

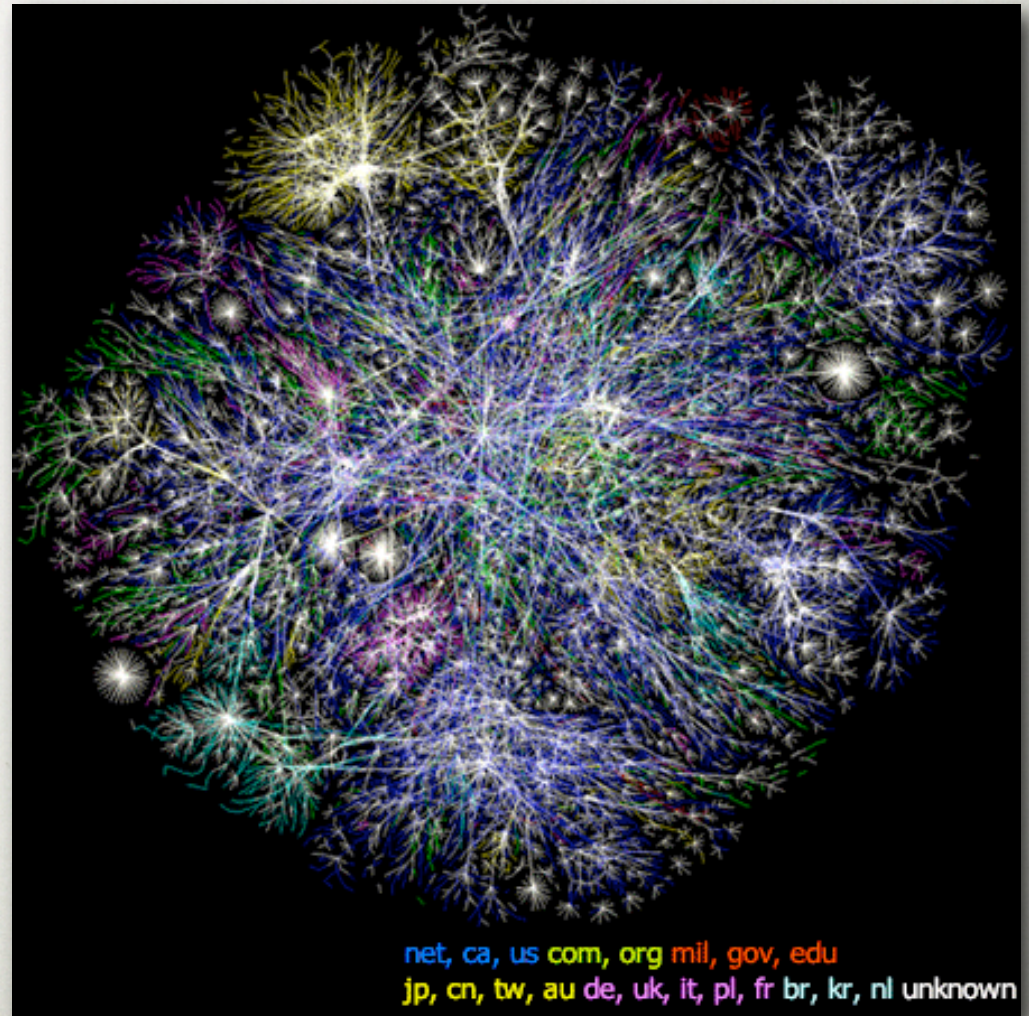
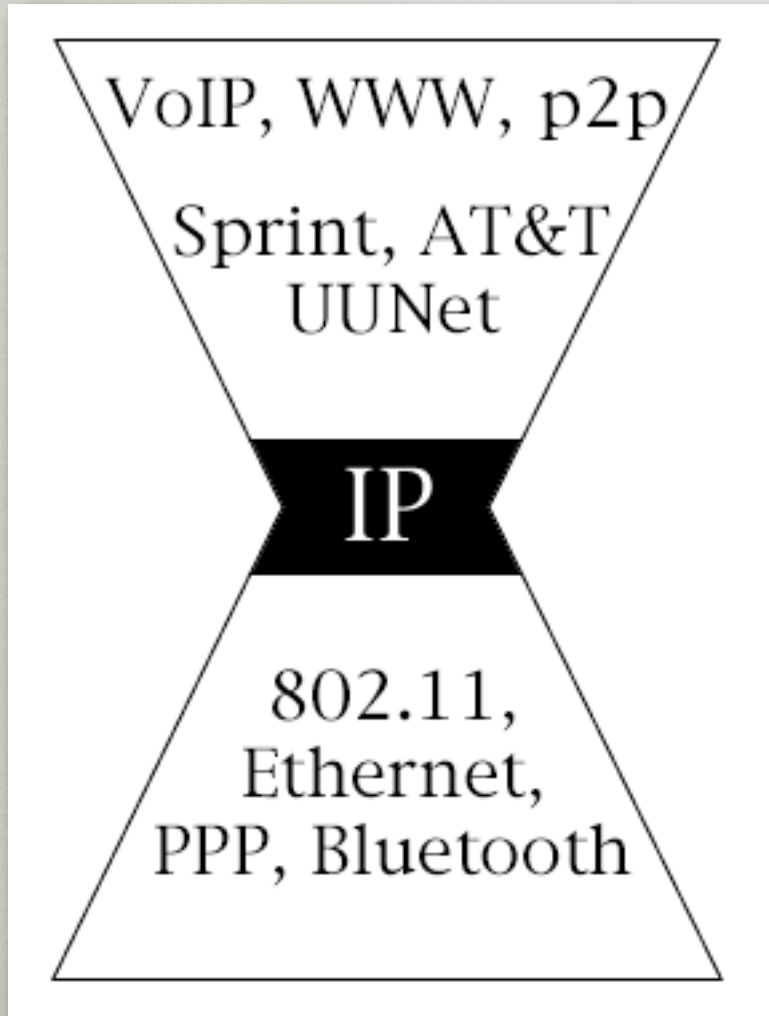
A network diagram with several circular nodes connected by lines, set against a dark gray background. The nodes are arranged in a roughly triangular pattern, with lines connecting them to form a network structure.

THE BIAS OF TRACEROUTE : ACCURACY & SCALING IN INTERNET MAPPING

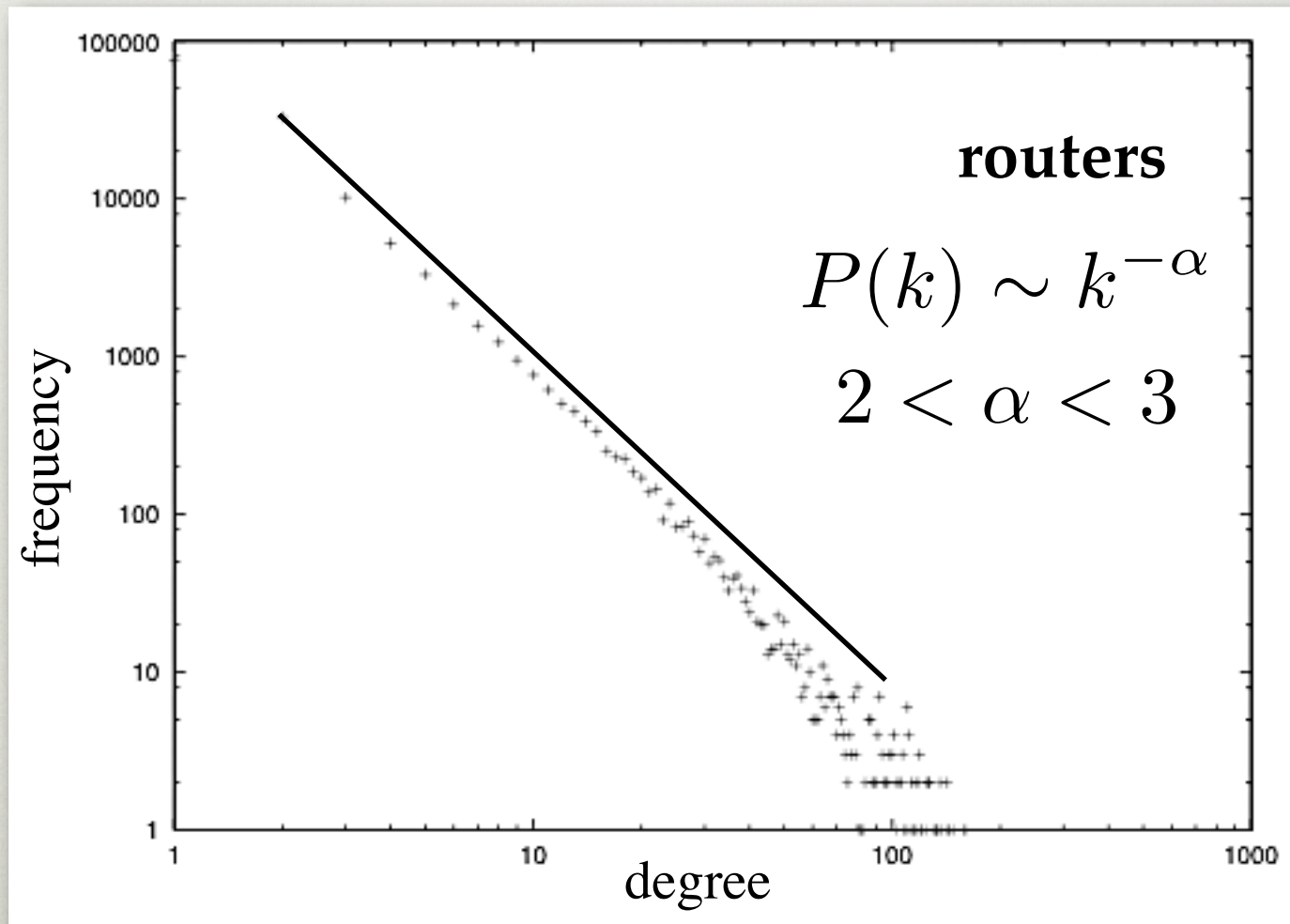
Aaron Clauset
UNM Computer Science at
CAIDA WIT
11 May 2006

with Christopher Moore
David Kempe, and
Dimitris Achlioptas

INTERNET MAPS (2)



INTERNET MAPS (3)

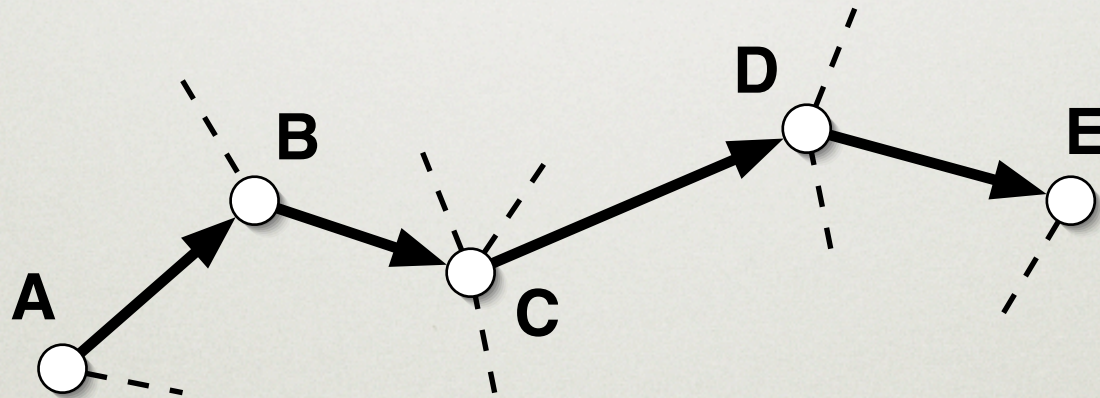


originally shown by [FFF99]

THE PROBLEM

Topology must be inferred

- router links not directly observable
- links queried indirectly by paths



- [FFF99] data based on single-source traceroute

INHERENT BIAS

Lakhina et al., INFOCOM 2003

- Traceroute sampling apparently *biased*
- Distant edges less likely to be observed, and nearby edges are over-represented
- Power-law degree distributions can appear when none exist...
- ... even in Erdős-Rényi random graphs!
- Our first result **confirms this analytically.**

RANDOM GRAPHS

E-R random graphs

$$G(n, p = c/n)$$

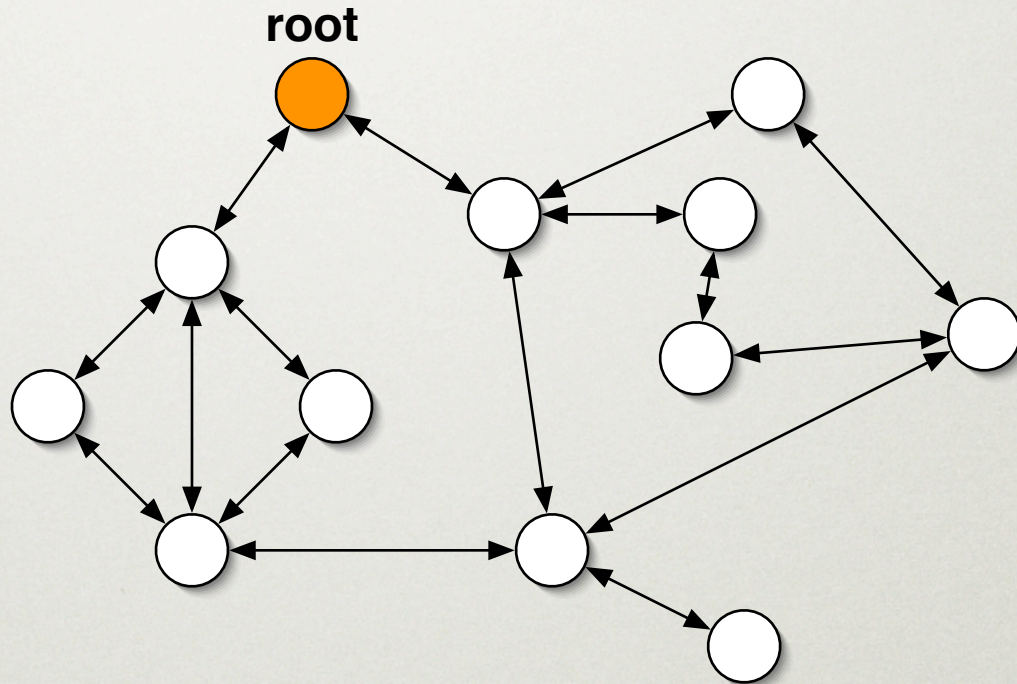
- n vertices; each pair connected with probability p
- Degree distribution $P(k)$ is Poisson with mean c
- For $c > 1$, $\Theta(n)$ vertices form a giant component.

SOME PERSPECTIVE

- Random graphs are totally unrealistic models of anything **real**
- But, random structures are good **null-models**
- If you think some property **P** is *interesting* but a totally random graph also exhibits **P**, then maybe not so interesting...
- Many power laws have trivial explanations, e.g., exponential sampling and certain multiplicative random processes

BACK TO TRACEROUTE

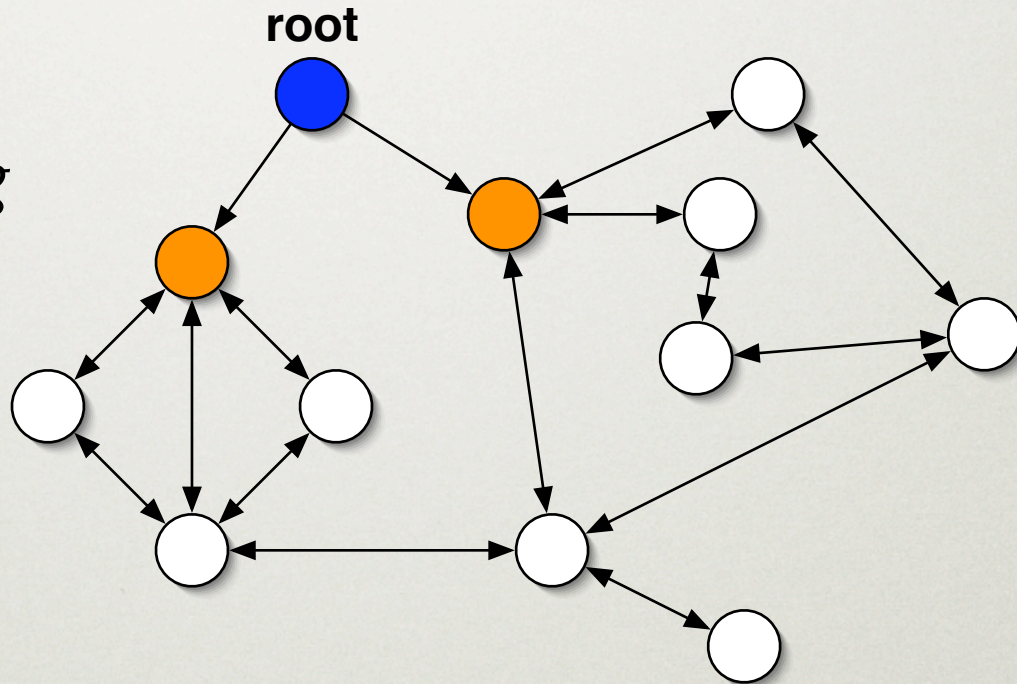
1. choose root



underlying network topology

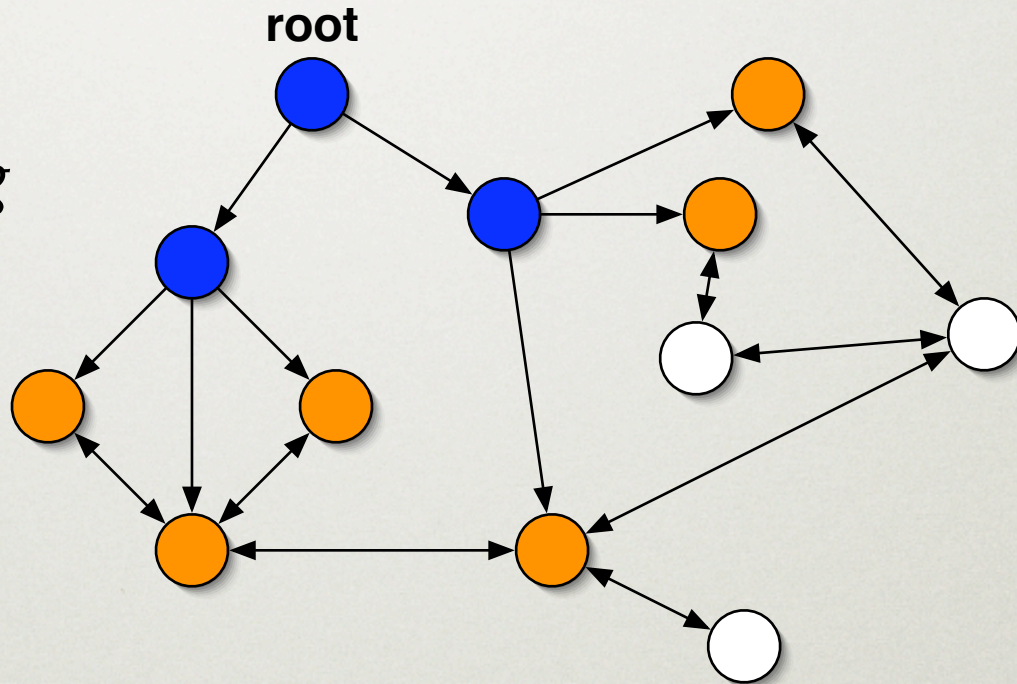
BUILDING A SPANNING TREE

1. choose root
2. add edges to *unknown* neighbors
they become *pending*



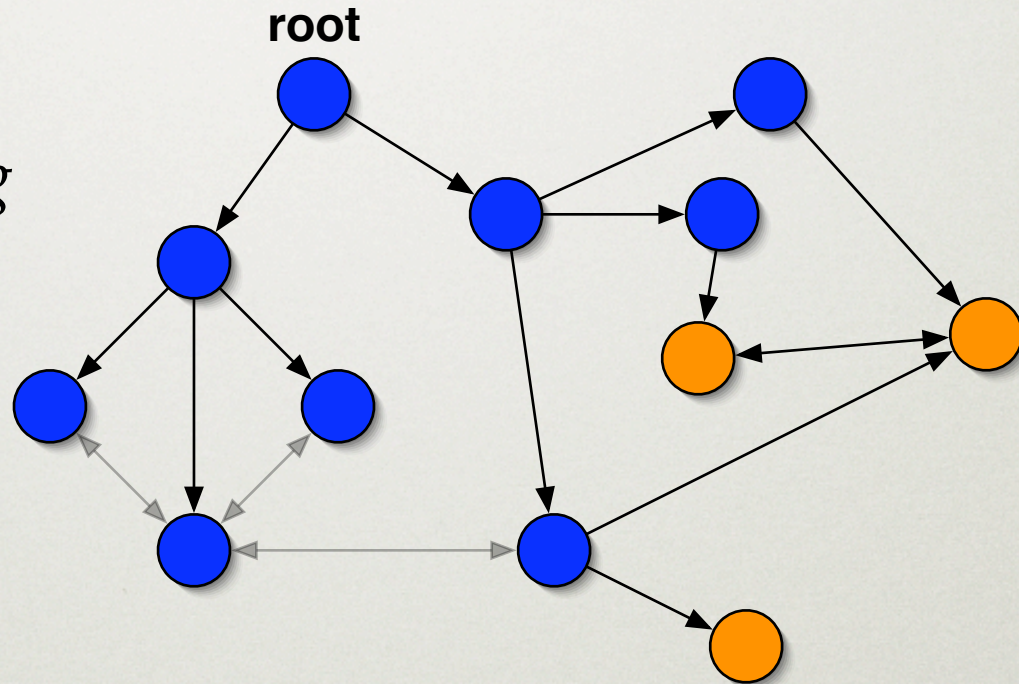
BUILDING A SPANNING TREE

1. choose root
2. add edges to *unknown* neighbors
they become *pending*
3. repeat 2 for each *pending* node;
it becomes *reached*



BUILDING A SPANNING TREE

1. choose root
2. add edges to *unknown* neighbors they become *pending*
3. repeat 2 for each *pending* node; it becomes *reached*



BUILDING A SPANNING TREE

label all vertices UNKNOWN, label one PENDING

while there are PENDING vertices

 choose* a PENDING vertex v

 label v REACHED

for every UNKNOWN neighbor u of v

 label u PENDING

 add (u, v) to TREE

* e.g., depth-, breadth- or random-first

ANALYSIS

- Let $S(T)$ be number of PENDING vertices at step T
- Let $U(T)$ be number of UNKNOWN vertices at step T
- The expected differences at each step are

$$E[U(T + 1) - U(T)] = -pU(T)$$

$$E[S(T + 1) - S(T)] = -1 + pU(T)$$

ANALYSIS

- We rescale to $t = T/n, u = U/n, s = S/n$ to obtain

$$\frac{du}{dt} = -cu$$

$$\frac{ds}{dt} = cu - 1$$

with solutions

$$u(t) = e^{-ct}$$

$$s(t) = 1 - t - e^{-ct}$$

- Wormald's theorem makes this rigorous; we can predict $S(T)$ and $U(T)$ w.h.p. to within $o(n)$.

ANALYSIS

- At each time t , degree of a PENDING vertex is a Poisson distribution

$$\text{Poisson}(m, k) = \frac{e^{-m} m^k}{k!}$$

with mean $m = cu(t)$

- So, integrate this up to time t_f , at which point every vertex in giant component is REACHED

ANALYSIS

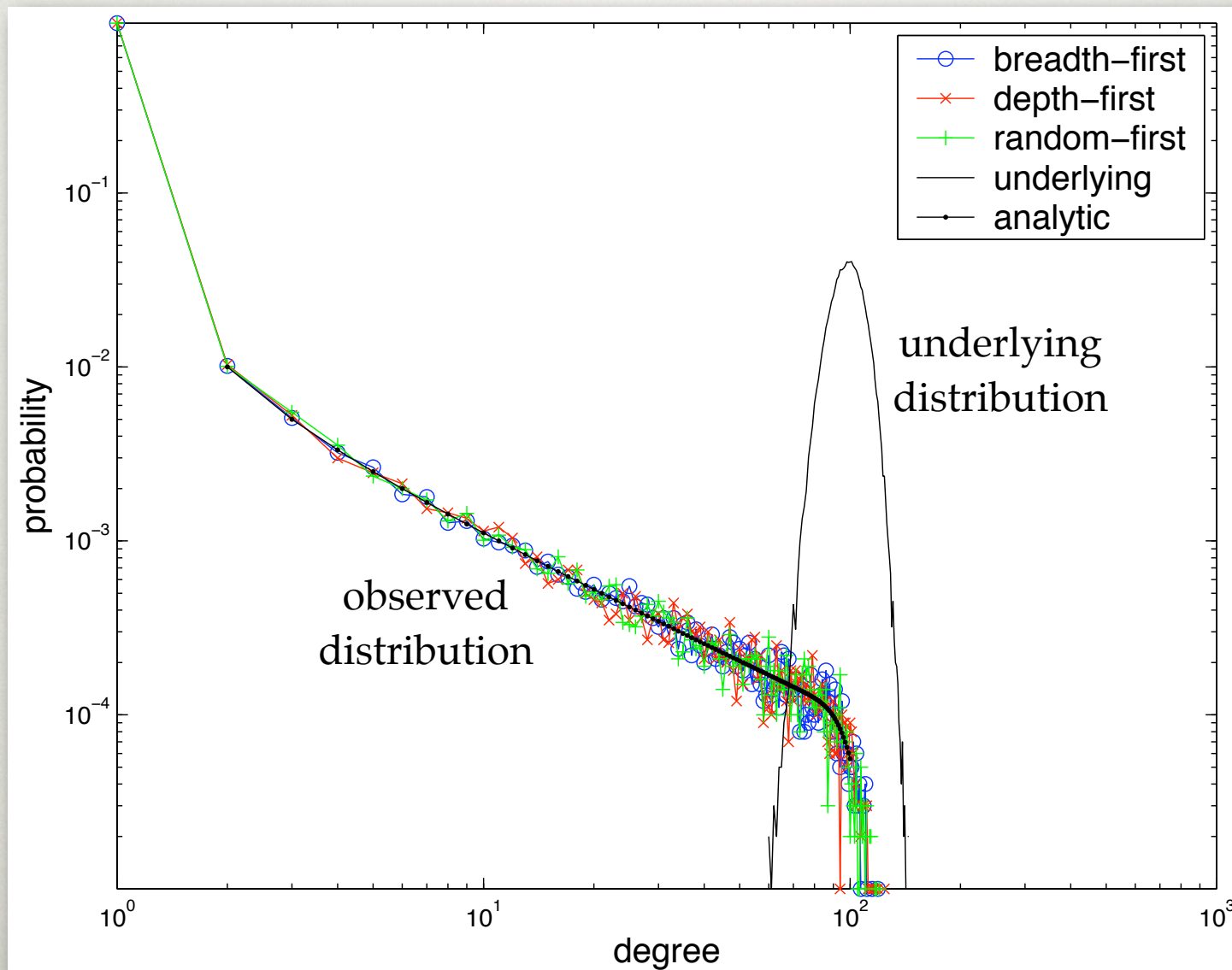
- Thus, averaging over all vertices in TREE

$$P(k + 1) = \frac{1}{t_f} \int_0^{t_f} dt \text{Poisson}(cu(t), k)$$

which is a power law with exponent $\alpha = -1$

$$P(k + 1) = (1 - o(1)) \frac{\Gamma(k)}{ck!} \sim \frac{1}{ck}$$

THEORY & EXPERIMENT



MORE GENERALLY

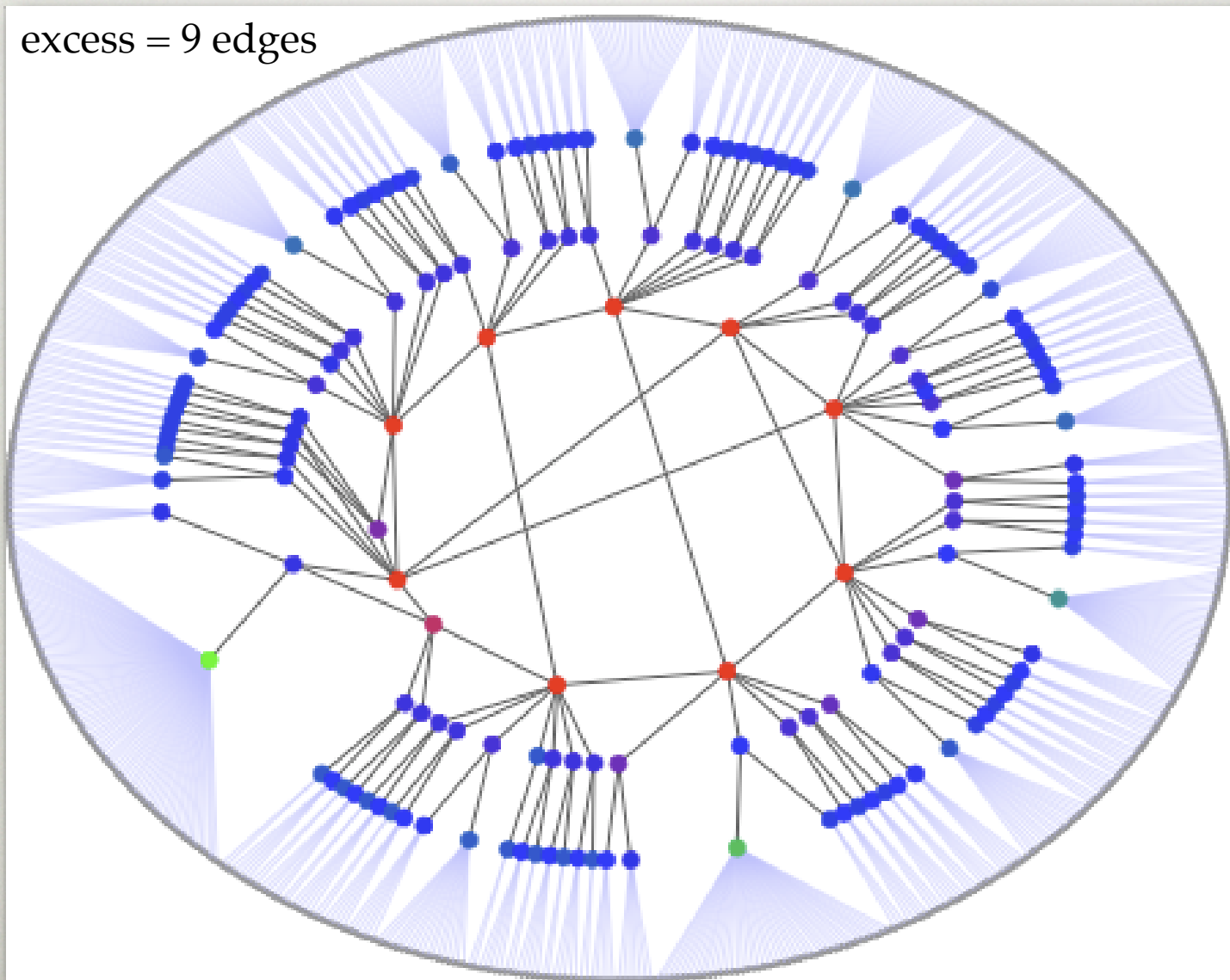
- Generalized process for random graphs with arbitrary underlying degree distributions, e.g.,
 - d-regular $\xrightarrow{\text{obs}}$ power law!
 - Poisson $\xrightarrow{\text{obs}}$ power law
 - power law $\xrightarrow{\text{obs}}$ different power law
- Mathematical crank to get observed distribution

$$g^{\text{obs}}(z) = z \int_0^1 dt g' \left[t - \frac{(1-z)}{g'(1)} g' \left(\frac{g'(t)}{g'(1)} \right) \right]$$

IT GETS WORSE

- Suppose underlying graph has a power law
- What is *observed* exponent?
- Severity of bias depends on number of **redundant edges** (the ones missed by traceroute)

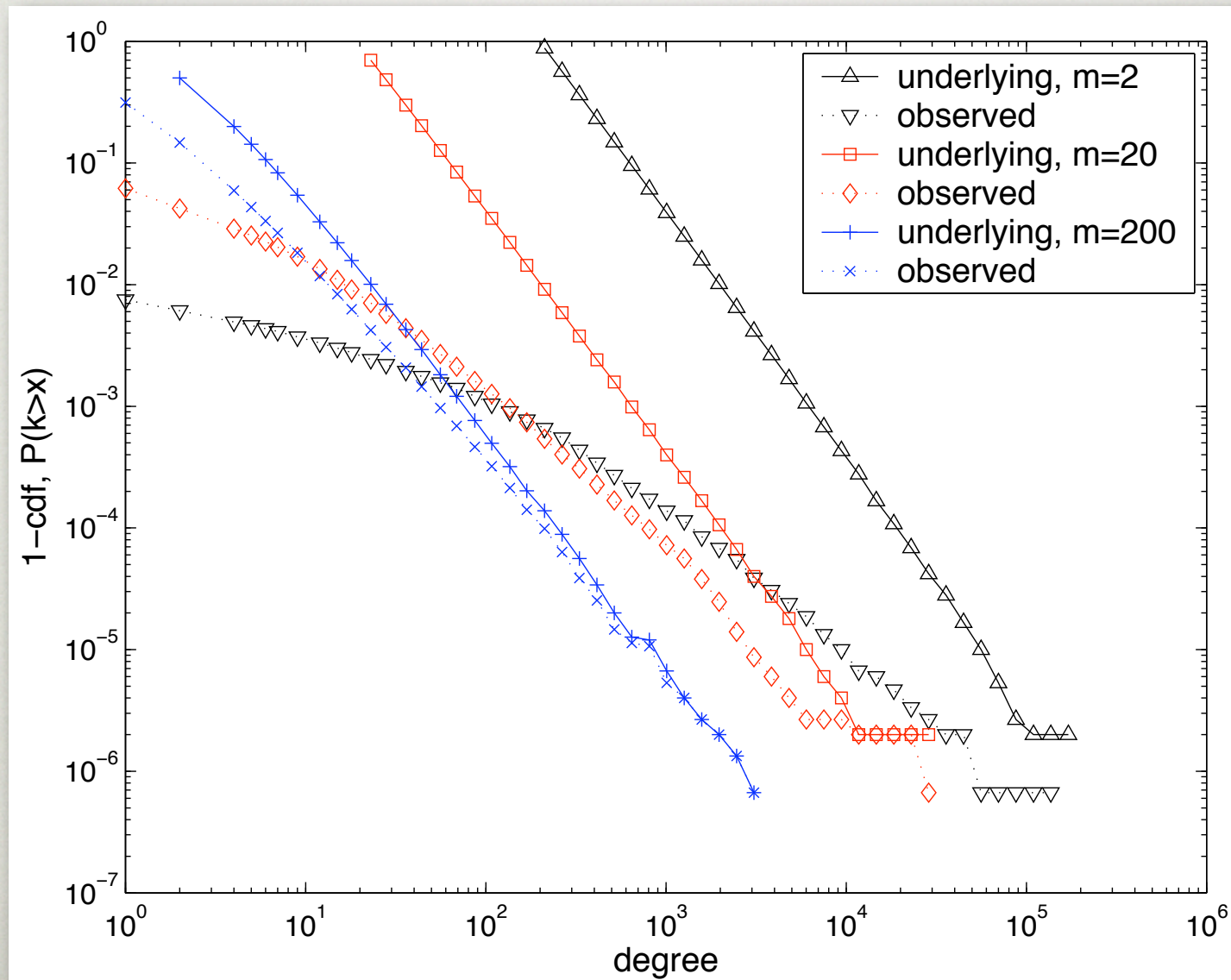
LOW EDGE-REDUNDANCY



IT GETS WORSE

- Suppose underlying graph has a power law
- What is *observed* exponent α_{obs} ?
- Severity of bias depends on number of **redundant edges** (the ones missed by traceroute)
 - You pick any underlying exponent α ,
 - You pick observed exponent $\alpha_{\text{obs}} < \alpha$
 - I can pick mean degree $\langle k \rangle$ so that traceroute gives you α_{obs} from α

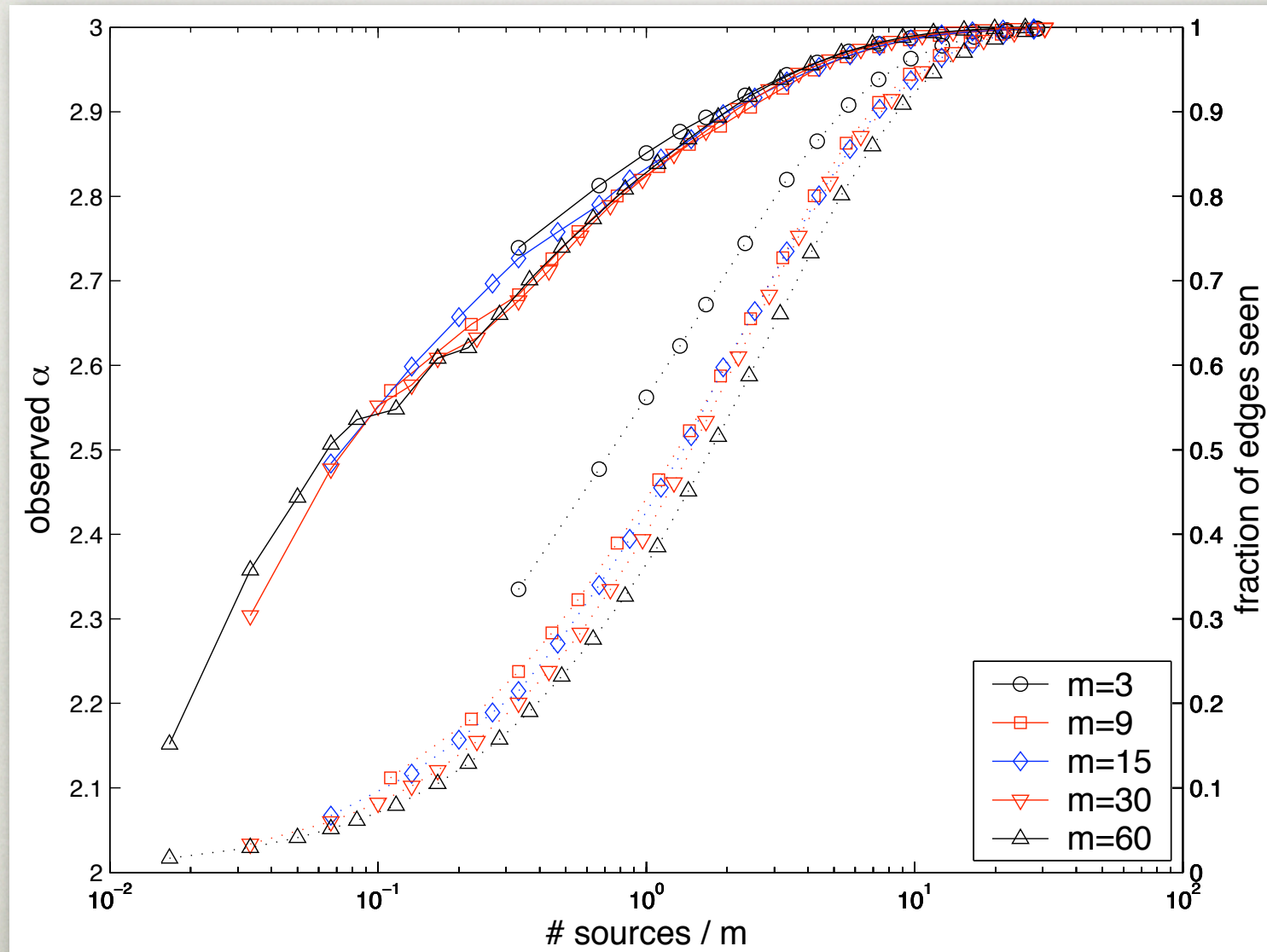
IT GETS WORSE



AND WORSE STILL...

- Okay, so just use more sources!
- Bias remains severe **until we see almost all edges**
- The marginal value of an additional source is low, but positive!

AND WORSE STILL...



CHALLENGES

- **Topological** inferences from traceroute **must** account for its bias (or be heavily caveated)
- But, how to do this? Not yet clear... some ideas:
 - Estimate real marginal value of new sources, for many sources (DIMES?)
 - Estimate mean degree $\langle k \rangle$
 - Use convergence rates to extrapolate
 - Use model of hierarchical organization to estimate the missing edges

FIN
