

# Finding Dataset Shortcuts with Grammar Induction

Dan Friedman, Alexander Wettig, and Danqi Chen

Department of Computer Science, Princeton University



## Introduction

Many NLP datasets contain **shortcuts**, simple decision rules that achieve **high in-domain accuracy** but **fail to generalize** to the intended test distribution. **If we can identify shortcuts, we can mitigate them** by collecting more training data or using robust optimization methods.

**Prior work** Earlier work on finding shortcuts relied on **intuition** and **manually designed probes**. Prior work on automatically identifying shortcuts has focused on **simple features that can be explicitly enumerated**, like unigrams [6, 1], or **qualitative interpretability methods**, like saliency maps [2, 5].

**Proposal:** Induce **dataset-specific grammars** to **formally characterize patterns** in sentence and sentence-pair classification datasets.

## Method

Given text classification dataset  $\mathcal{D} \subseteq \mathcal{X} \times \mathcal{Y}$ , where  $y \in \mathcal{Y}$  is a categorical label and  $x \in \mathcal{X}$  is either a sentence or a pair of sentences,  $x = (x^a, x^b)$ .

### 1. Grammar induction

Induce a grammar for (unlabeled) training instances  $x_1, \dots, x_N$  and get maximum likelihood trees  $t_1, \dots, t_N$ .

• Sentence datasets: **context-free grammar** (CFG)

• Sentence-pair datasets: **synchronous CFG** [7]

Given a grammar  $\mathcal{G}$ , find parameters  $\theta^*$ :

$$\theta^* = \arg \max_{\theta} \sum_{i=1}^N \log \sum_{t: \text{yield}(t)=x_i} p(t \mid \mathcal{G}, \theta).$$

### 2. Finding features

Define a boolean indicator feature  $Z_s \in \{0, 1\}$  for each complete subtree  $s$  and calculate the mutual information:

$$I(Z_s; Y) = \sum_{z_s \in \{0, 1\}} \sum_{y \in \mathcal{Y}} p(y, z_s) \log \frac{p(y, z_s)}{p(y)p(z_s)}.$$

Group subtrees according to root label and majority class.

## Finding Shortcuts

### IMDb movie review dataset

Root	Desc.	Patterns	% Maj.
5	Negative actors	ed wood, steven seagal, uwe boll	95.5
29	Negative ratings	4 / 10, 3 / 10, 1 / 10, 2 / 10	96.8
8	Negative durations	30 minutes, 10 minutes, five minutes	76.7
5	Positive actors	walter matthau, jon voight, james stewart	88.6
29	Positive ratings	10 / 10, 8 / 10, 7 / 10	98.8
8	Positive durations	many years	69.1

### Subjectivity dataset

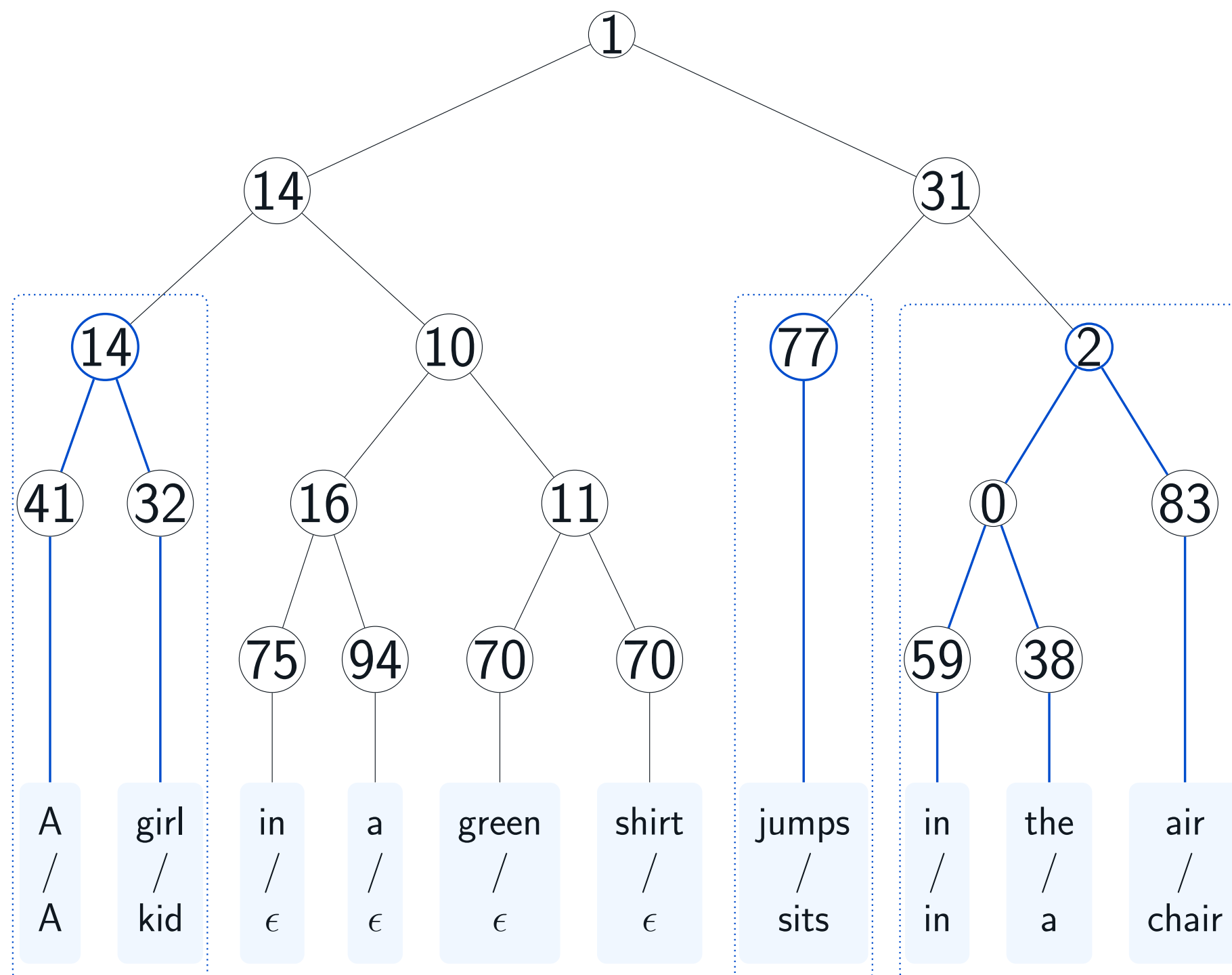
Root	Desc.	Patterns	% Maj.
27	Review NP	a movie, the film, the movie, this movie	86.3
3	Review VB	comes off, 's hard, makes up, 'd expect	87.4
27	Summary NP	his life, his wife, his father, his mother	80.1
3	Summary VB	finds himself, finds out, falls in love	85.5

### SNLI

Root	Desc.	Patterns	% Maj.
44	Copy verb	walking/walking, running/running	68.4
14	Subj. hypernym	a woman/a person, a man/a human	45.7
4	Expletive	a /there is, /there are, two /there are	63.0
14	Subj. antonym	a man/a woman, a boy/a girl, a dog/a cat	82.5
78	Verb antonym	standing/sitting, walking/sitting	92.6
85	Adj. antonym	black/white, red/blue, /empty, /living	76.9
35	Add object	/[UNK], /work, /get, /friends, /park	59.8
85	Add adj.	/tall, /sad, /[UNK], /new, /big	72.1
17	Add PP phrase	/to work, /to get, /to buy, /the park	71.4

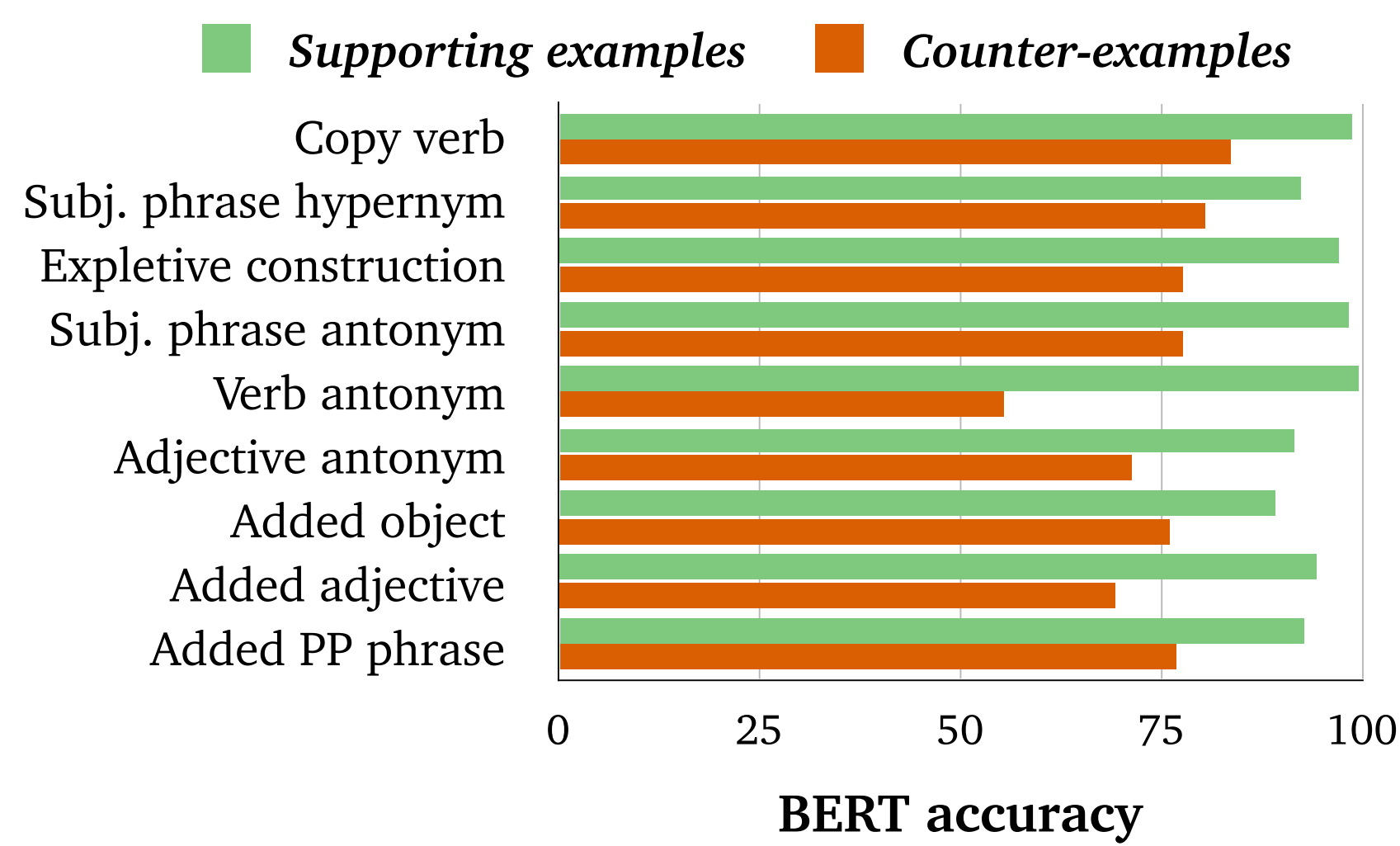
### Quora Question Pairs

Root	Desc.	Patterns	% Maj.
70	Additions	/[UNK], /in, /a, /-, /for	60.7
49	Deletions	[UNK]/, in/, a/, like/, of/	61.6
59	Change Q	why/how, why/what, how/why, why/can	70.8
14	How-to	how can/how can, how do/how can	66.2
25	Topics	new year/new year, world war/world war	82.9
59	Same Q	how/how, why/why, when/when	60.3

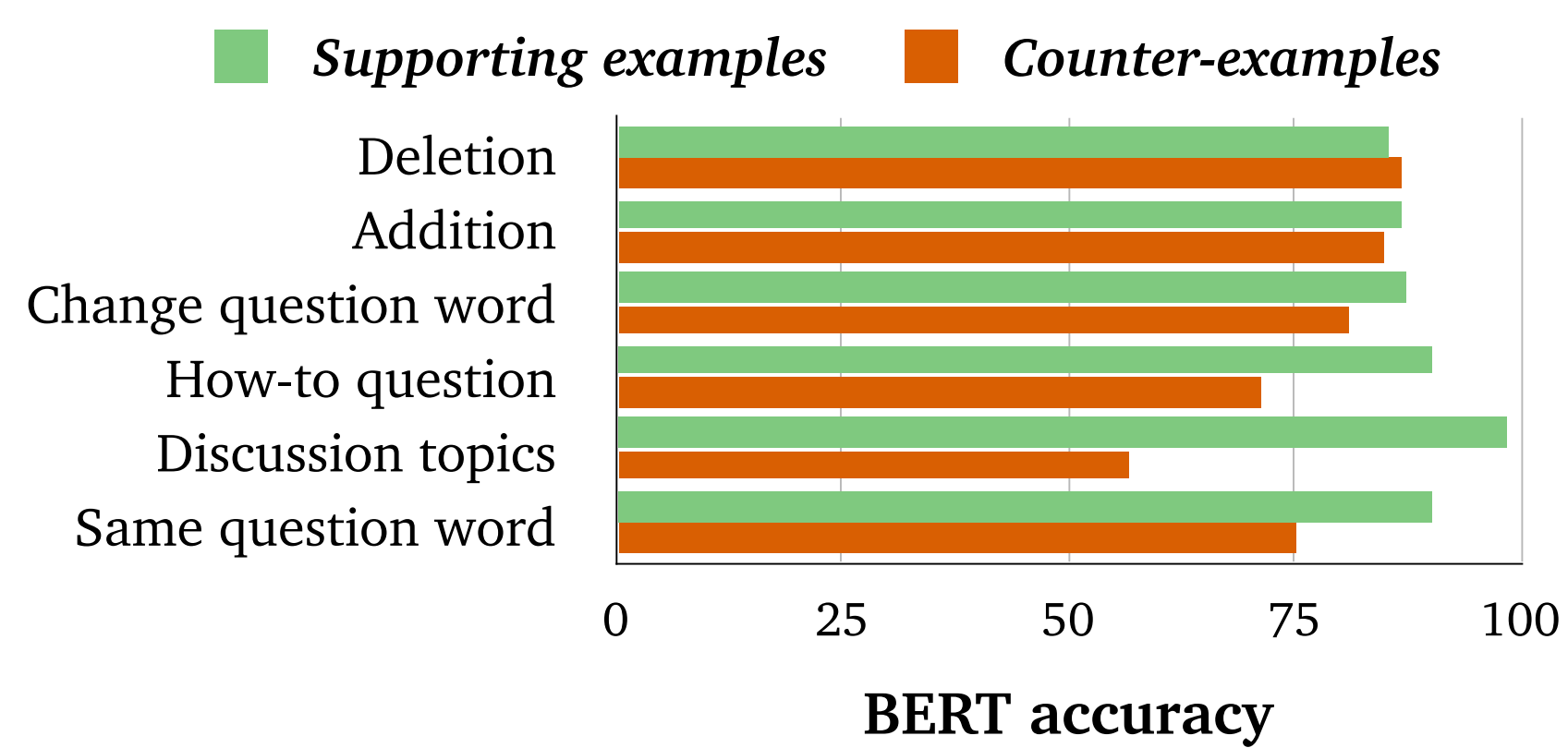


## Do Models Use these Shortcuts?

### SNLI



### QQP

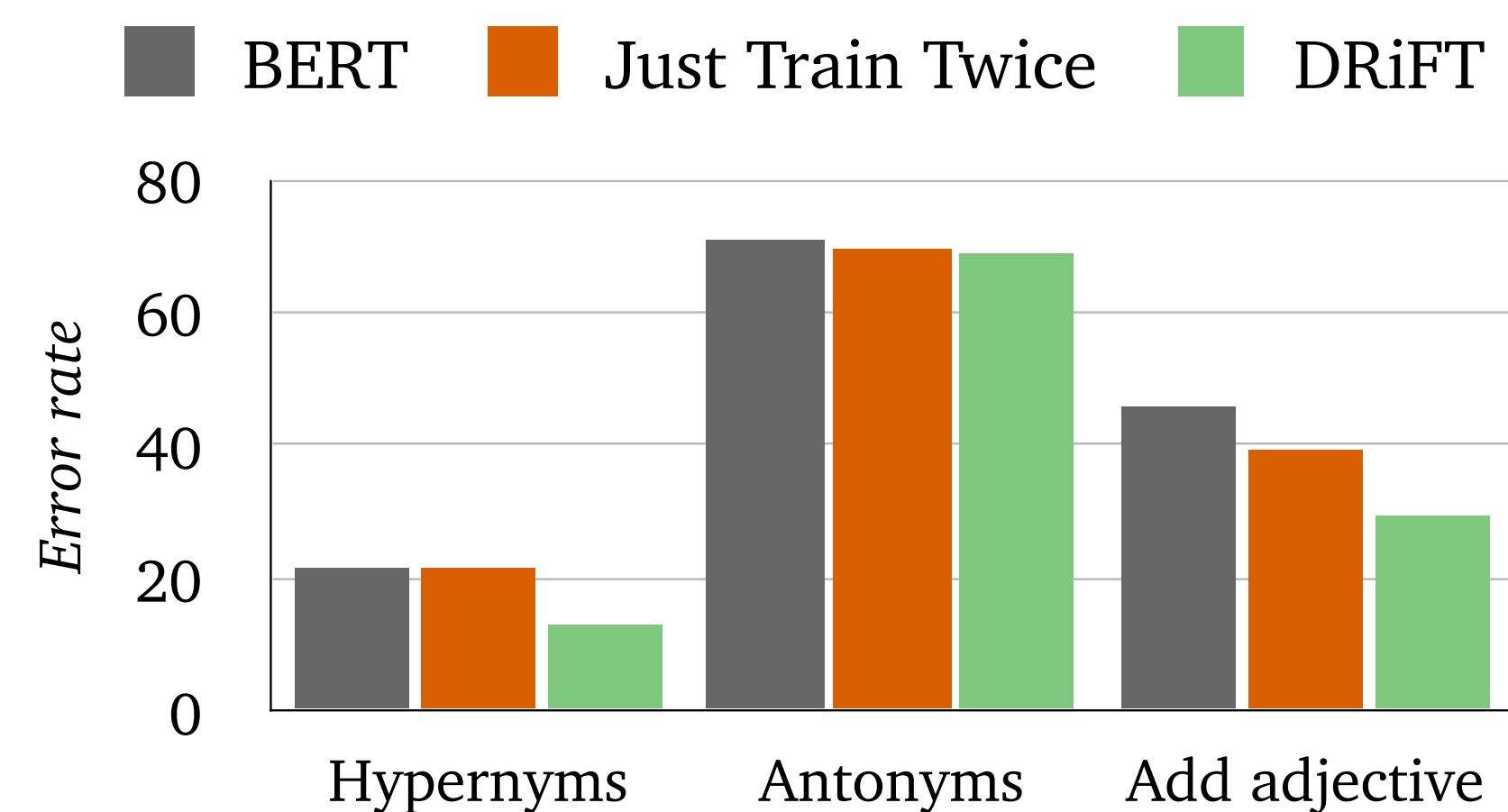


## Generating counter-examples

### SNLI

Edit	# Sets	Error
<b>Hypernyms</b> A man is smoking at sunset. A <b>man</b> +person smoking a cigarette.	389	21.8±0.8
<b>Antonyms</b> Two black dogs splash around on the beach. The dogs are playing with a <b>+white</b> ball.	281	71.1±3.8
<b>Add adjective</b> A man taking photos of nature. A <b>+sad</b> man is taking photos of a wedding.	1,470	45.6±8.4

## Mitigating Shortcut Learning

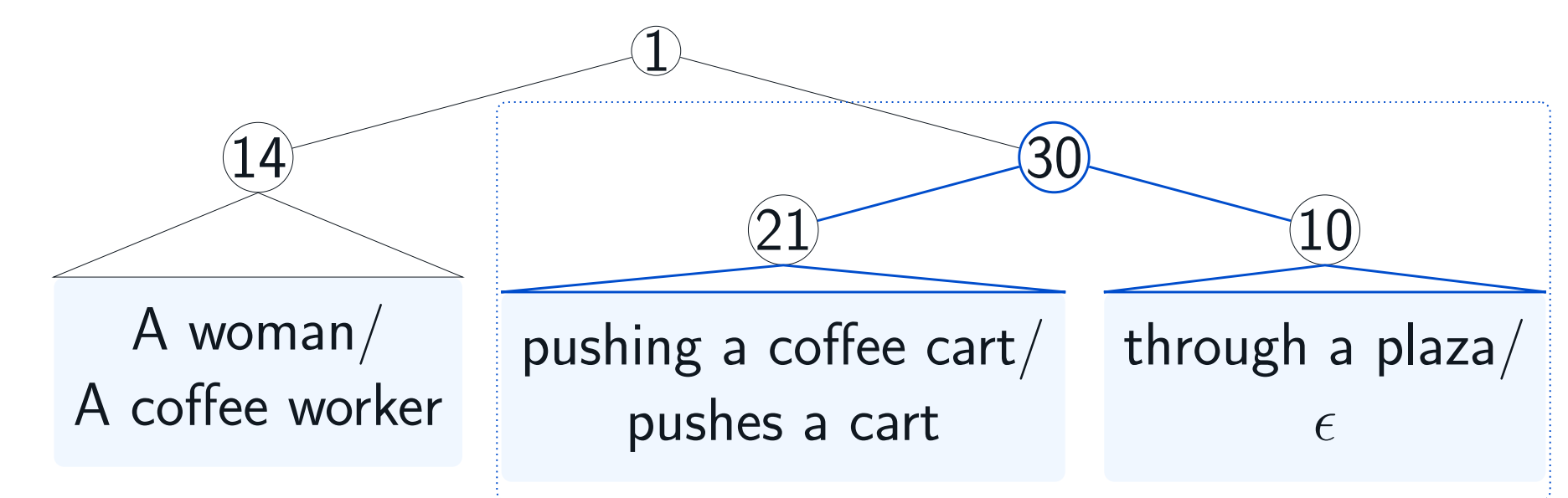


Robust optimization methods on SNLI. Just Train Twice [4] assumes **shortcuts** are **unknown**. DRIFT [3] incorporates information about the **known shortcuts** we discovered.

<b>14</b> A girl/A kid	<b>77</b> jumps/sits	<b>2</b> in the air/in a chair
Top subtrees for <i>entailment</i>		
A man/A person A man/A human	walking/walking walk/walking	in the grass/outside down the street/outside
Top subtrees for <i>contradiction</i>		
A man/A woman A woman/A man	/sitting stand/sit	at night/during the day in the air/down
Top subtrees for <i>neutral</i>		
A man/A tall human A man/An old man	/competes running/running	down the street/home in the sand/on the beach

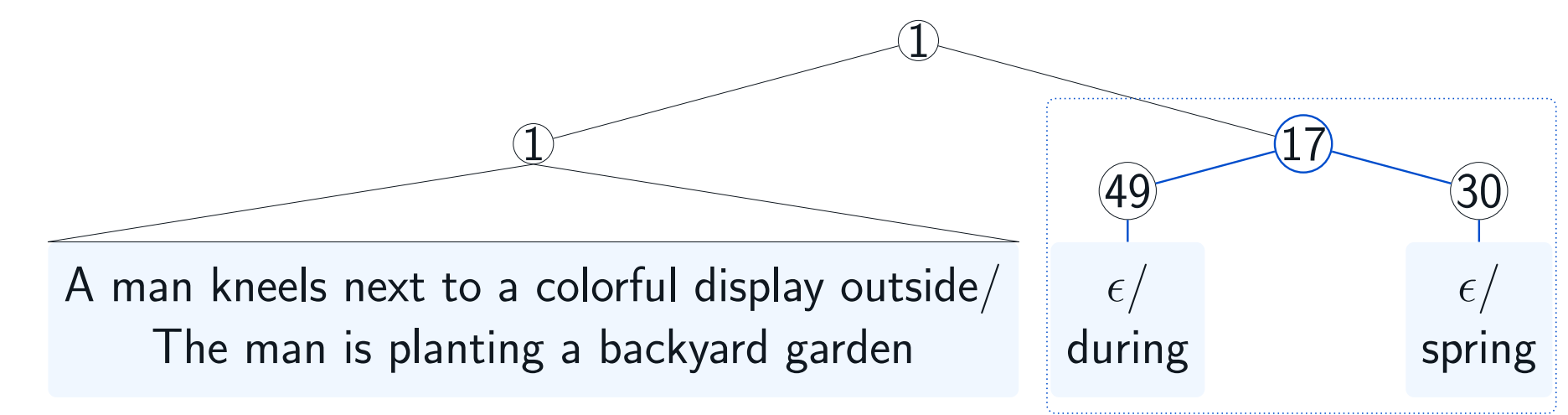
## Comparing grammar features and saliency maps

A woman **pushing** a **coffee** cart through a plaza. A **coffee** worker **pushes** a cart.



Label: Neutral, BERT prediction: Entailment

A man kneels **next** to a colorful display **outside**. The man is planting a backyard **garden** during **spring**.

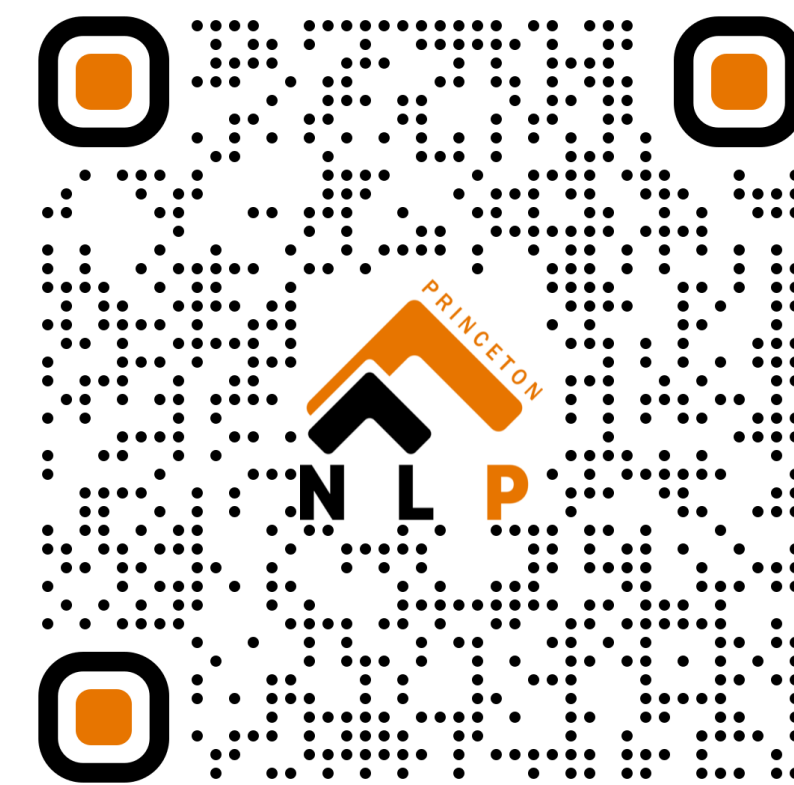


Label: Contradiction, BERT prediction: Neutral

## Conclusions

An approach for automatically **finding dataset shortcuts** by inducing **dataset-specific grammars**. The grammar can be used to:

- **Discover** interesting shortcut features
- **Diagnose** classifier errors
- **Mitigate** shortcut learning



## References

- [1] Matt Gardner et al. "Competency Problems: On Finding and Removing Artifacts in Language Data". In: *EMNLP*. 2021.
- [2] Xiaochuang Han, Byron C Wallace, and Yulia Tsvetkov. "Explaining Black Box Predictions and Unveiling Data Artifacts through Influence Functions". In: *ACL*. 2020.
- [3] He He, Sheng Zha, and Haohan Wang. "Unlearn Dataset Bias in Natural Language Inference by Fitting the Residual". In: *DeepLo*. 2019.
- [4] Evan Z Liu et al. "Just Train Twice: Improving group robustness without training group information". In: *ICML*. 2021.
- [5] Pouya Pezeshkpour et al. "Combining Feature and Instance Attribution to Detect Artifacts". In: *ACL Findings*. 2022.
- [6] Zhao Wang and Aron Culotta. "Identifying Spurious Correlations for Robust Text Classification". In: *EMNLP Findings*. 2020.
- [7] Dekai Wu. "Stochastic inversion transduction grammars and bilingual parsing of parallel corpora". In: *Computational Linguistics* (1997).