# Stealing Links from Graph Neural Networks
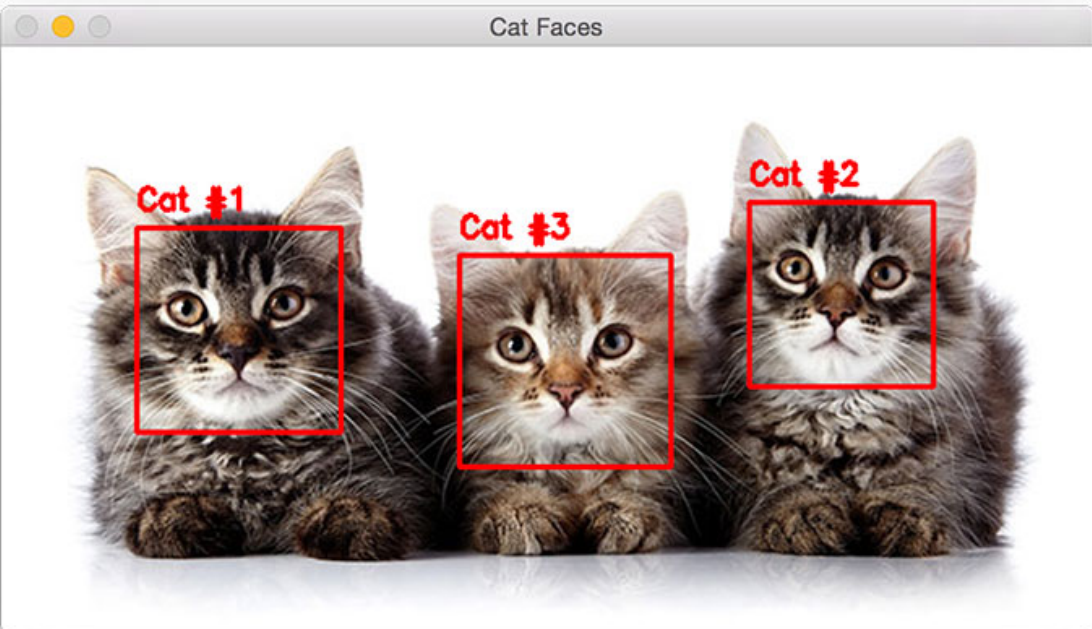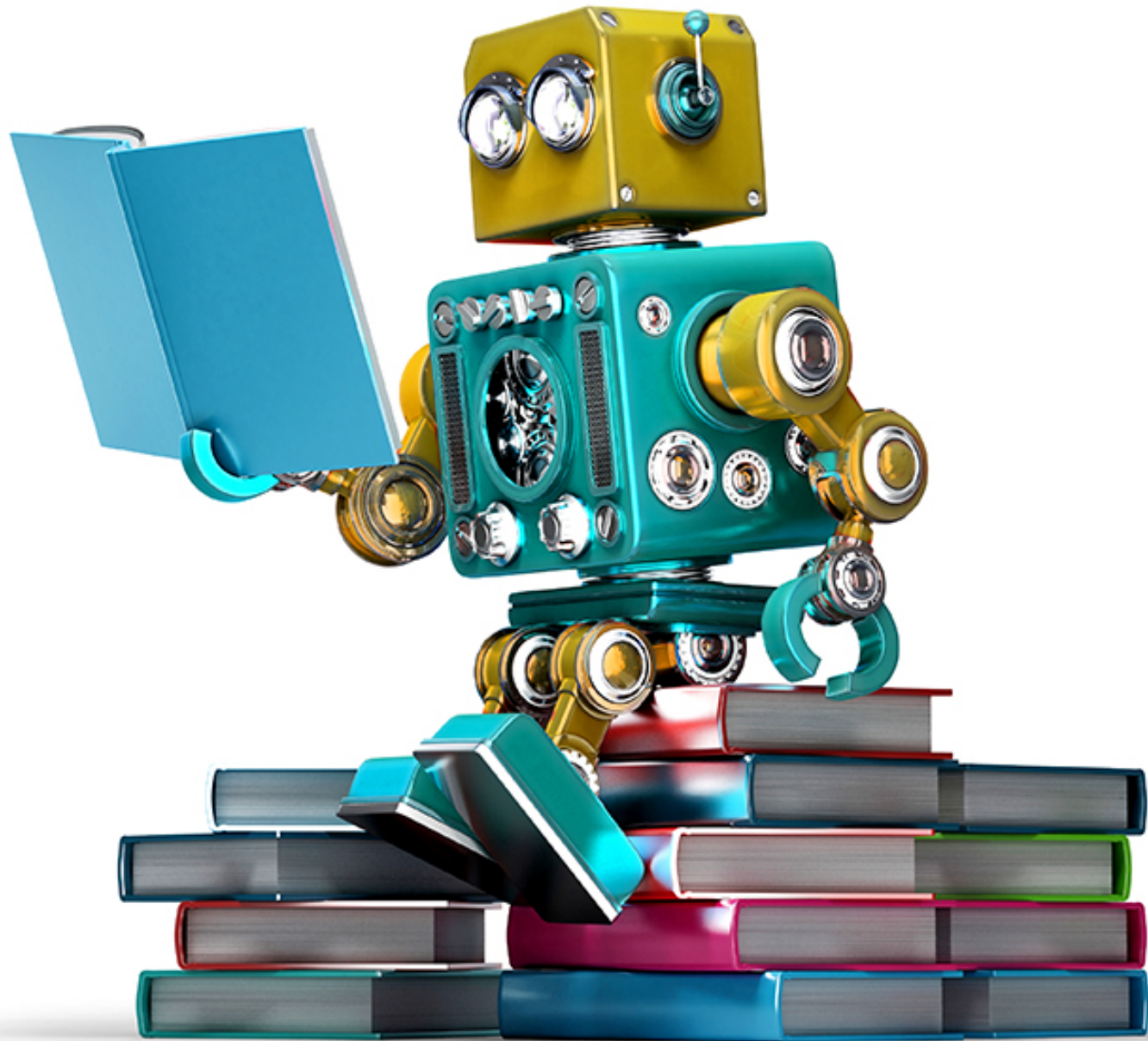
**Xinlei He**[1], Jinyuan Jia[2], Michael Backes[1],
Neil Zhenqiang Gong[2], Yang Zhang[1]

[1]CISPA Helmholtz Center for Information Security
[2]Duke University

# Era of Machine Learning
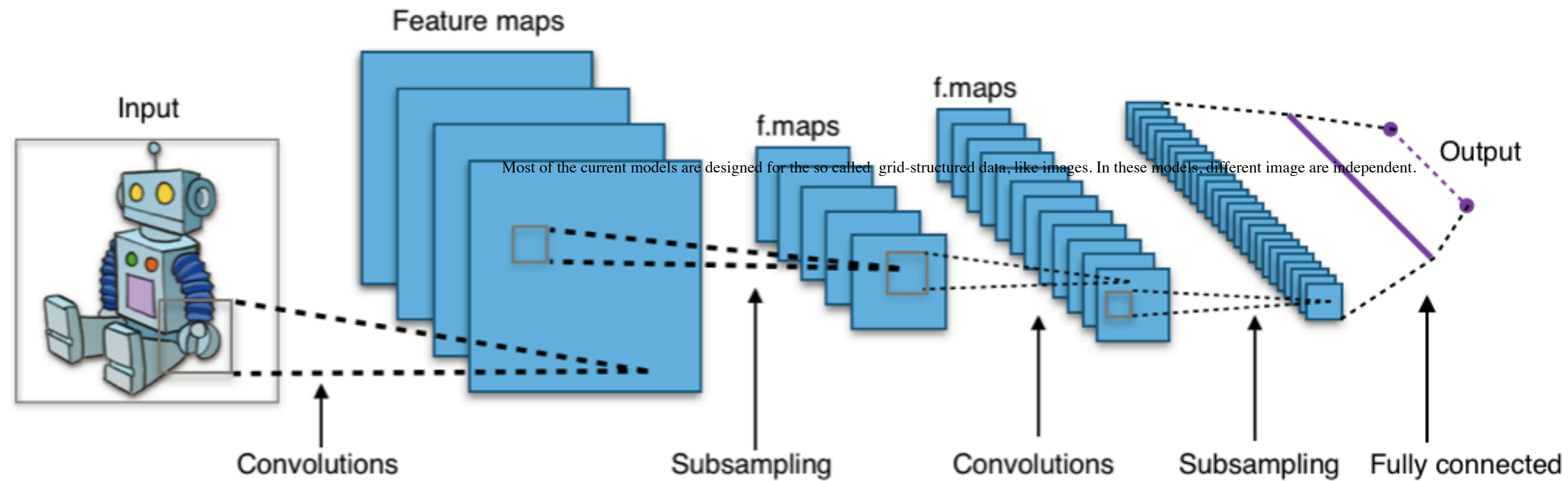


Data

# Machine Learning Pipeline

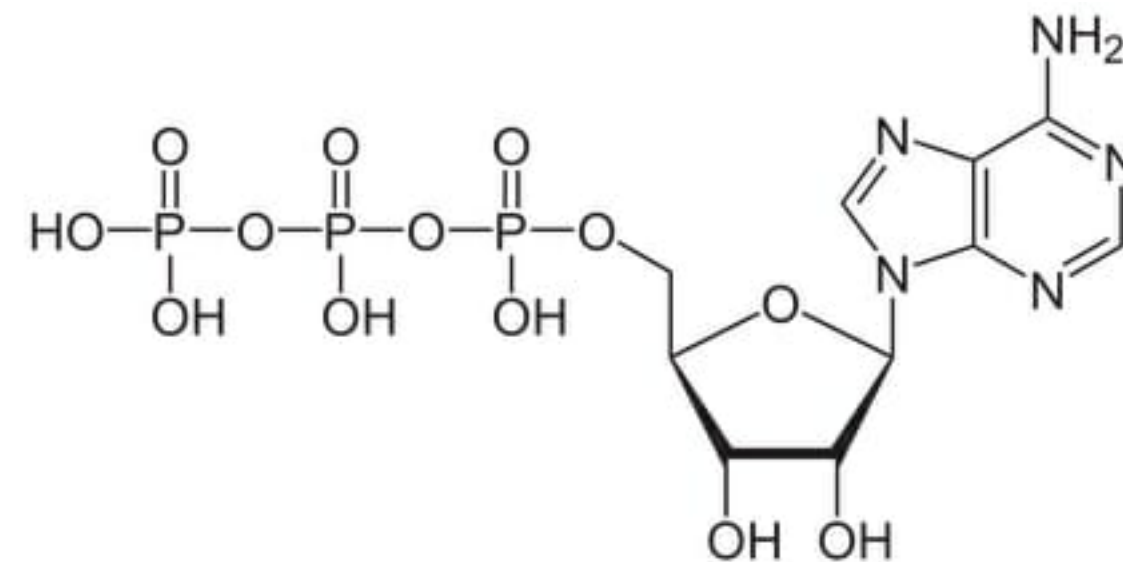Modern machine learning excels at exploiting grid-structured data

# Many Data are Graphs

Graphs are combinatorial structures, have arbitrary sizes, and contain multi-modal information
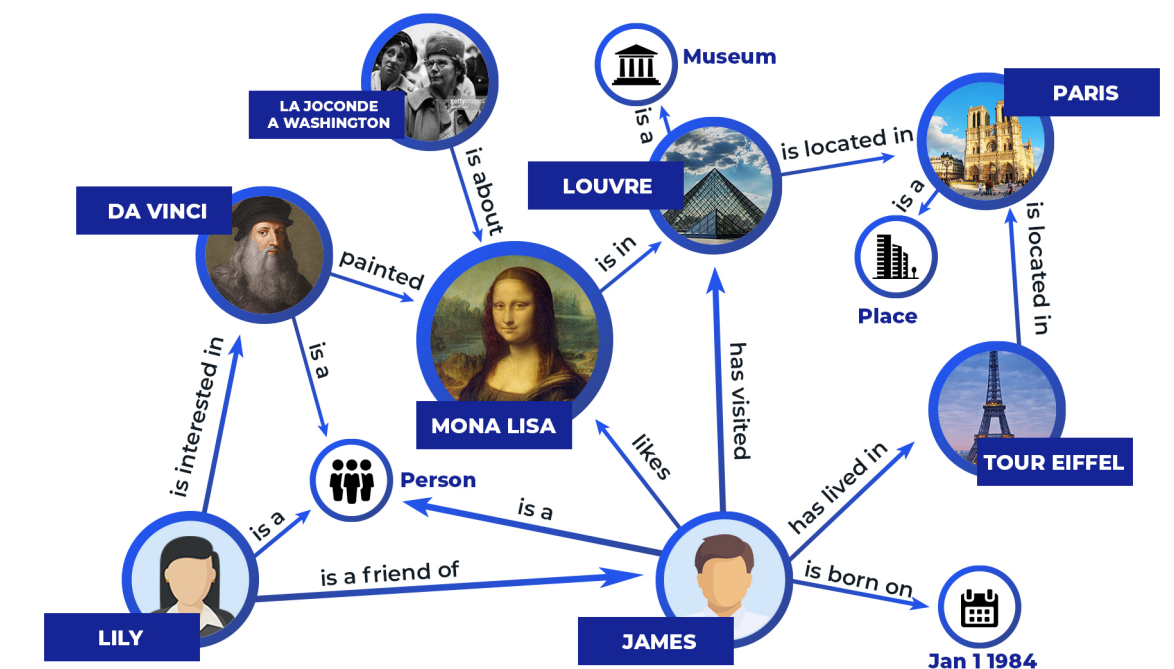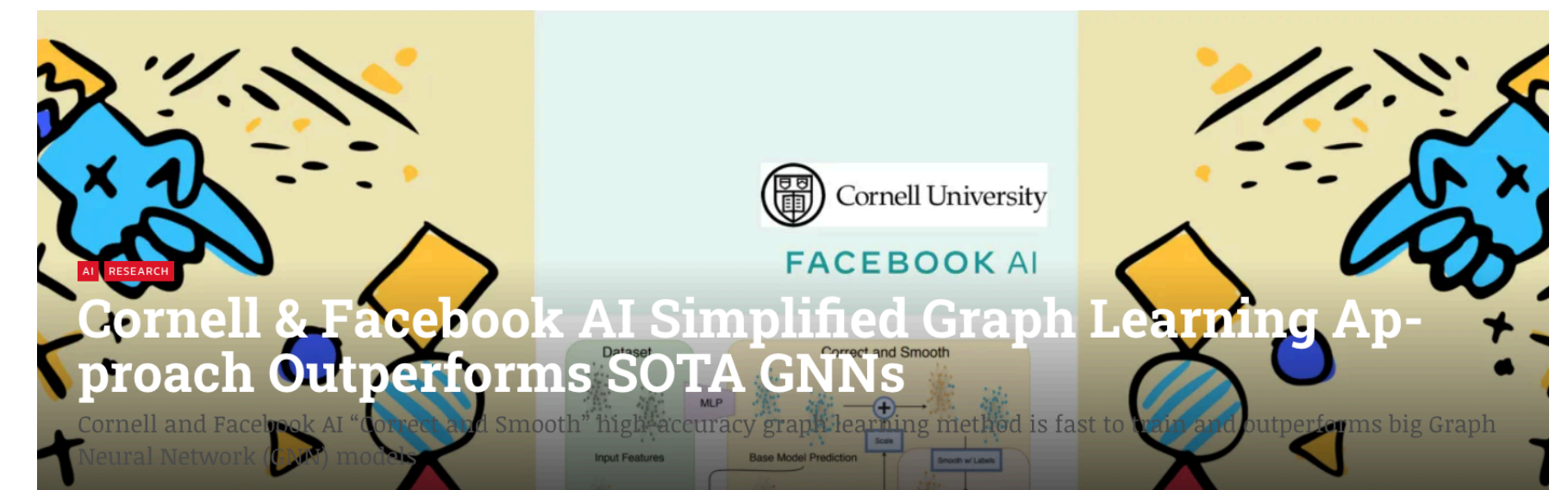
Social Networks

Molecules

Knowledge Graphs

# Graph Neural Networks

# Graph Neural Networks (Transductive)

[12, 32, 6, 0.3]

[61, 2, 13, 7.2]

[22, 78, 5, 9.1]

[14, 10, 9, 1.2]

[15, 32, 9, 4.1]

panda

cat

?

?

dog

GNN

?

?

GNN

70
35
0

panda  dog  cat

70
35
0

panda  dog  cat

**Research question: Given two nodes used to train a black-box GNN, can we predict whether they are linked?**

# Attack Taxonomy

[12, 32, 6, 0.3]

[61, 2, 13, 7.2]

[22, 78, 5, 9.1]

[14, 10, 9, 1.2]

[15, 32, 9, 4.1]

GNN

- Attacker can have either of these 3 knowledge

- Totally 8 different attack models

**Node Features**

[12, 32, 6, 0.3]

[14, 10, 9, 1.2]

[22, 78, 5, 9.1]

[15, 32, 9, 4.1]

[61, 2, 13, 7.2]

**Partial Graph**

**Shadow Dataset**

[14, 6, 9]

[10, 5, 8]

[5, 3, 12]

[8, 5, 13]

[12, 7, 8]

# Attack 0

# Attack 0

**Correlation performs the best!**



Figure 1: AUC for Attack-0 on all the 8 datasets with all the 8 distance metrics. The x-axis represents the dataset and the y-axis represents the AUC score.

Table 15: Prediction results for Attack-0 on all the 8 datasets with Correlation distance.

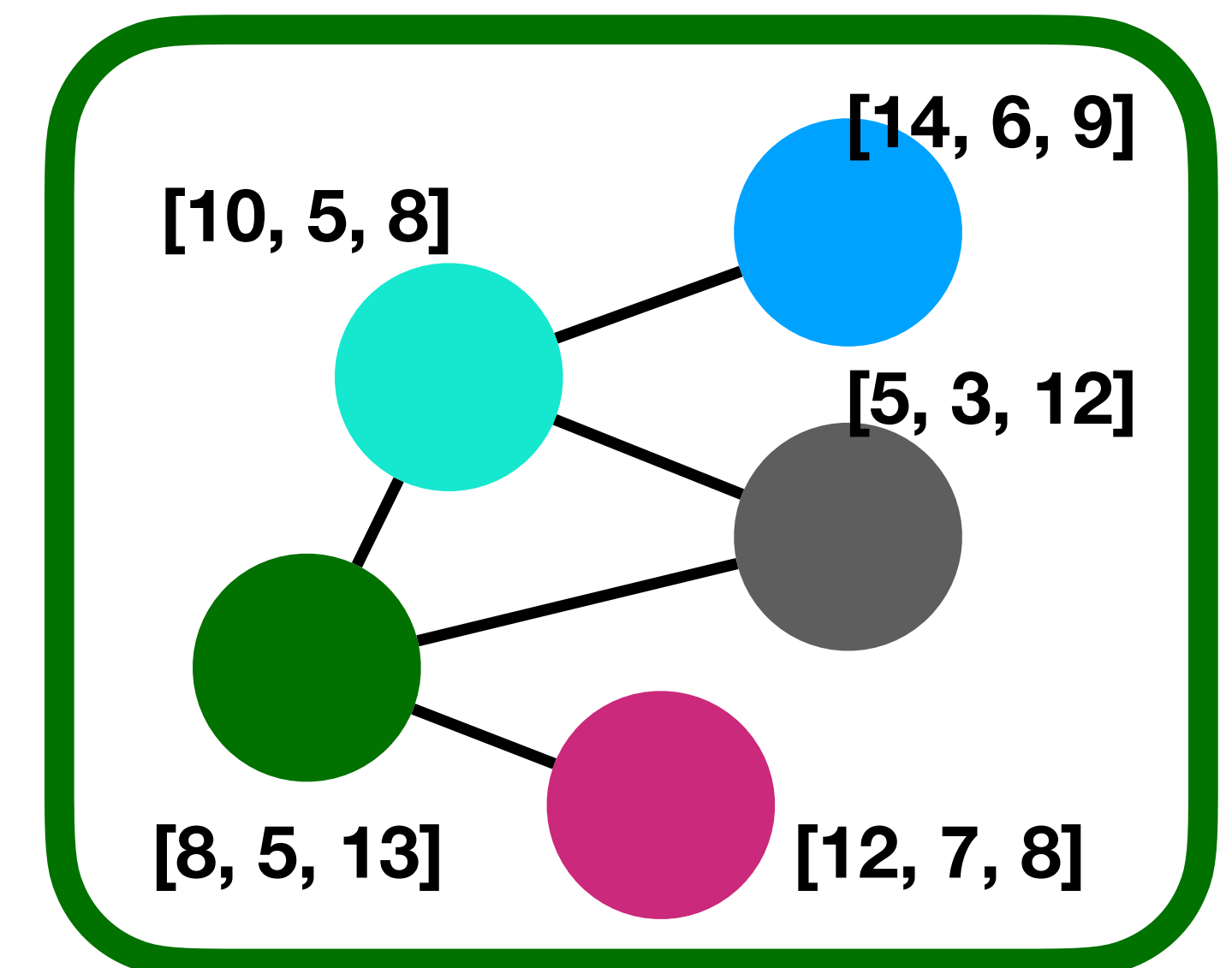| Dataset | Precision | Recall | F1-Score | AUC |
|---|---|---|---|---|
| AIDS | 0.524 | 0.996 | 0.687 | 0.691 |
| COX2 | 0.523 | 0.987 | 0.684 | 0.867 |
| DHFR | 0.555 | 0.977 | 0.708 | 0.765 |
| ENZYMES | 0.501 | 1.000 | 0.667 | 0.630 |
| PROTEINS_full | 0.540 | 0.998 | 0.701 | 0.815 |
| Citeseer | 0.788 | 0.991 | 0.878 | 0.959 |
| Cora | 0.777 | 0.966 | 0.861 | 0.929 |
| Pubmed | 0.691 | 0.965 | 0.806 | 0.874 |

**Use KMeans to give a concrete prediction**

# Attack 1

# Attack 1



**Shadow**

50
25
0
Cook  Actor  Barber  Coach

70
35
0
Cook  Actor  Barber  Coach

**Target**

GNN

70
35
0
panda  dog  cat

70
35
0
panda  dog  cat

**Distance**

**Dimension mismat**

**Entropy**

| Operator | Definition | Operator | Definition |
|---|---|---|---|
| Average | $\dfrac{f_i(u) + f_i(v)}{2}$ | Weighted-L1 | $\lvert f_i(u) - f_i(v) \rvert$ |
| Hadamard | $f_i(u) \cdot f_i(v)$ | Weighted-L2 | $\lvert f_i(u) - f_i(v) \rvert^2$ |

Aditya Grover and Jure Leskovec. node2vec: Scalable Feature Learning for Networks. In KDD 2016.

# Attack 1

**Table 4: Average AUC with standard deviation for Attack-1 on all the 8 datasets. Best results are highlighted in bold.**

| Target Dataset | AIDS | COX2 | DHFR | ENZYMES | PROTEINS_full | Citeseer | Cora | Pubmed |
|---|---|---|---|---|---|---|---|---|
| | | | | | Shadow Dataset | | | |
| AIDS | - | $0.720 \pm 0.009$ | $0.690 \pm 0.005$ | $\mathbf{0.730 \pm 0.010}$ | $0.720 \pm 0.005$ | $0.689 \pm 0.019$ | $0.650 \pm 0.025$ | $0.667 \pm 0.014$ |
| COX2 | $0.755 \pm 0.032$ | - | $0.831 \pm 0.005$ | $0.739 \pm 0.116$ | $\mathbf{0.832 \pm 0.009}$ | $0.762 \pm 0.009$ | $0.773 \pm 0.008$ | $0.722 \pm 0.024$ |
| DHFR | $0.689 \pm 0.004$ | $\mathbf{0.771 \pm 0.004}$ | - | $0.577 \pm 0.044$ | $0.701 \pm 0.010$ | $0.736 \pm 0.005$ | $0.740 \pm 0.003$ | $0.663 \pm 0.010$ |
| ENZYMES | $\mathbf{0.747 \pm 0.014}$ | $0.695 \pm 0.023$ | $0.514 \pm 0.041$ | - | $0.691 \pm 0.030$ | $0.680 \pm 0.012$ | $0.663 \pm 0.009$ | $0.637 \pm 0.018$ |
| PROTEINS_full | $0.775 \pm 0.020$ | $0.821 \pm 0.016$ | $0.528 \pm 0.038$ | $0.822 \pm 0.020$ | - | $\mathbf{0.823 \pm 0.004}$ | $0.809 \pm 0.015$ | $0.809 \pm 0.013$ |
| Citeseer | $0.801 \pm 0.040$ | $0.920 \pm 0.006$ | $0.842 \pm 0.036$ | $0.846 \pm 0.042$ | $0.848 \pm 0.015$ | - | $\mathbf{0.965 \pm 0.001}$ | $0.942 \pm 0.003$ |
| Cora | $0.791 \pm 0.019$ | $0.884 \pm 0.005$ | $0.811 \pm 0.024$ | $0.804 \pm 0.048$ | $0.869 \pm 0.012$ | $\mathbf{0.942 \pm 0.001}$ | - | $0.917 \pm 0.002$ |
| Pubmed | $0.705 \pm 0.039$ | $0.796 \pm 0.007$ | $0.704 \pm 0.042$ | $0.708 \pm 0.067$ | $0.752 \pm 0.014$ | $0.883 \pm 0.006$ | $\mathbf{0.885 \pm 0.005}$ | - |

**For all best performing shadow datasets, attack 1 is better than attack 0**
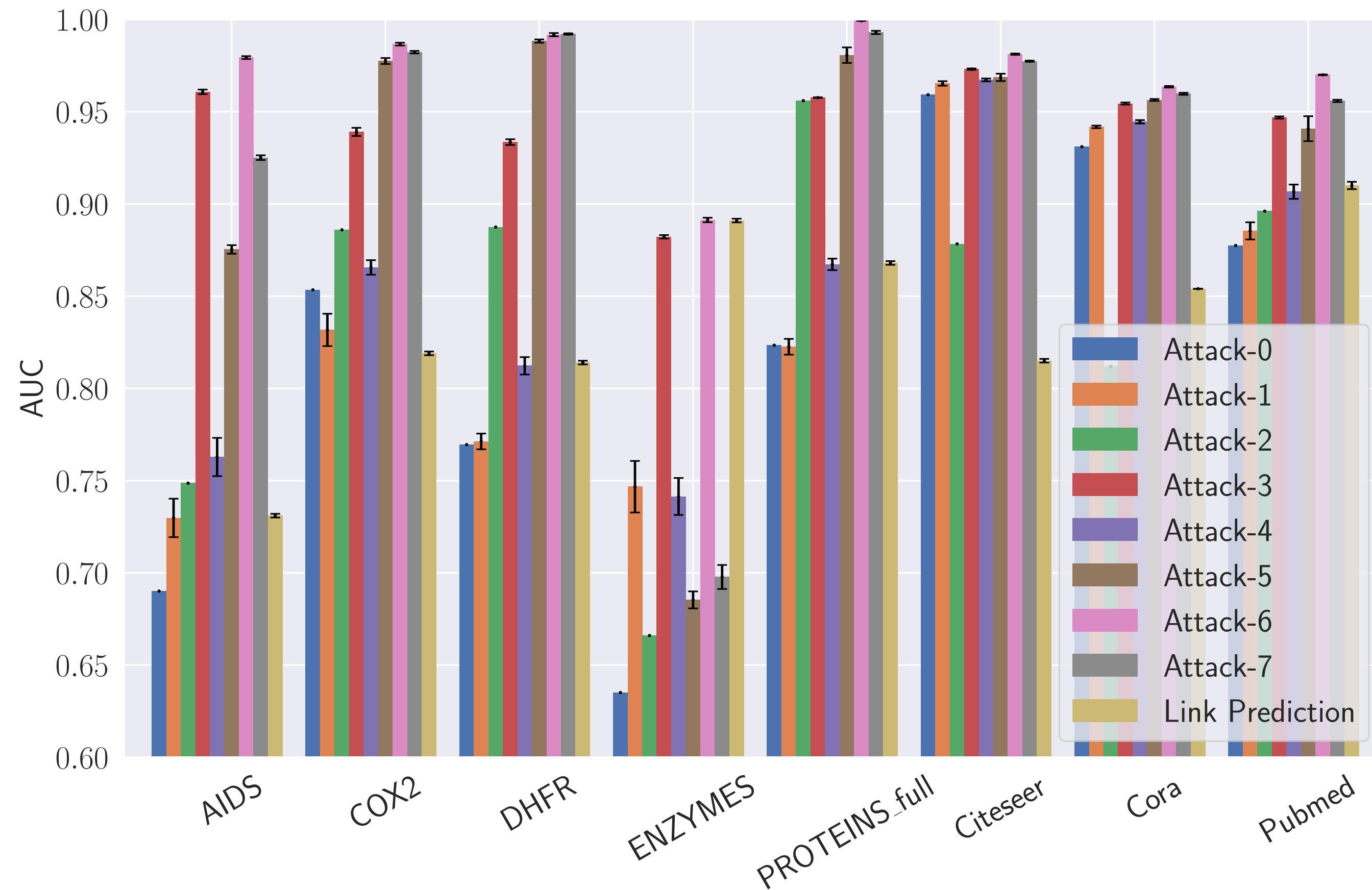
13

# Attack 1

Figure 3: The last hidden layer's output from the attack model of Attack-1 for 200 randomly sampled positive node pairs and 200 randomly sampled negative node pairs projected into a 2-dimension space using t-SNE. (a) Cora as the shadow dataset and Citeseer as the target dataset, (b) Cora as the shadow dataset and ENZYMES as the target dataset.

# Evaluation of All Attacks



- More knowledge leads to better attack performance

- Partial graph contains the strongest signal

- Shadow dataset is the weakest

- Better performance than traditional link prediction, this means GNN indeed leaks graph information

# Conclusion

➡️ **We are the first to propose link stealing attack against GNNs**

➡️ **Our attacks can effectively ste** Questions?

➡️ **More information leads to bett**

➡️ **Transferring attack can achieve good performance**

**Code is available at https://github.com/xinleihe/link_stealing_attack**

**Xinlei He**
**CISPA Helmholtz Center for Information Security**
**@AllenXinleiHe**
**http://www.xinlei.info/**

# Thanks!