# Application Overlays (P2P); Network Transport Layer: Overview; UDP; Stop-and-Wait ARQ

**Qiao Xiang**, Congming Gao

https://sngroup.org.cn/courses/cnns-xmuf23/index.shtml
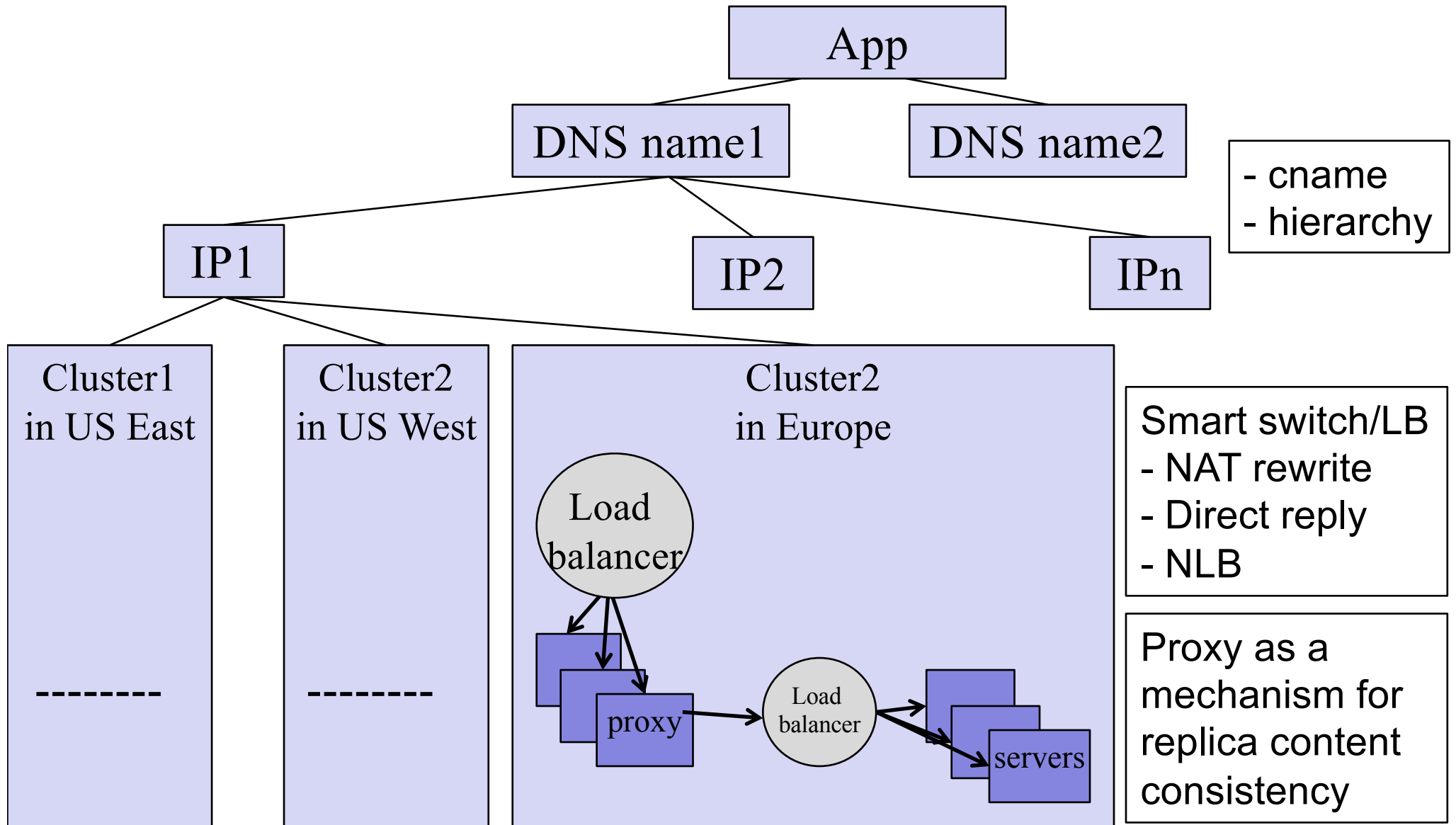
10/31/2023

# Outline

- ❑ Admin and recap
- ❑ Application overlays
- ❑ Overview of transport layer
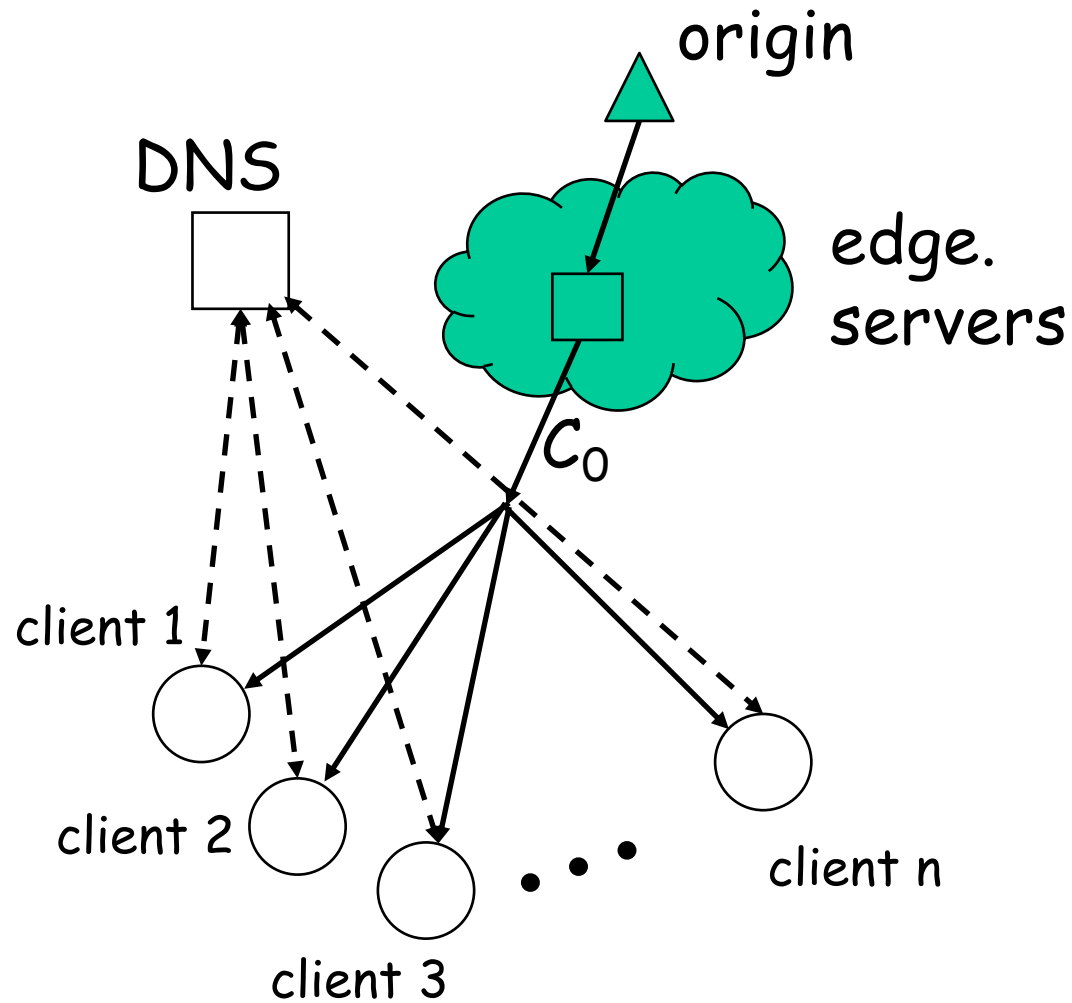- ❑ UDP
- ❑ Reliable data transfer, the stop-and-go protocols

# Admin

❑ Lab assignment 3 due on Nov. 19

❑ Midterm exam on Nov. 9 (during lab class)
  - cover from introduction to application layer
  - 15-16 subjective questions over 100 minutes
  - 1-page cheat sheet allowed
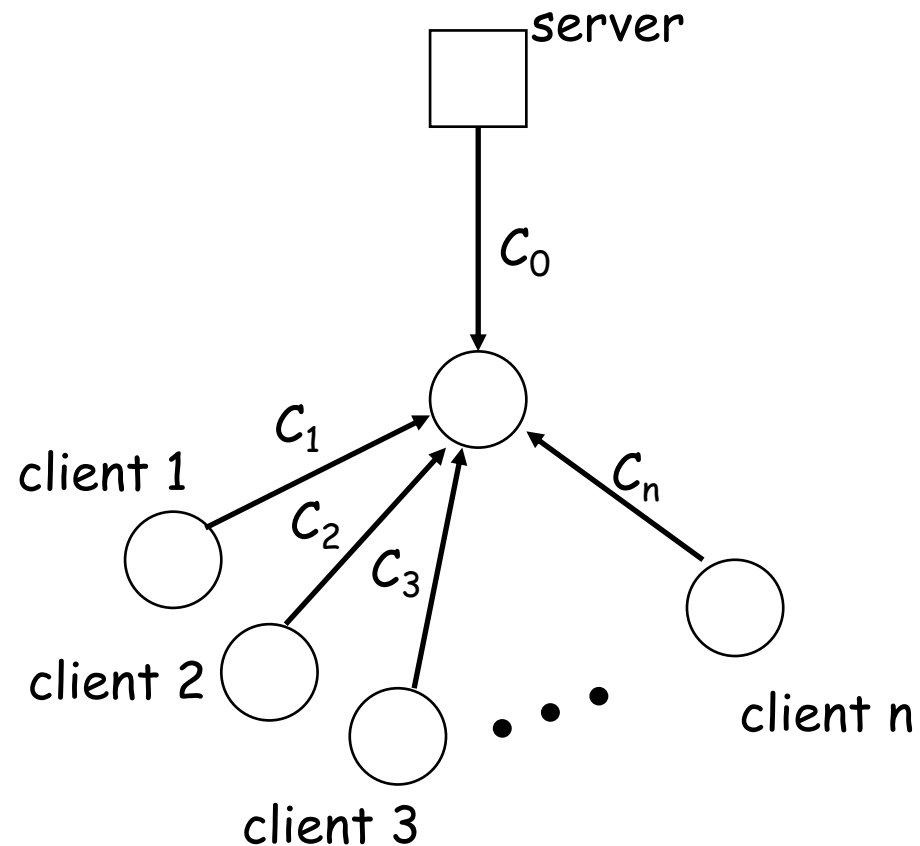
# Recap: Direction Mechanisms

App

DNS name1     DNS name2

- cname
- hierarchy

IP1     IP2     IPn

Cluster1
in US East

Cluster2
in US West

Cluster2
in Europe

Load balancer

proxy

Load balancer

servers

-------

-------

Smart switch/LB
- NAT rewrite
- Direct reply
- NLB

Proxy as a mechanism for replica content consistency

# Scalability of Server-Only Approaches



origin

DNS

edge. servers

$C_0$

client 1

client 2

client 3

client n

# An Upper Bound on Scalability

□ Idea: use resources from both clients and the server

□ Assume

    o need to achieve same rate to all clients

    o only uplinks can be bottlenecks
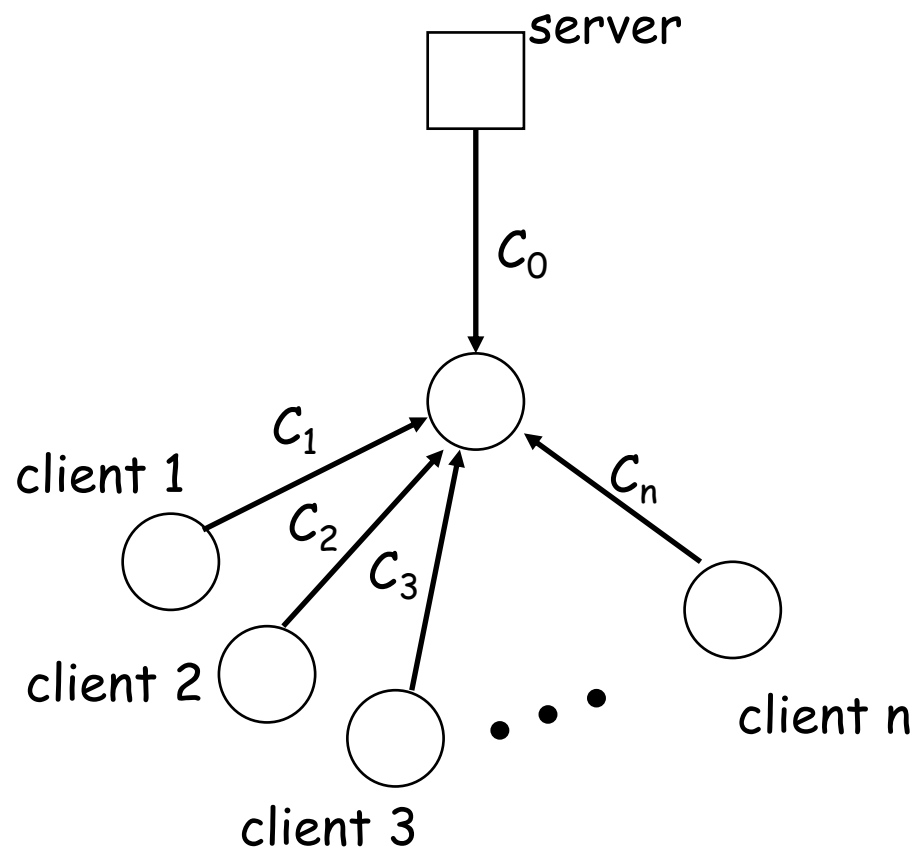
□ What is an upper bound on scalability?

server

$C_0$

client 1   $C_1$

$C_2$

$C_3$

$C_n$

client 2

client 3

client n

# The Scalability Problem

□ Maximum throughput

**R = min{$C_0$, ($C_0$+$\Sigma C_i$)/n}**

□ The bound is theoretically approachable

server

$C_0$

$C_1$

client 1

$C_2$

client 2

$C_3$

client 3

$C_n$

client n
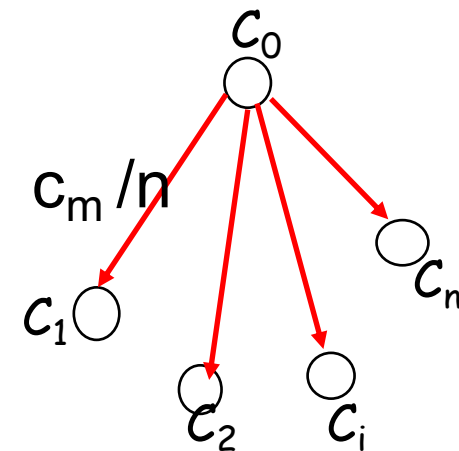
• • •

# Theoretical Capacity: upload is bottleneck

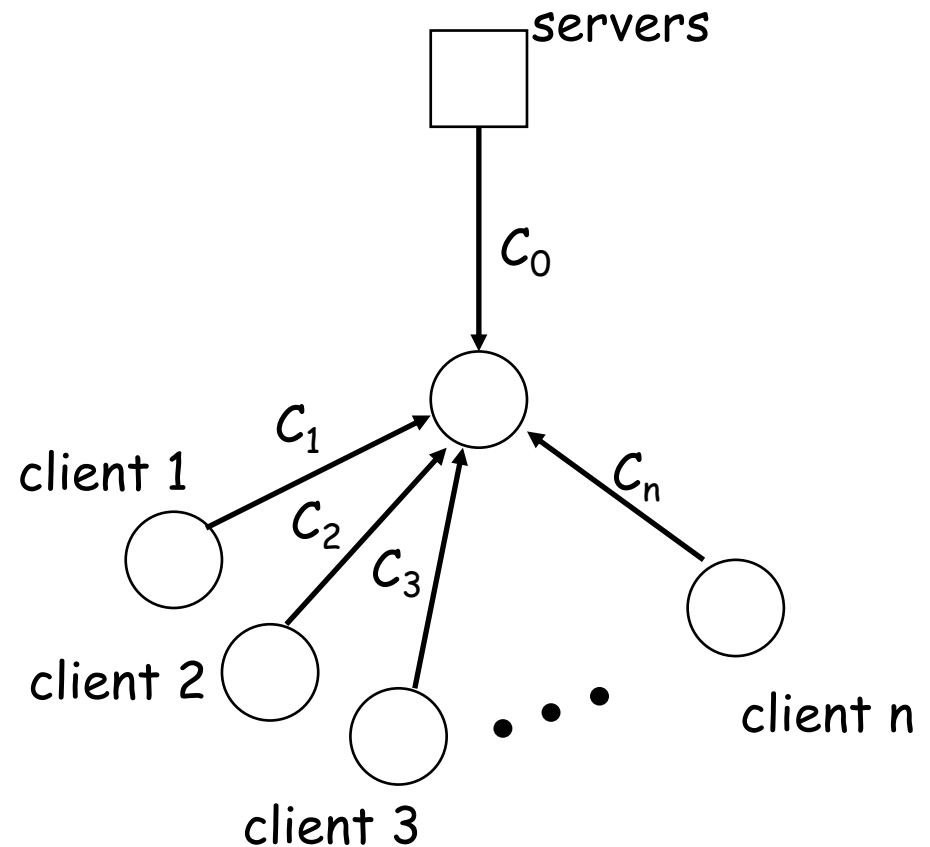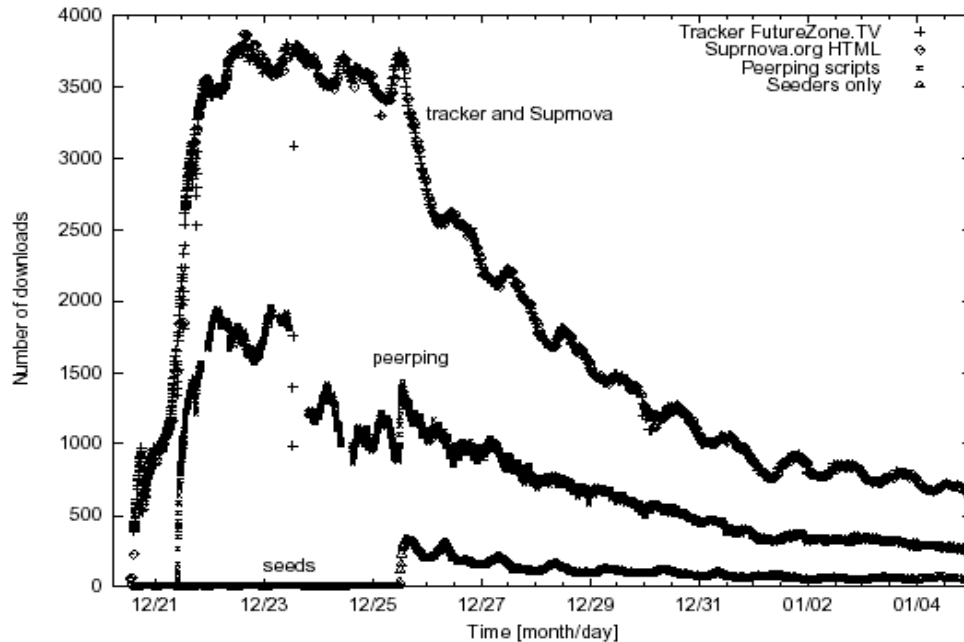$$R = \min\{C_0, (C_0 + \Sigma C_i)/n\}$$

- Assume $c_0 > (C_0 + \Sigma C_i)/n$

- Tree i:

  server $\rightarrow$ client i: $c_i /(n-1)$
  client i $\rightarrow$ other n-1 clients

- Tree 0:
  server has remaining
  $c_m = c0 - (c1 + c2 + \dots cn)/(n-1)$
  send to client i: $c_m/n$

$c_0$

$c_i /(n-1)$

$c_n$

$c_i$

$c_1$   $c_2$
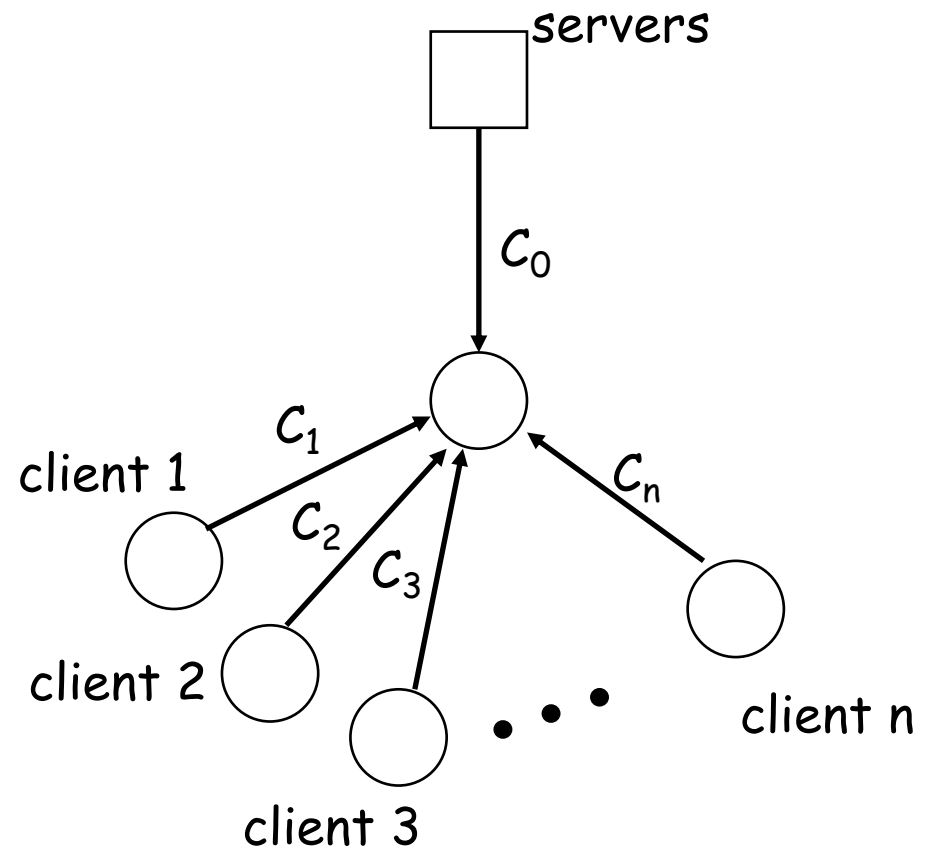
$c_0$

$c_m /n$

$c_n$

$c_1$

$c_2$   $c_i$

# Why not Building the Trees?



❑ Clients come and go (churns): maintaining the trees is too expensive

❑ Each client needs N connections

# Server+Host (P2P) Content Distribution: Key Design Issues

□ **Robustness**
- ○ Resistant to churns and failures

□ **Efficiency**
- ○ A client has content that others need; otherwise, its upload capacity may not be utilized

□ **Incentive: clients are willing to upload**
- ○ Some real systems nearly 50% of all responses are returned by the top 1% of sharing hosts

servers

$C_0$
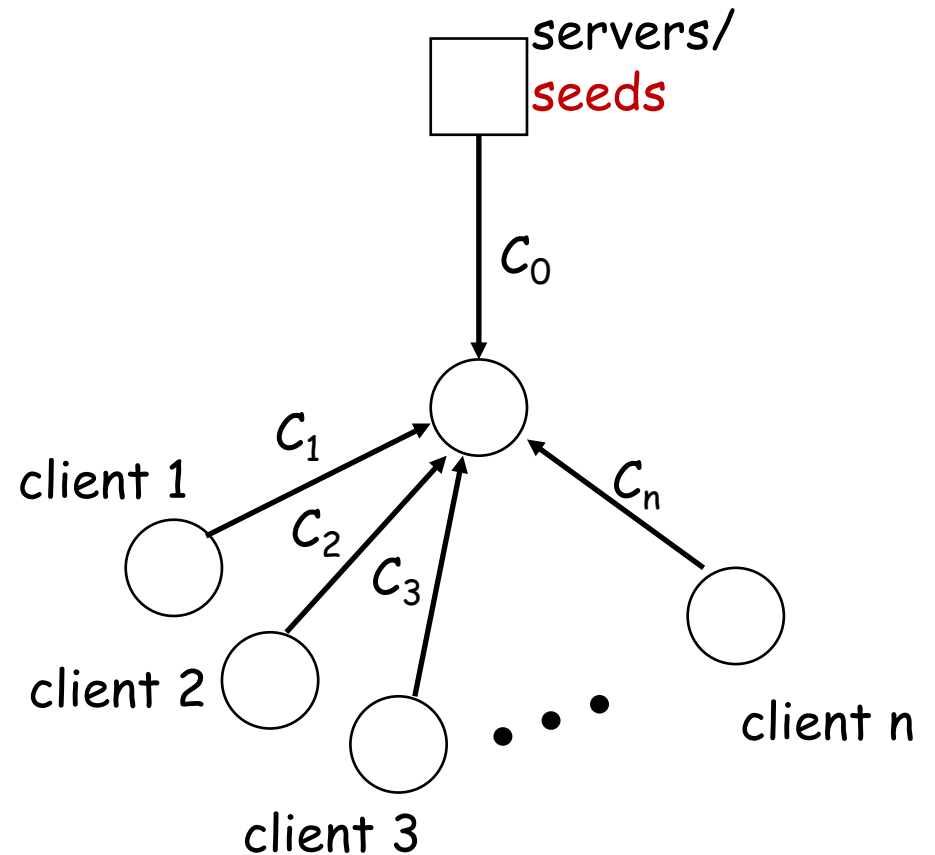
$C_1$

client 1

$C_2$

client 2

$C_3$

client 3

$C_n$

client n

# Discussion: How to handle the issues?

- Robustness

- Efficiency

- Incentive



servers/seeds

$c_0$

$c_1$

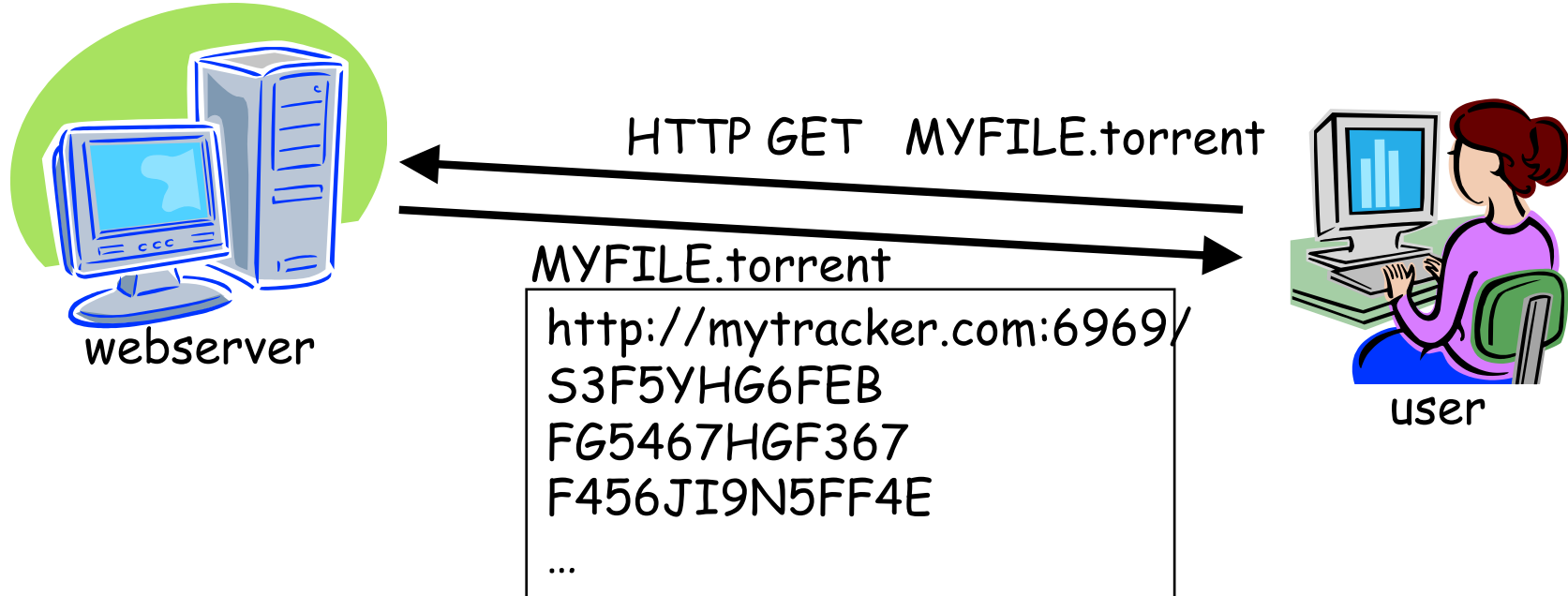client 1

$c_2$

client 2

$c_3$

client 3

$c_n$

client n

# Example: BitTorrent

❑ A P2P file sharing protocol

❑ Created by Bram Cohen in 2004

    o Spec at bep_0003:
       http://www.bittorrent.org/beps/bep_0003.html

# BitTorrent: Lookup



HTTP GET   MYFILE.torrent

webserver

MYFILE.torrent

http://mytracker.com:6969/
S3F5YHG6FEB
FG5467HGF367
F456JI9N5FF4E
...

user

13

# Metadata (.torrent) File Structure

❏ Meta info contains information necessary to contact the tracker and describes the files in the torrent

- o URL of tracker
- o file name
- o file length
- o piece length (typically 256KB)
- o SHA-1 hashes of pieces for verification
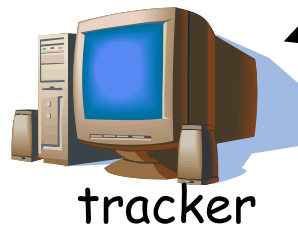- o also creation date, comment, creator, …

# Tracker Protocol

❑ Communicates with clients via HTTP/HTTPS

❑ Client GET request
  o info_hash: uniquely identifies the file
  o peer_id: chosen by and uniquely identifies the client
  o client IP and port
  o numwant: how many peers to return (defaults to 50)
  o stats: e.g., bytes uploaded, downloaded

❑ Tracker GET response
  o interval: how often to contact the tracker
  o list of peers, containing peer id, IP and port
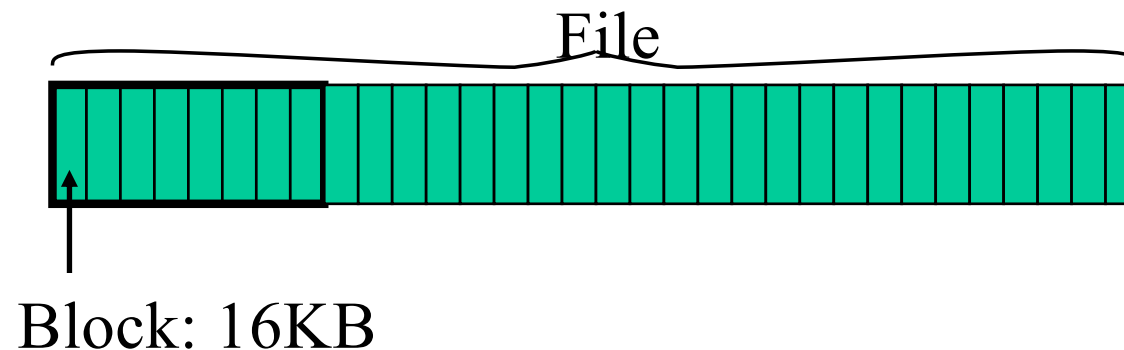  o stats

# Tracker Protocol

webserver

user

"register"

list of peers

```
ID1    169.237.234.1:6881
ID2    190.50.34.6:5692
ID3    34.275.89.143:4545
...
ID50 23        :68
```

tracker

Peer 50

Peer 2

Peer 1

# Robustness and efficiency: Piece-based Swarming

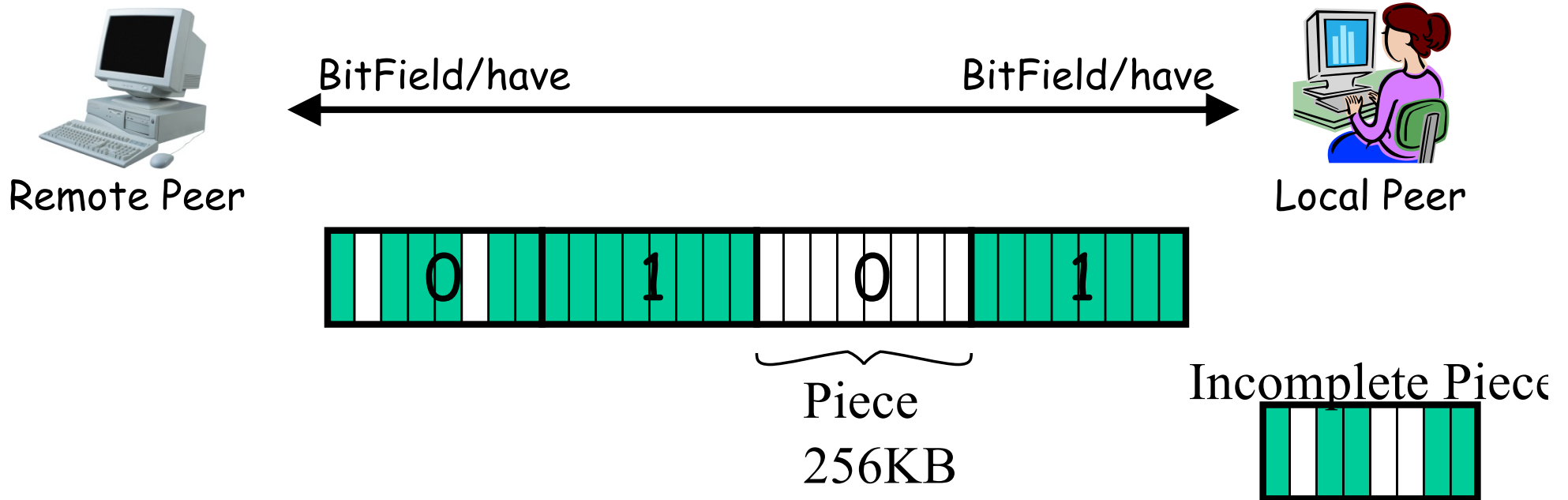□ Divide a large file into small blocks and request block-size content from different peers (why?)

Block: unit of download

File



Block: 16KB

□ If do not finish downloading a block from one peer within timeout (say due to churns), switch to requesting the block from another peer

# Detail: Peer Protocol

(Over TCP)

BitField/have ← → BitField/have

Remote Peer

Local Peer

```
0      1      0      1
```
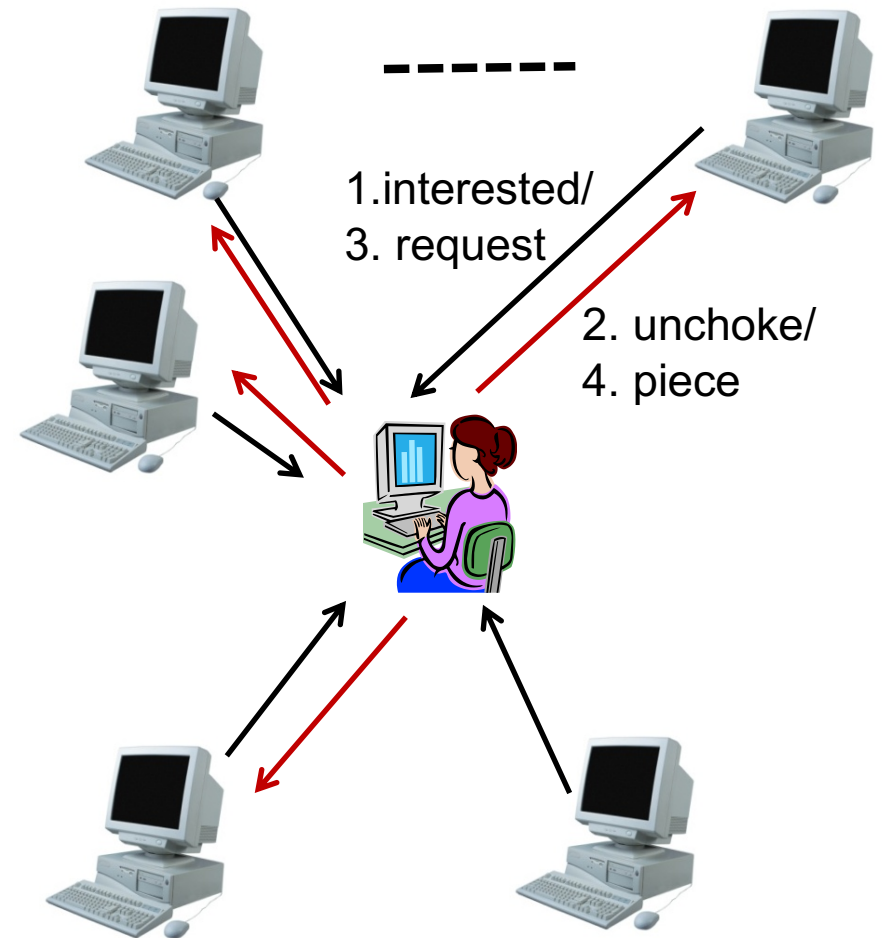
Piece
256KB

Incomplete Piece

❑ **Peers exchange bitmap representing content availability**
- o `bitfield` msg during initial connection
- o `have` msg to notify updates to bitmap
- o to reduce bitmap size, aggregate multiple blocks as a piece

# Peer Request

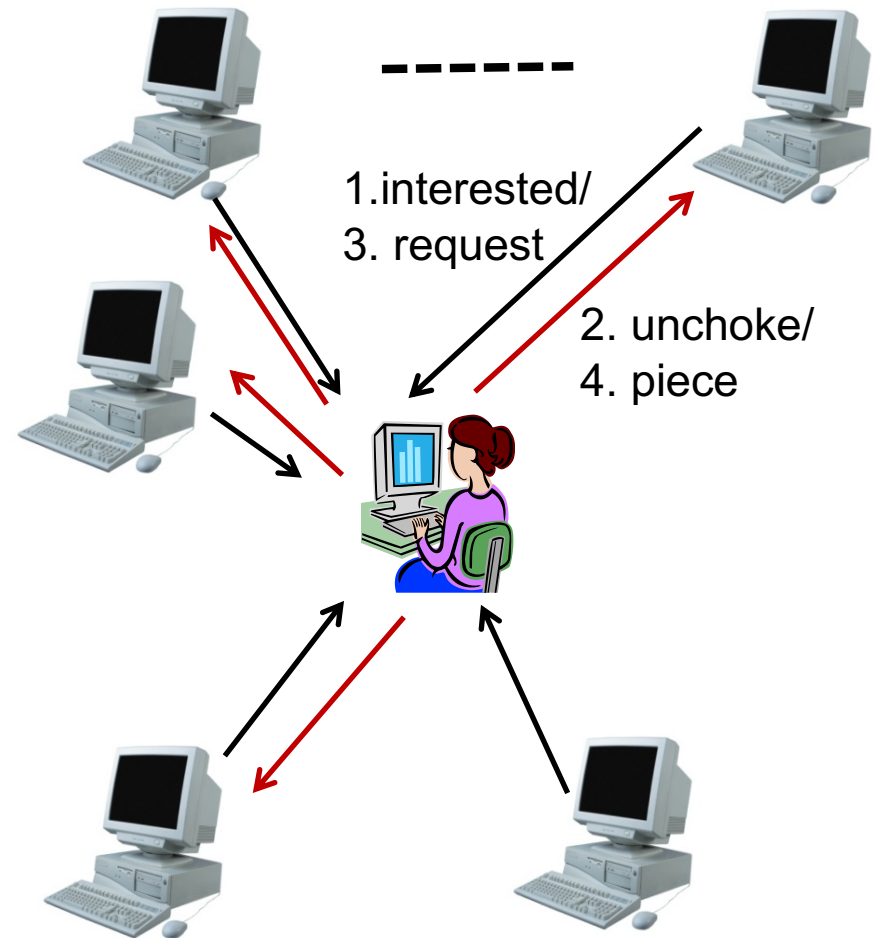- ❑ If peer A has a piece that peer B needs, peer B sends `interested` to A

- ❑ `unchoke`: indicate that A allows B to request

- ❑ `request`: B requests a specific block from A

- ❑ `piece`: specific data

1.interested/
3. request

2. unchoke/
4. piece

# Key Design Points

❑ `request`:

   ○ which data blocks to request?

❑ `unchoke`:

   ○ which peers to serve?

1.interested/
3. request

2. unchoke/
4. piece

# Request: Block Availability

❑ Request (local) rarest first
   o achieves the fastest replication of rare pieces
   o obtain something of value

# Block Availability: Revisions

- ❑ **When downloading starts (first 4 pieces): choose at random and request them from the peers**
  - o get pieces as quickly as possible
  - o obtain something to offer to others

- ❑ **Endgame mode**
  - o defense against the "last-block problem": cannot finish because missing a few last pieces
  - o send requests for missing pieces to all peers in our peer list
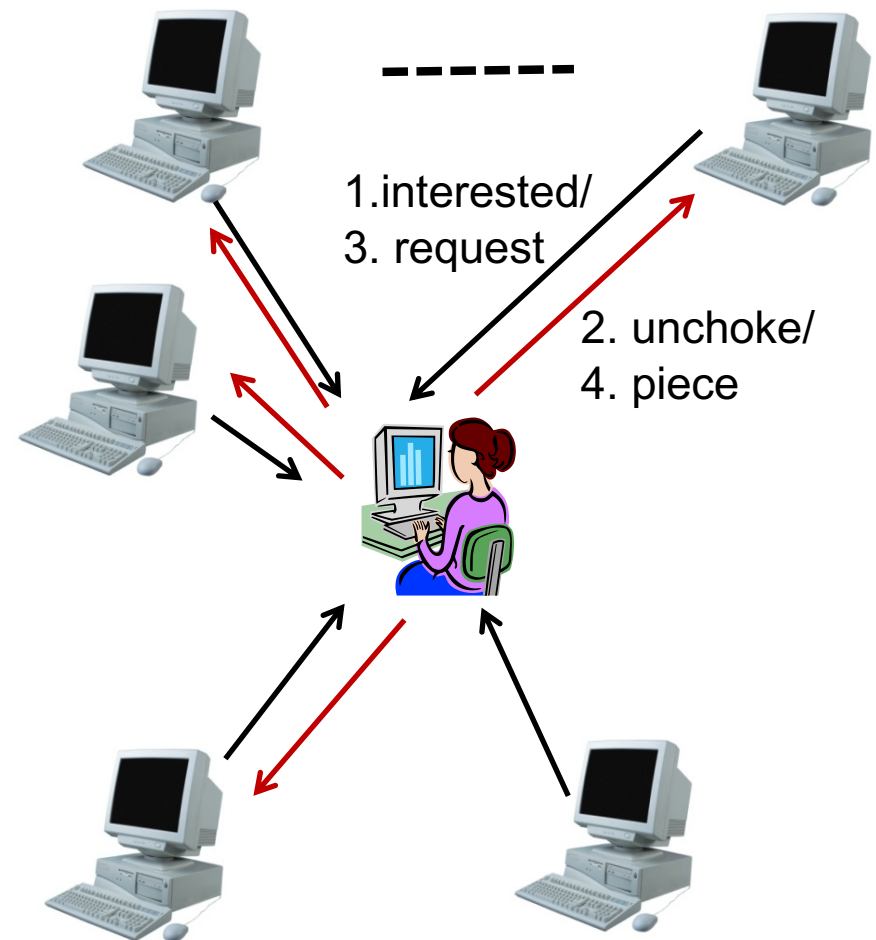  - o send `cancel` messages upon receipt of a piece

# BitTorrent: Unchoke

❑ Periodically (typically every 10 seconds) calculate data-receiving rates from all peers

❑ Upload to (*unchoke*) the fastest

- constant number (4) of unchoking slots

- partition upload bw equally among unchoked



1.interested/
3. request

2. unchoke/
4. piece

commonly referred to as "tit-for-tat" strategy

# Optimistic Unchoking

❑ Periodically select a peer at random and upload to it
  - o typically every 3 unchoking rounds (30 seconds)

❑ Multi-purpose mechanism
  - o allow bootstrapping of new clients
  - o continuously look for the fastest peers (exploitation vs exploration)

# BitTorrent Fluid Analysis

- Normalize file size to 1
- x(t): number of downloaders (also known as leechers) who do not have all pieces at time t.
- y(t): number of seeds in the system at time t.
- $\lambda$: the arrival rate of new requests.
- $\mu$: the uploading bandwidth of a given peer.
- c: the downloading bandwidth of a given peer, assume $c \geq \mu$.
- $\theta$: the rate at which downloaders abort download.
- $\gamma$: the rate at which seeds leave the system.
- $\eta$: indicates the effectiveness of downloader sharing, $\eta$ takes values in [0, 1].

# System Evolution

$$\frac{dx}{dt} = \lambda - \theta x(t) - \min\{cx(t), \mu(\eta x(t) + y(t))\},$$

$$\frac{dy}{dt} = \min\{cx(t), \mu(\eta x(t) + y(t))\} - \gamma y(t),$$

Solving steady state: $\frac{dx(t)}{dt} = \frac{dy(t)}{dt} = 0$

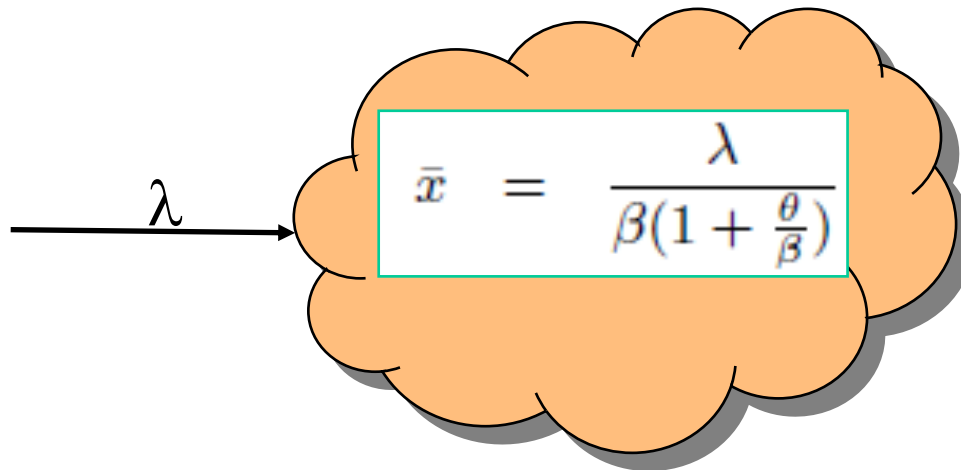Define $\quad \frac{1}{\beta} = \max\{\frac{1}{c}, \frac{1}{\eta}(\frac{1}{\mu} - \frac{1}{\gamma})\}$

$$\bar{x} = \frac{\lambda}{\beta(1 + \frac{\theta}{\beta})}$$

$$\bar{y} = \frac{\lambda}{\gamma(1 + \frac{\theta}{\beta})}.$$

"Modeling and Performance Analysis of BitTorrent-Like Peer-to-Peer Networks", SIGCOMM'04
https://conferences.sigcomm.org/sigcomm/2004/papers/p444-qiu1.pdf

# System State

$$\bar{x} = \frac{\lambda}{\beta(1 + \frac{\theta}{\beta})}$$

$$\bar{y} = \frac{\lambda}{\gamma(1 + \frac{\theta}{\beta})}.$$

Q: How long does each downloader stay as a downloader?

$$\lambda \qquad \bar{x} = \frac{\lambda}{\beta(1 + \frac{\theta}{\beta})}$$

$$T = \frac{1}{\theta + \beta}$$

$$\frac{1}{\beta} = \max\{\frac{1}{c}, \frac{1}{\eta}(\frac{1}{\mu} - \frac{1}{\gamma})\}$$

Key take-away: not scaling inverse with system size (x)

- New requests comes, new bandwidth also comes

27

# Recap

- ❏ Applications
  - ❏ Client-server applications
    - Single server
    - Multiple servers load balancing
  - ❏ Application overlays (distributed network applications) to
    - scale bandwidth/resource (BitTorrent)
    - distribute content lookup (Freenet, DHT, Chord) [optional]
    - distribute content verification (Block chain) [optional]
    - achieve anonymity (Tor) [optional]