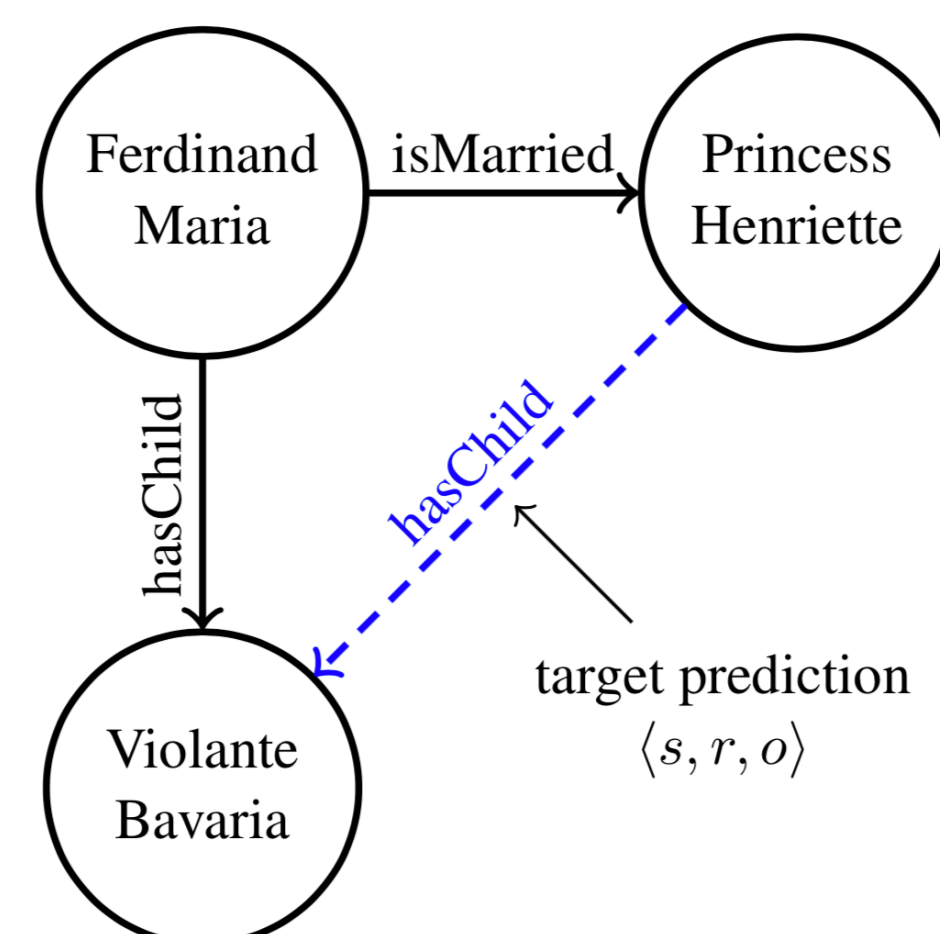


## Graph Embeddings for Link Prediction

In this work, we propose efficient adversarial modifications for link prediction models to evaluate robustness, and study interpretability and error correcting.

**Completing Knowledge Graphs:**  
Predicting a missing link from observed graph structure.



- Existing models:
  - Embed  $s, r$ , and  $o$
  - Maximize score  $\psi(s, r, o)$  for observed facts
    - DistMult:  $e_s R e_o$
    - ConvE:  $f(\text{vec}(f([\bar{e}_s; \bar{r}_r * w]))) W e_o$
- Embeddings are inscrutable...
  - Are these embeddings robust to small changes?
  - Can we explain why a fact/link was predicted?

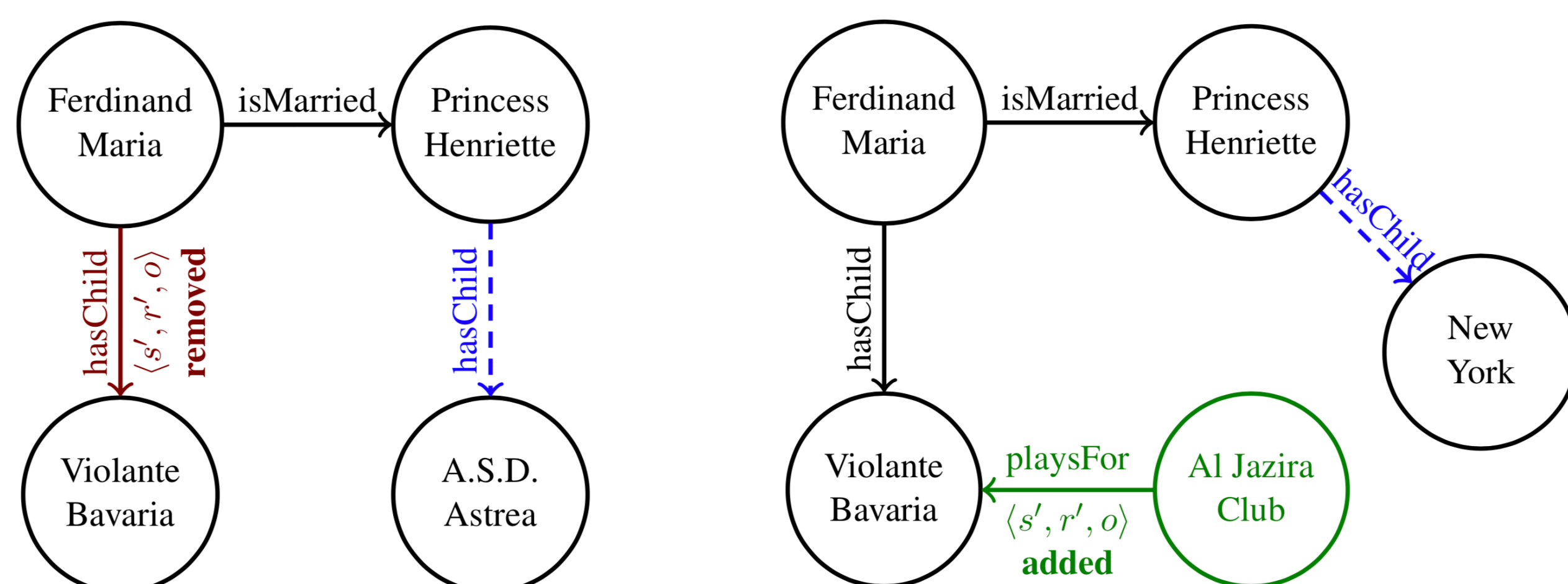
## Adversarial Modifications (CRIAGE)

- Completion Robustness and Interpretability via Adversarial Graph Edits (CRIAGE)

Minimally change the graph so that target fact prediction changes the most after embeddings are relearned.

Removing an existing link

Adding a fake link



## Efficiently Identifying the Modification

For target triple  $\langle s, r, o \rangle$  and graph  $G$ , we identify:

- Removing / Adding**: Find  $(s', r', o)$  such that score  $\psi(s, r, o)$  trained on  $G$  is maximally different from score  $\bar{\psi}(s, r, o)$  trained after removing or adding  $(s', r', o)$ :

$$\text{argmax}_{(s', r')} \psi(s, r, o) - \bar{\psi}(s, r, o)$$

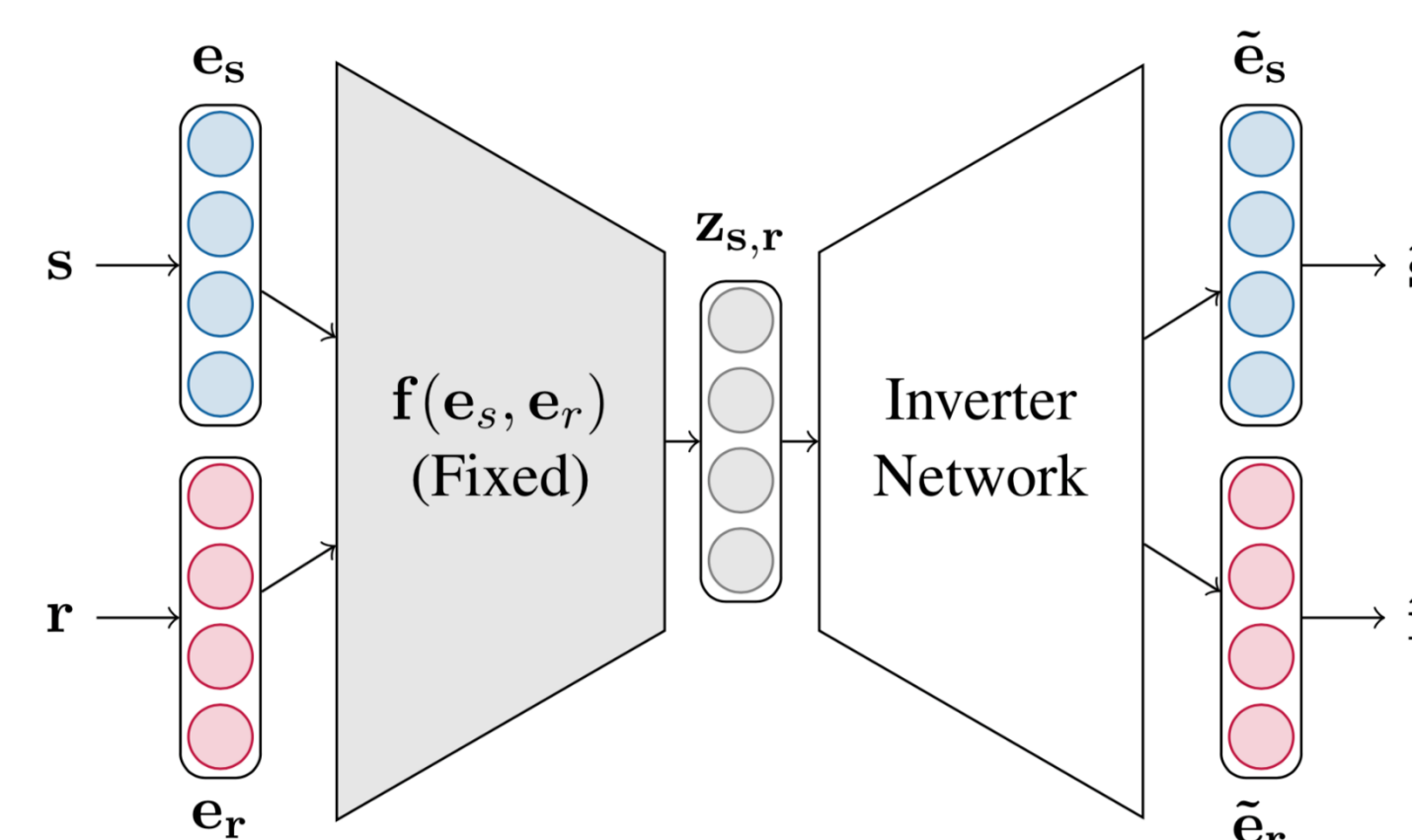
Two Primary Challenges:

- Retraining is too expensive: Taylor approximation on gradient of loss and utilizing graph structure.

$$\Delta(\nabla_e \text{loss}(e)) = H_e(\text{loss}) \times (e - \bar{e})$$

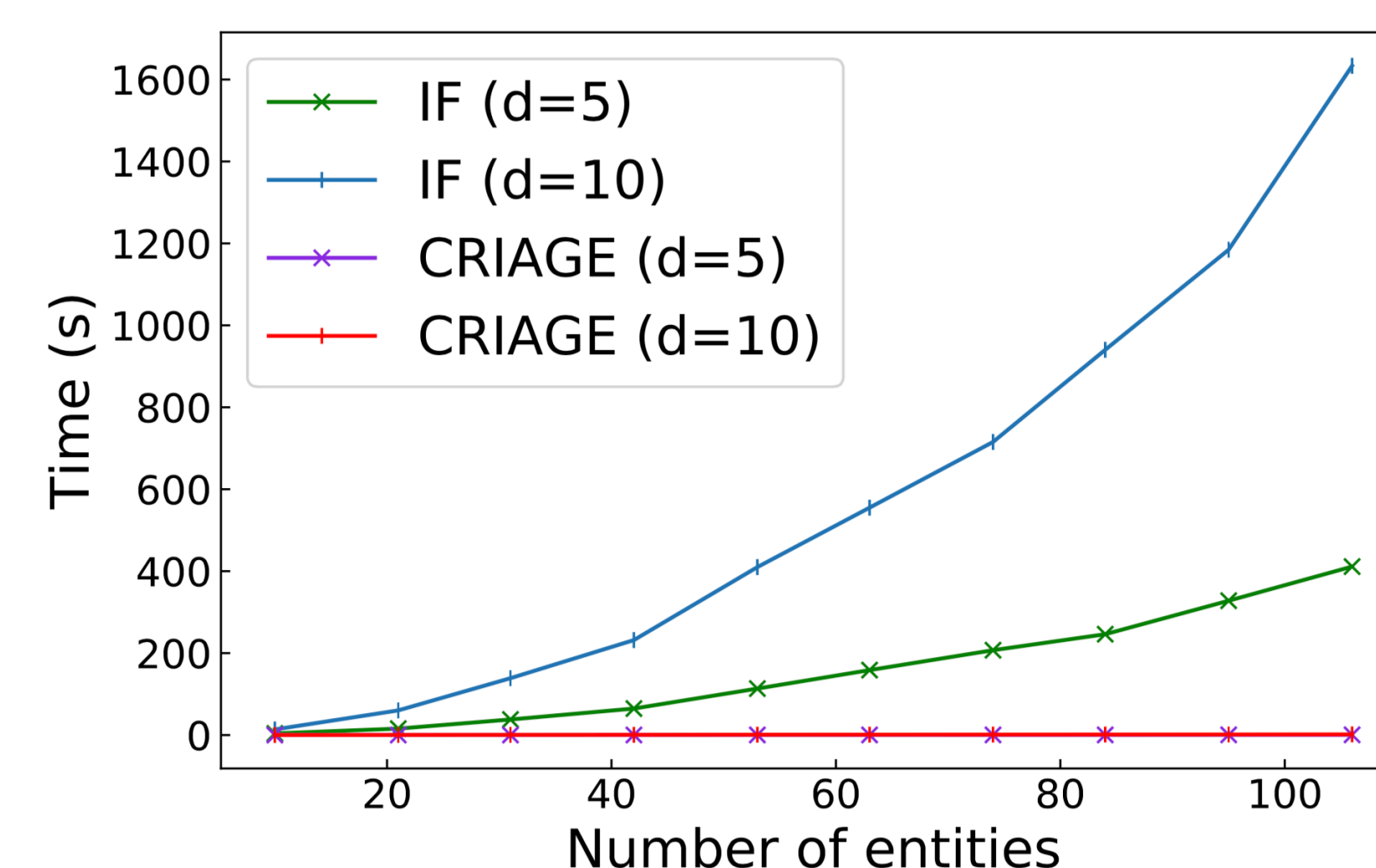
$e, \bar{e}$  = optimum embedding &  $H$  = Hessian  
 $\Rightarrow \bar{e} = e - H_e(\text{loss})^{-1} \times \Delta(\nabla_e \text{loss}(e))$

- Too many links to search: Learn a continuous space of links using an inverter, and use gradient descent.



## CRIAGE vs Influence Functions

- Influence Functions (IF)\*:
  - Similar motivation, but doesn't exploit graph structure

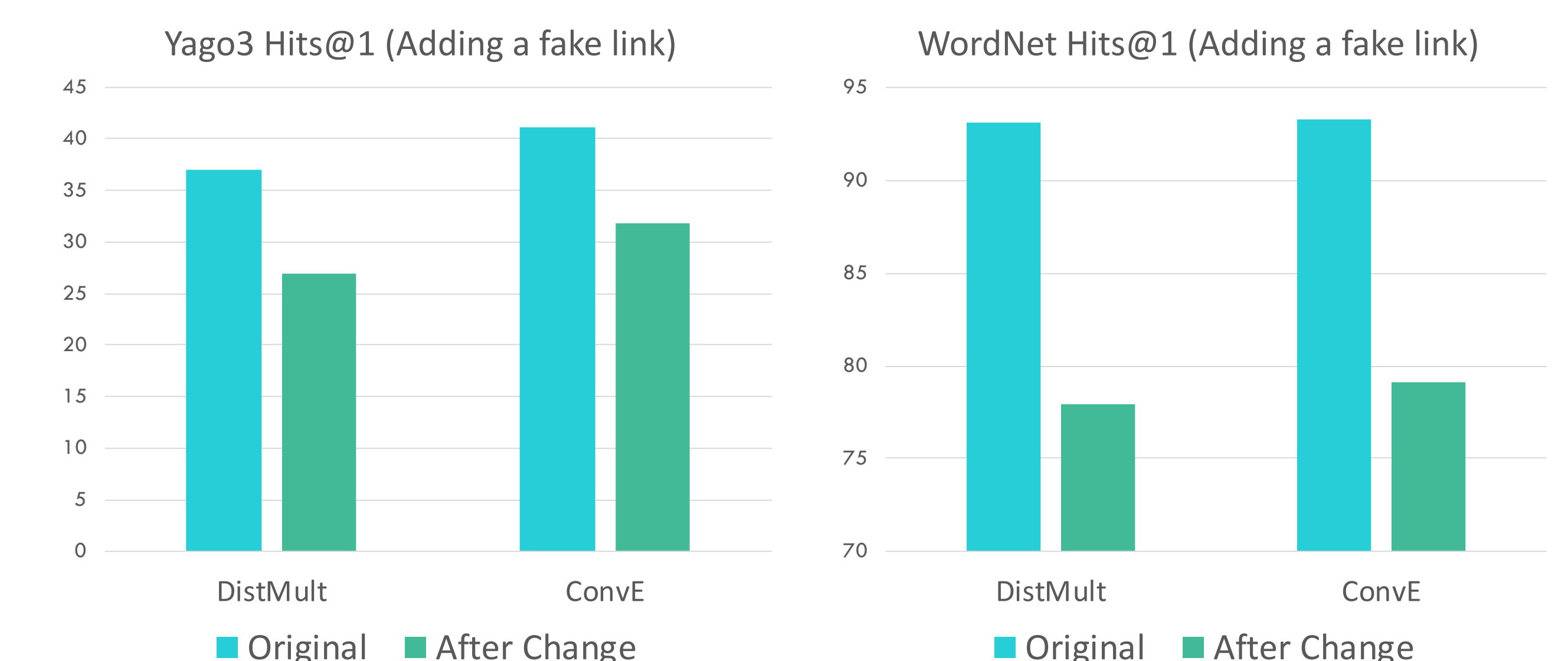


\* Koh, Pang Wei, and Percy Liang. "Understanding black-box predictions via influence functions."

## Robustness and Interpretability

### Robustness

Does adding a fake link affect performance?



### Interpretability

Which link, when removed, changes the prediction?  
Find common patterns in removed link  $R(a, b)$

DistMult and ConvE:  $\text{isMarriedTo}(a, c) \wedge \text{hasChild}(c, b) \Rightarrow \text{hasChild}(a, b)$

Only in DistMult:  $\text{playsFor}(a, c) \wedge \text{isLocatedIn}(c, b) \Rightarrow \text{wasBornIn}(a, b)^*$   
 $\text{isAffiliatedTo}(a, c) \wedge \text{isLocatedIn}(c, b) \Rightarrow \text{diedIn}(a, b)^*$

Only in ConvE:  $\text{hasAdvisor}(a, c) \wedge \text{graduatedFrom}(c, b) \Rightarrow \text{graduatedFrom}(a, b)$   
 $\text{influences}(a, c) \wedge \text{influences}(c, b) \Rightarrow \text{influences}(a, b)$

\* Identified as rules by [Yang et. al. 2015]

### Error Correcting

Introduce errors and see if we can detect it.  
Choose neighbor w/ least  $\psi(s, r, o) - \bar{\psi}(s, r, o)$  as incorrect.

