Anmol Agarwal*, Shrey Gupta*, Vamshi Bonagiri, Manas Gaur, Joseph Reagle, and Ponnurangam Kumaraguru

IIIT Hyderabad, KAI2 Lab @ UMBC, Northeastern University

## Motivation and Problem Statement

Researchers dealing with public user-generated content often need to paraphrase content related to sensitive topics like health, violence, drug use, etc, before making it public.

- Existing AI-based automated word spinners (e.g., SpinRewriter, WordAI) are often ineffective as their paraphrased content is still locatable on search engines.
- **Introducing:** an *unsupervised black-box adversarial framework* to paraphrase content such that querying snippets of text from it on search engines does not lead back to the original content on the web.

Given a sentence 's,' we paraphrase the sentence with the aim:

- **Non-locatability**: the sentence's source is not in the top-K results when the sentence is queried on search engines
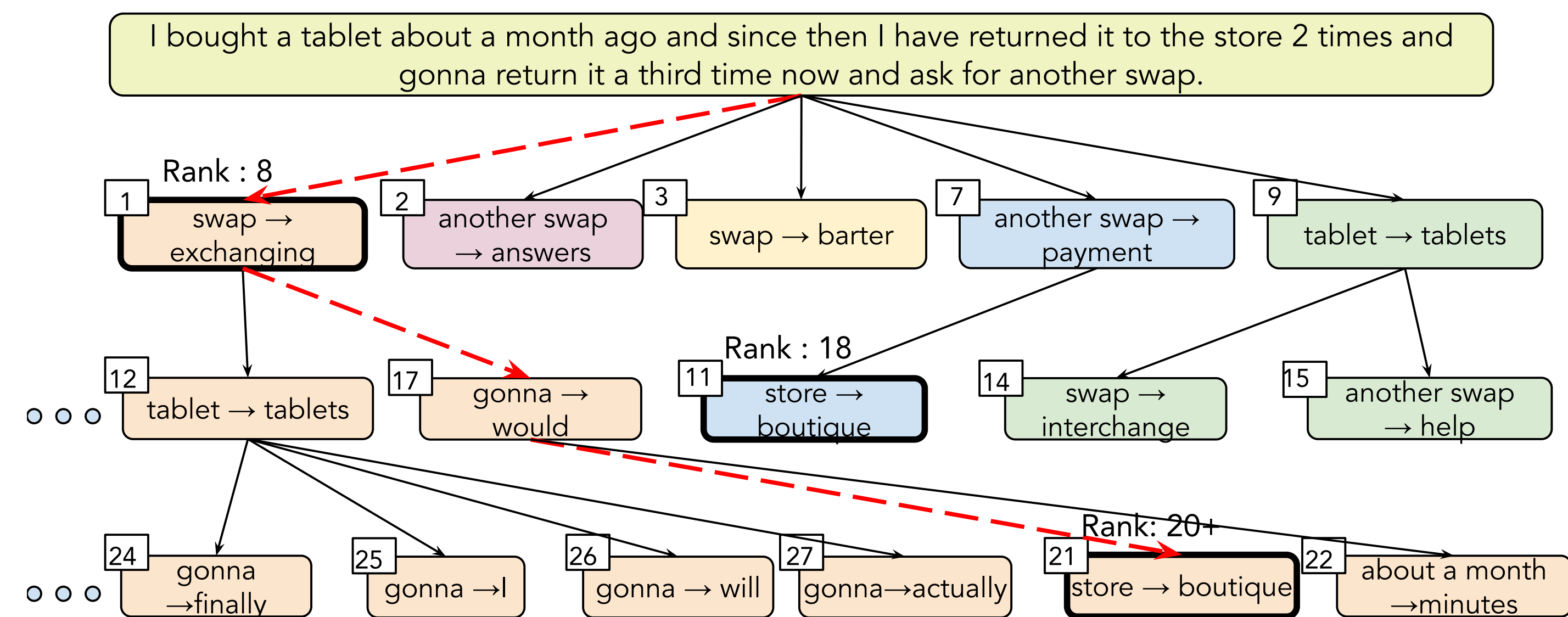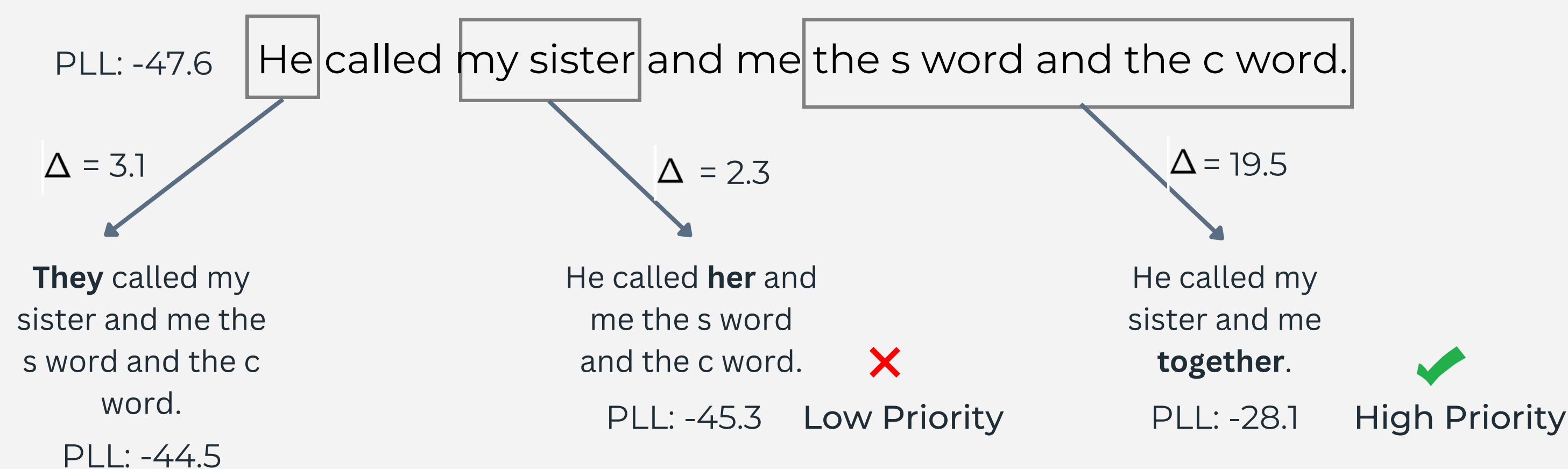- **Fidelity**: the semantic meaning of the sentence is preserved

## Which part of the sentence to attack first ?

**STEP 1:** Create a constituency-based parse tree for the sentence.
**STEP 2:** Prioritise parse tree nodes to attack based on PLL scores.
**STEP 3:** Rank candidates based on *PLL score difference* after replacing masks with BERT suggestions.
**PLL:** Probability of a sentence from BERT, by iteratively masking every word in the sentence and then summing the log probabilities.

PLL: -47.6 | He | called my | sister | and me | the s word and the c word.

$\Delta$ = 3.1

$\Delta$ = 2.3

$\Delta$ = 19.5

**They** called my sister and me the s word and the c word.
PLL: -44.5

He called **her** and me the s word and the c word. ✗ Low Priority
PLL: -45.3

He called my sister and me **together**. ✓ High Priority
PLL: -28.1



Constituency Parse Tree

Original Ranking Results

Attack Locations

Query

Feedback Results

Ranked Documents

top-K ... K ... irrelevant

1 2 ... K ... N-1 N

Rank



I bought a tablet about a month ago and since then I have returned it to the store 2 times and gonna return it a third time now and ask for another swap.

Rank : 8
1 swap → exchanging
2 another swap → answers
3 swap → barter
7 another swap → payment
9 tablet → tablets

12 tablet → tablets
17 gonna → would
Rank : 18
11 store → boutique
14 swap → interchange
15 another swap → help

24 gonna →finally
25 gonna →I
26 gonna → will
27 gonna→actually
Rank: 20+
21 store → boutique
22 about a month →minutes

## Types of Attack

Attacking by generating replacements using a combination of:
1) **BERT masked language model** : maintains grammar; independent of the phrase being replaced
2) **Synonyms in Counter-fitting vector space**: depends on phrase being replaced; decreases grammar quaility

## Multi-level attack for multi-level perturbations

Expanding single-phrase attack to multiple levels using Beam Search.

$$f(s_{paraphrased}) = (1 - \alpha) * \underbrace{Sim(s_{org}, s_{paraphrased})}_{\substack{\text{semantic similarity} \\ \text{(distance from origin)}}} + \alpha * \underbrace{\frac{(Rank(s_{paraphrased}, D_{source}) - 1)}{20}}_{\substack{\text{non-locatablity of source} \\ \text{(estimated distance to target)}}}$$

## Results

We succeed in disguising 82% of the queries when there are 3 beam levels and 5 nodes per parse tree are expanded.

### References

- Reagle, J. and Gaur, M. 2022. Spinning words as disguise: Shady services for ethical research?. First Monday, vol. 27, no. 1, Jan. 2022
- Jin Yong Yoo and Yanjun Qi. 2021. Towards Improving Adversarial Training of NLP Models. EMNLP 2021

R&D SHOWCASE 2023

TECH ALCHEMY
Ag > Au

INTERNATIONAL INSTITUTE OF INFORMATION TECHNOLOGY
HYDERABAD

IIIT Hyderabad 25