# Link Prediction in Networks with Core-Fringe Data
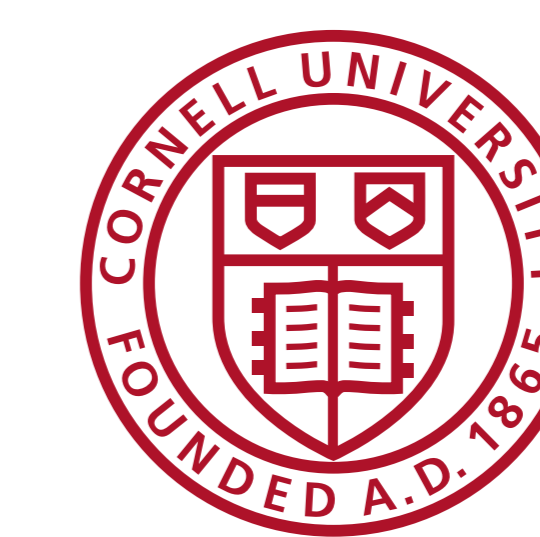
## Austin R. Benson and Jon Kleinberg

arb@cs.cornell.edu, kleinber@cs.cornell.edu

Code & data → https://github.com/arbenson/cflp
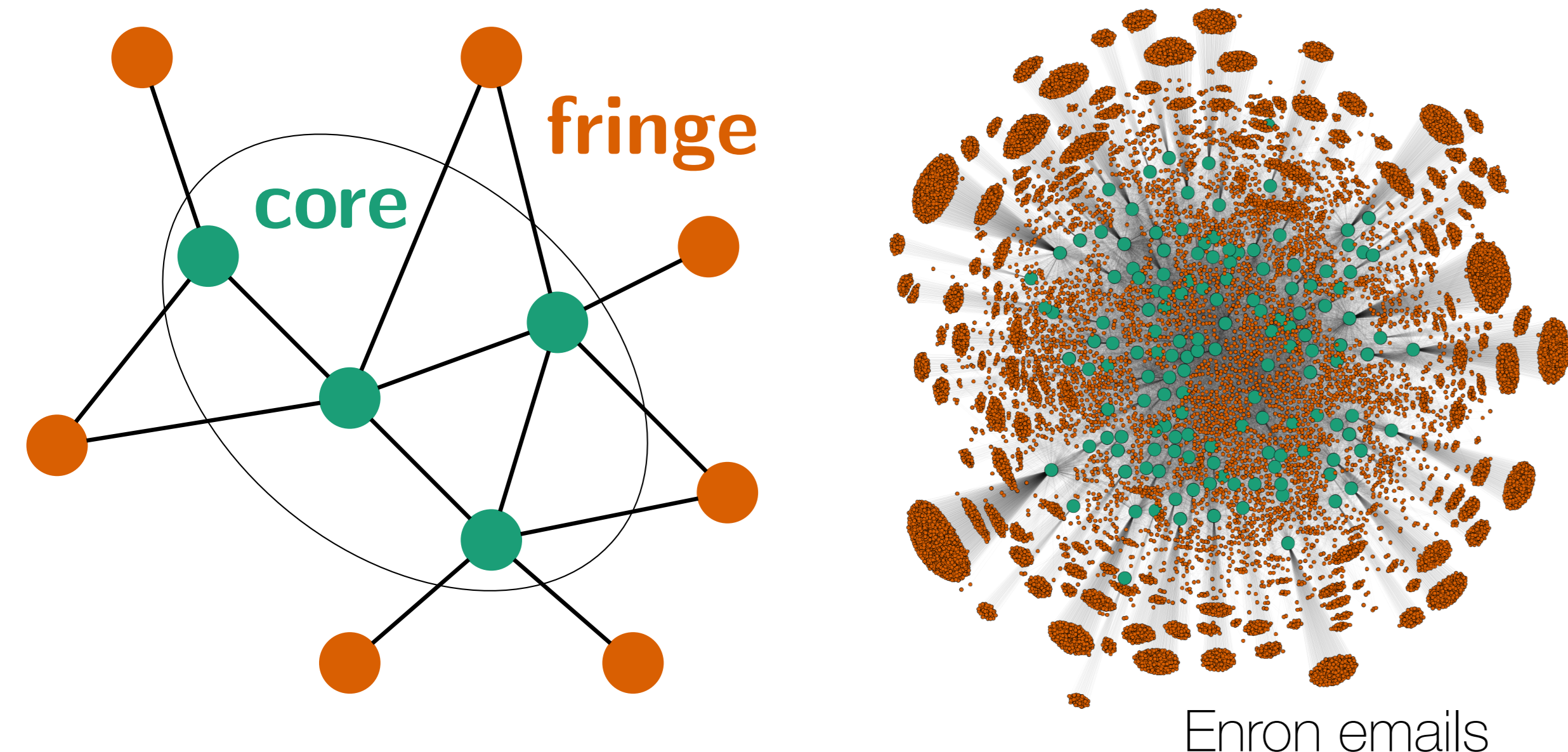
Cornell University

## Partially measured graphs & core-fringe data.

We often measure graph data by recording interactions involving a *core* set of nodes:
- Email of company employees
- Phone calls of all customers of a service provider

We end up with a dataset that includes the core along with a potentially much larger set of *fringe* nodes.



fringe / core



Enron emails

## Does knowing the fringe help with link prediction in the core?
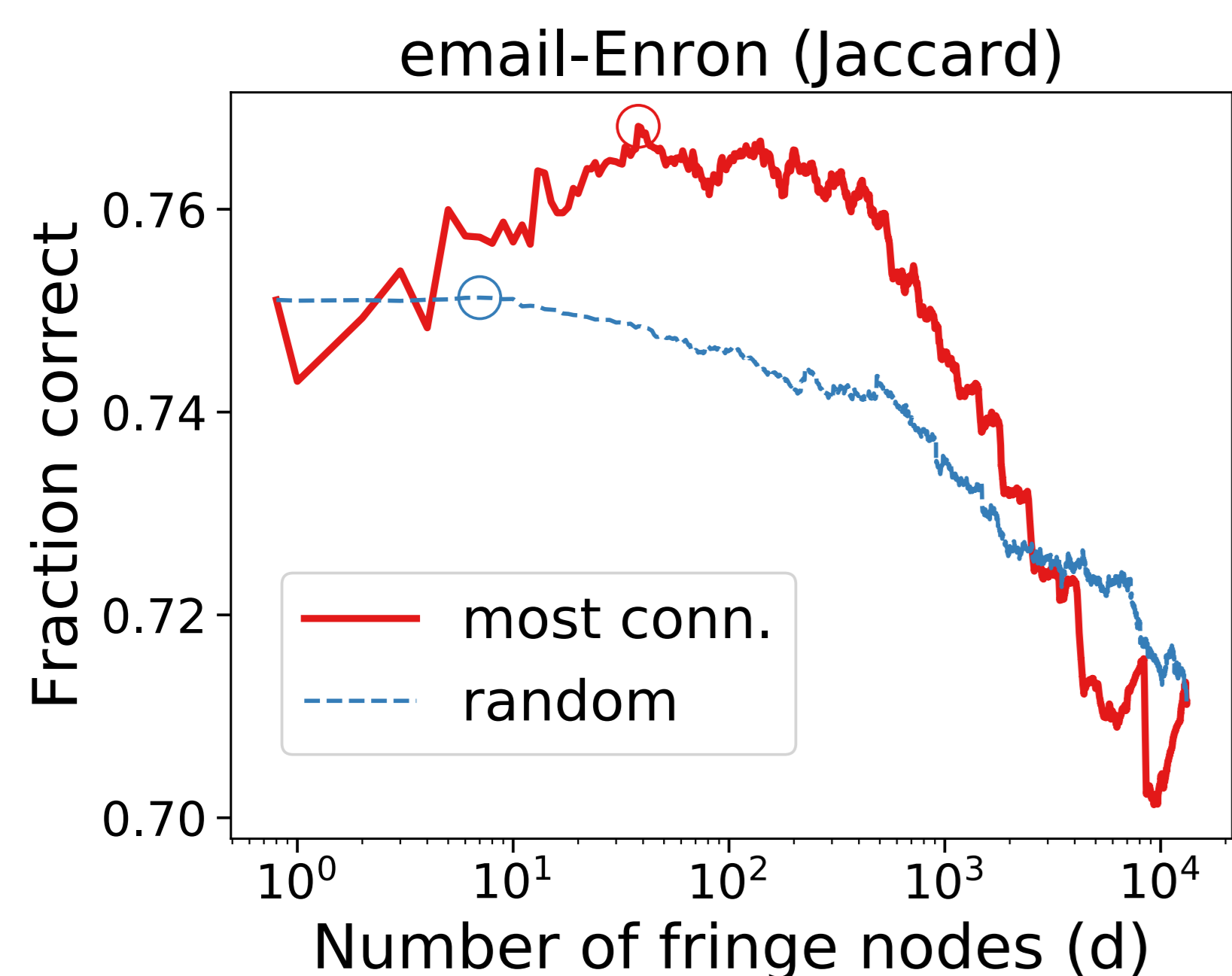
1. Fix a link prediction algorithm

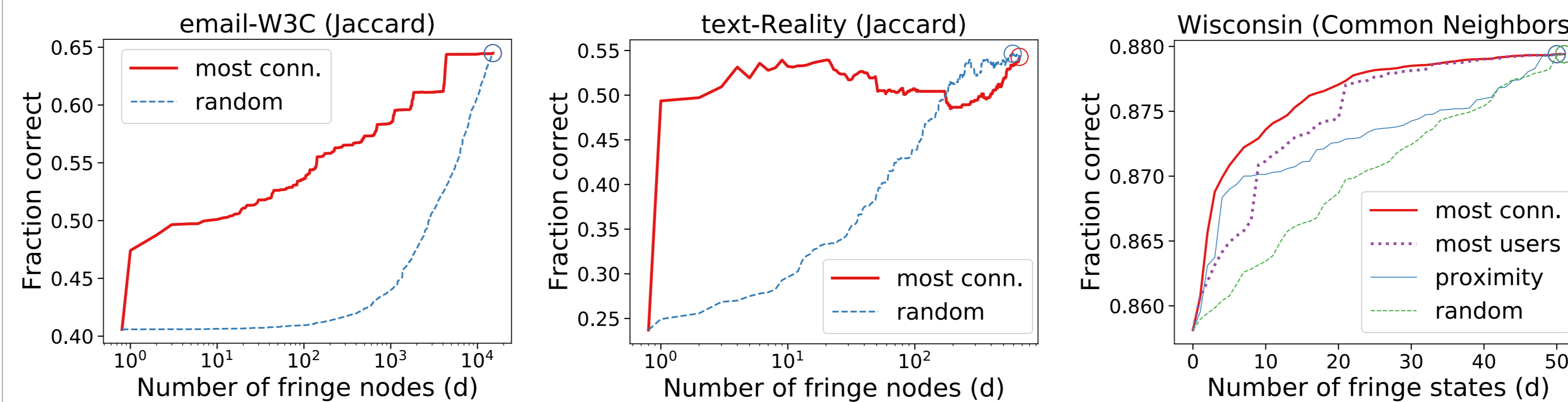   $\text{CommonNeighbors}(u, v) = |N(u) \cap N(v)|$

   $\text{Jaccard}(u, v) = |N(u) \cap N(v)| / |N(u) \cup N(v)|$

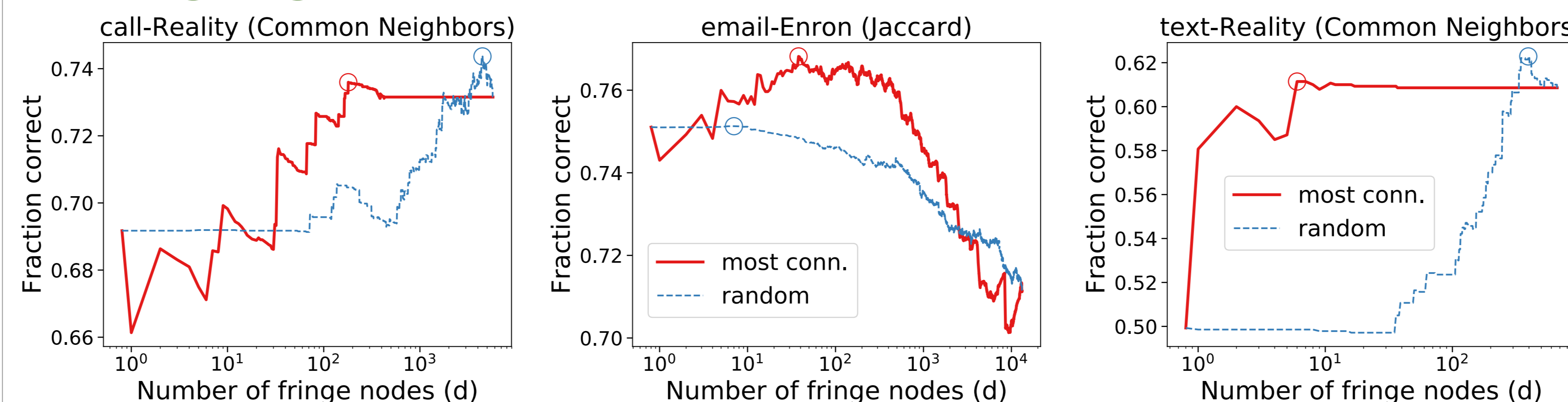2. Include fringe nodes and connections in some order, *possibly changing the algorithm predictions*.

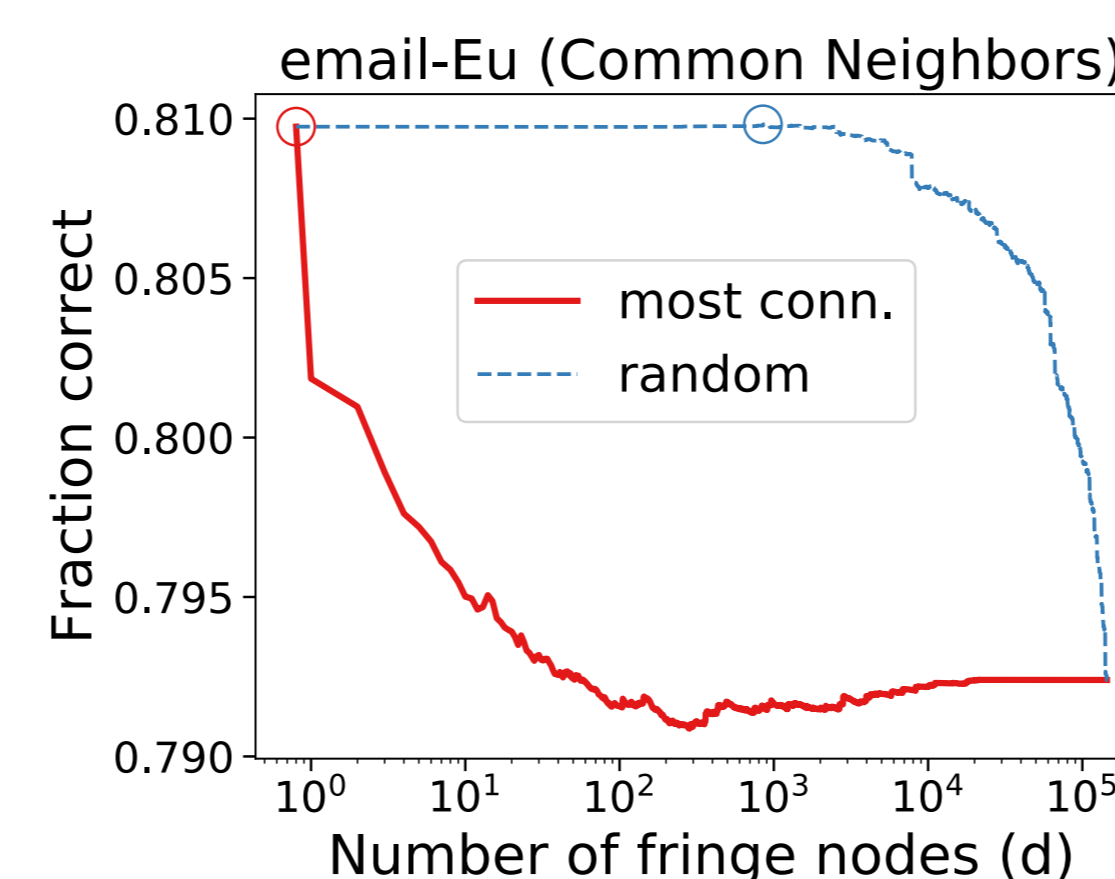3. Measure link prediction accuracy as a function of the number of fringe nodes included.

email-Enron (Jaccard)



most conn. / random

Fraction correct vs. Number of fringe nodes (d)

## Sometimes, we want all of the information from the fringe.

email-W3C (Jaccard)



most conn. / random

text-Reality (Jaccard)



most conn. / random

Wisconsin (Common Neighbors)



most conn. / most users / proximity / random

Fraction correct vs. Number of fringe nodes (d) / Number of fringe states (d)

## Sometimes, an intermediate amount of fringe gives the best performance.

call-Reality (Common Neighbors)



email-Enron (Jaccard)



most conn. / random

text-Reality (Common Neighbors)



most conn. / random

Fraction correct vs. Number of fringe nodes (d)

## Sometimes, any fringe information hurts.

email-Eu (Common Neighbors)



most conn. / random

Fraction correct vs. Number of fringe nodes (d)

## Sometimes, performance saturates with more fringe information.

email-Avocado (Common Neighbors)



most conn. / random

Dane (Common Neighbors)



most conn. / most users / proximity / random

Dane (Jaccard)



most conn. / most users / proximity / random

Fraction correct vs. Number of fringe nodes (d) / Number of fringe counties (d)

## Random graph models can explain the diversity of behaviors.

- Fix algorithm as number of common neighbors.
- Random graph model where edge {u, v} is more likely than edge {w, z} (latent random variables).
- Algorithm can use fringe info., parameterized by $d$.

   $X_d = \text{CommonNeighbors}(u, v)$

   $Y_d = \text{CommonNeighbors}(w, z)$

- Goal is to maximize $\text{Prob}(X_d > Y_d)$.
- Solution is

$$\max_d \text{SNR}(Z_d) = \frac{\mathbb{E}(Z_d)}{\sqrt{\mathbb{V}(Z_d)}}, \; Z_d = X_d - Y_d$$

## Core-fringe stochastic block model

$$\begin{array}{c} \\ 1 \\ 2 \\ 3 \\ 4 \end{array} \begin{bmatrix} p & q & r & s \\ q & p & s & r \\ r & s & 0 & 0 \\ s & r & 0 & 0 \end{bmatrix}$$
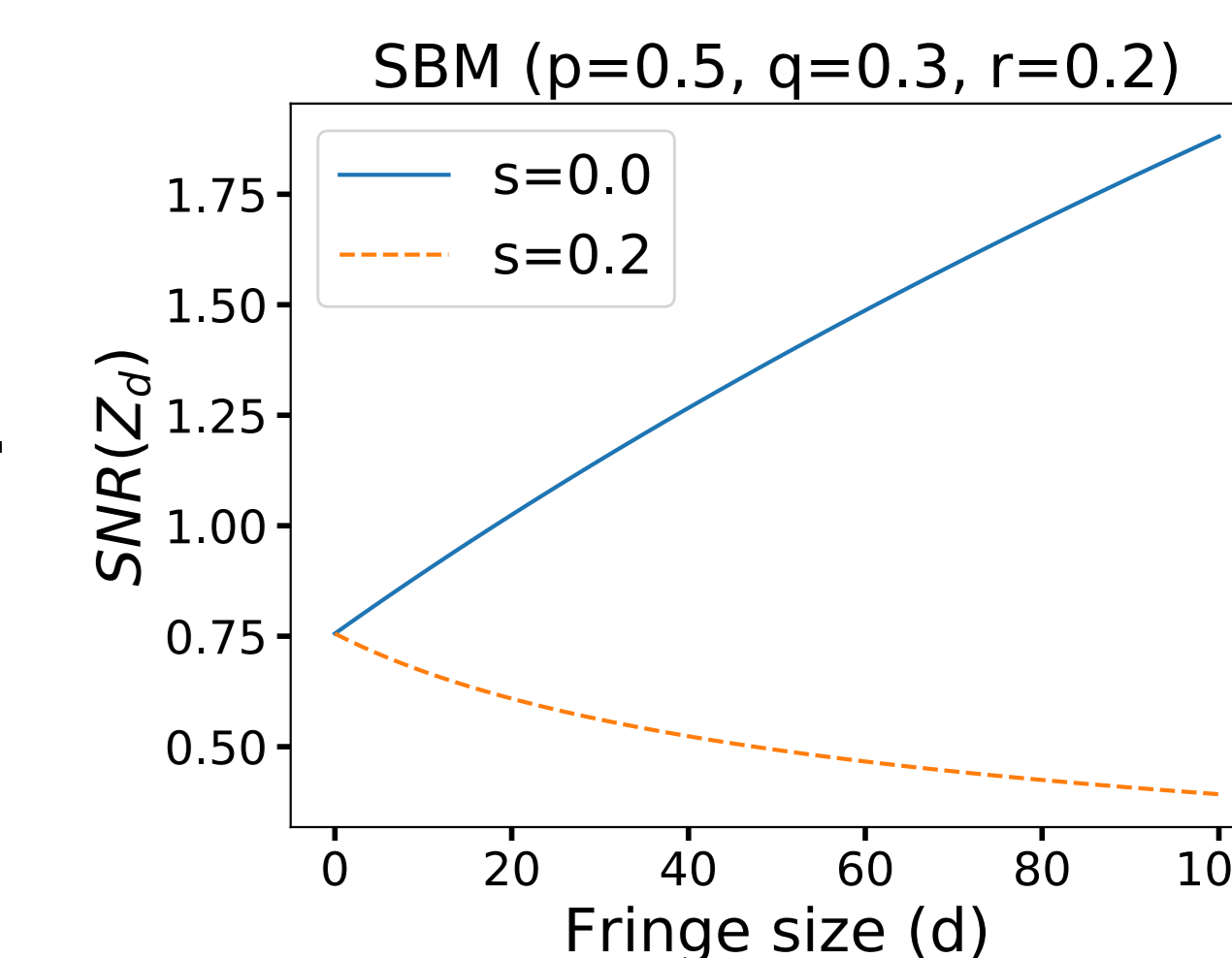
- $p > q$; core is blocks 1 & 2
- $u, v, w$ in block 1; $z$ in block 2
- $d$ is number of fringe nodes included from blocks 3 and 4.

SBM (p=0.5, q=0.3, r=0.2)



s=0.0 / s=0.2

$SNR(Z_d)$ vs. Fringe size (d)

**Lemma (no-fringe optimality).**
If $r = s$, then $\text{SNR}(Z_d)$ decreases monotonically in $d$.

**Lemma (all-fringe optimality).**
If $r > 0$, $s = 0$, then $\text{SNR}(Z_d)$ increases monotonically in $d$.

## Core-fringe small-world model

- 1-D lattice of nodes, $\text{Prob}(\text{edge } (i, j)) = 1 / |i - j|$
- Core is $\{-c, -c + 1, \ldots, c - 1, c\}$
- $d$ includes fringe $\{-(c + d), \ldots, -(c + 1), c + 1, \ldots c + d\}$

Lattice model (c = 10, z-w = 9)



v-u = 2 / v-u = 3 / v-u = 4 / v-u = 5 / v-u = 6

$SNR(Z_d)$ vs. Fringe size (d)

**Theorem (saturation).**
$$\lim_{d \to \infty} \text{SNR}(Z_d) = S^* > 0$$

**Theorem (intermediate-fringe optimality).**
If $\text{SNR}(Z_0) < \text{SNR}(Z_1)$, then $d^* = \arg\max_d \text{SNR}(Z_d)$ satisfies $0 < d^* < \infty$.