# On the Performance of Percolation Graph Matching

Lyudmila Yartseva
Matthias Grossglauser
EPFL


COSN
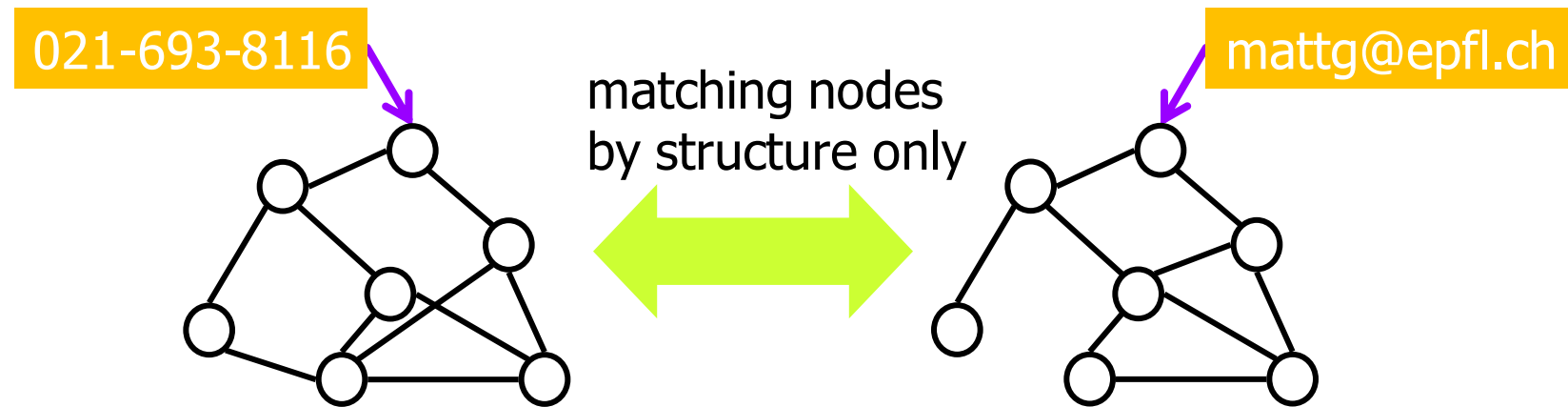Oct 2013

EPFL
ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

# Privacy of Networks

- ## Adversary has:
  - Anonymized network: unlabeled graph
  - Side information: labeled graph – similar but not identical



anonymized social network

side information

matching nodes by structure only

Adam

Barbara

Carlos

# Graph Matching Applications

- ## Social networks:
  - Correlating different domains



021-693-8116    mattg@epfl.ch

matching nodes by structure only

- ## Security:
  - Identifying computer viruses by function-call patterns
- ## Computer vision:
  - Segment adjacency graph to find similar images
- ...

# $G(n, p; s)$ Sampling Model
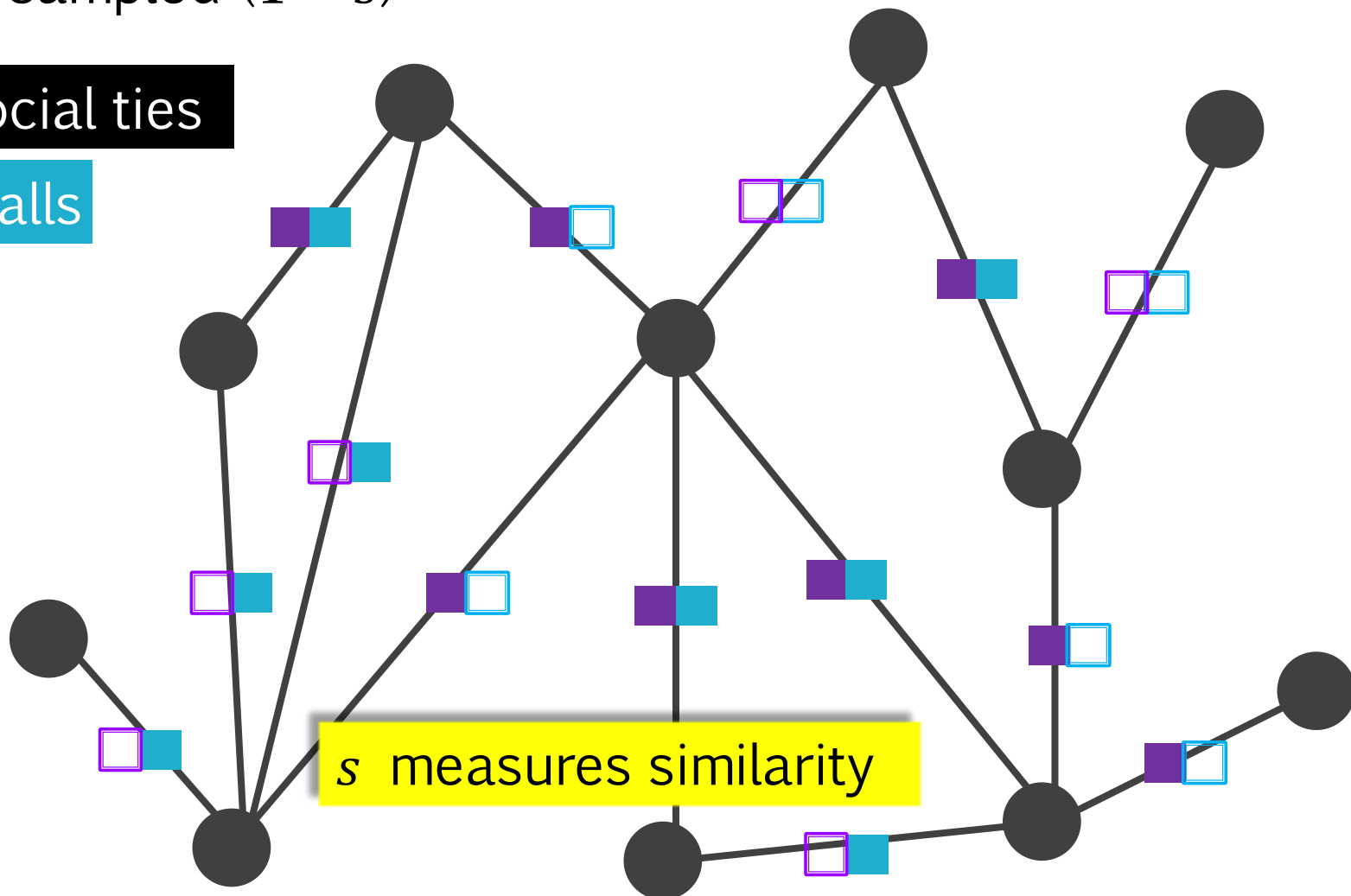
Generator $G = G(n, p)$
- ■ sampled $(s)$
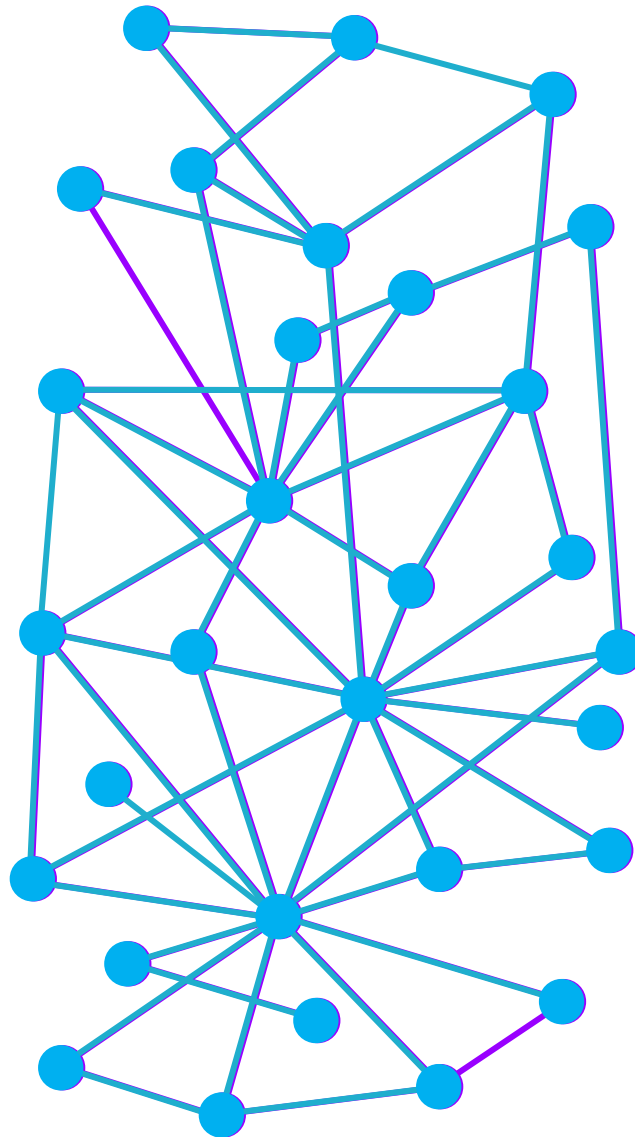- □ not sampled $(1 - s)$

"real" social ties

phone calls

emails

$s$ measures similarity

# Result 1: GM is Easy with ∞ Resources

- ## Theorem [PG11]:
  - For the $G(n, p; s)$ matching problem, if

$$nps \frac{s^2}{2-s} = 8 \log n + \omega(1)$$

threshold for $aug(G) = 1$
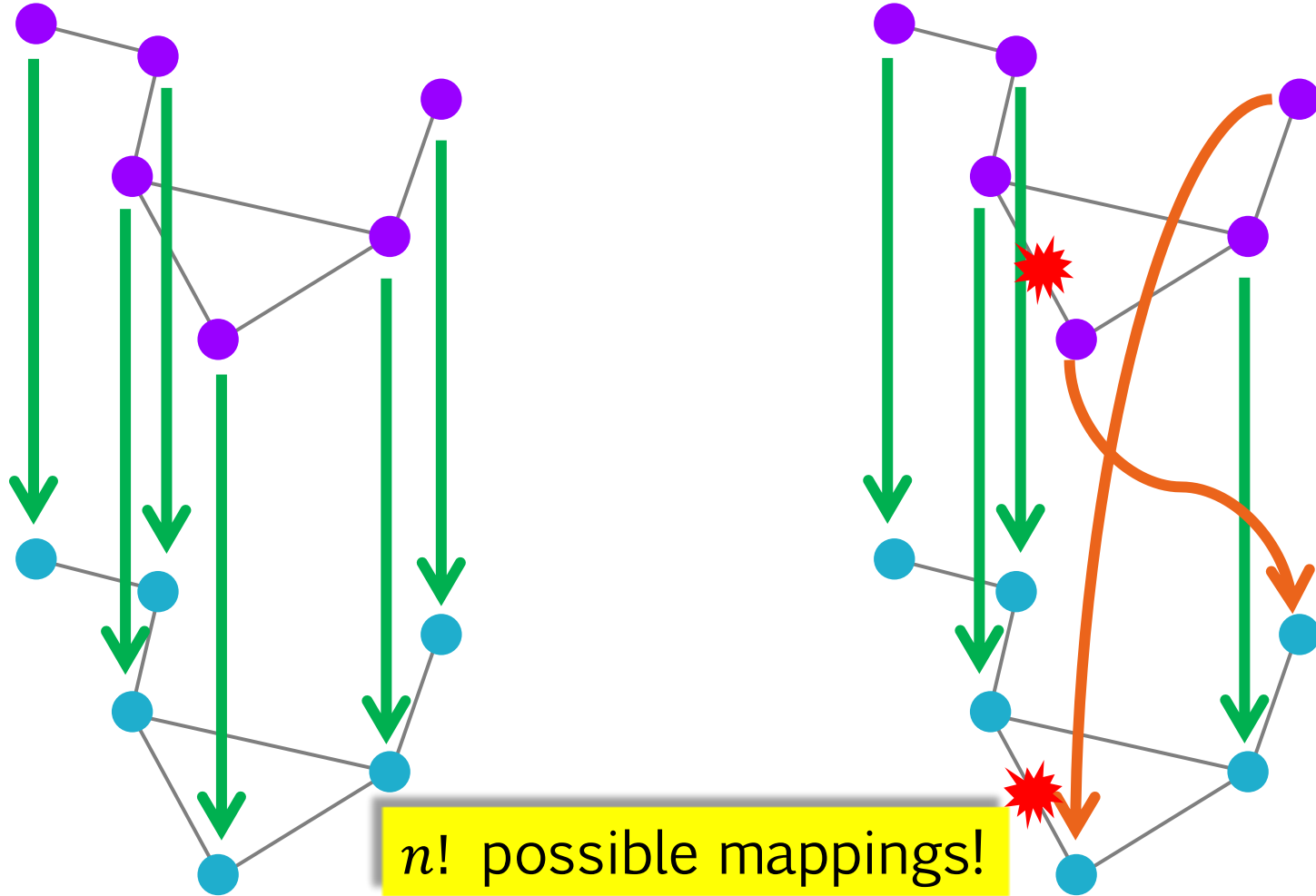
$nps$: $E[\text{degree}]$ of $G_{1,2}$    penalty for dissim of $G_{1,2}$

then $G_{1,2}$ can be perfectly matched a.a.s.

- ## Interpretation:
  - Surprisingly weak condition: degree growing faster than ~$\log n$ enough to break anonymity
  - Decrease with $s$ only quadratic

# Mappings and Edge Mismatch



$n!$ possible mappings!

$$\Delta(\pi_0) = 0$$

$$\Delta(\pi) = 2$$

# Approach

- ## Assumption:
  - Attacker has infinite computational power
  - Can try all possible mappings π and compute edge mismatch function $\Delta(\pi)$

- ## Question:
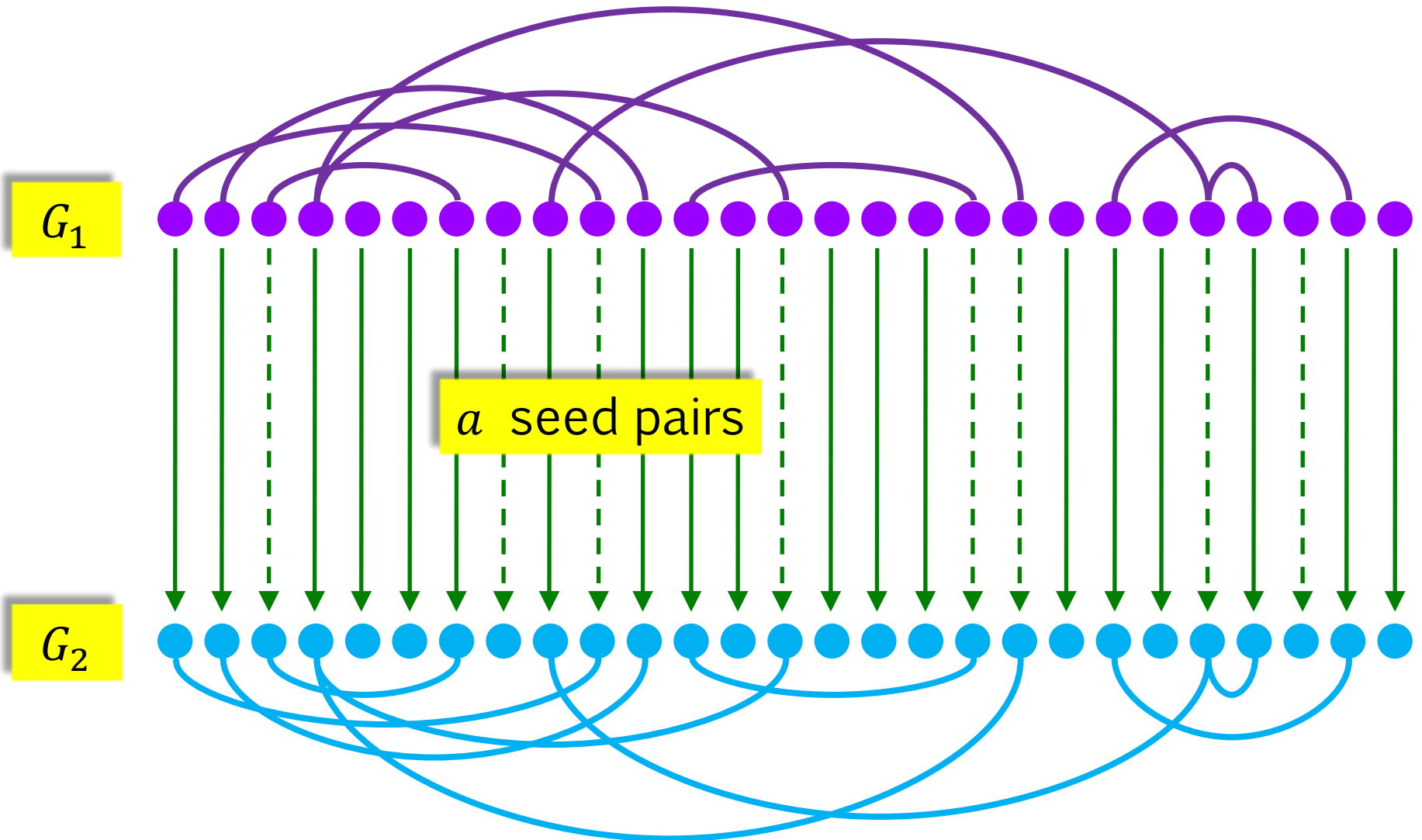  - Are there conditions on $p, s$ such that

  $$P\left\{\pi_0 \ \text{unique} \quad \min \quad \text{of} \quad \Delta(\pi)\right\} \to 1$$

  - If yes: adversary would be able to match vertex sets only through the structure of the two networks!

- ## Note:
  - $G(n, p; s)$ model: statistically uniform, low clustering, degree distribution not skewed -> conjecture: harder than real networks

# Result 2: Graph Matching with Seeds



$G_1$

$a$ seed pairs

$G_2$

# Questions

- **How many seeds are needed?**
- **Is there a phase transition?**
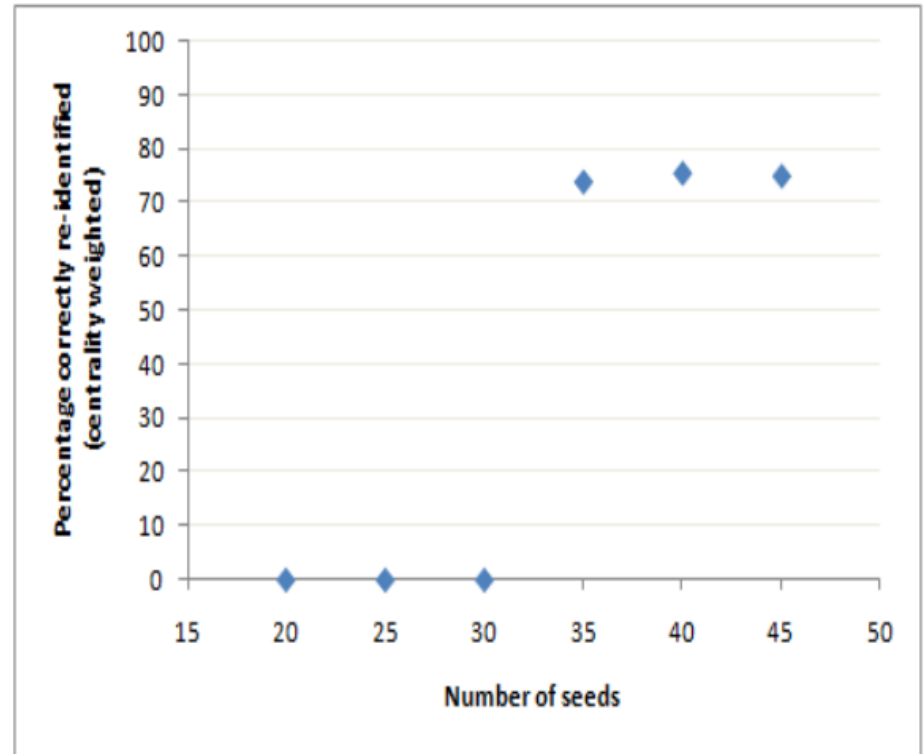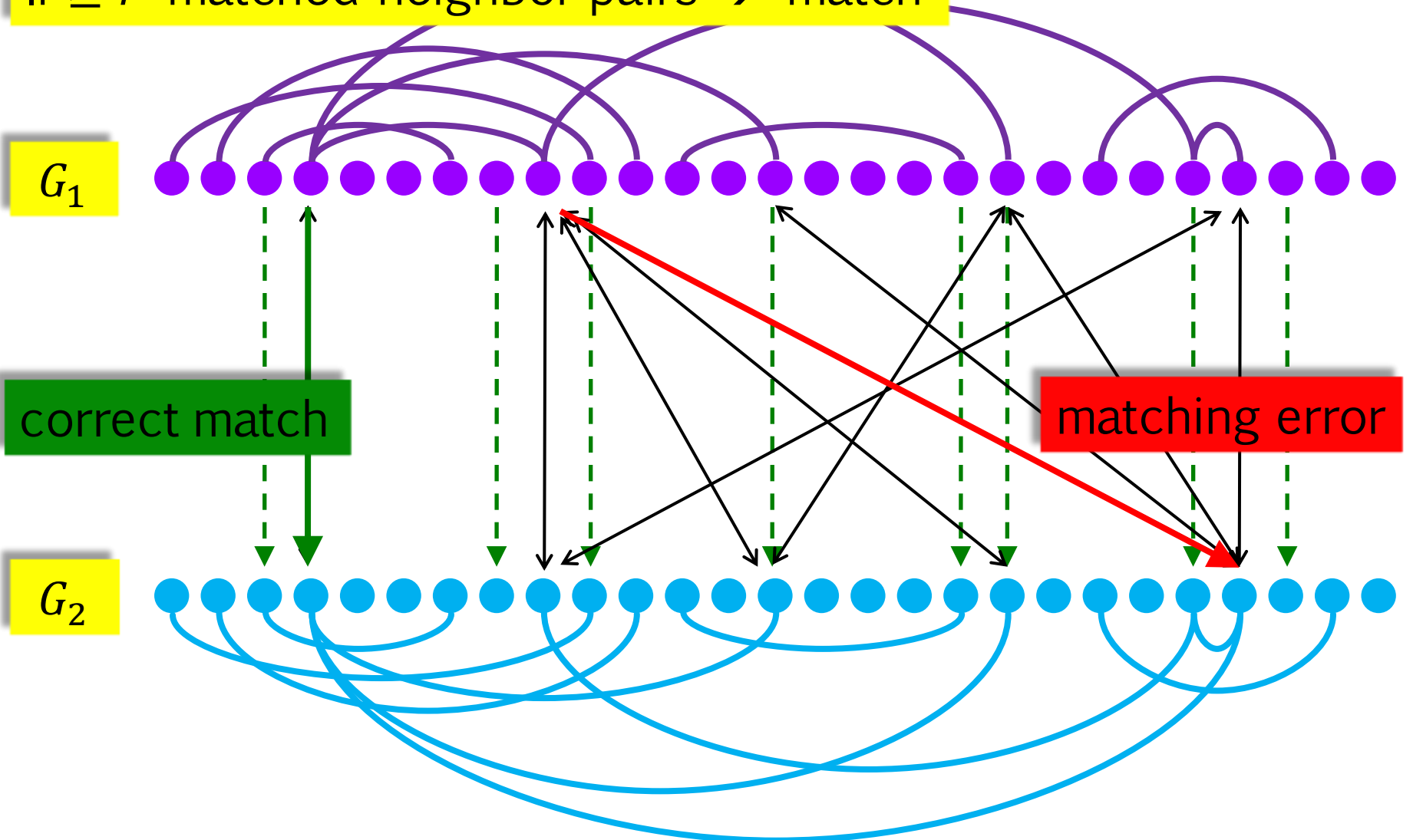- **How efficiently can we match?**
- **Tuning parameters?**



Figure 2. The fraction of nodes re-identified depend sharply on the number of seeds. Node overlap: 25% Edge overlap: 50%

From [A. Narayanan, V. Shmatikov, "De-anonymizing social networks", IEEE Symp. on Security and Privacy, 2009]

# Percolation Graph Matching Algorithm

If $\geq r$ matched neighbor pairs → match

$G_1$

$G_2$

correct match

matching error

# Result: Seed Set Size Threshold for $G(n, p; s)$

- Theorem 2: phase transition in # seeds $a$

  - For $n^{-1} \ll ps^2 \ll s^2 n^{-\frac{4}{r}}$:

    - If $\dfrac{a}{a_c} \to \alpha < 1$, final map is $o(n)$ w.h.p.

    - If $\dfrac{a}{a_c} > \alpha > 1$, final map is $n - o(n)$ w.h.p.

- Seed set size threshold:

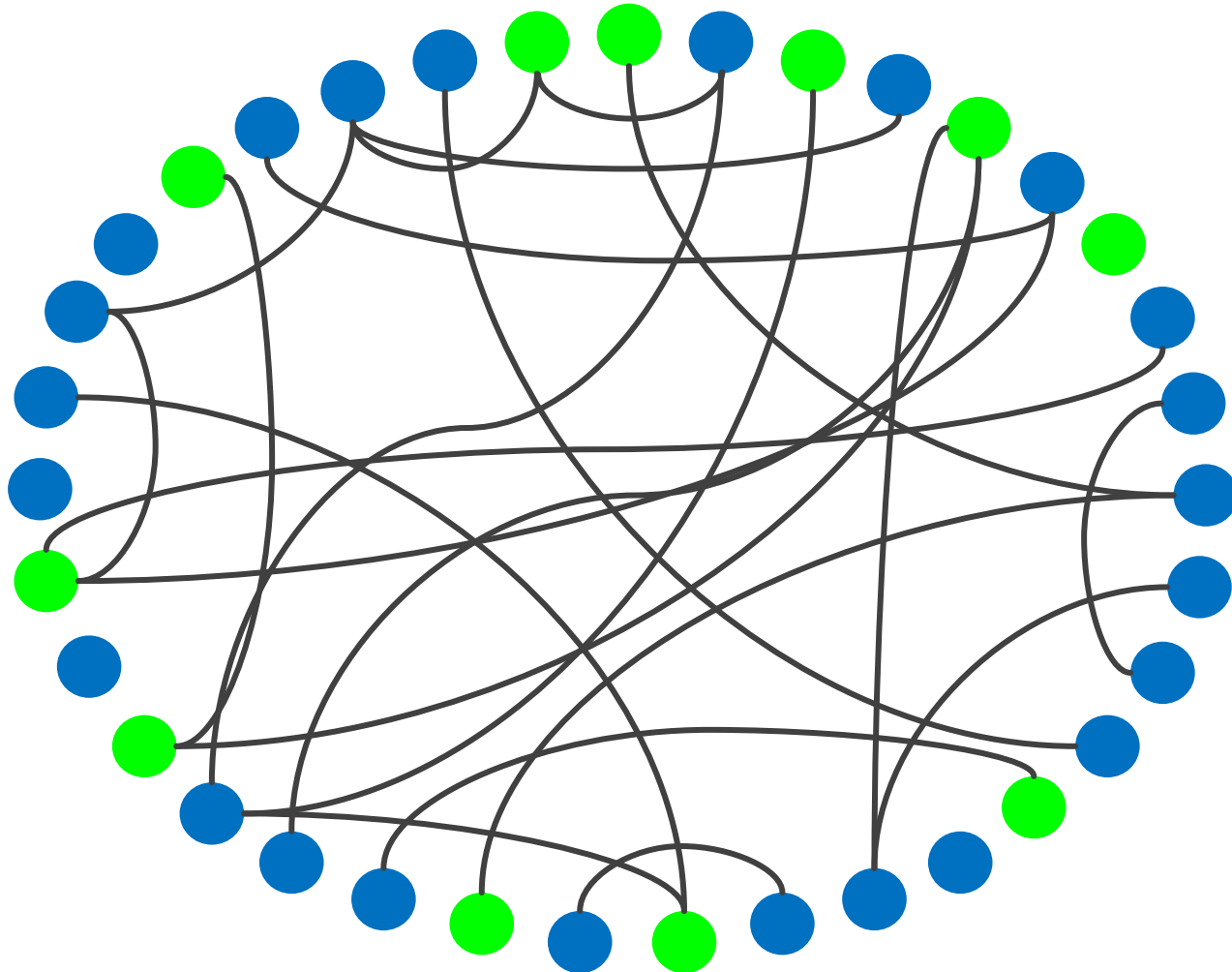  - $a_c = (1 - r^{-1}) \left( \dfrac{(r-1)!}{n(ps^2)^r} \right)^{1/(r-1)}$

  - Slowly densifying network: constant $r$
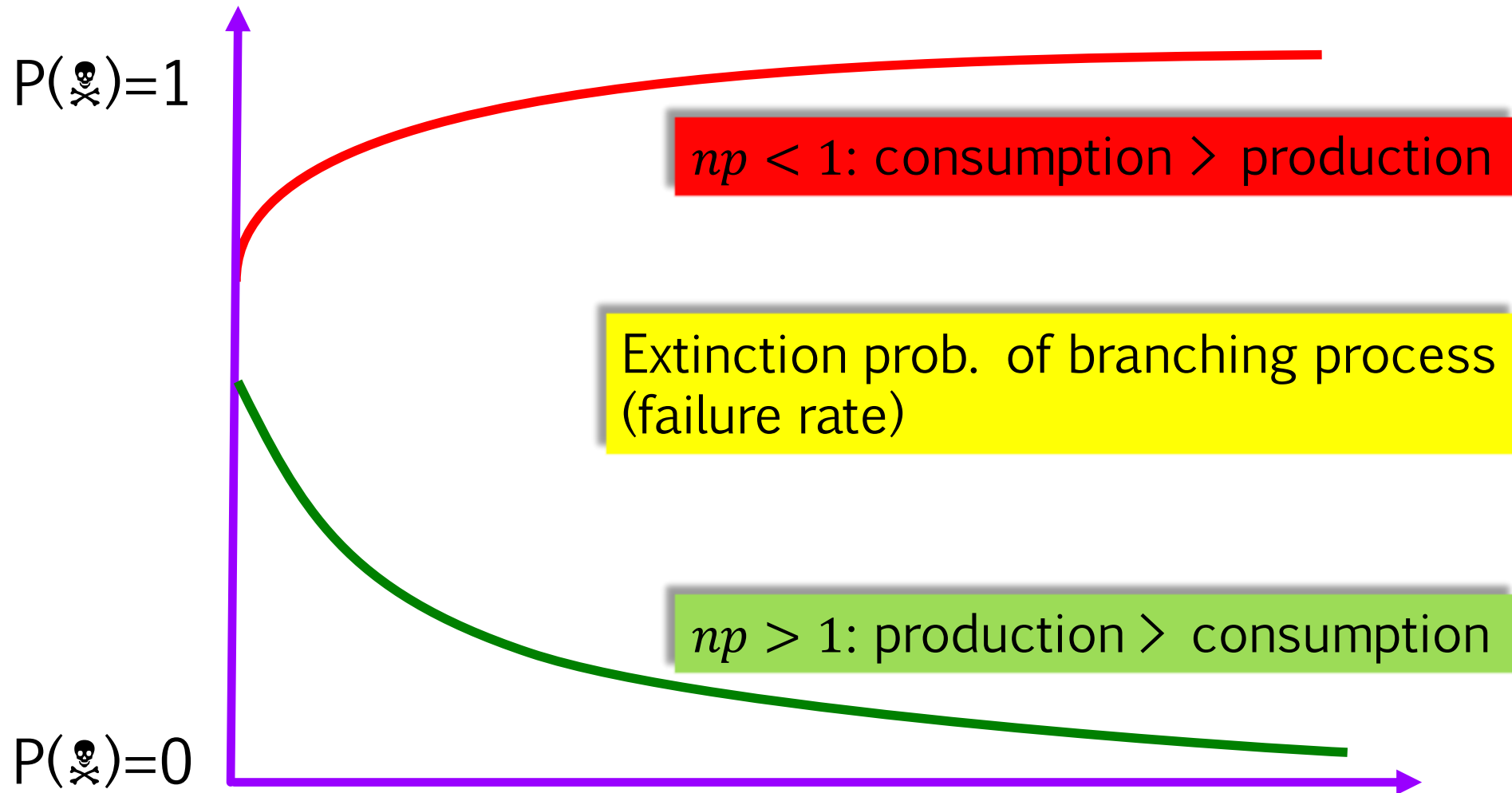  - Growth of $a_c$: a bit less than linear
    - $p = \log n / n$, $s$ fixed $\rightarrow$ $a_c \propto n (\log n)^{-r/(r-1)}$

# Bootstrap Percolation for $G(n,p)$

Activation from $r$ neighbors



[S. Janson, T. Luczak, T. Turova, T. Vallier, Bootstrap Percolation on the Random Graph $G(n,p)$, Annals Applied Prob., 22(5), 2012]

# Giant Component: Branching Process

$P(\skull)=1$

$P(\skull)=0$

$np < 1$: consumption $>$ production

Extinction prob. of branching process (failure rate)

$np > 1$: production $>$ consumption
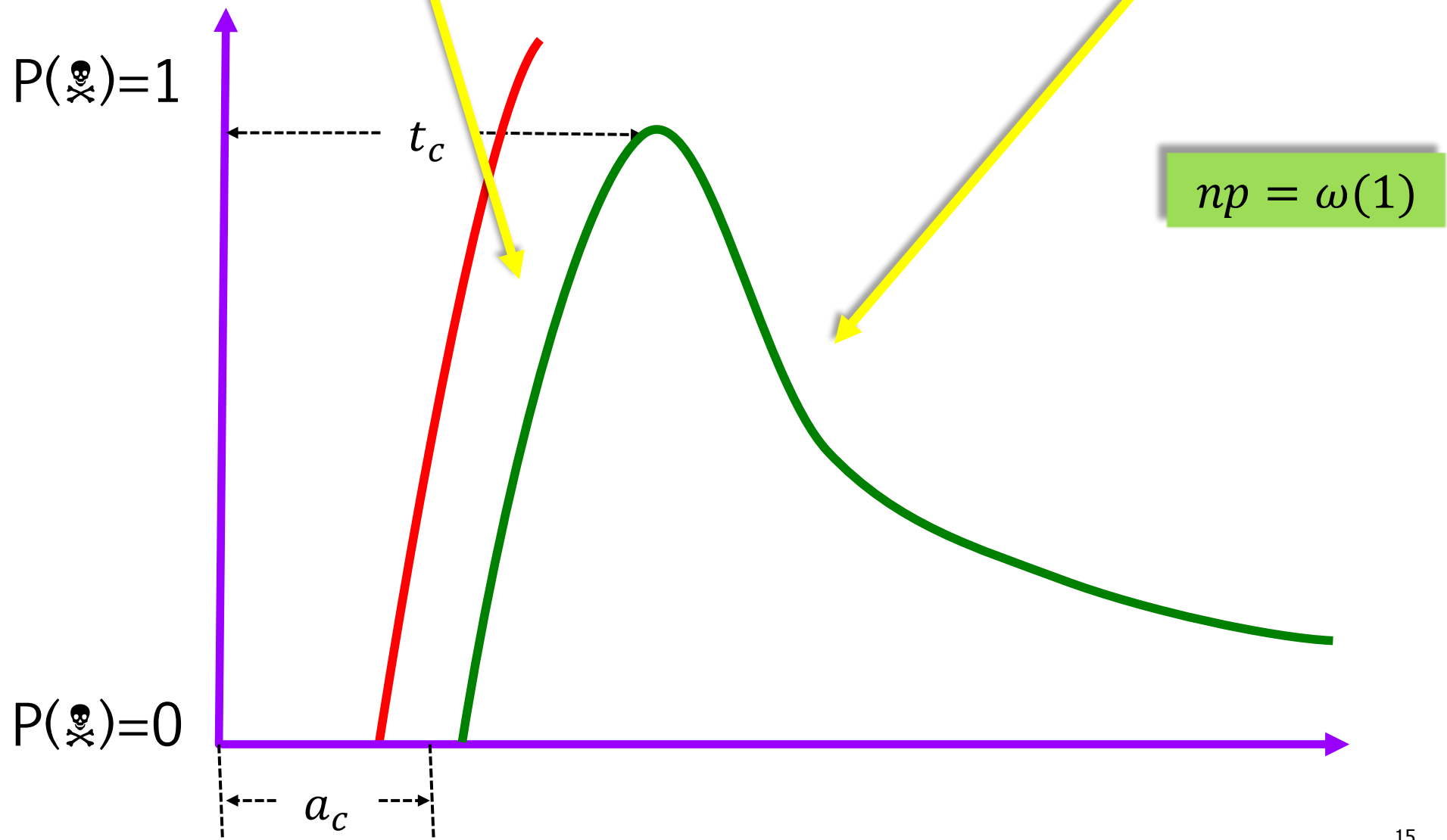
# Bottleneck in Bootstrap Percolation

consumption > production

production > consumption

$P(\text{☠})=1$

$t_c$

$np = \omega(1)$

$P(\text{☠})=0$

$a_c$
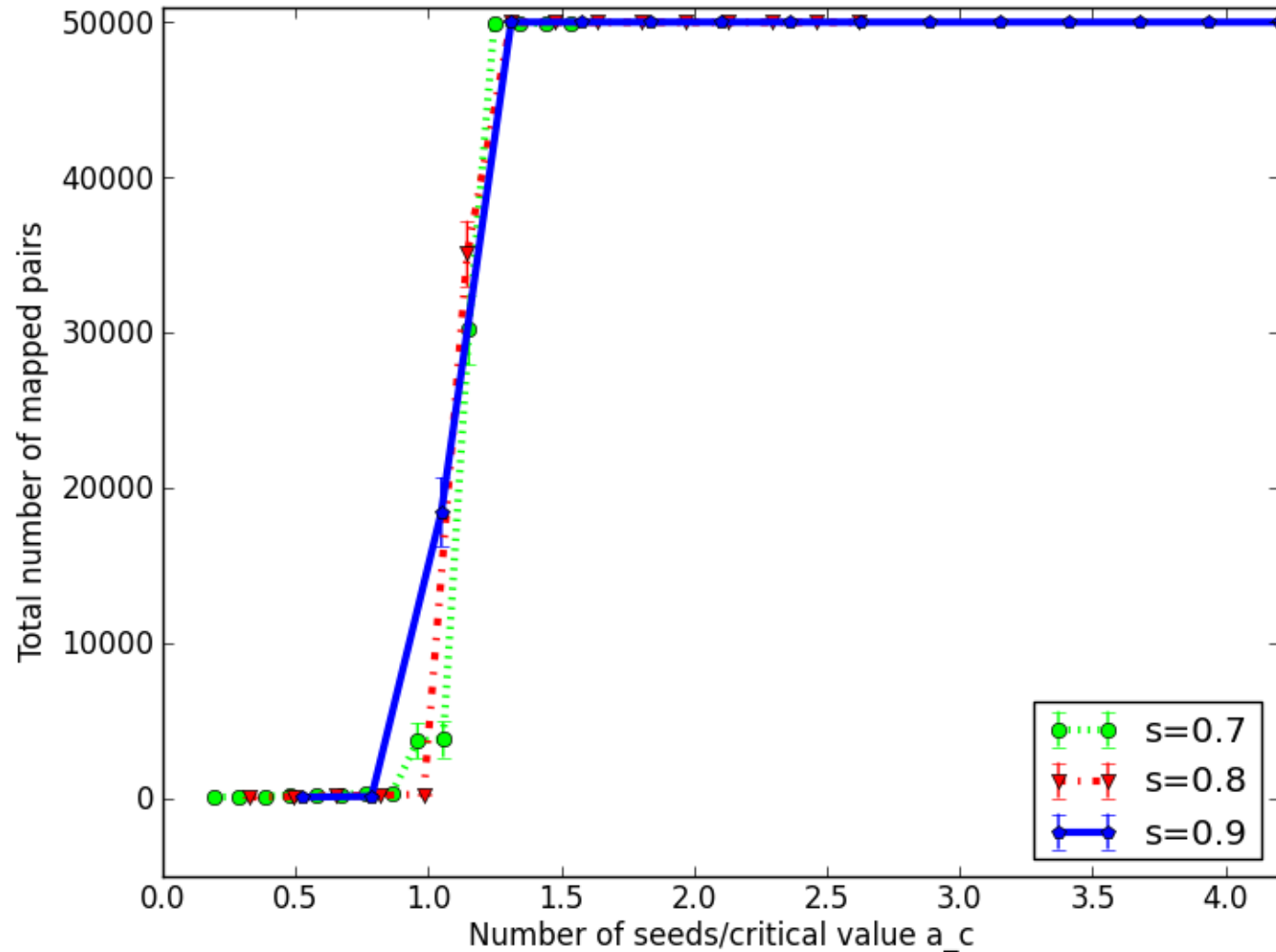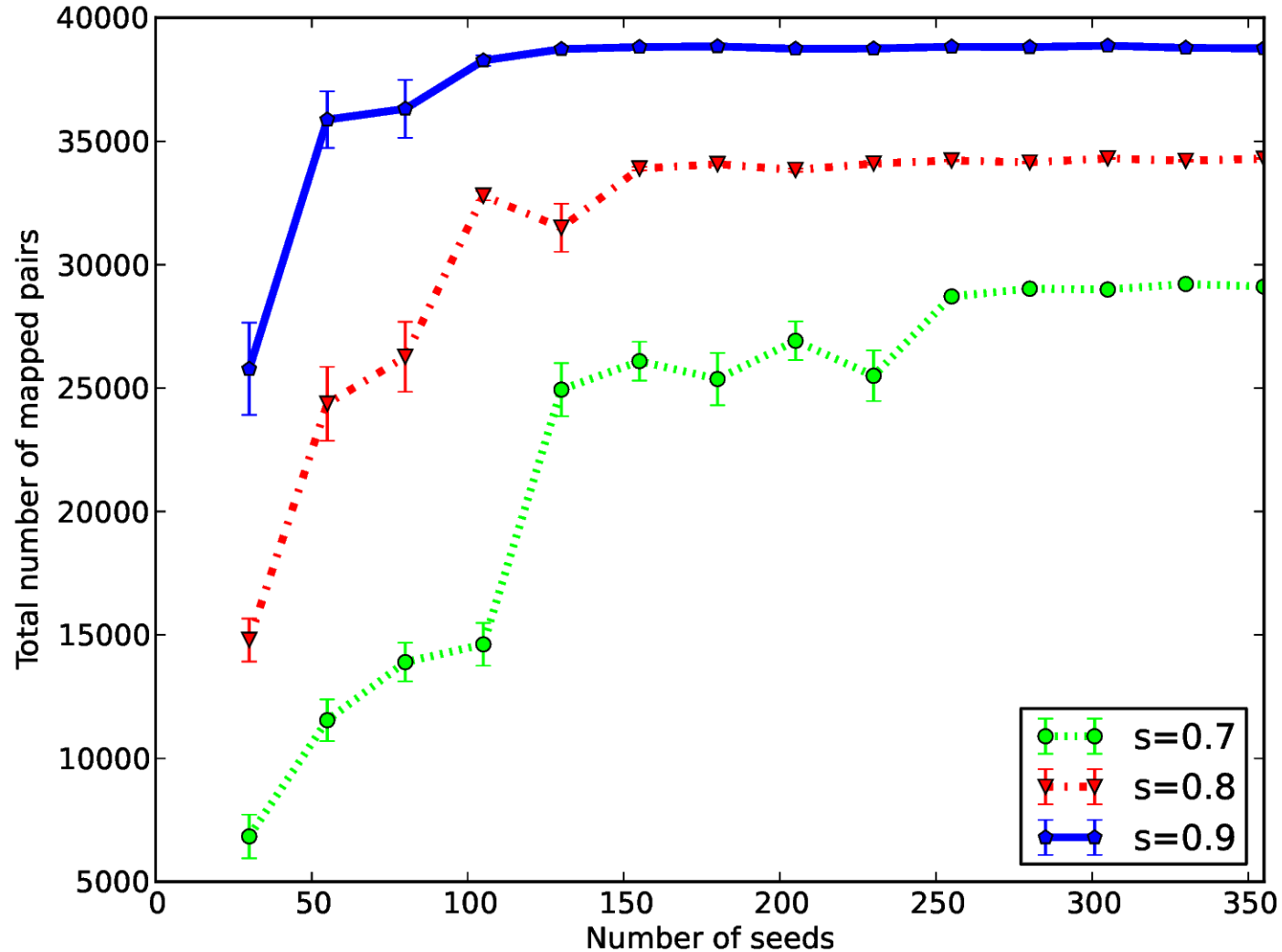
15

# Simulation of PGM with $G(n, p; s)$ Network

# Simulation of PGM with $G(n, p; s)$ Network

# Slashdot Social Network

# Result 3: Getting Started

- ## How to find seeds? [PFG13]
  - Efficient (polynomial) algorithm to generate seed set
  - Does not work for $G(n, p)$
- ## Real graphs:
  - More heterogeneous than $G(n, p)$: degree skew, transitivity
  - Provides features for nodes

# Finding Seeds: Bayesian Framework



Fingerprint:
$(X_1 = 4,$
$D_{1,1} = 1,$
$D_{1,2} = 3)$

Fingerprint:
$(X_1 = 1,$
$D_{1,1} = 4,$
$D_{1,2} = 2)$

seed1

seed2

Fingerprint:
$(X_2 = 3,$
$D_{2,1} = 3,$
$D_{2,2} = 1)$

Fingerprint:
$(X_2 = 3,$
$D_{2,1} = 1,$
$D_{2,2} = 3)$

# Seed: Bayesian Framework

**Node:** $U_1$ ←——— **=?** ———→ **Node:** $U_2$

**Fingerprint:**
$(X_1 = 4,$
$D_{1,1} = 1,$
$D_{1,2} = 3)$

**Fingerprint:**
$(X_1 = 1,$
$D_{1,1} = 4,$
$D_{1,2} = 2)$

Network sampling model:
$P(\text{fp1}, \text{fp2} \mid U_1 = U_2),$
$P(\text{fp1}, \text{fp2} \mid U_1 \neq U_2)$

Single-pair posterior:
$P(U_1 = U_2 \mid \text{fp1}, \text{fp2})$

**Fingerprint:**
$(X_2 = 3,$
$D_{2,1} = 3,$
$D_{2,2} = 1)$

**Fingerprint:**
$(X_2 = 3,$
$D_{2,1} = 1,$
$D_{2,2} = 3)$

Jointly MAP matching:
Best bipartite matching $\pi$ s.t. max
P(all matched correctly | all fingerprints)

# Conclusion

- **Graph matching problem:**
  - Social networks: privacy; merging
  - Model as noisy graph isomorphism problem
  - How much information in network structure?
- $G(n,p;s)$ **random graph model:**
  - Parsimonious: density $(p)$, similarity $(s)$
  - Information-theoretic characterization of feasible region – condition is quite mild
- **Percolation Graph Matching algorithm:**
  - Simple algorithm, propagating evidence over node pairs
  - Actually works very well in practice; parsimonious $(r)$
- **Analysis:**
  - Sharp phase transition in seed set size $(a)$, confirms empirical observation