



*Fork me on GitHub*

Luca Costabello  
Sumit Pai  
Chan Le Van  
Rory McGrath  
Nicholas McCarthy

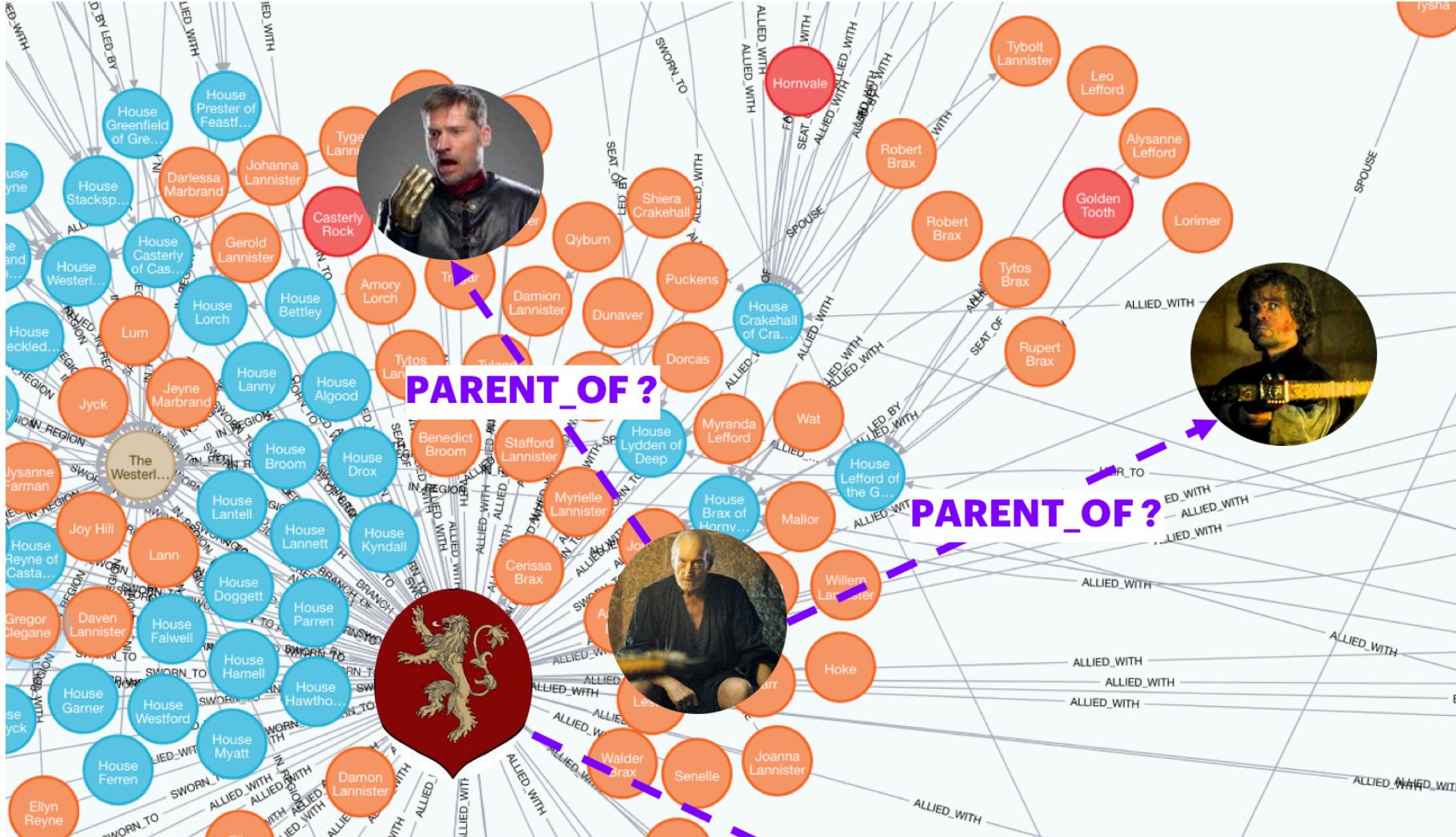
# AmpliGraph

## A Library for Graph Representation Learning

[ampligraph.org](http://ampligraph.org)

Open source Python library that predicts links between concepts in a knowledge graph.

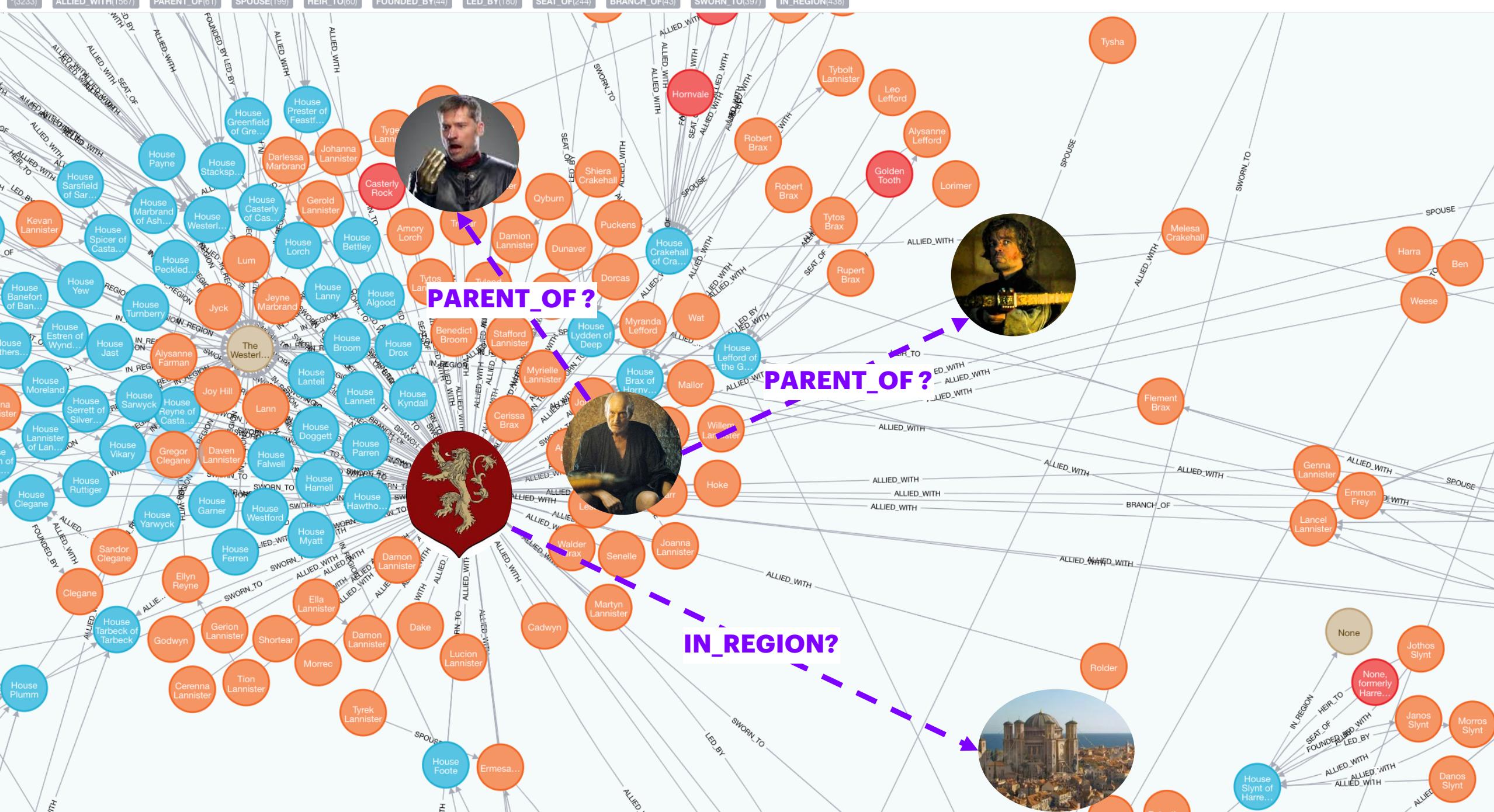
AmpliGraph is a suite of neural machine learning models for relational Learning, a branch of machine learning that deals with supervised learning on knowledge graphs.

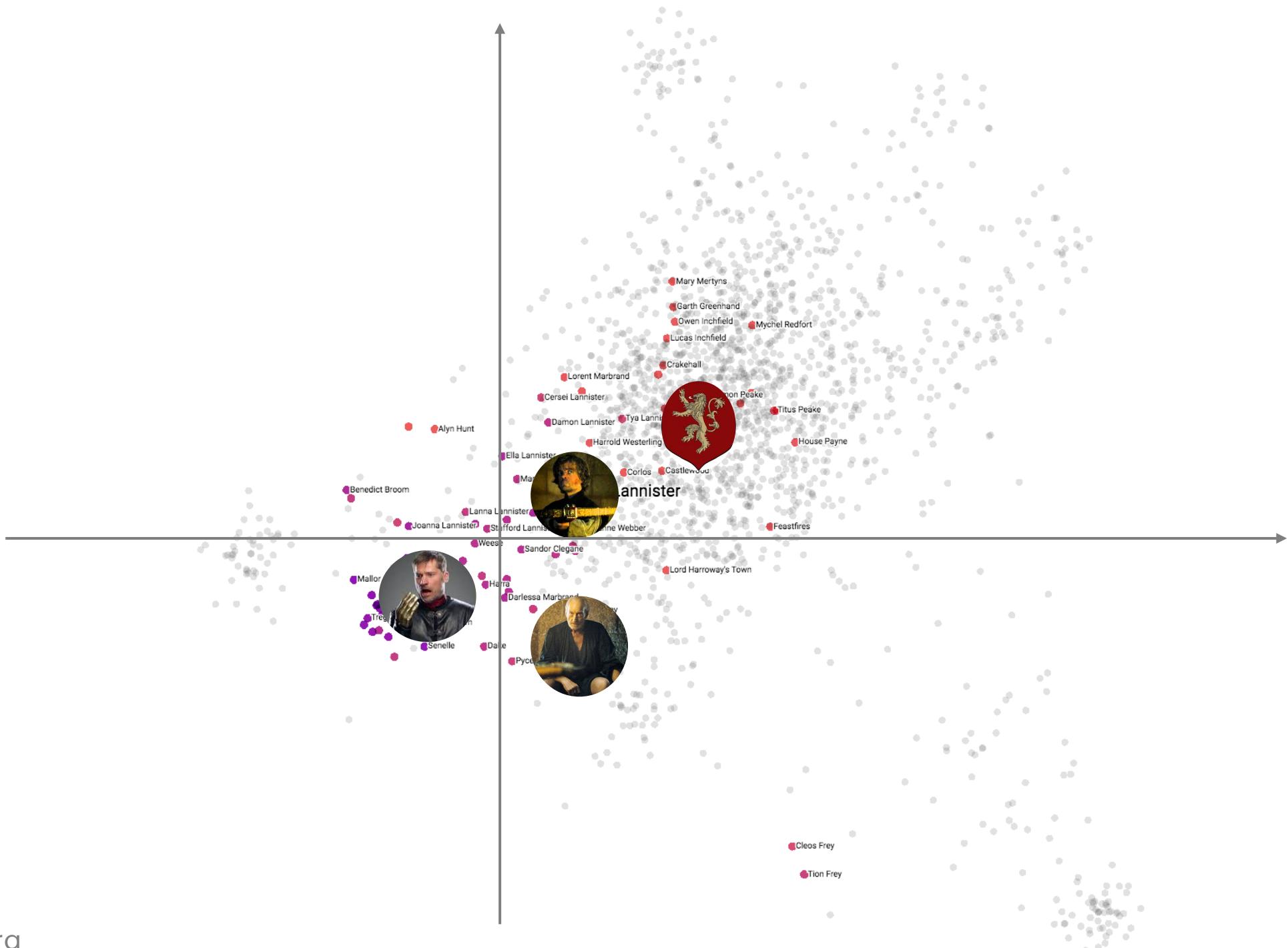


# **Use AmpliGraph if you need to:**

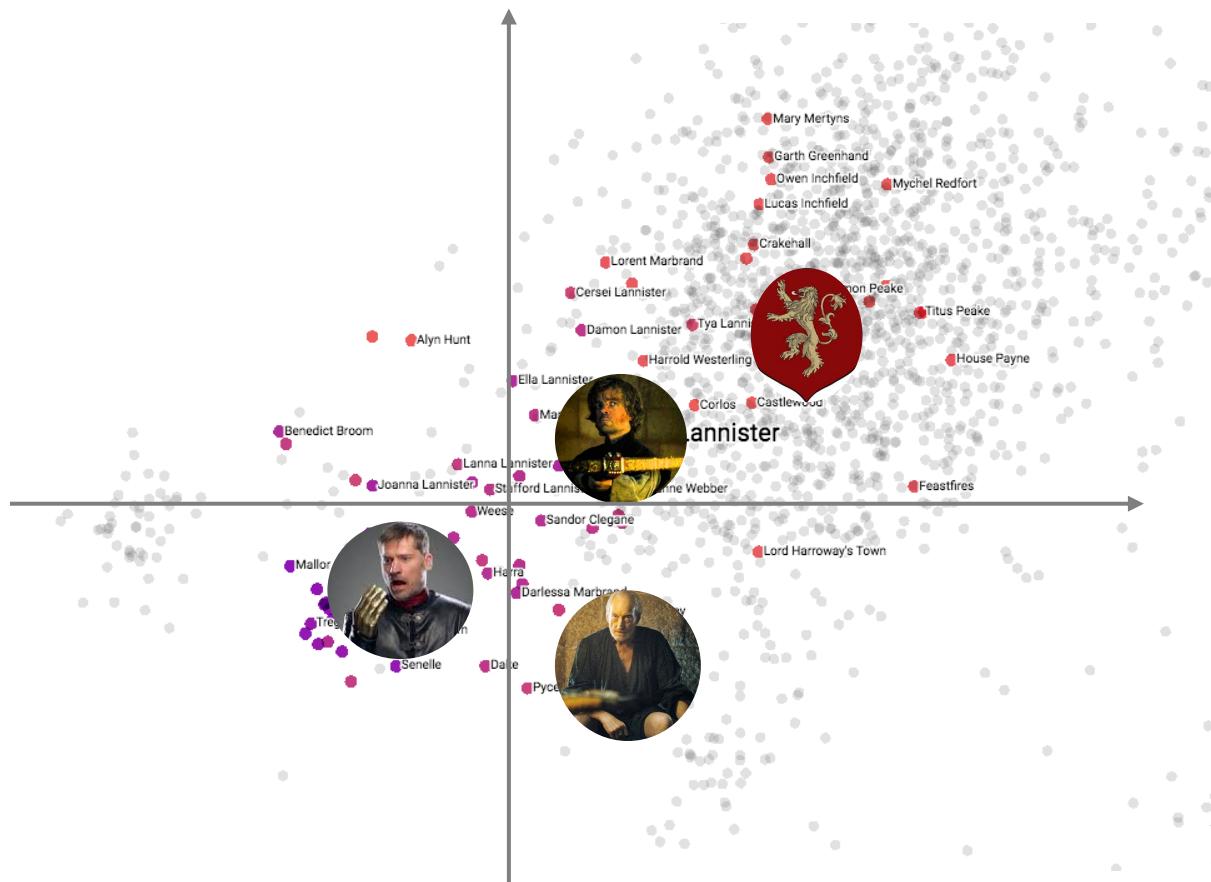
- Discover new knowledge from an existing graph
- Complete a large and incomplete knowledge graph
- Cluster similar concepts according to their relations
- Identify near-duplicate concepts
- Use graph embeddings in downstream tasks

AmpliGraph's machine learning models generate **knowledge graph embeddings**, vector representations of concepts in a metric space:

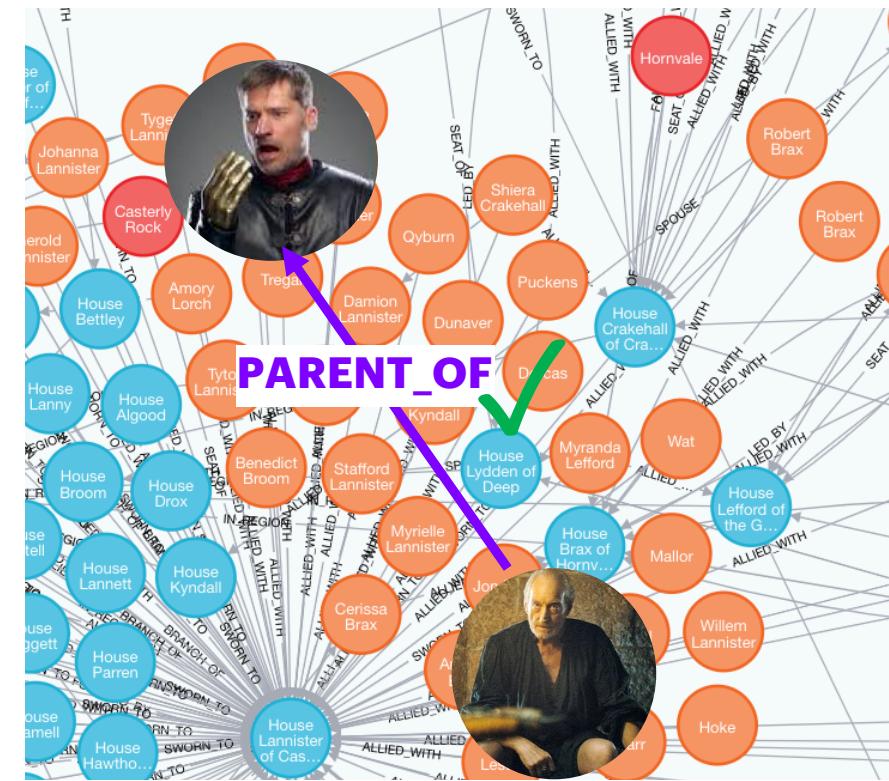




It then combines embeddings with model-specific scoring functions to predict unseen and novel links:



$$f(\text{Tywin Lannister}, \text{PARENT\_OF}, \text{Jaime Lannister}) = 0.85$$



# Anatomy of a KG Embedding Model

- Knowledge Graph (KG)  $\mathcal{G}$
- Scoring function for a triple  $f(t)$
- Loss function  $\mathcal{L}$
- Optimization algorithm
- Negatives generation strategy

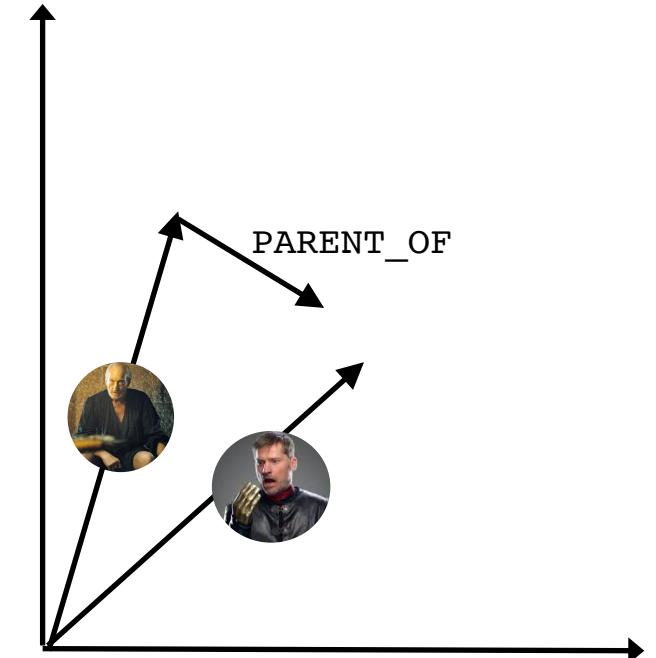
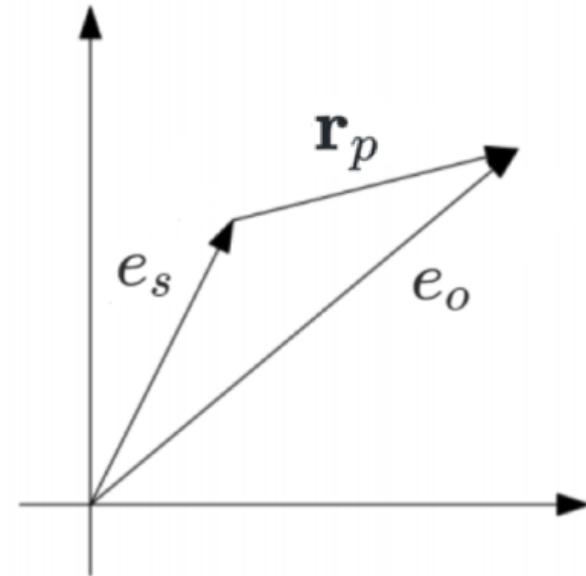
# Scoring function $f$

$f$  assigns a score to a triple  $(s, p, o)$ :

- **TransE:** Translating Embeddings [Bordes13]

$$f_{TransE} = -\|(\mathbf{e}_s + \mathbf{r}_p) - \mathbf{e}_o\|_n$$

( , PARENT\_OF,  )



# Loss function $\mathcal{L}$

- Pairwise margin-based:

$$\mathcal{L}(\Theta) = \sum_{t^+ \in \mathcal{G}} \sum_{t^- \in \mathcal{C}} \max(0, [\gamma + f(t^-; \Theta) - f(t^+; \Theta)])$$

- Absolute Margin
- Negative log-likelihood
- Self-adversarial
- ...

# Optimization Procedure

- Estimate the optimal parameters  $\Theta$  with (same variant of) SGD:

$$\min_{\Theta} \mathcal{L}(\Theta)$$

- On-the-fly negatives generation with the negatives strategy

# Negatives generation strategies

- **Uniform sampling:** Generate and uniformly sample corruptions  $\mathcal{C}$ :

$$\mathcal{C}((s, p, o)) = \{(\hat{s}, p, o) | \hat{s} \in \mathcal{E}\} \cup \{(s, p, \hat{o}) | \hat{o} \in \mathcal{E}\}$$

- Based on **generative adversarial network** [Cai17]
- **Self-adversarial negative sampling** [Sun19]

## Negatives Generation with Uniform Sampling

$$\mathcal{E} = \{Tyrion, Jamie, Tywin, House\_Lannister\}$$

$$\mathcal{R} = \{PARENT\_OF\}$$

$$t_i \in \mathcal{G} = (\text{Tywin PARENT_OF Tyrion})$$

$$\mathcal{C}_{t_i} =$$

Tywin	PARENT_OF	House_Lannister
Tywin	PARENT_OF	Jamie
Jamie	PARENT_OF	Tyrion
House_Lannister	PARENT_OF	Tyrion

# Evaluation Metrics

## Learning to Rank metrics

*How well are positive triples ranked against their corruptions?*

$$Hits@N = \frac{1}{|Q|} \sum_{i=1}^{|Q|} 1 \text{ if } rank_{(s,p,o)_i} \leq N$$

$$MR = \frac{1}{|Q|} \sum_{i=1}^{|Q|} rank_{(s,p,o)_i} \text{ [Mean Rank]}$$

$$MRR = \frac{1}{|Q|} \sum_{i=1}^{|Q|} \frac{1}{rank_{(s,p,o)_i}} \text{ [Mean Reciprocal Rank]}$$

## Example

s	p	o	score	rank	*
Tywin	PARENT_OF	House_lannister	0.789	1	
Tywin	PARENT_OF	Tyrion	0.753	2	*
Tywin	PARENT_OF	House_Lannister	0.695	3	
Sansa	PARENT_OF	Tyrion	0.695	3	
Cersei	PARENT_OF	Tyrion	0.234	5	

s	p	o	score	rank	*
Tywin	PARENT_OF	Jamie	0.852	1	*
Sansa	PARENT_OF	Jamie	0.345	2	
Cersei	PARENT_OF	Jamie	0.293	3	
Tywin	PARENT_OF	Sansa	0.201	4	
Tywin	PARENT_OF	Cersei	0.156	5	

$$Hits@3 = \frac{1}{2} \sum_{i=1}^2 1 \text{ if } rank_{(s,p,o)_i} \leq 3 = 1$$

$$Hits@1 = \frac{1}{2} \sum_{i=1}^2 1 \text{ if } rank_{(s,p,o)_i} \leq 1 = 0.5$$

Import a dataset and split into train/test:

```
X = load_from_csv("datasets/", "GoT_graph.csv")
X_train, X_test = train_test_split_no_unseen(X, test_size=10)
```

Initialize a graph embedding model:

```
model = ComplEx(k=150, epochs=500, eta=5, loss='nll')
```

Generate the embeddings:

```
model.fit(X_train)
```

Visualize the embeddings:

```
create_tensorboard_projector_files(model, 'GoT_embeddings.tsv')
```



Import a dataset and split into train/test:

```
X = load_from_csv("datasets/", "GoT_graph.csv")
X_train, X_test = train_test_split_no_unseen(X, test_size=10)
```

Initialize a graph embedding model:

```
model = ComplEx(k=150, epochs=500, eta=5, loss='nll')
```

Generate the embeddings:

```
model.fit(X_train)
```

Visualize the embeddings:

```
create_tensorboard_projector_files(model, 'GoT_embeddings.tsv')
```

Assess predicting power:

```
ranks = evaluate_performance(X_test, model=model)
mrr_score(ranks)
hits_at_n_score(ranks, n=10)
```

MRR: 0.393462

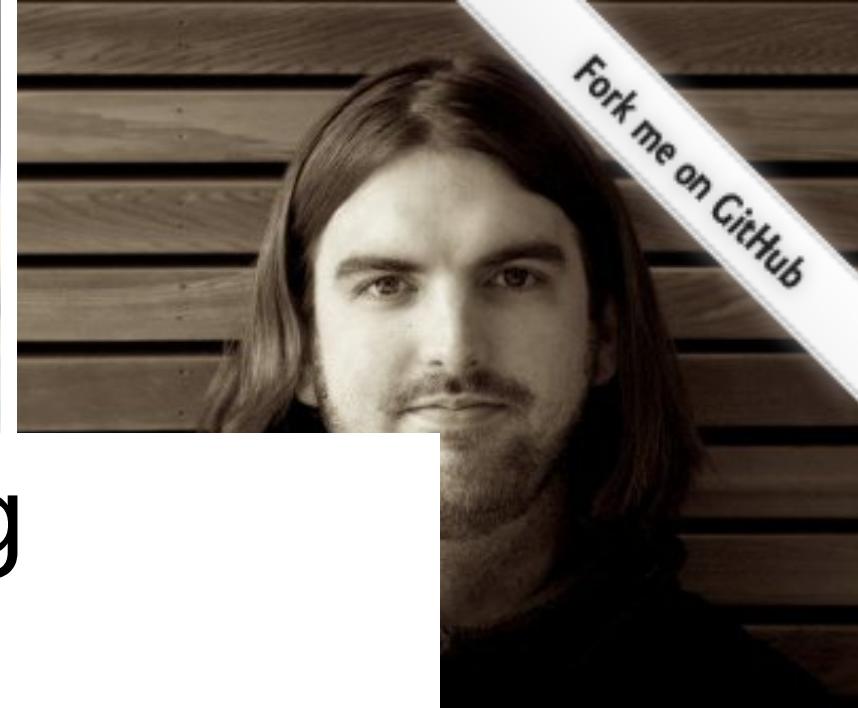
Hits@10: 0.495000

```
evaluate_performance(np.array([[ 'Daenerys Targaryen', 'SPOUSE', 'Jon Snow']]), model=model)
```

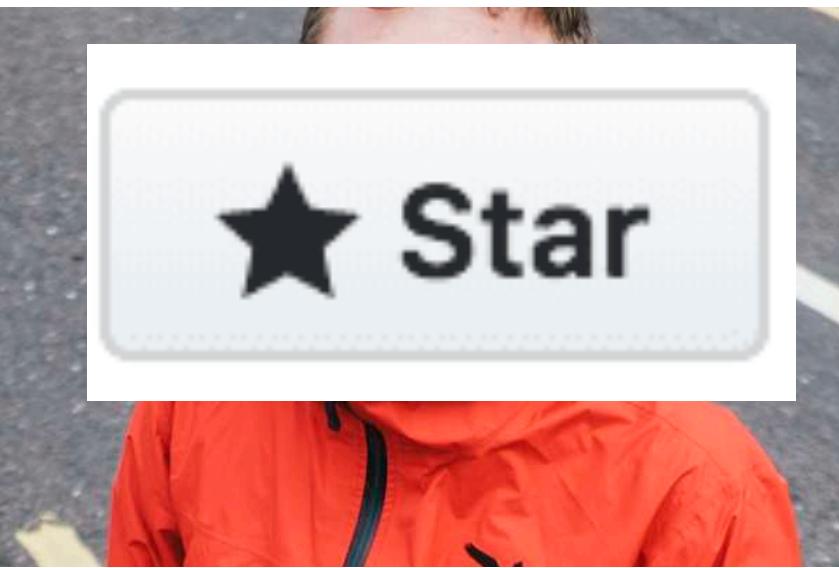
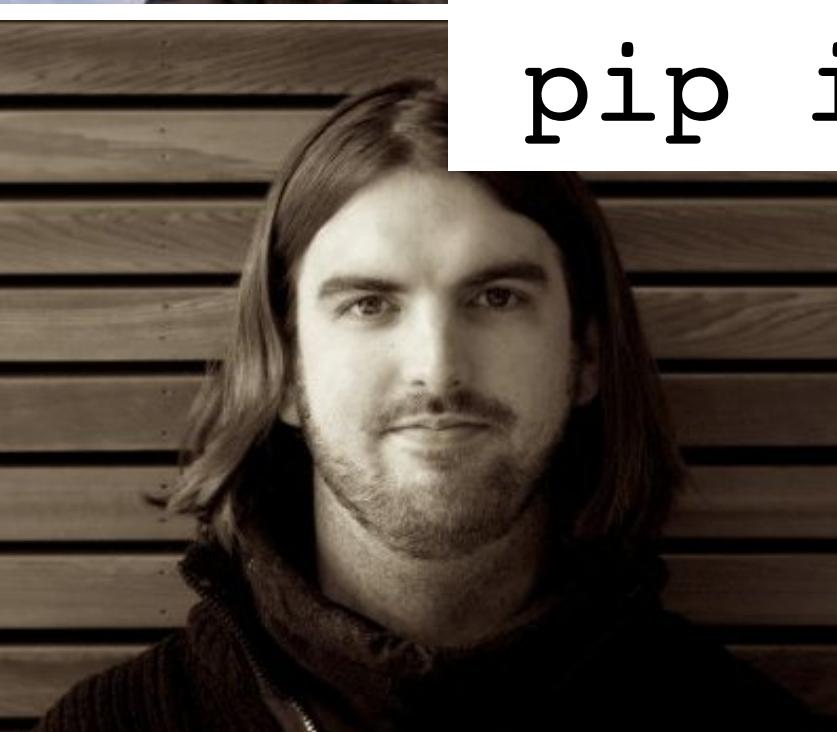
1844/6350

```
evaluate_performance(np.array([["King's Landing", 'SEAT_OF', 'Brandon Stark']]), model=model)
```

385/6350



# ampligraph.org



pip install ampligraph