13], XLNET [14] and BERT [9]):
 1. are more likely to “draw” an ITA when the adjective under negation is positive rather than negative (**H1**);
 2. show a higher discrepancy in ITA strength for T-antonyms as opposed to O-antonyms (**H2**).

	Task 1		Task 2		H1	H2
	H1	H2	H1	H2		
< .001						
< .01						
< .05						
< .1						
n.s.						
	GPT-2	.14 (S) T+O	n.s.	.59 (L) T+O	.19 (S) T+O	.30 (S)
	XLNet	.03 (N) T+O	.37 (M)	.42 (M) T+O	.35 (M) T+O	.17 (S)
	BERT	.15 (S) T	.14 (S)	n.s.	.22 (S) T+O	.31 (S)

Table 1: Test results for Tasks 1 and 2.¹ **Table 2:** Test results for Task 3.

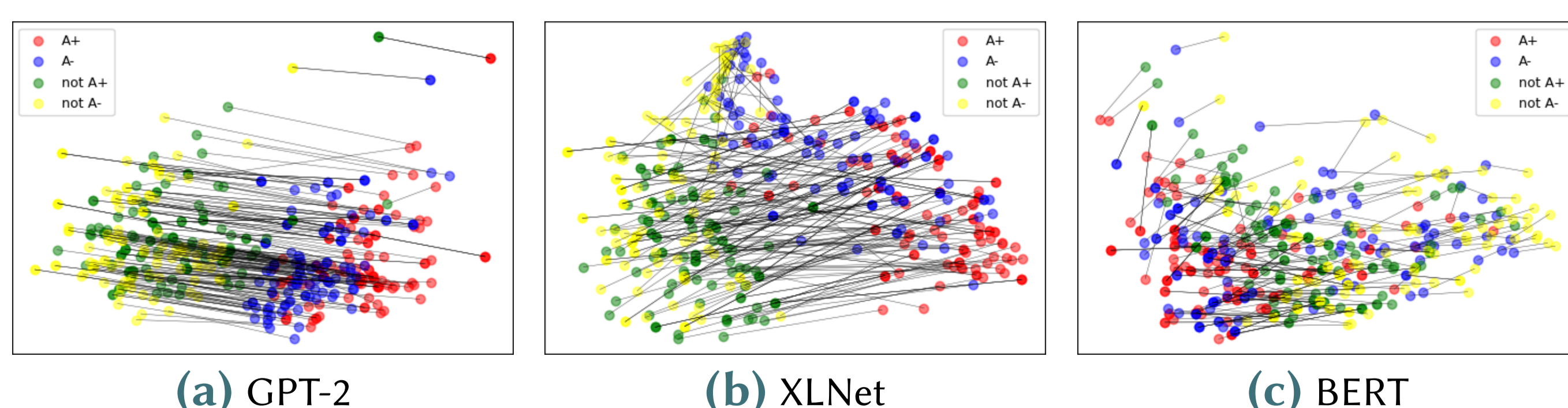


Figure 1: 2D-PCA-reduced embeddings of the 3 models. Lines represent the effect of negation for any given adjective.

Task 1: Sentence surprisal (S_{sent})

Hypothesis 1: $\Delta S_{sent} \triangleq S_{sent}(1b) - S_{sent}(1a) > 0$

- ▶ 111 pairs of antonyms (48 T, 63 O) tested using a one-sided paired-sample Wilcoxon test.
- ▶ **3/3 LLMs captured the contrast** but with small effect sizes.
- ▶ Group-by-group testing (T vs. O) reveals that only for GPT-2 and BERT, *both* groups verify H1 (after corrections).

Hypothesis 2: $(\Delta S_{sent})_{O-group} < (\Delta S_{sent})_{T-group}$

- ▶ We compared the T- and O-groups using a Mann Whitney U-Test on the surprisal contrasts between (1a) and (1b).
- ▶ **2/3 LLMs captured the effect of transparency.**

Task 2: Target-word surprisal (S_w)

H1: $\Delta S_A \triangleq (S_{A^+}(1b) - S_{A^+}(2b)) - (S_{A^-}(1a) - S_{A^-}(2a)) > 0$

- ▶ We compared the individual surprisals associated with the second adjective in (1a)/(1b).
- ▶ Because LLMs tend to assign morphologically complex words higher surprisals, we used (2a)/(2b) as baselines.
 (2) a. He is **not tall**_(A⁺). She is **short**_(A⁻).
 b. # He is **not short**_(A⁻). She is **tall**_(A⁺).
- ▶ **2/3 LLMs captured the contrast.** For both models, individual groups (T/O) were also linked to a significant effect after corrections.

H2: $(\Delta S_A)_{O-group} < (\Delta S_A)_{T-group}$

- ▶ **p-values for H2 in this task were not significant**, i.e., no effect of morphological transparency was detected...

Task 3: Comparison in word embeddings

- ▶ We focus on the representation that LLMs assign to A^+ , A^- , and their negations: \vec{A}^+ , \vec{A}^- , $\vec{not A^+}$, $\vec{not A^-}$.
- ▶ Semantic closeness between those vectors is measured by cosine similarity (=angle between 2 vectors).

H1: $\Delta Cos \triangleq Cos(\vec{A}^-, \vec{not A^+}) - Cos(\vec{A}^+, \vec{not A^-}) > 0$

- ▶ We test if $\vec{not A^+}$ is “closer” to \vec{A}^- than $\vec{not A^-}$ is to \vec{A}^+ .
- ▶ **3/3 embeddings captured the contrast**, suggesting that ITA strength translates into topological distance.

H2: $(\Delta Cos)_{O-group} < (\Delta Cos)_{T-group}$

- ▶ **3/3 embeddings captured the effect of transparency**, although the effect sizes were small...

Discussion & outlook

- ▶ LLMs seem to distinguish positive and negative adjectives w.r.t their semantic closeness to their antonym (H1) and somewhat differentiate between T- and O-antonyms (H2).
- ▶ The relative weakness of this last result may be due to LLMs’ tokenization strategy not aligning with human-like morphological segmentation.
- ▶ This should also be contrasted with two other results:
 - As already noted for older models [11], **LLMs represent $not A^+$ as closer to $not A^-$ than to A^- , and vice-versa** (cf. Fig. 1)!
 - A refinement of BERT (RoBERTa-MNLI [12]) fine-tuned for Natural Language Inference was **more likely to conclude ($He not A^- \implies He is A^+$) than ($He not A^+ \implies He is A^-$)**, *contra* H1(!)

¹p-values color-coded. Effect sizes are Cliff’s Δ . N, S, M, L resp. mean ‘negligible’, ‘small’, ‘medium’, ‘large’. Each cell also lists which subgroup(s) (T, O) drive(s) the effect in Task 1.

