

A day of Internet life

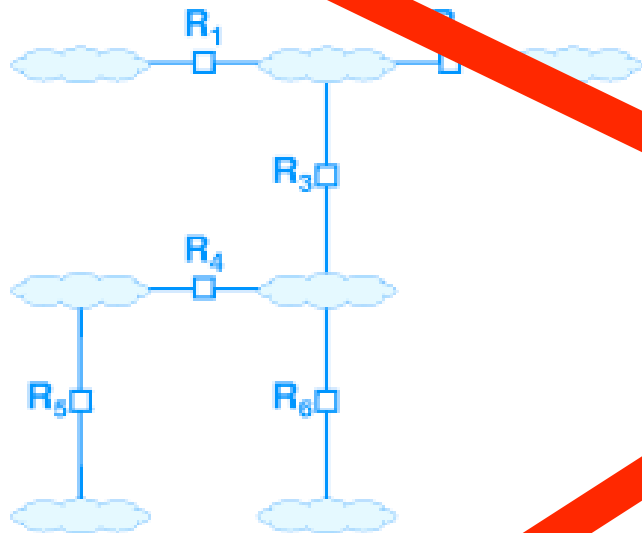
Mark Meiss^{1,3}, Filippo Menczer^{1,2}, Alessandro Vespignani²



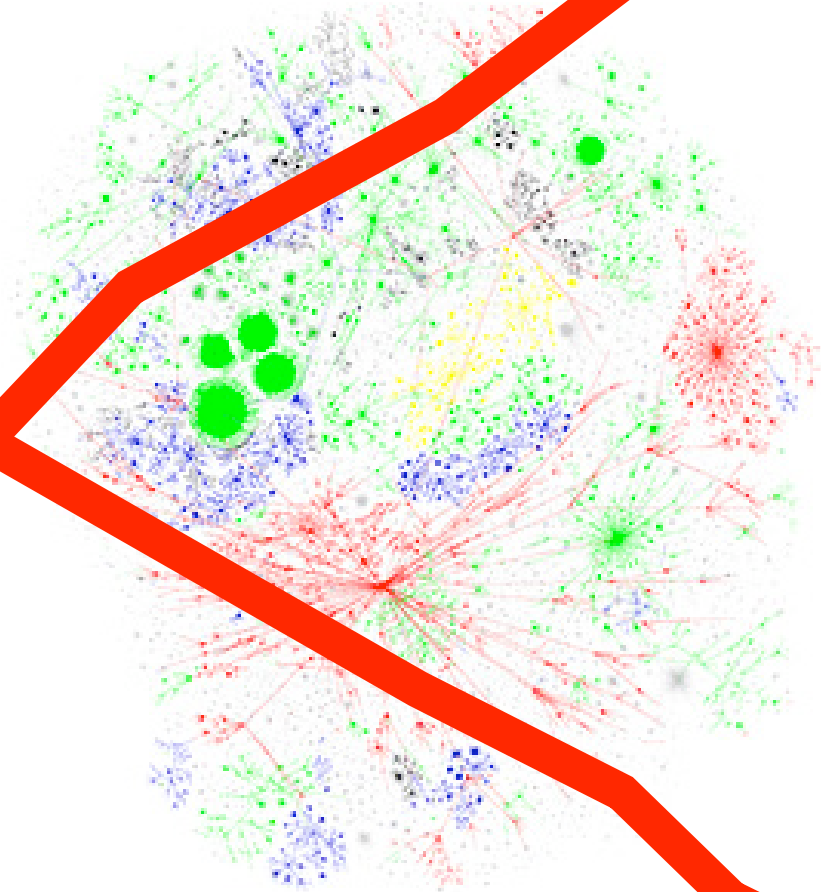
(1) Department of Computer Science
(2) Department of Informatics
(3) Advanced Network Management Lab
Indiana University, Bloomington



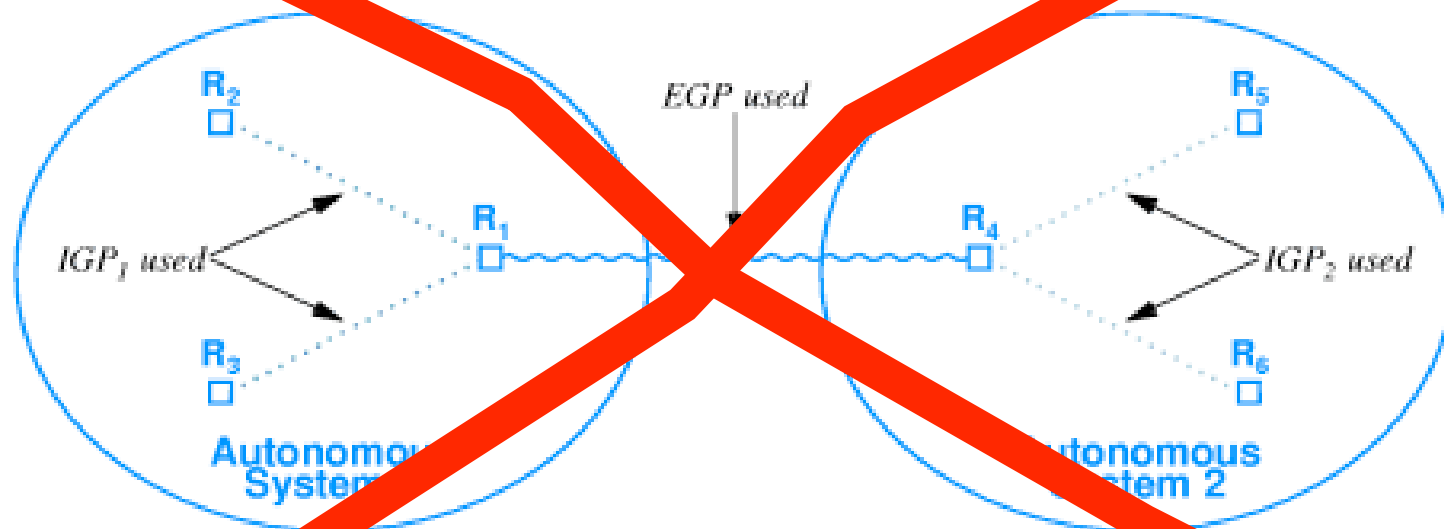
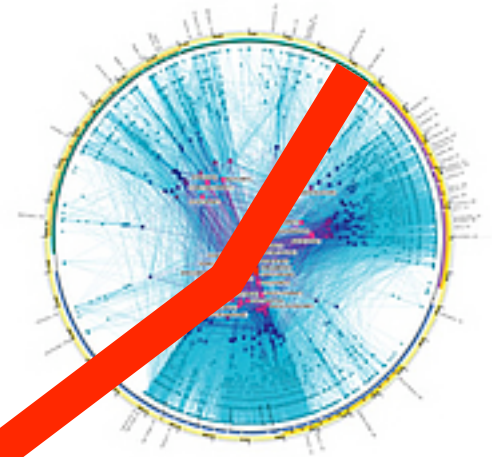
“The Internet”?



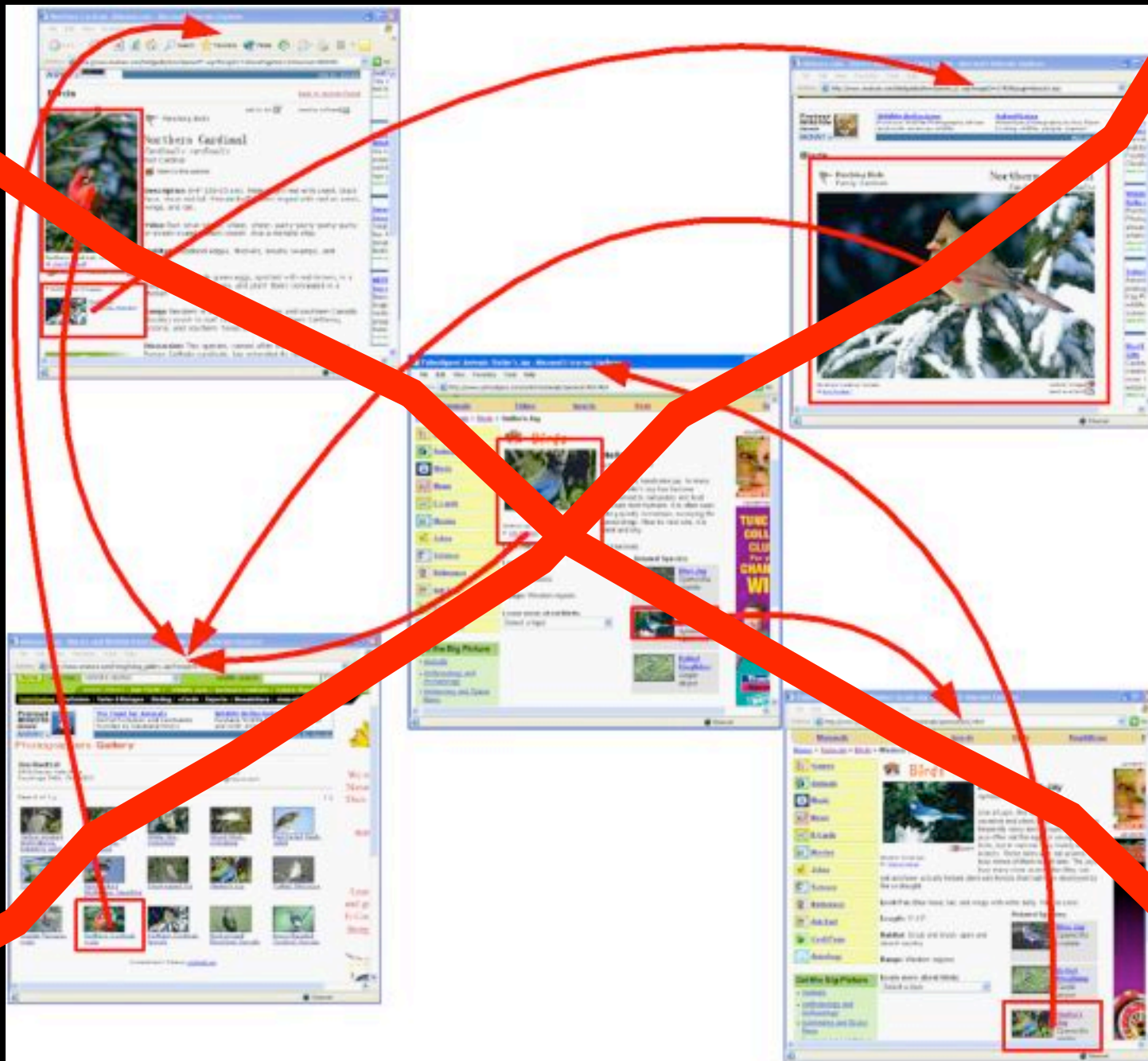
(a)

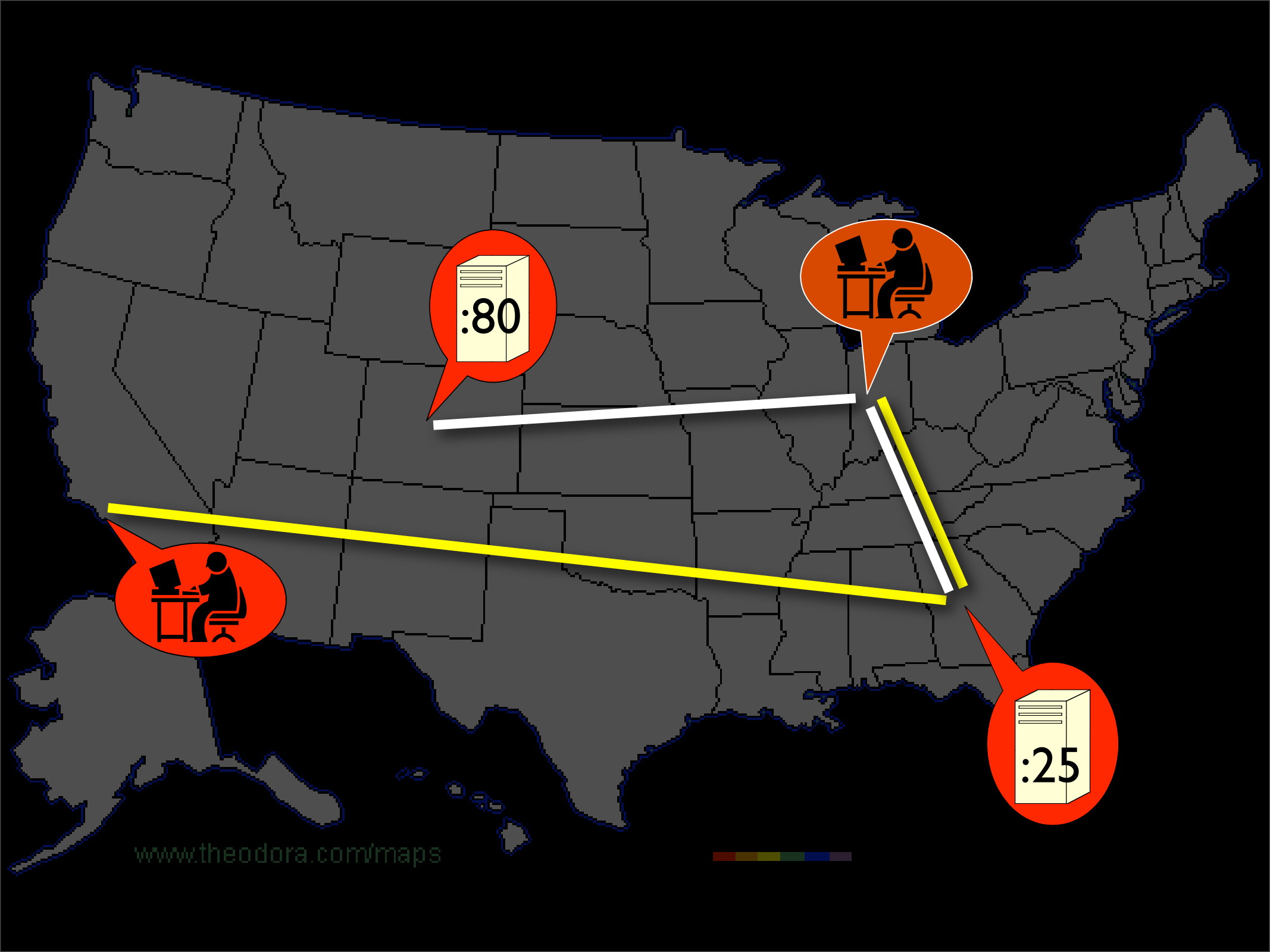


“The Internet”?



“The Internet”?

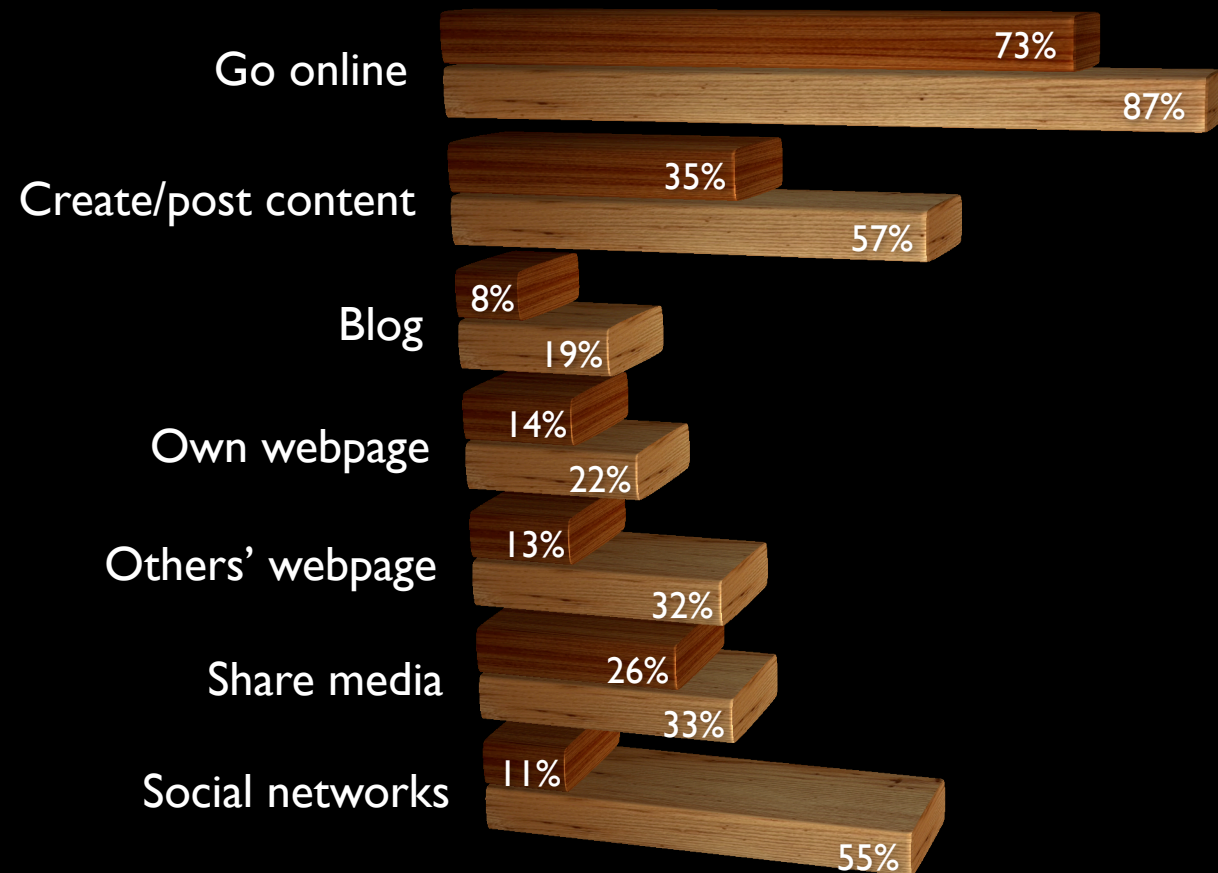




What do we do online?

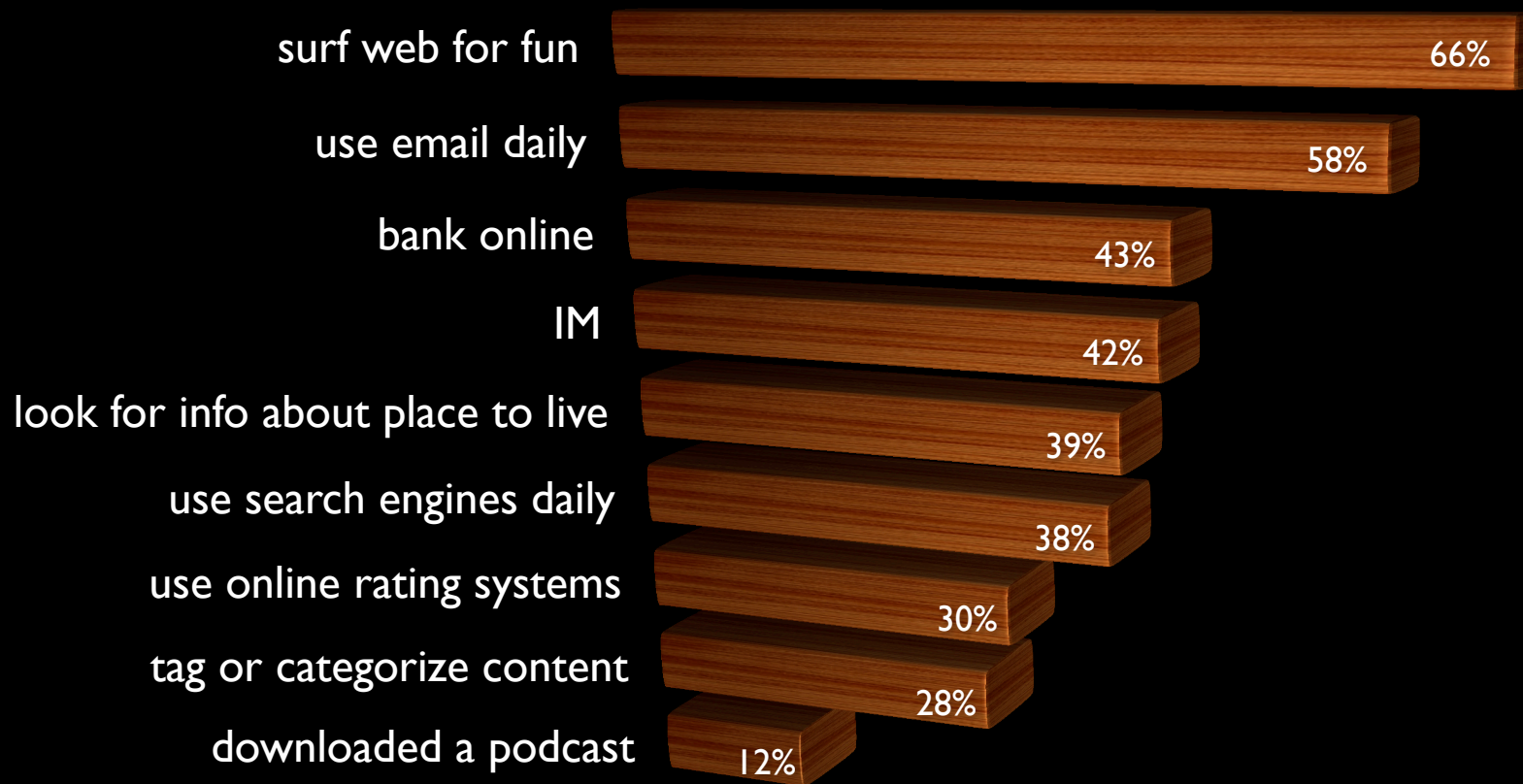
■ Adults ■ Teens

- 60% of home internet users have broadband access
- 34% of internet users have wireless access



Source: PIP (pewinternet.org)

More stats...



Source: PIP (pewinternet.org)

What's missing

- Distributional information, correlations, dependencies
- Who talks to whom, how, how much...?
- How do these activities impact Internet traffic?
- Can we characterize patterns of traffic?
- Can we deduct what people are doing from looking at the traffic?
 - Covert activities?
 - Without looking at IP addresses, or inspecting payloads, or looking at all packets?

Traffic networks

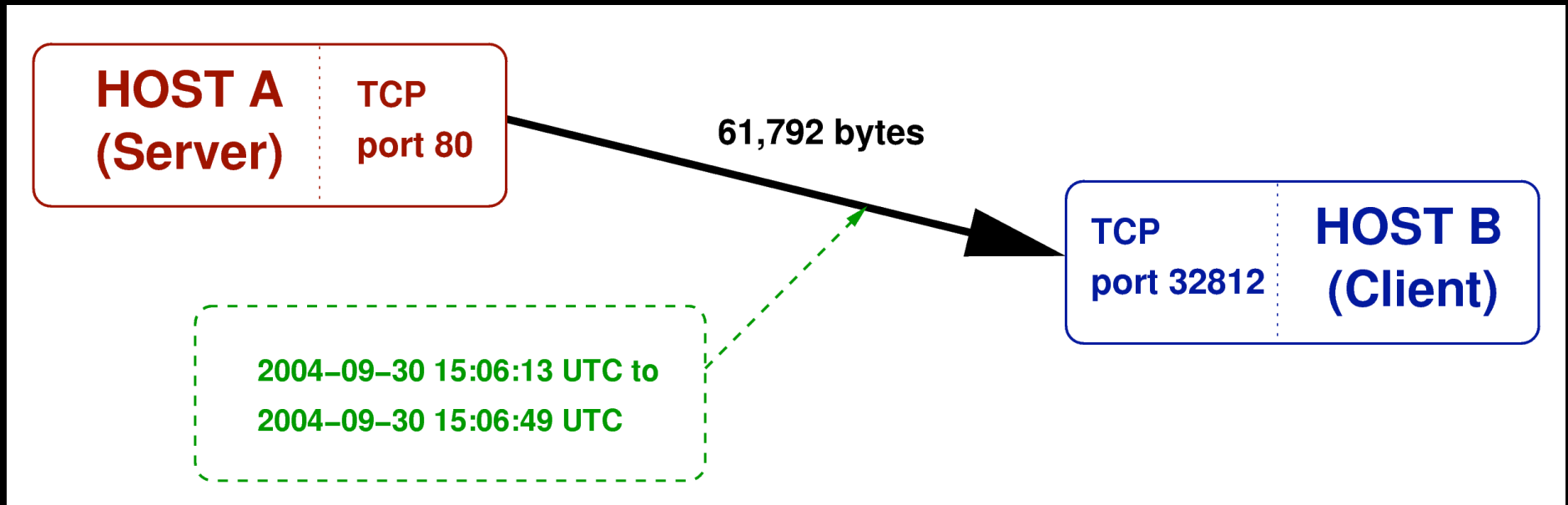
- Build networks from Internet traffic data
- Two proofs of concept:
 - Behavioral networks (host to host)
 - Application networks (app to app)

- TCP/IP network connecting research and educational institutions in the U.S.
- Over 200 universities and corporate research labs
- Hundreds of thousands of undergraduates
- Also provides transit service between Pacific Rim and European networks
- High capacity (never congested)
- Now peered with commodity Internet

Internet2/Abilene

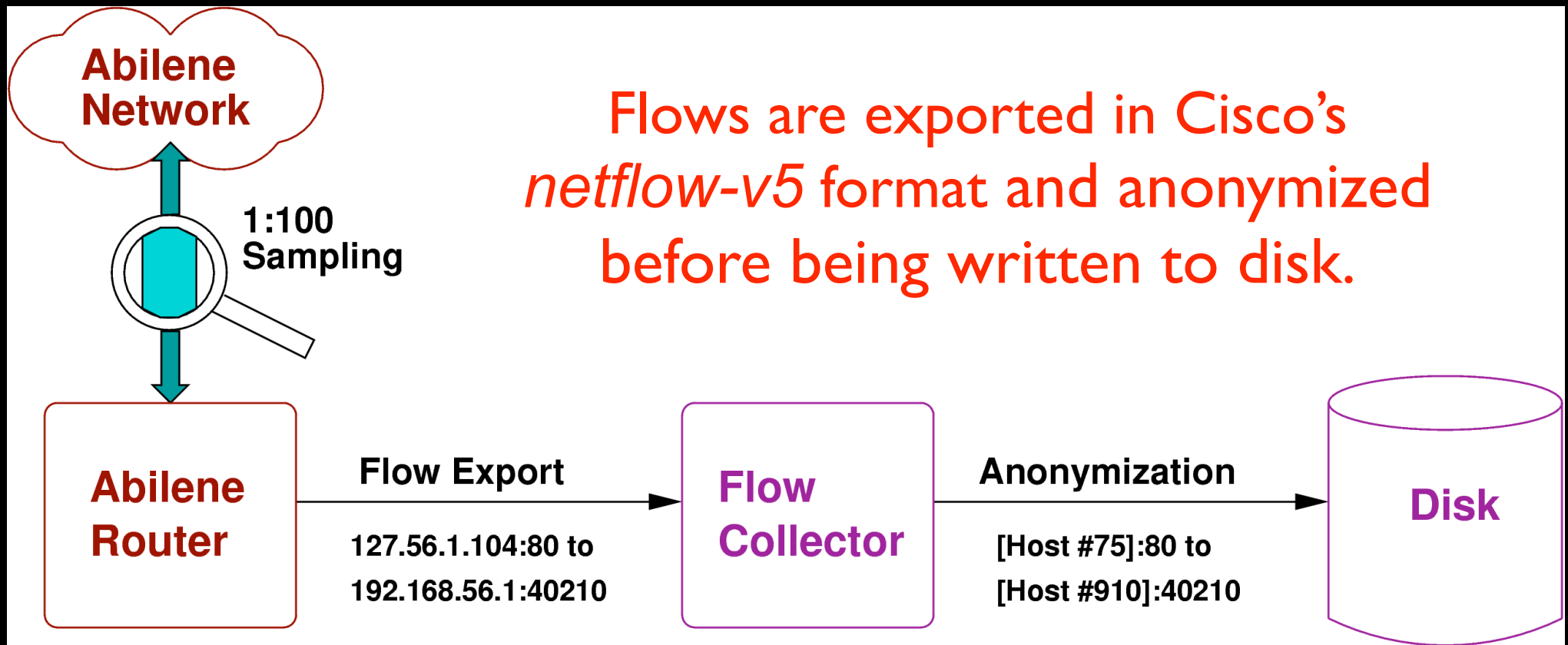


Network flow data



A successful TCP session contains ***two*** flows

Flow collection

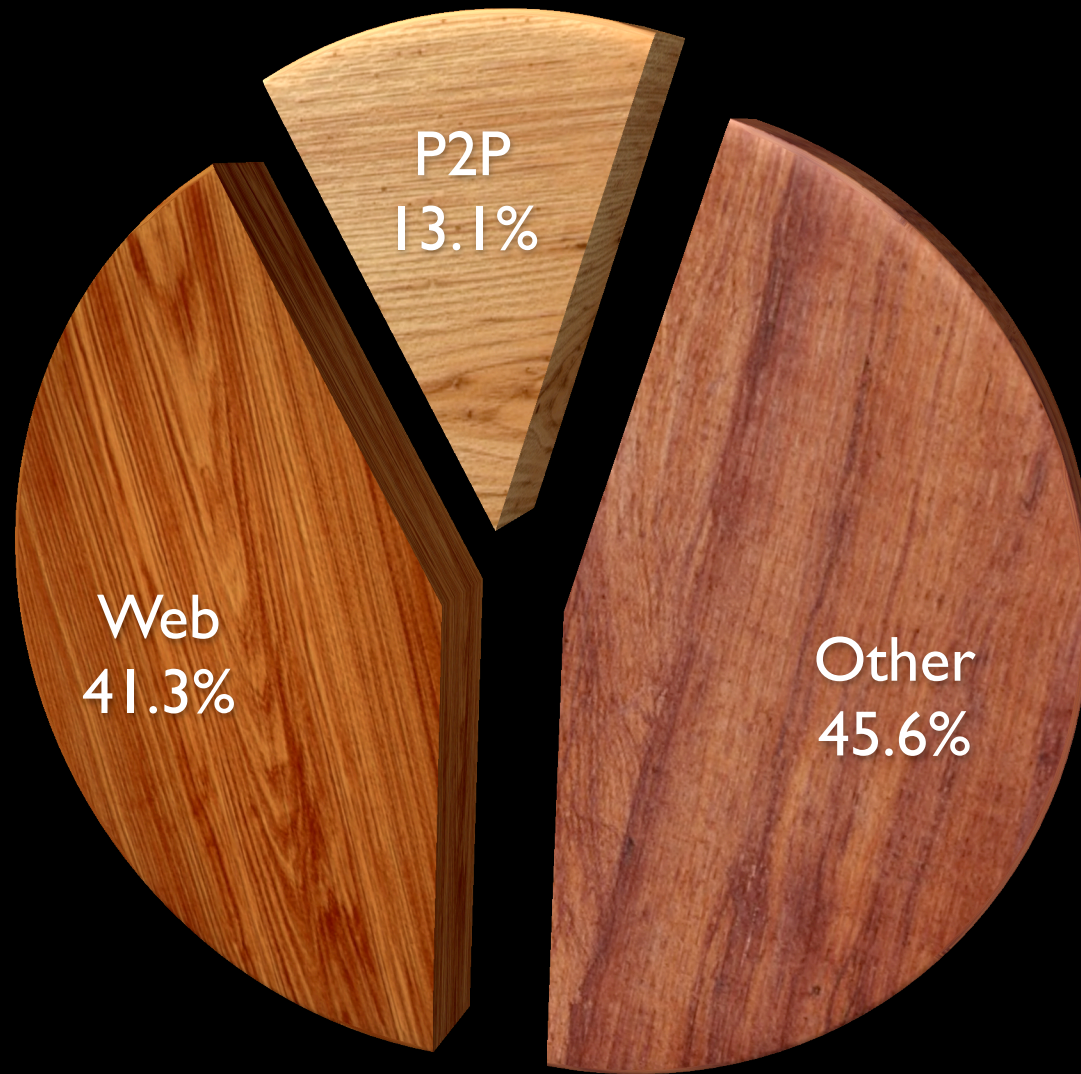


In a typical day: 35GB raw data at 3.1 Mbps

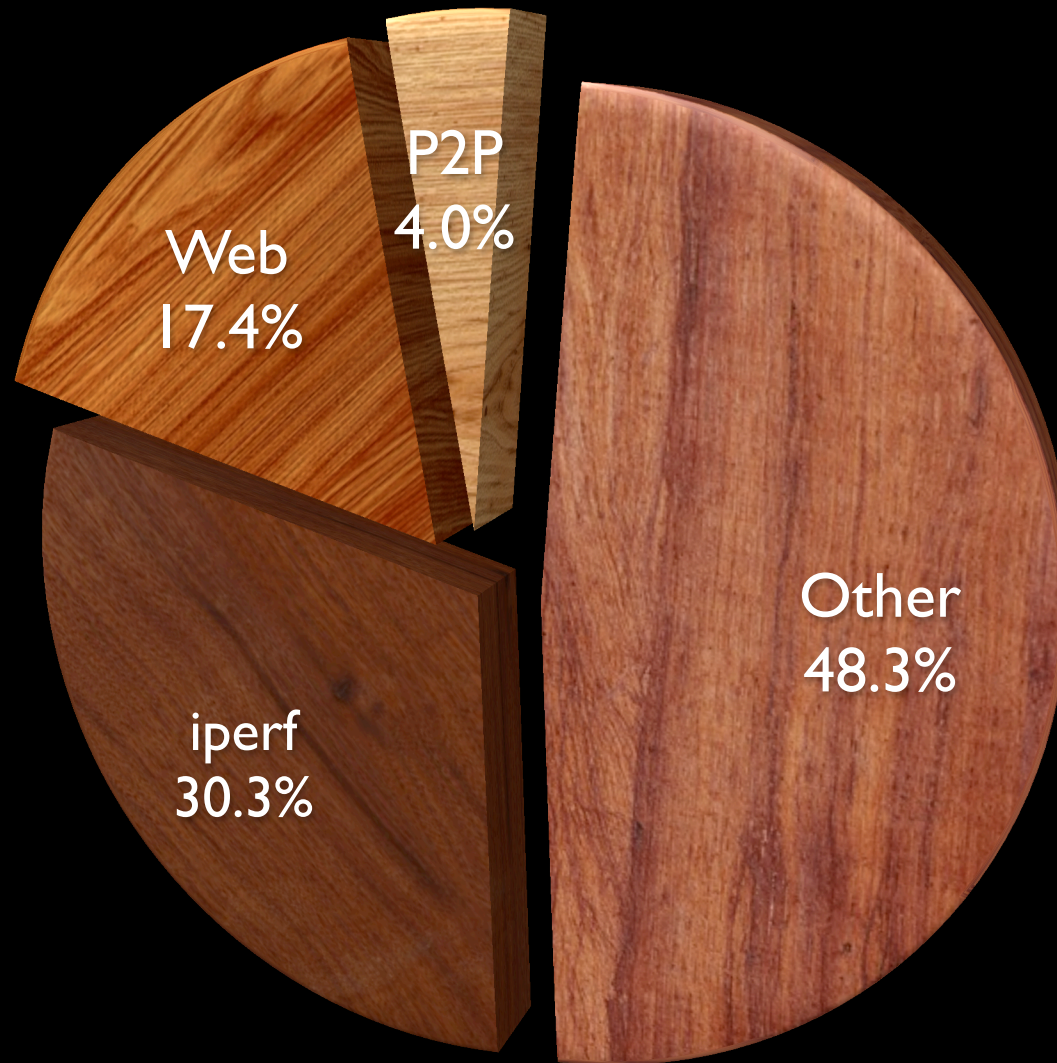
Data set for analysis

- Full 24-hour day of network flow data starting at 2005-04-14 05:00:00 UTC
 - a typical day
 - 600M flows
 - 15M unique hosts
 - 2TB (1% sampling)

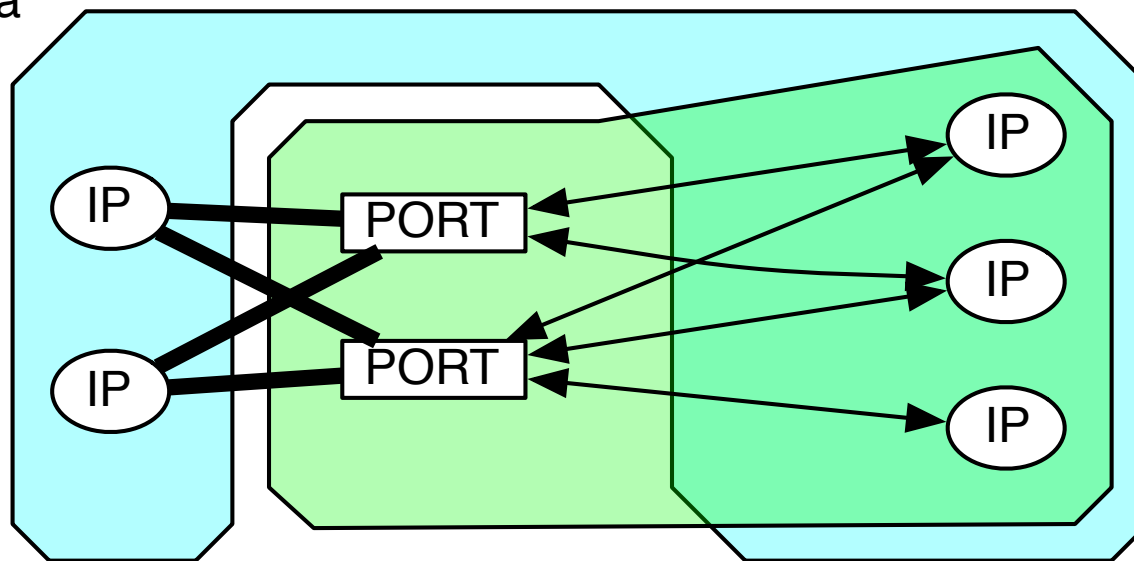
Flows



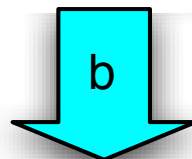
Traffic



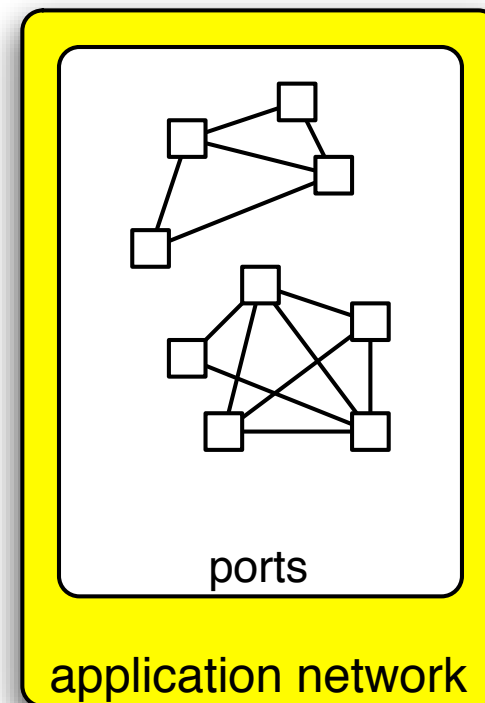
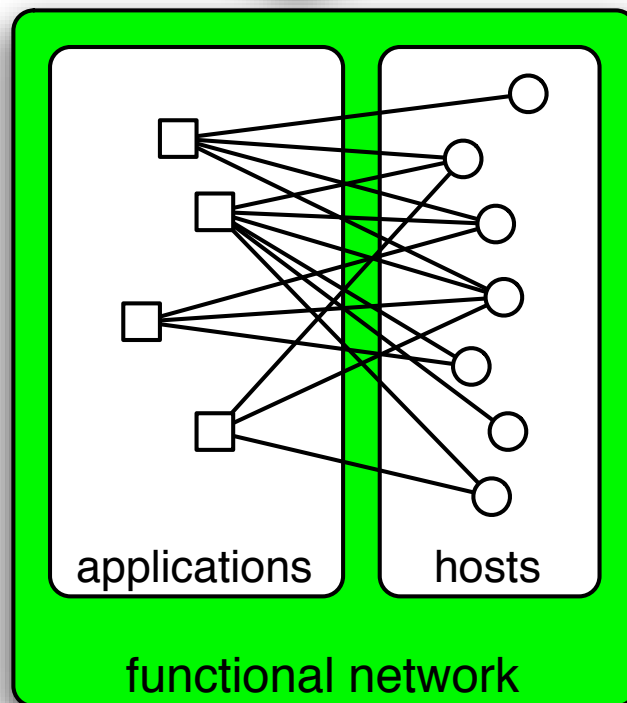
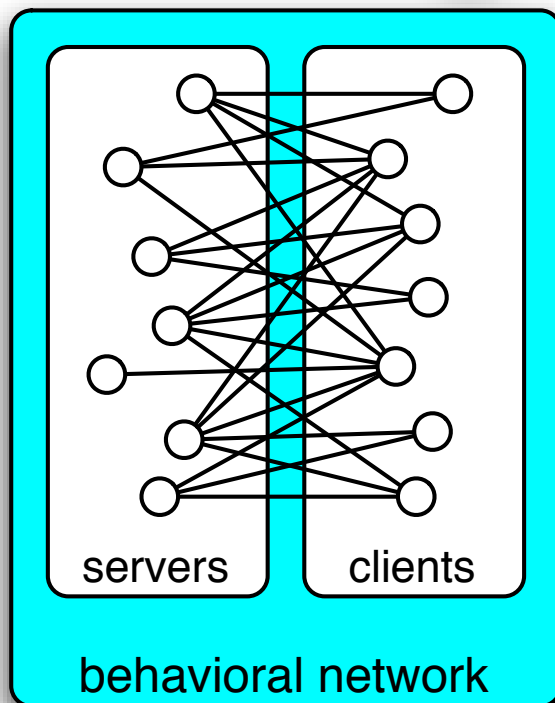
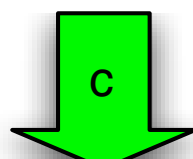
a



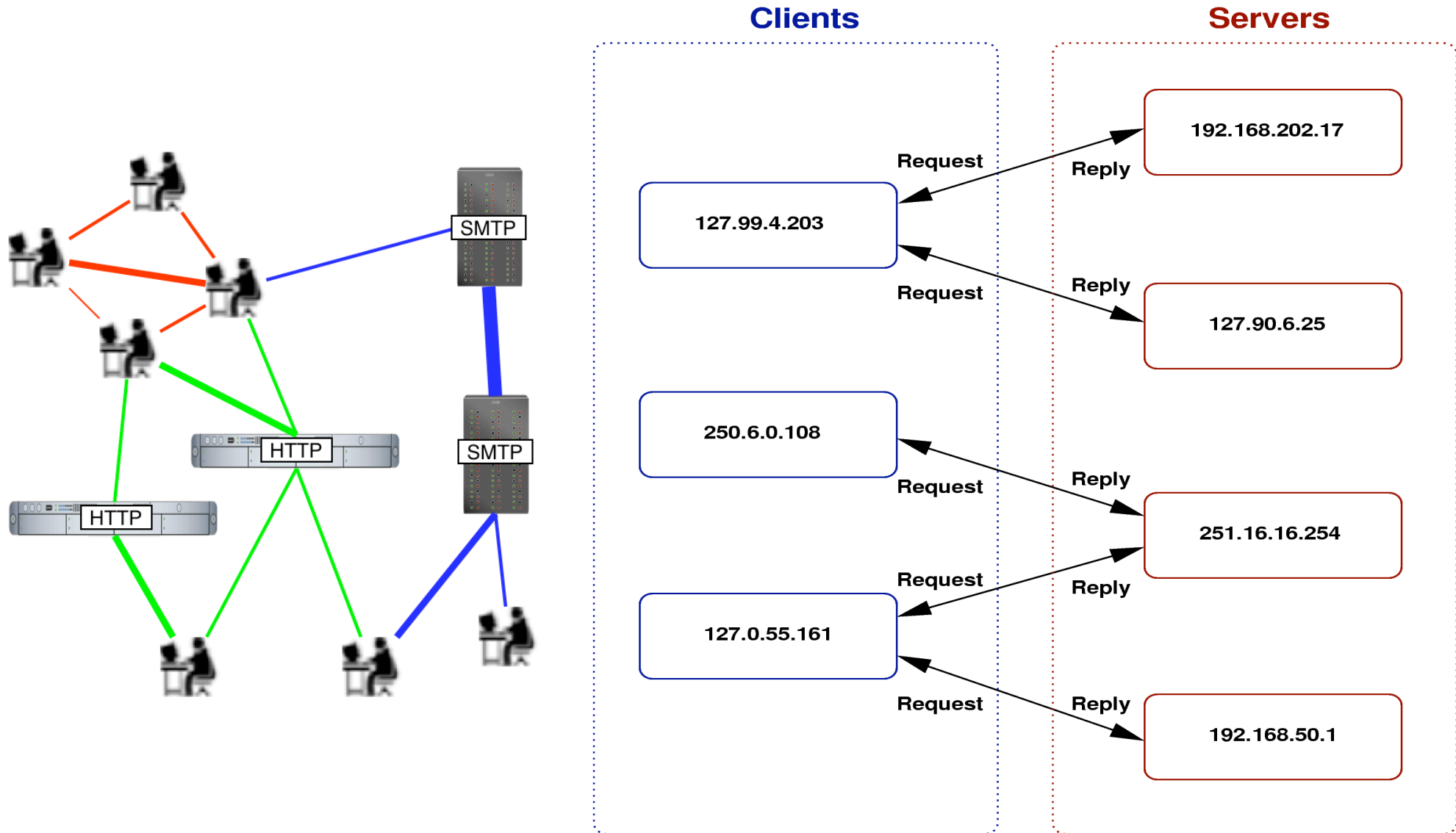
b



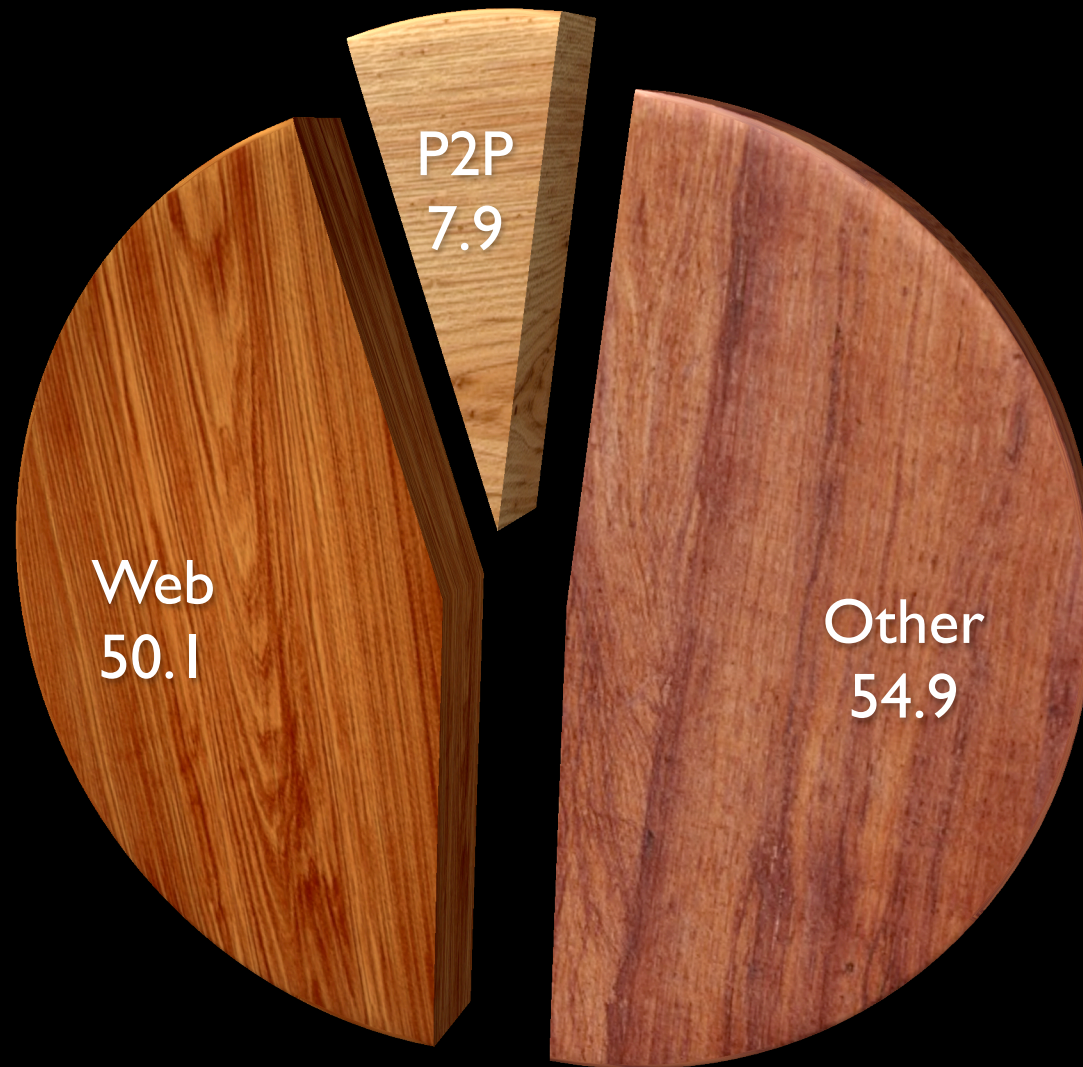
c



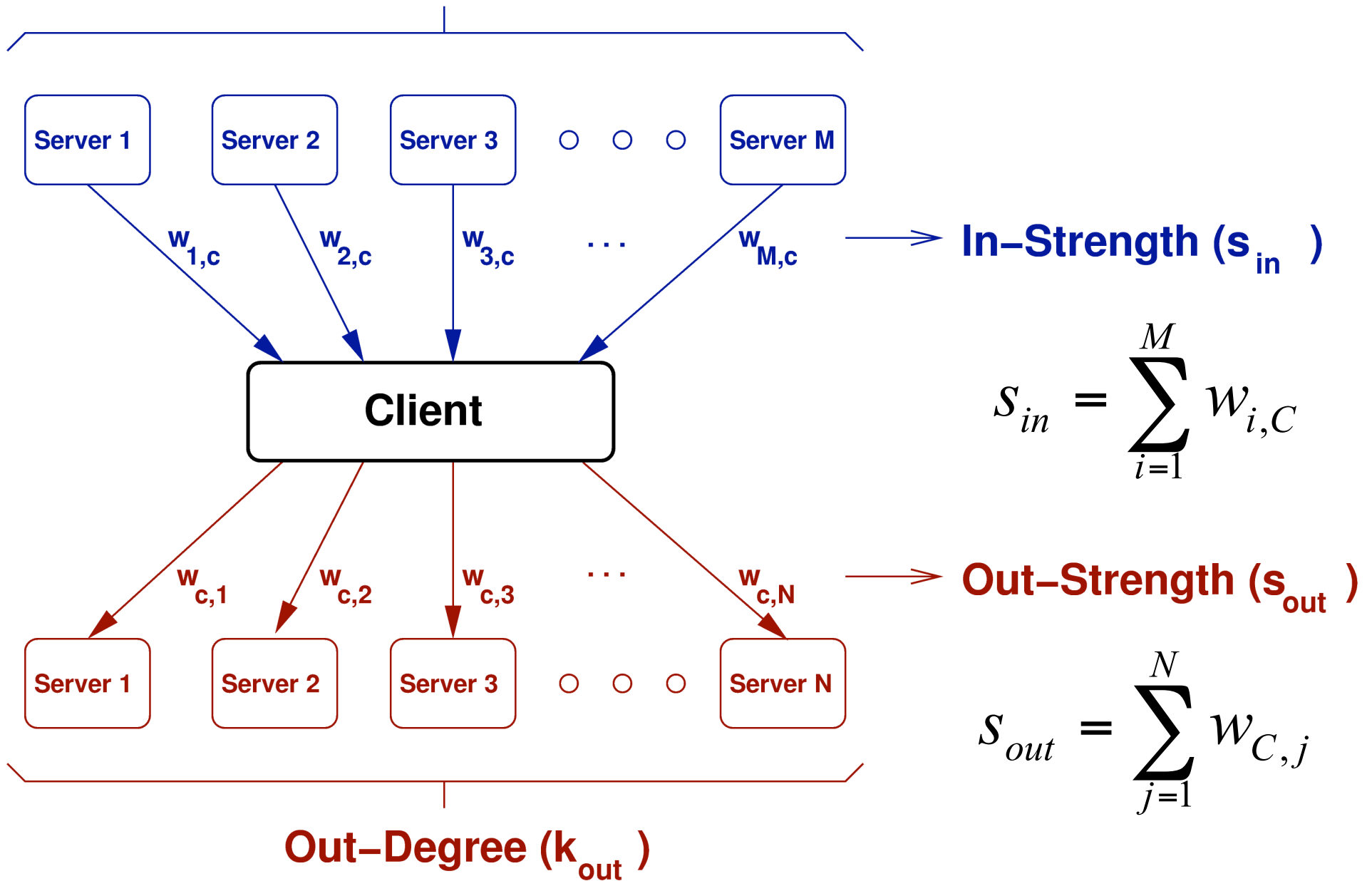
Behavioral networks



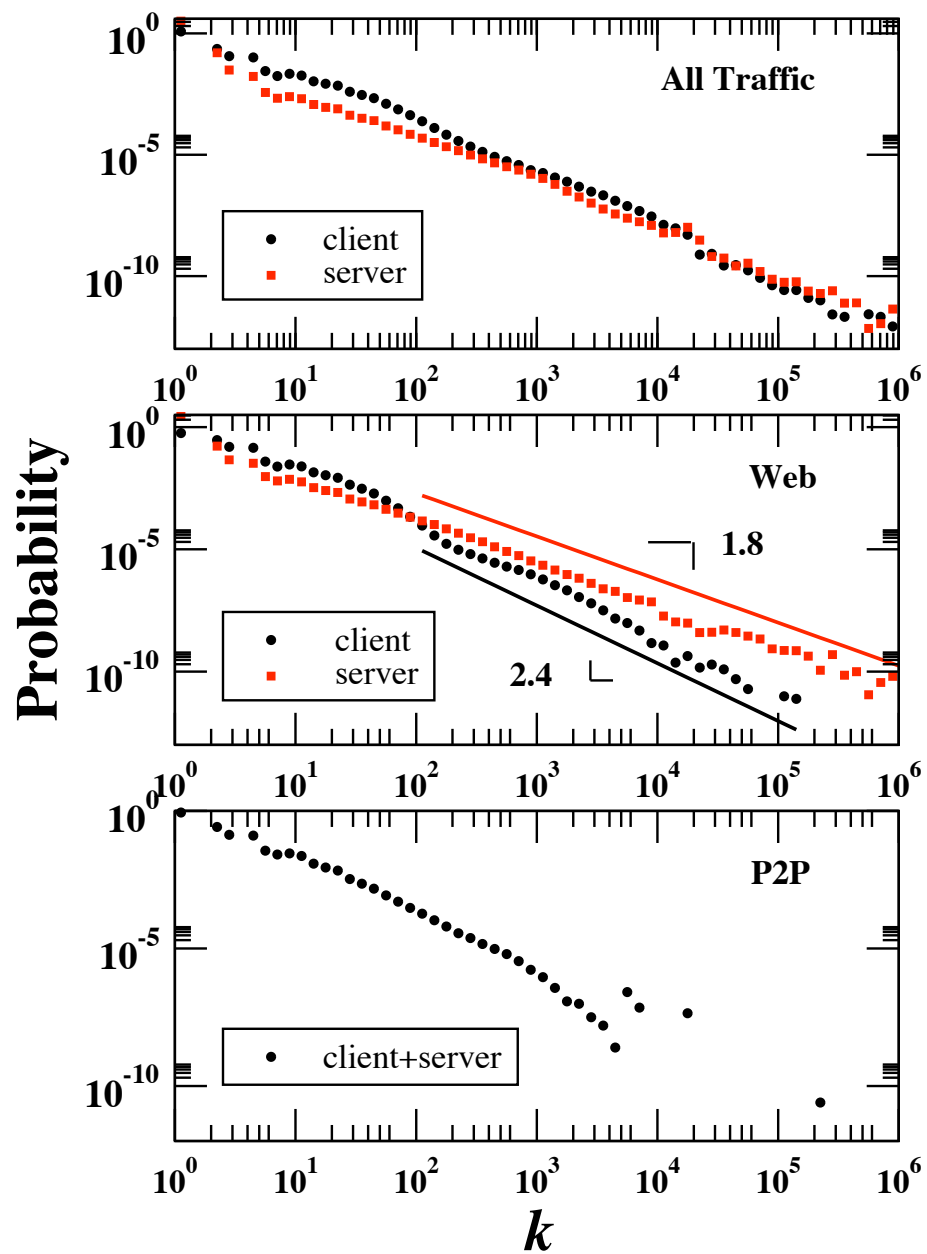
Edges ($\times 10^6$)



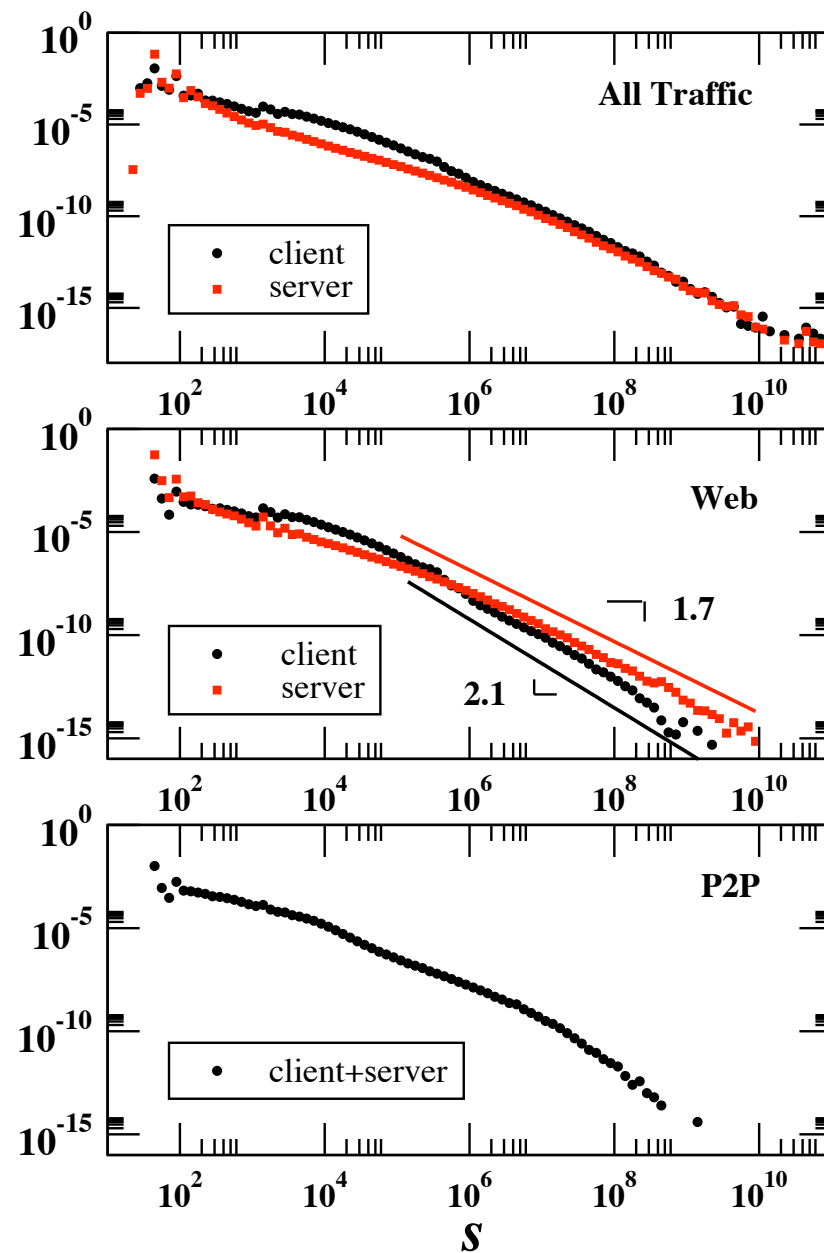
In-Degree (k_{in})



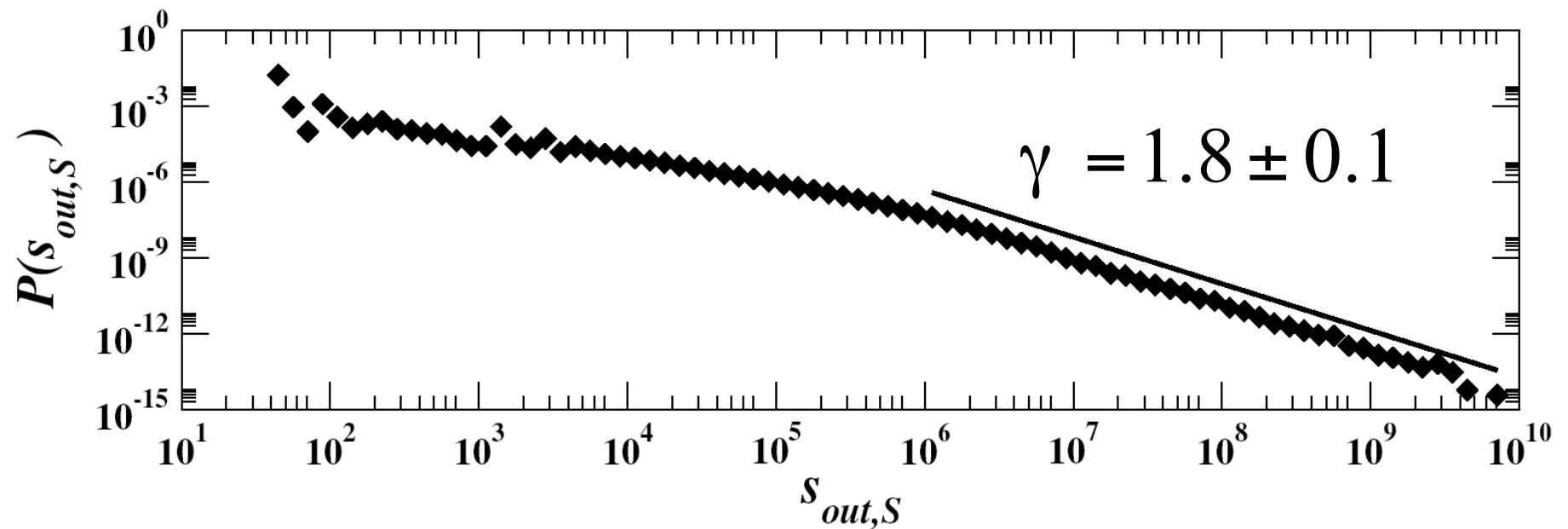
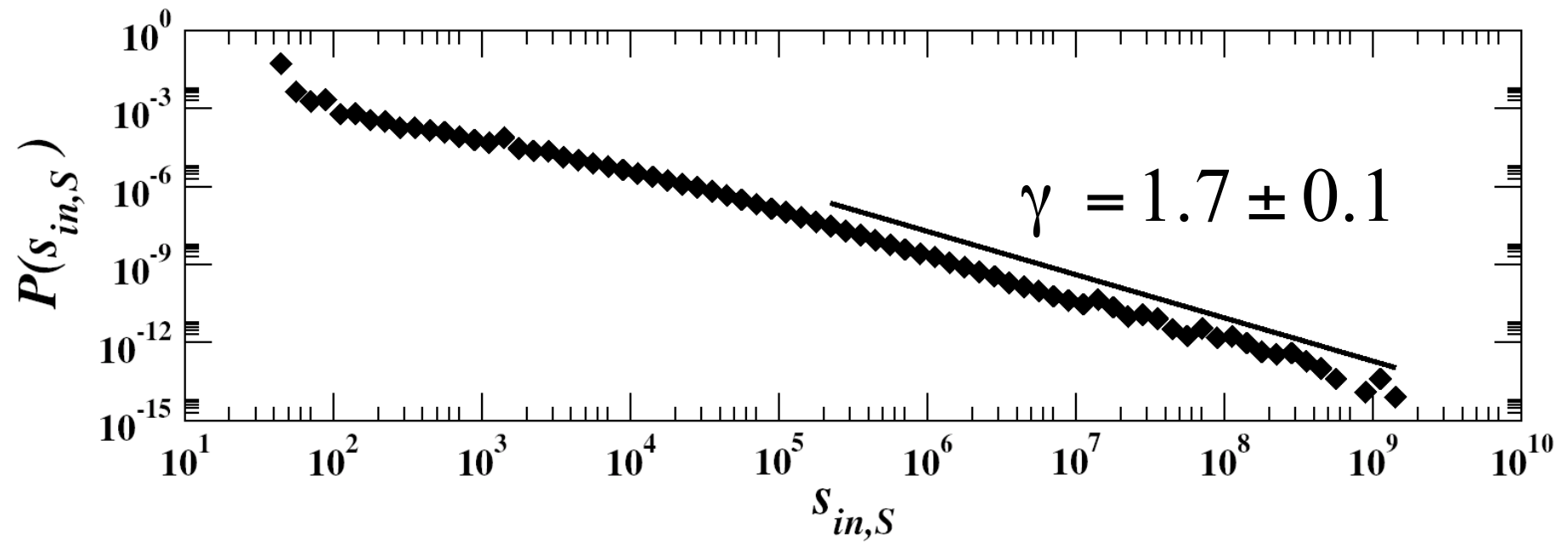
Degree

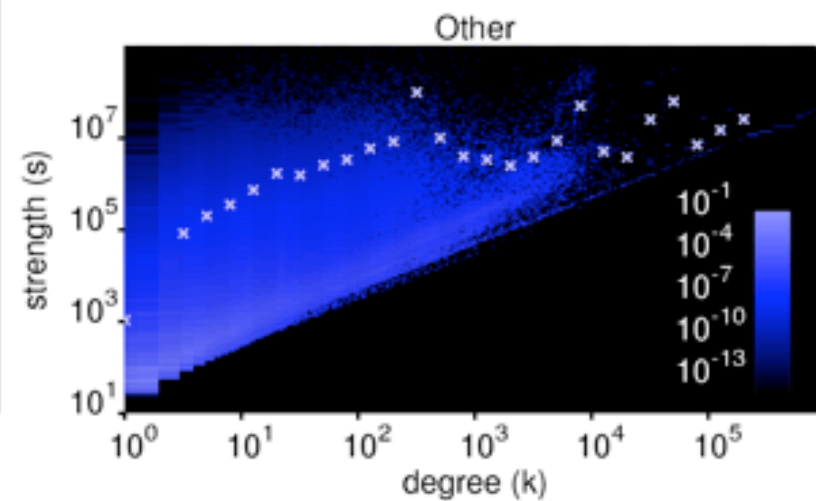
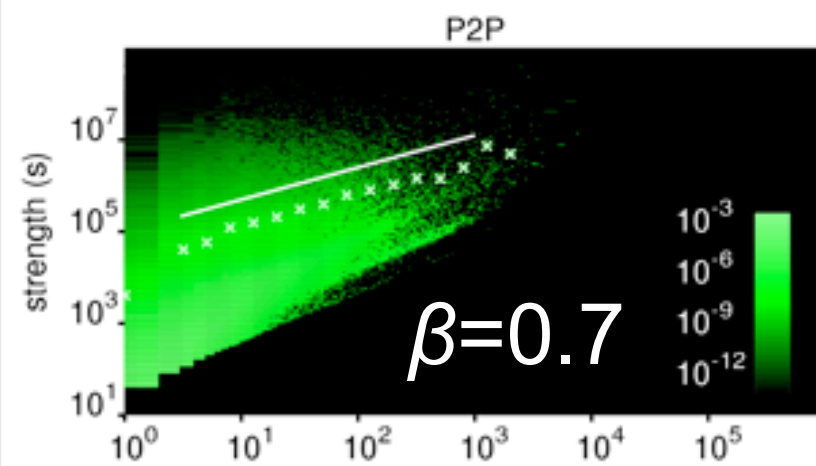
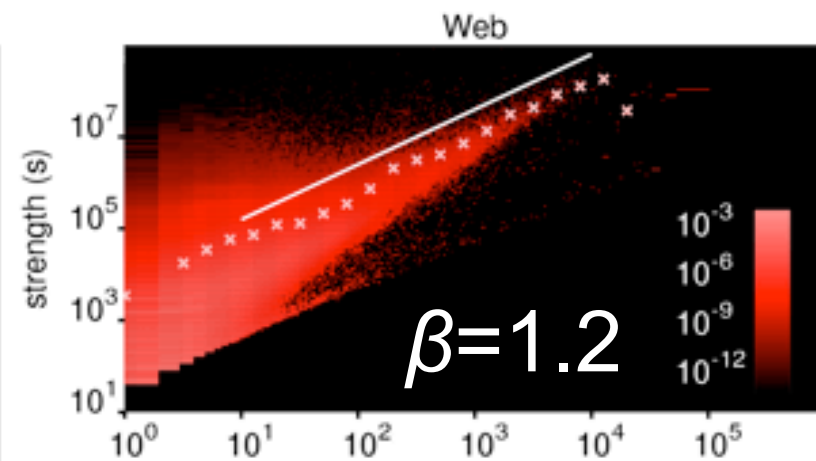
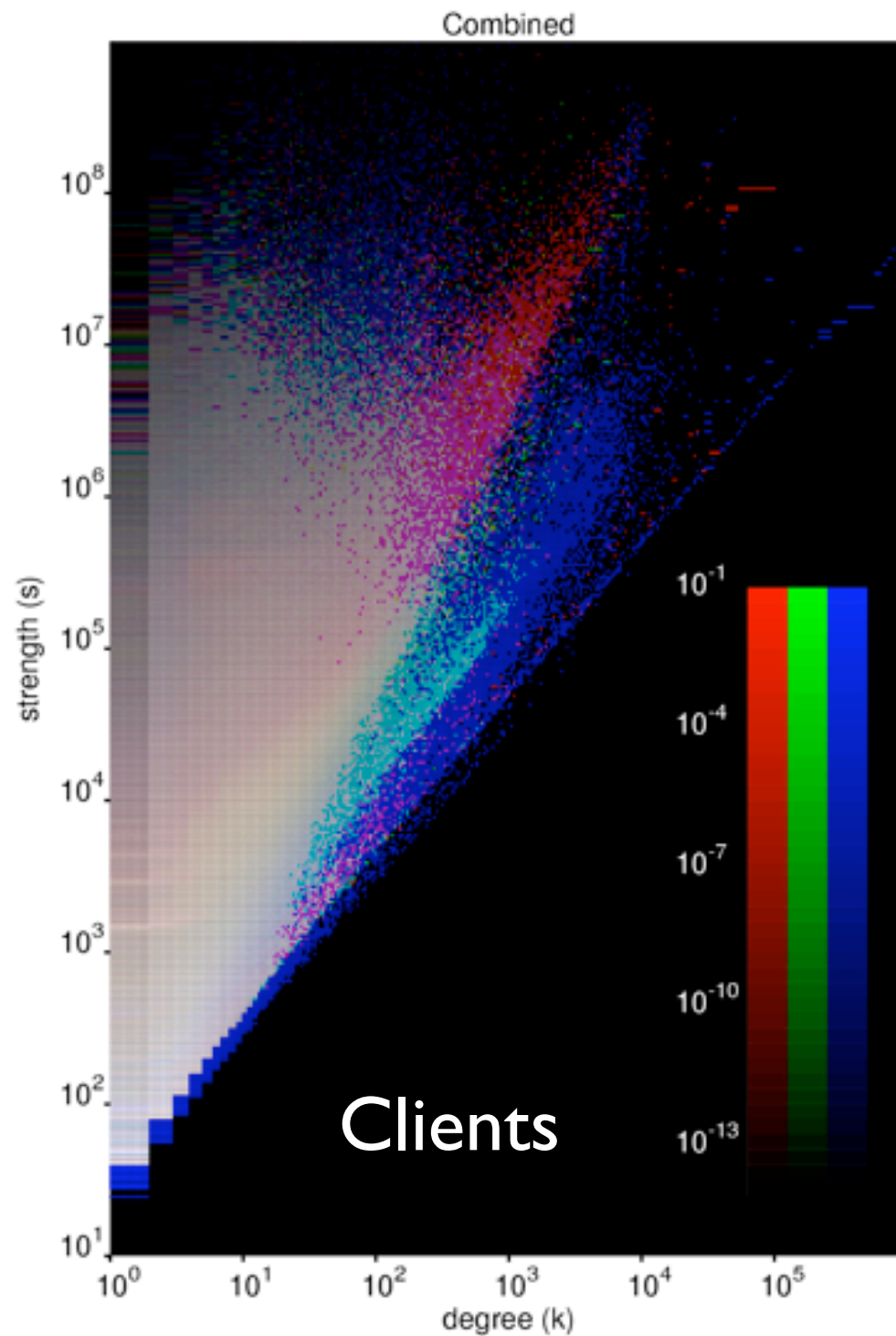


Strength

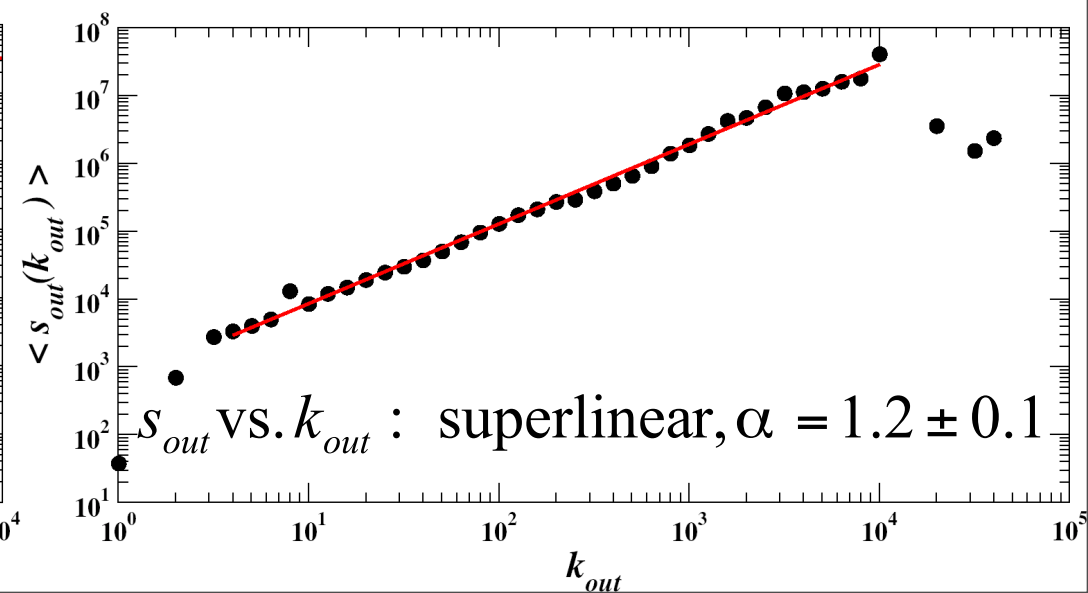
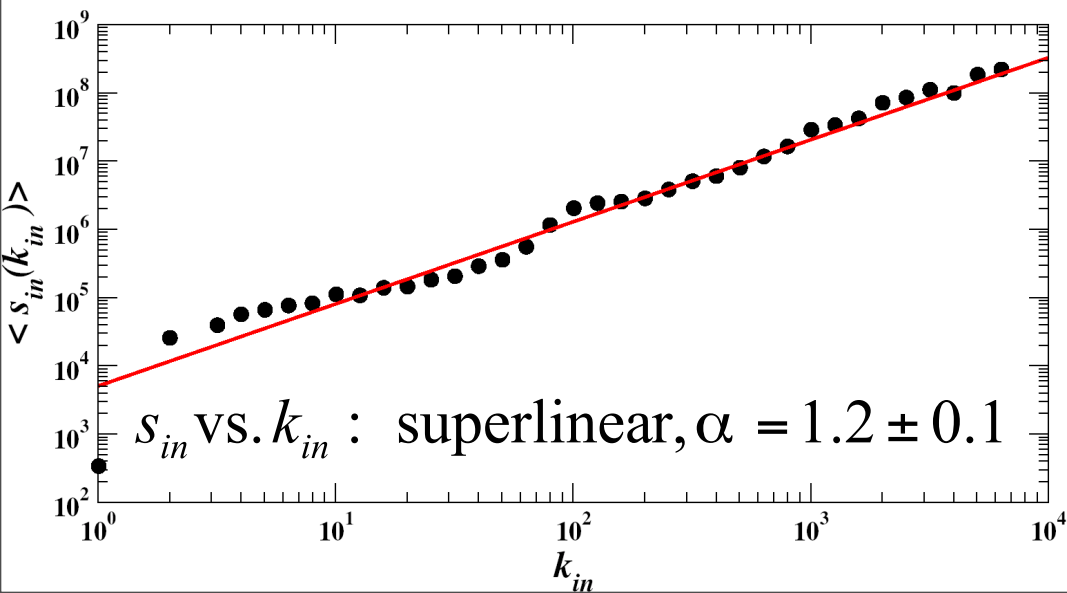
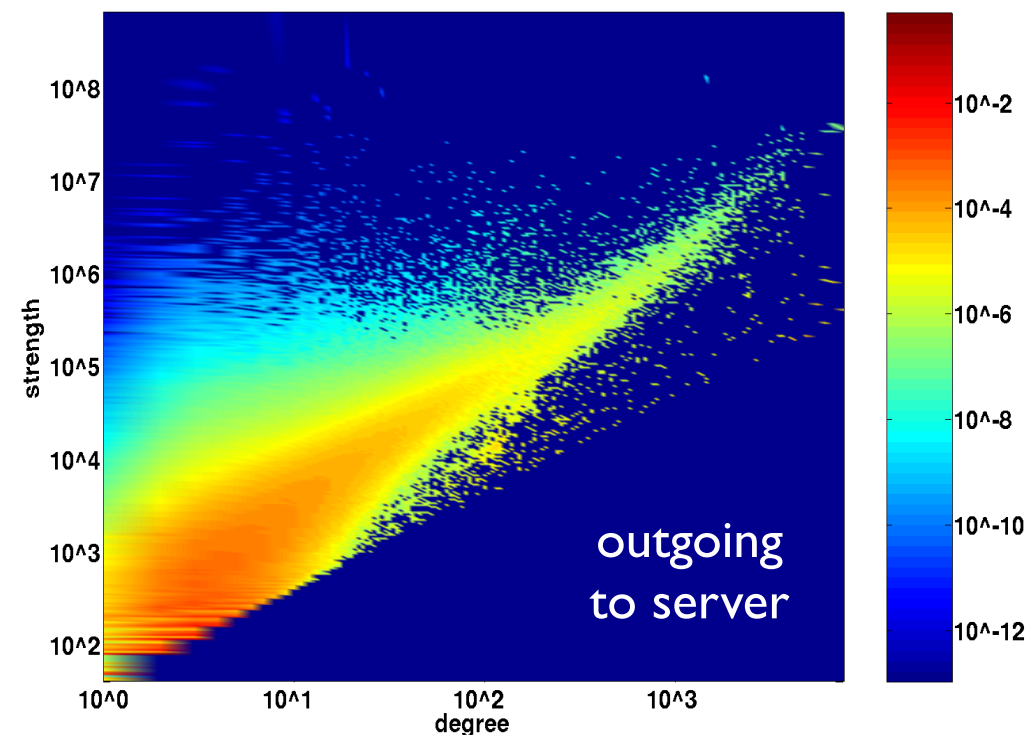
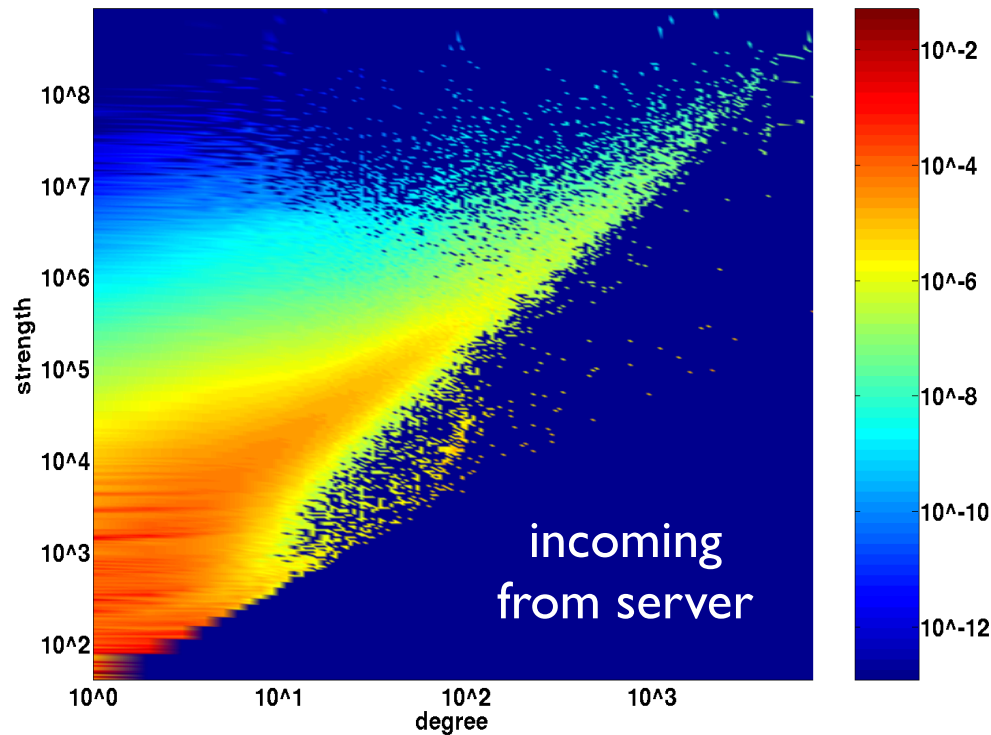


Web servers: Strength distributions





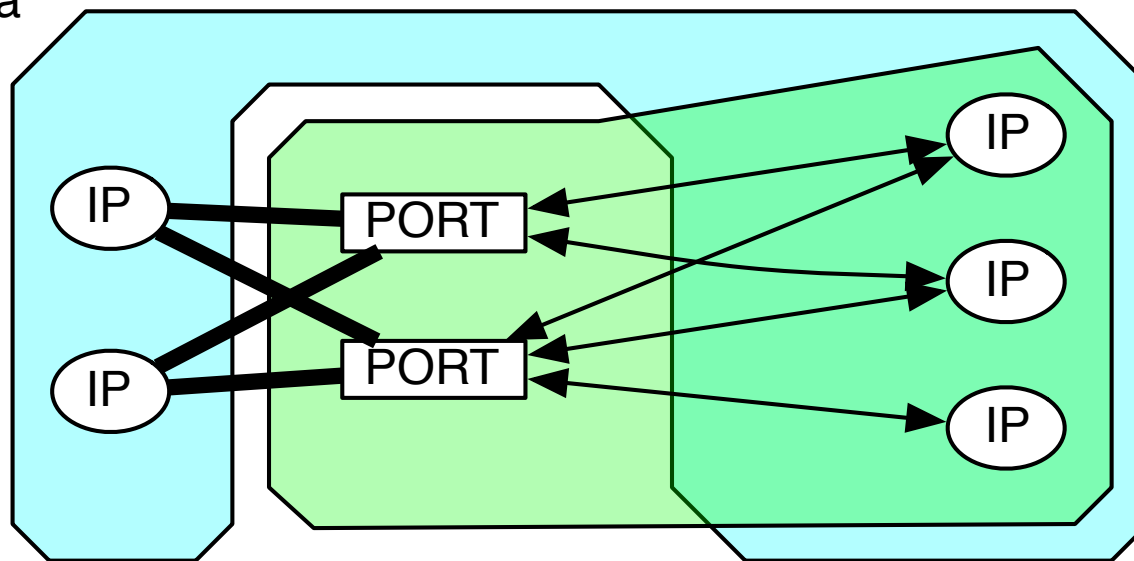
Web scalability problem: The more servers a client contacts, the more data it exchanges with **each** server!



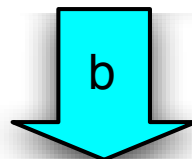
Summary

- **Power-law distributions** are found in all aspects of traffic networks: degree, strength, and weight distributions for both clients and servers
 - The **strength** distribution for Web **servers** lacks any **mean value**
 - The relationship between **degree** and **strength** for Web clients is **super-linear**
 - **Models** must be able to account for these heavy-tailed distributions and non-linear coupling
 - Classes of traffic can be characterized by their strength-degree **signatures**

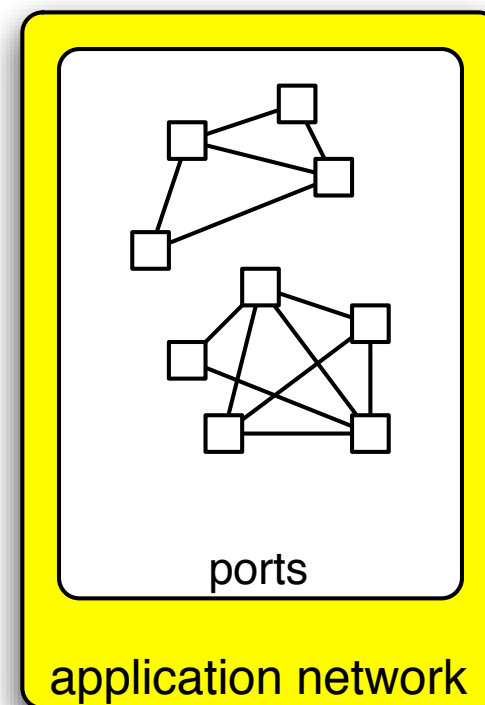
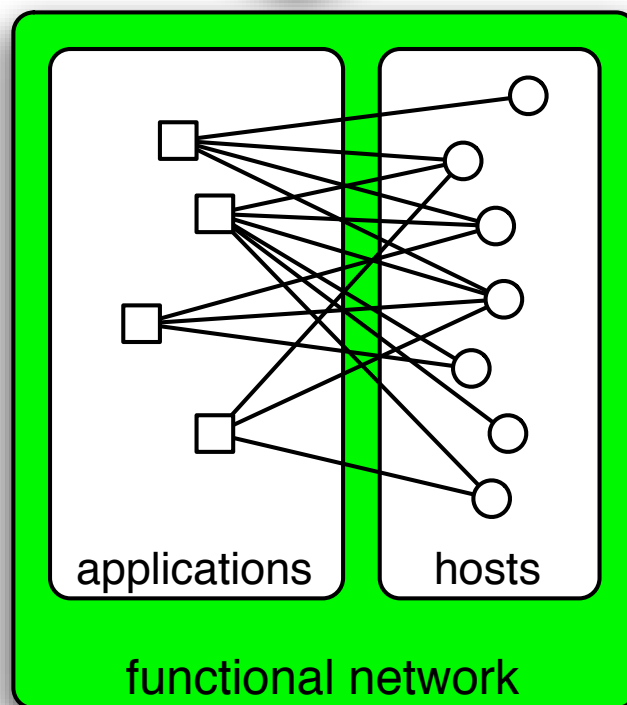
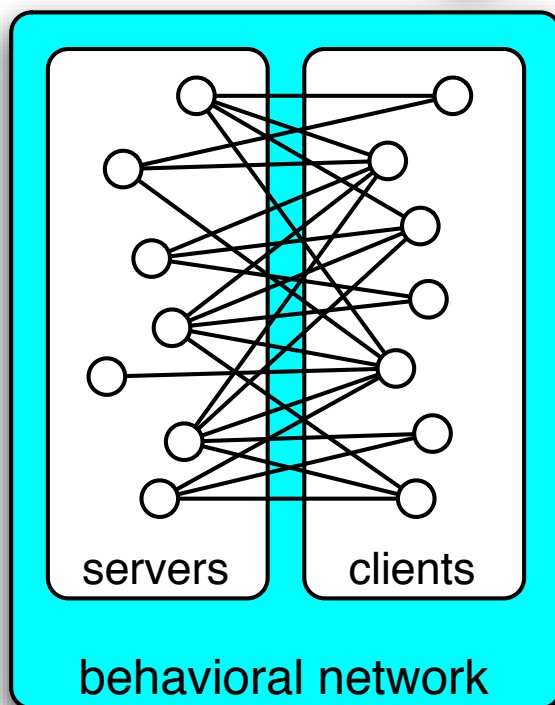
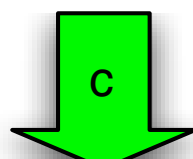
a



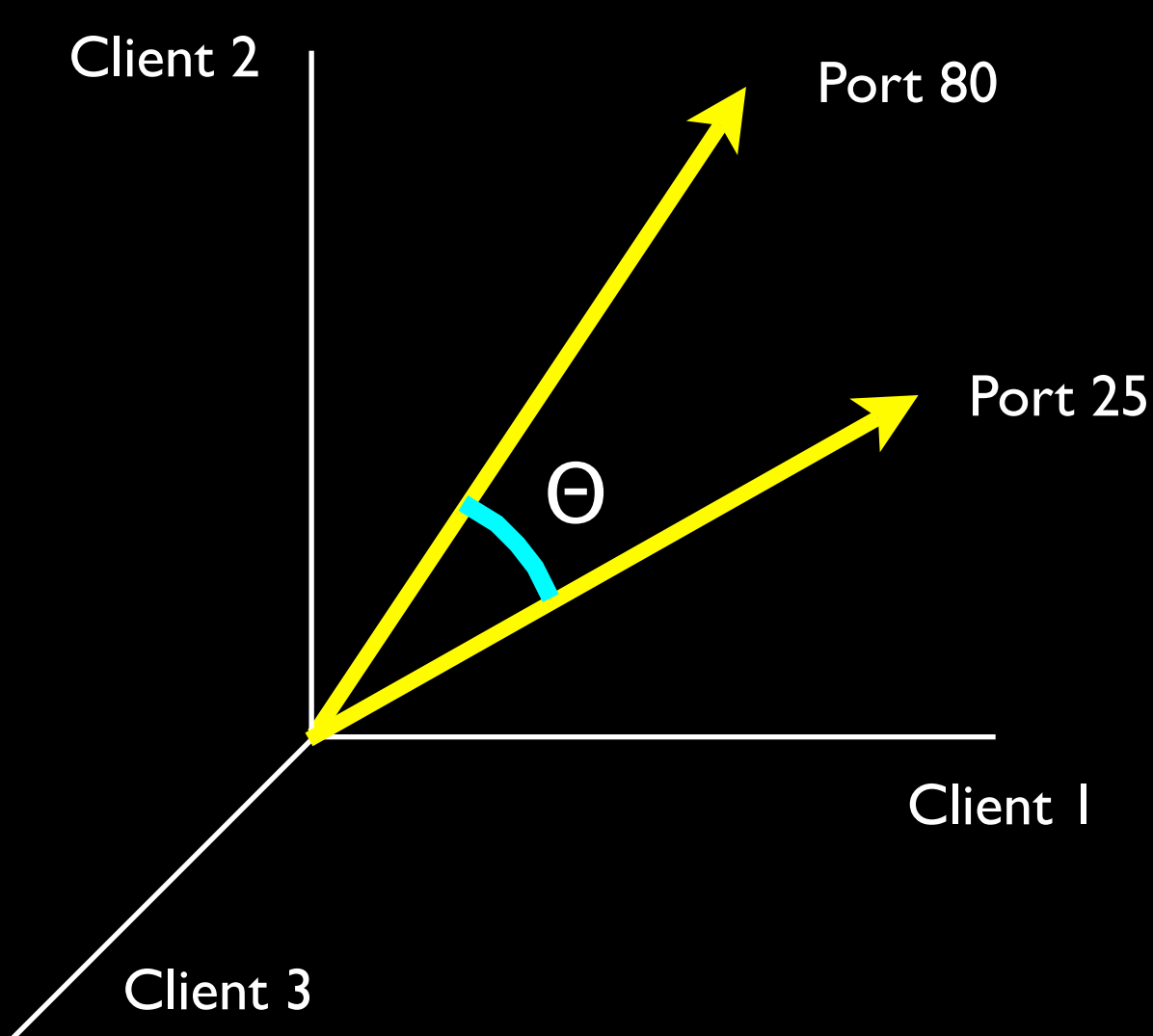
b



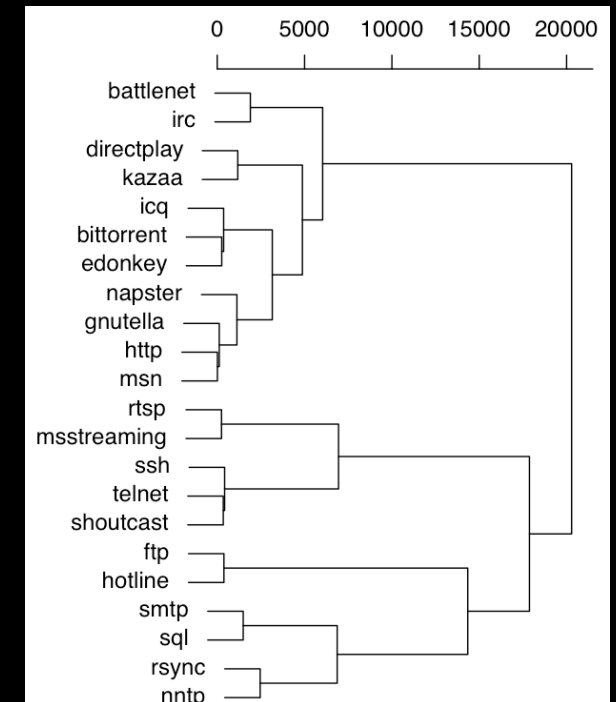
c

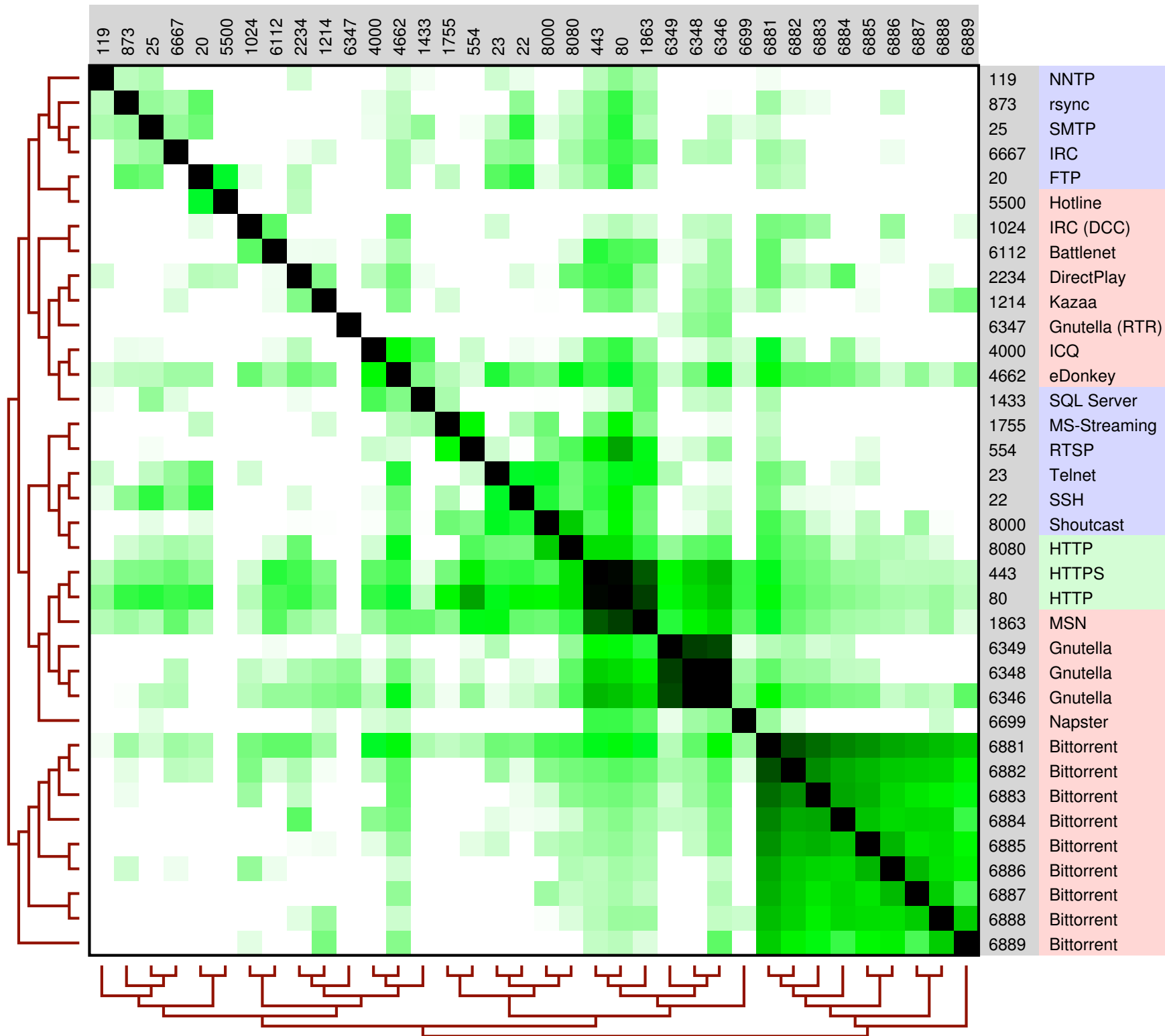


From functional to application clusters



$$d = \frac{1}{\cos \theta} - 1$$





Unknown port	Application
388	Unidata/LDM
19101	Clubbox

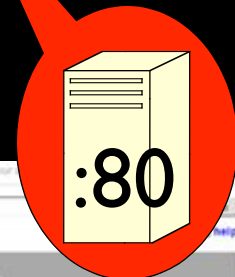
Port	App	Match?
388	weather data transfer	yes
19101	individual file shares	yes
9080	team collaboration	yes
8090	Weblog server	yes
5020	BBFTP file transfer	partial
42899	<i>unknown</i>	?
8301	several trojans	partial
1025	many different trojans	yes
20000	BitTorrent	yes
59174	<i>unknown</i>	?
20001	several trojans	partial
15002	biology collaboration tool	partial
16881	BitTorrent	yes
9000	several trojans	partial
3124	Web proxy (VWindows)	yes
39281	grid-based computing	partial

Summary

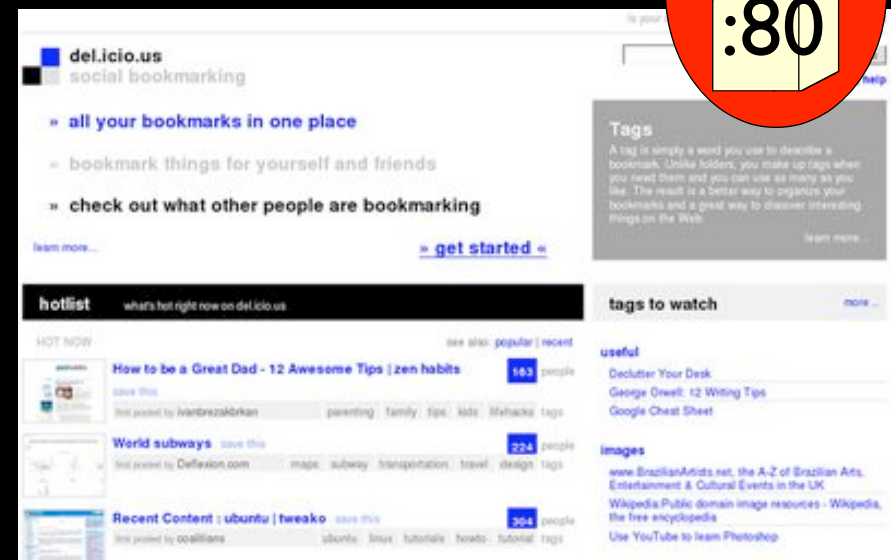
- Use application networks to guess the **function** of **unknown or covert** ports from topological information alone
 - No need to look at payload — preserve **privacy**
 - No need to inspect all packets — **efficient** (possibly real-time)
 - Detect potentially **malicious** traffic



Clicks?



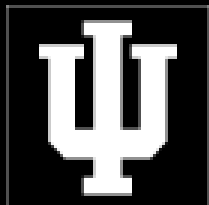
Not all dynamics are captured



Future applications

- Anomaly detection
- Application fingerprinting
- Classification of Web clients according to their purpose: *browsers, crawlers, scanners, proxies*, etc.
 - This may provide insight into scalable design
- Identify unknown or covert activities in real time
- Using traffic data to improve the performance of search engines

Thanks



Indiana University School of
informatics

