# Adversarial Attacks on Knowledge Graph Embeddings via Instance Attribution Methods
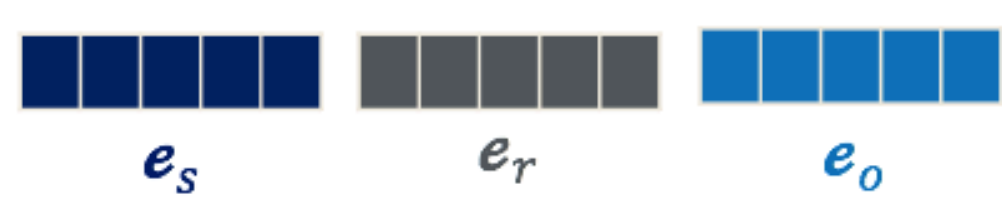
**Peru Bhardwaj, John Kelleher, Luca Costabello, Declan O'Sullivan**

Minimize $\mathcal{L}$ by updating $e_s, e_r, e_o$

$e_s$  $e_r$  $e_o$

Scores for positive triples are higher than scores for negative triples

Scoring Function
$f(e_s, e_r, e_o)$

$y_{s,r,o} \in [0,1]$
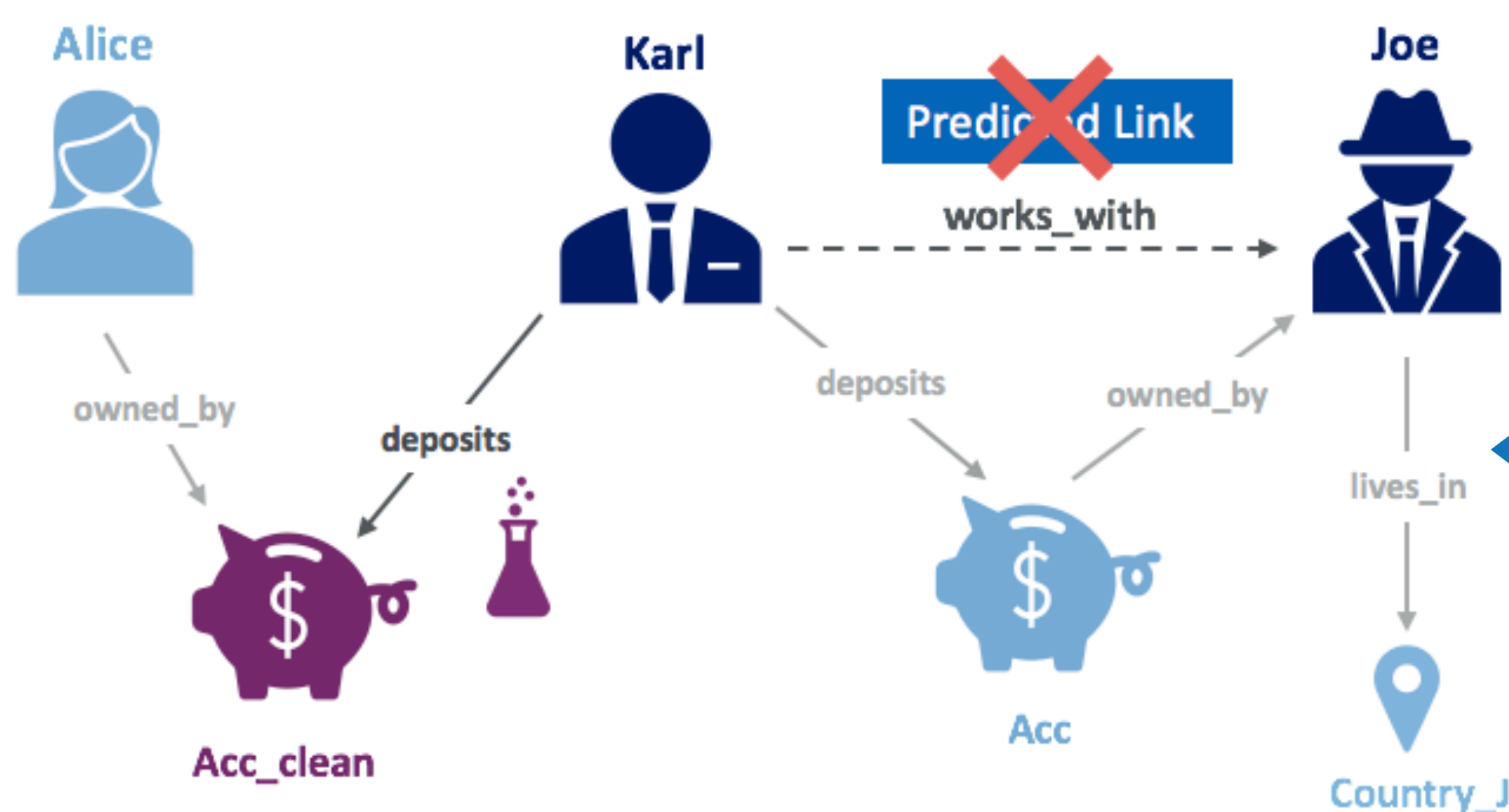
Loss Function
$\mathcal{L}(\hat{y}_{s,r,o}, y_{s,r,o})$

$\hat{y}_{s,r,o} \in \{0,1\}$

Generate negatives by corrupting s/o

| s | r | o |
|---|---|---|
| Karl | credit_card | Card_X |
| Karl | credit_card | Card_Y |
| Karl | credit_card | Card_K |
| Person_X | credit_card | Card_K |
| Person_Y | credit_card | Card_K |

## What are Knowledge Graph Embeddings?

## Motivation

Knowledge graph embedding models enable representation learning on multi-relational graphs and are used in security sensitive domains. But their security vulnerabilities have received little attention.

## What are Instance Attribution Methods?

Methods to identify the training triple that is most influential to model's prediction on target triple

1. Similarity between feature representations of instances (Instance Similarity)
2. Similarity between gradients due to instances (Gradient Similarity)
3. Influence Functions

## Adversarial Deletions against Knowledge Graph Embeddings

Target triple becomes False by _deleting_ training triple
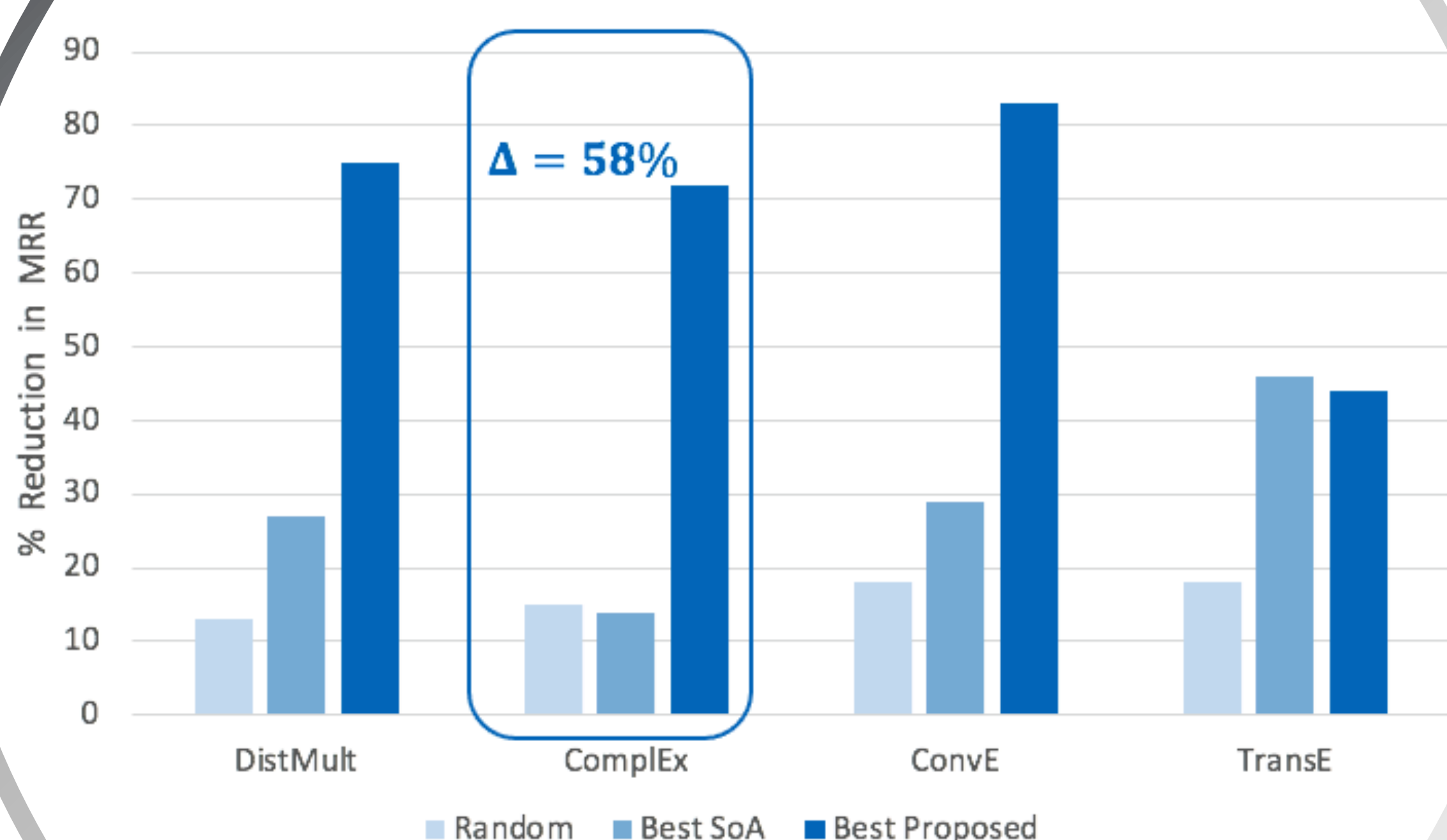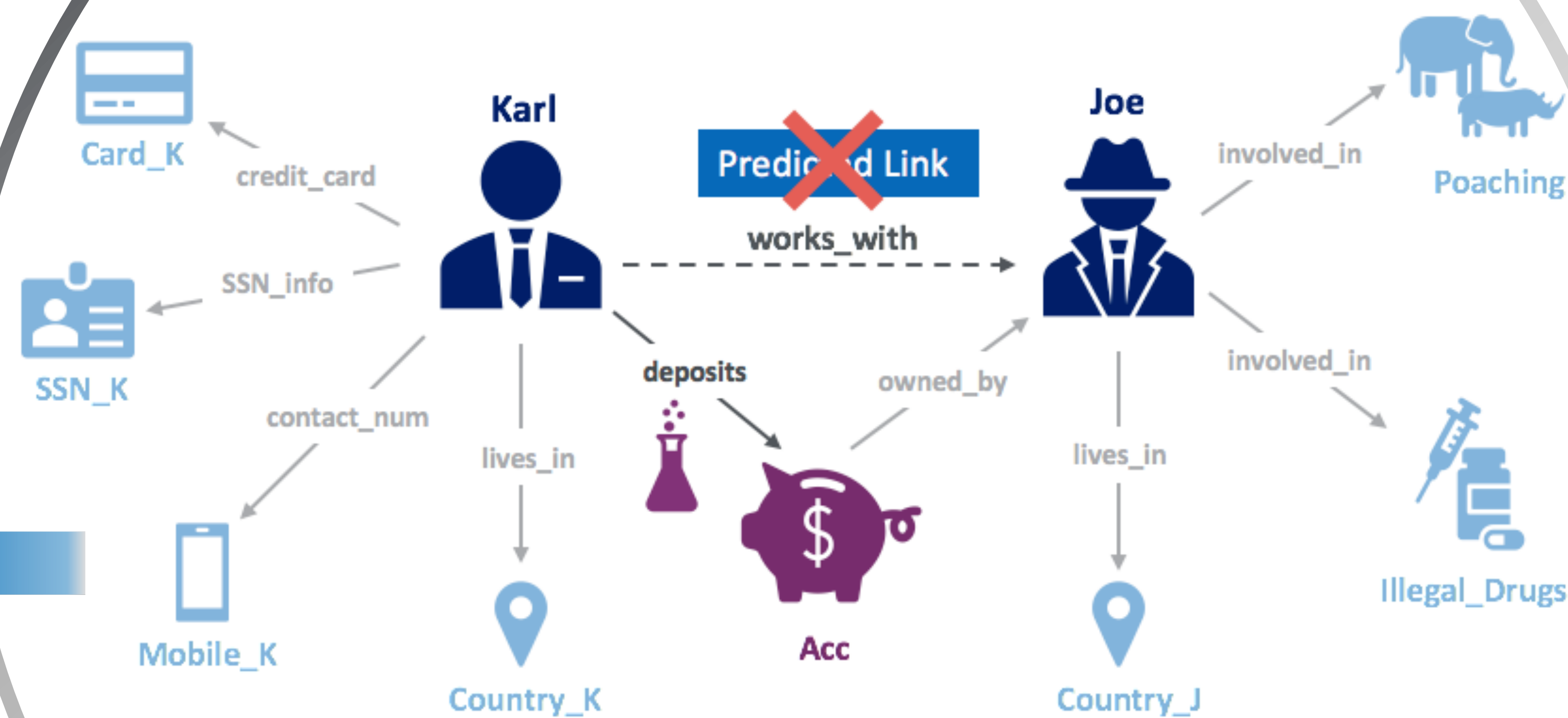
## Adversarial Additions

Target triple becomes False by _adding_ triples to training set

The malicious attacker selects adversarial addition by replacing an entity of the influential triple with the most dissimilar entity in embedding space.
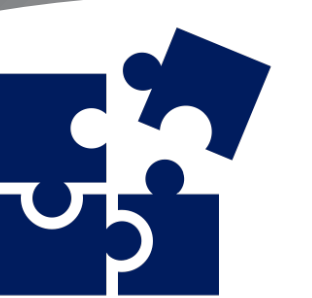
Example scenario for adversarial attacks against KGE models for fraud detection. The missing triple to predict is (Karl, works_with, Joe). Original KGE model predicts this triple as True.

But a malicious attacker uses the instance attribution methods to identify the influential triple and deletes it.

Now the KGE model predicts the target triple as False.

## Future Work

1. Influence of training sub-graph instead of individual triples
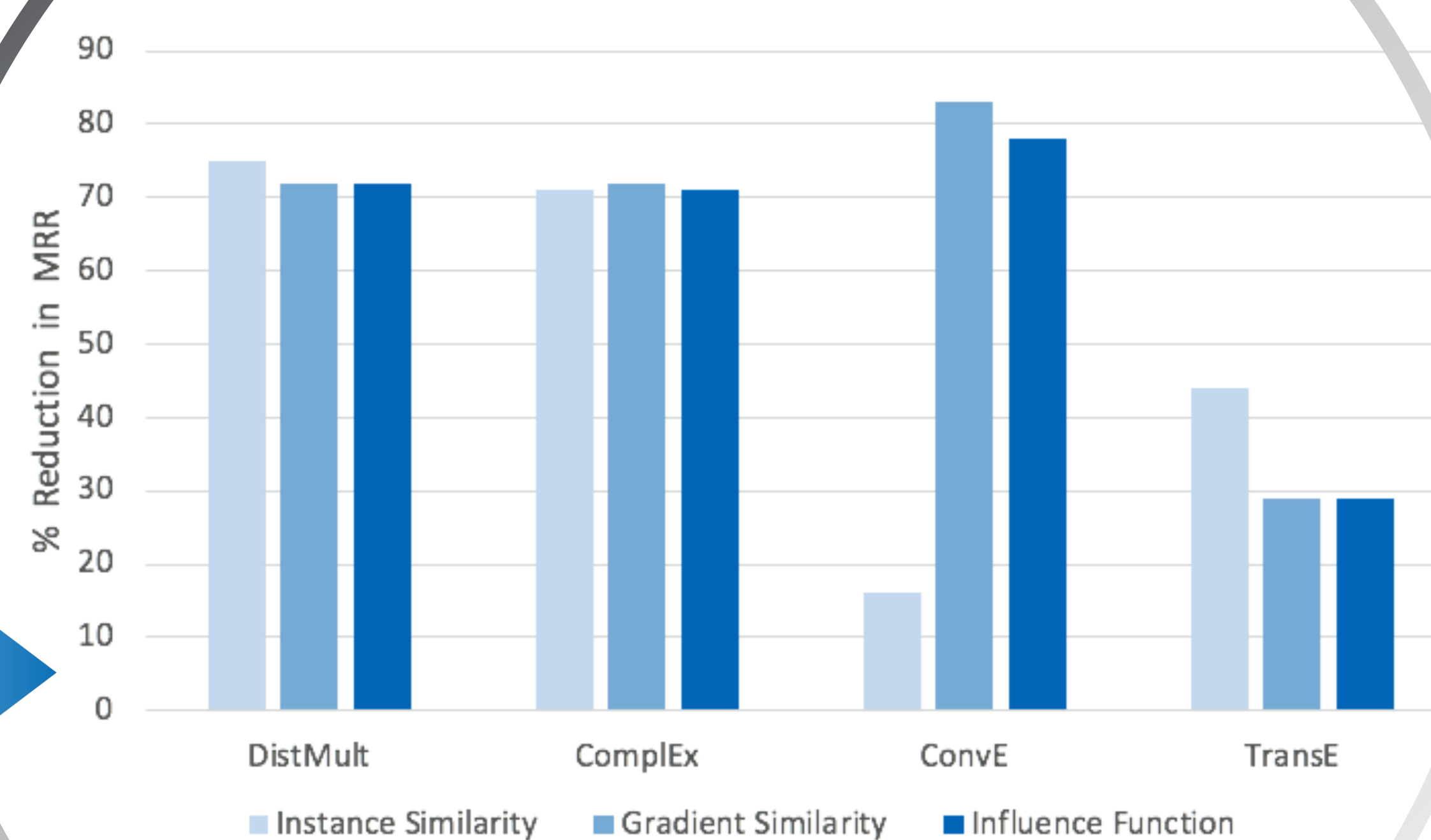2. Adversarial robustness of KGE models

## Challenges

How to measure impact of candidate edit on target prediction?
✓ Use instance attribution methods

How to search through combinatorial space of candidate adversarial additions?
✓ Replace entity in influential triple

## Proposed Vs State-of-Art

Adversarial Deletions – WN18RR

% Reduction in MRR

$\Delta = 58\%$

DistMult  ComplEx  ConvE  TransE

■ Random  ■ Best SoA  ■ Best Proposed

Adversarial Deletions – WN18RR

% Reduction in MRR

DistMult  ComplEx  ConvE  TransE

■ Instance Similarity  ■ Gradient Similarity  ■ Influence Function

## Instance Attribution Methods