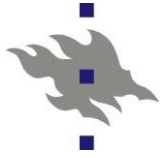HELSINGIN YLIOPISTO
HELSINGFORS UNIVERSITET
UNIVERSITY OF HELSINKI

# Peer-to-Peer Networks

Chapter 6: P2P Content Distribution
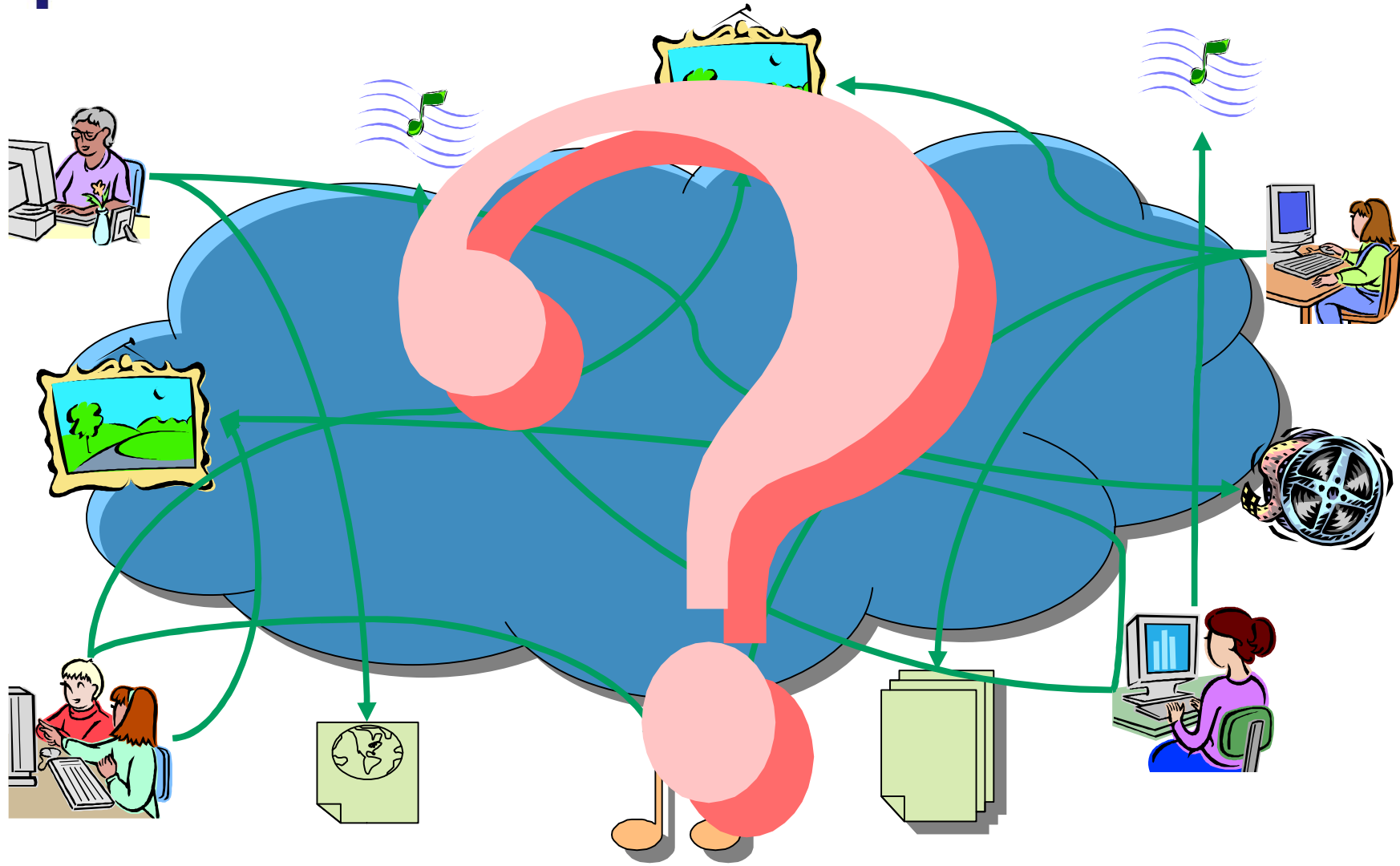
# Chapter Outline

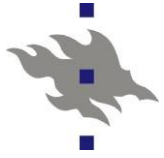n Content distribution overview

n Why P2P content distribution?
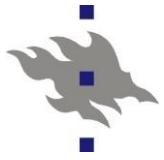
n Network coding

n Peer-to-peer multicast

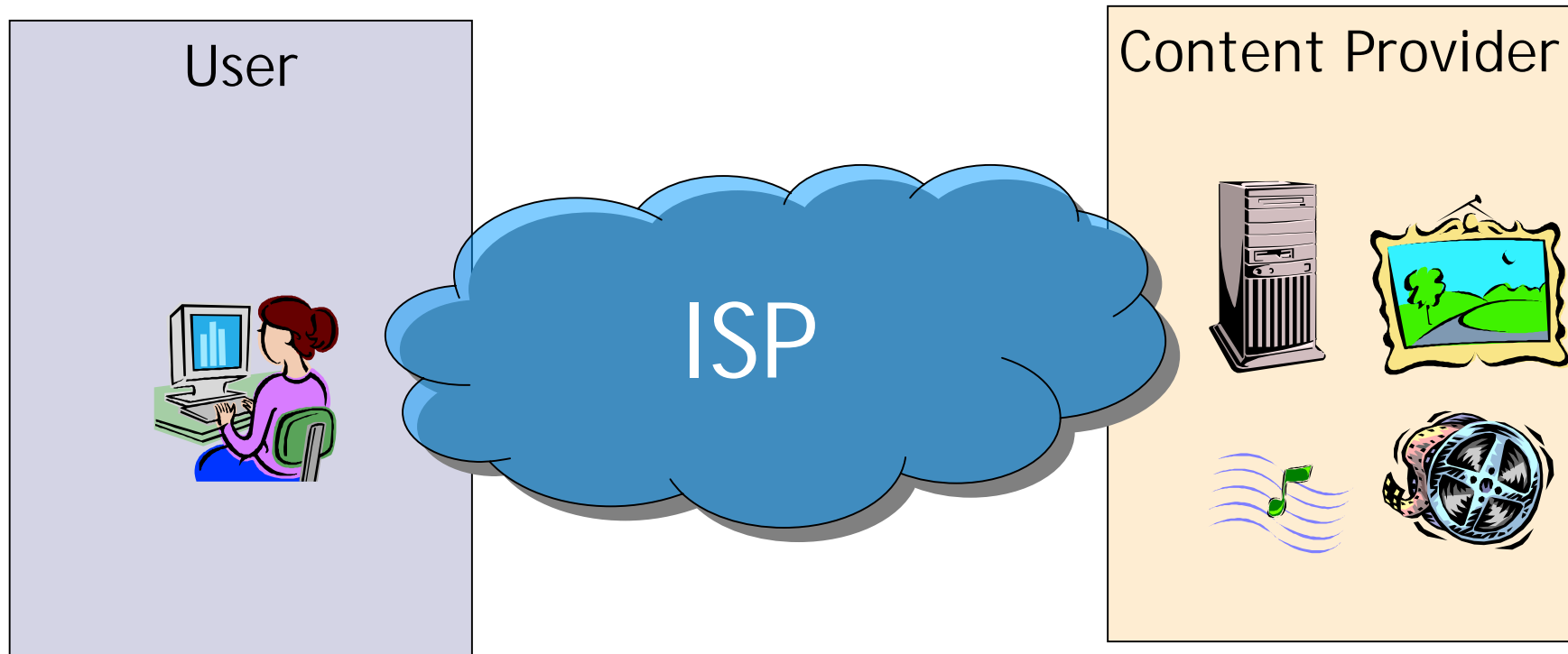# What is the Content Distribution Problem?

## Problem Definition
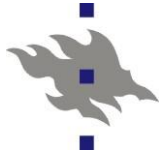
n    How to distribute (popular) content efficiently on Internet to a large number of geographically distributed people?

n    …while keeping your costs manageable! :-)

n    Who is "you"?

n    So, we need to answer the following questions…

1.    Who are the players?

2.    What is the content?

3.    How are costs determined?

4.    How is performance measured?

# The Players



User

ISP

Content Provider

- n ISP represents the network path between the user and the content provider
- n Typically, such a path contains several ISPs
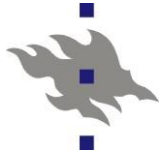
# Roles of the Players

n User

    n Wants to access content provided by content provider

    n Buys access to network from ISP

    n Possibly pays content provider for content

n Content provider

    n Produces content

    n Runs a server (or outsources this)

    n Buys access from an ISP

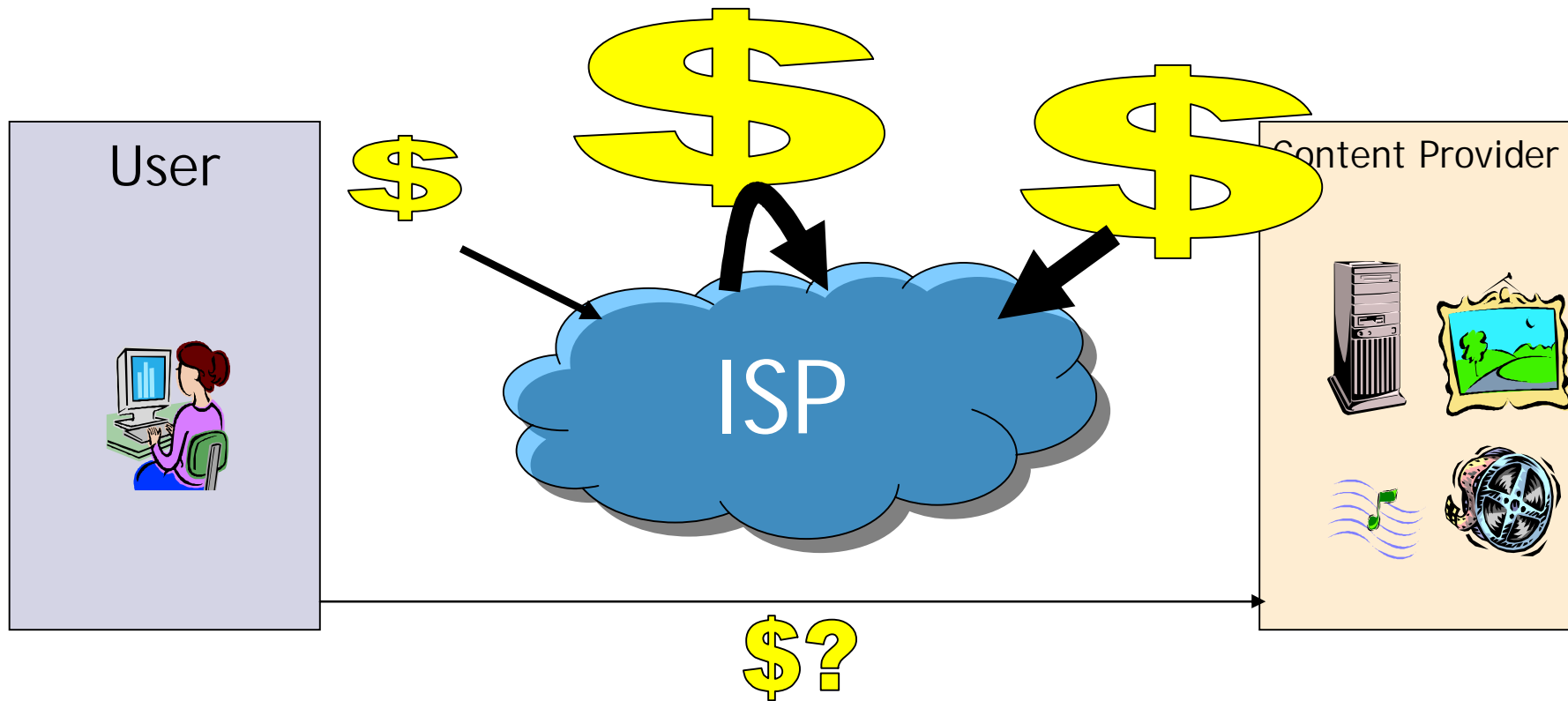n ISP

    n Provides the network connection for the others

# What Is Content?

n Content is any digital content the users want

n Examples:

- n Web pages
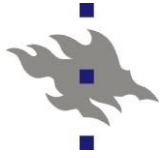- n Music files
- n Videos
- n Streaming media
- n and so on…

# Follow the Money

n Where does the money come from?

**User**

**ISP**

**Content Provider**

$

$ $

$ $

$?

# Business Relationships

n User

- n Pays money to ISP for connectivity
- n These days often flat-rate, can be per-byte
- n Might rarely pay content provider

n Content provider

- n Pays ISP for connectivity
- n Pays for running a server
- n Needs to get money from somewhere
- n *May or may not need to make a profit from content*

n ISP

- n Gets money from others
- n Uses money to run network
- n *Wants to make a profit*

# Performance Targets

- User
  - Wants content fast
  - Does not want to pay (too much)
- Content provider
  - Wants to get as many users as possible
  - While maintaining costs as low as possible
- ISP
  - Wants to make as much money as possible
- Which goal is the most important?

- In real world, user needs often not "relevant"
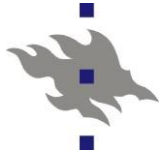  - Users do not contribute enough money
  - Users might not have a choice, e.g., only 1 ISP possible
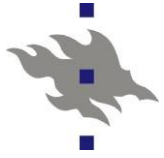- ISP is the most important in many cases
- Content providers have some say in some cases
  - Compare proxy caching and CDNs

# Content Distribution Technology

n Content distribution is an "old" problem

    n Read: Originates from the Web

n Technology currently in use:

    n Client-side web proxies/caches

    n Server-side load balancing

    n Mirror servers

    n Content Distribution Networks (CDN)

    n Peer-to-peer content distribution

    n Multicast (at least in theory…)

# File Sharing?!?

n Is file sharing content distribution?

n YES:

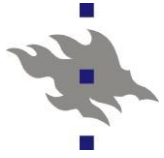    n There is content being distributed

n NO:

    n Goal is to *make content available*, not the actual distribution

    n Main focus is on searching and locating content

    n Actual distribution just a simple download with nothing fancy

n We exclude file sharing from content distribution

# Peer-to-Peer Content Distribution

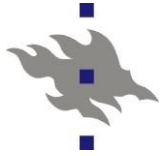n Definition:

*Peer-to-peer content distribution is content distribution which is implemented in a peer-to-peer fashion*

n Best/only current example: BitTorrent

n One file, replicate to everyone as fast as possible
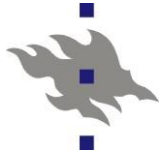
n Why P2P content distribution?

# P2P Content Distribution: Pros and Cons

**Advantages**

n Cheap for content provider

- n Need to seed file, that's it!
- n Server farms, CDNs, etc. are very, very expensive

n Capacity increases with increasing number of users

**Disadvantages**

n Users must want to help

- n Incentives?

n Need to trust users

n No (easy) way to optimize

- n Topological locality

n P2P has "bad reputation"

# Current State

n Real world:

n BitTorrent and nothing much else

n Research world:

n BitTorrent under active research
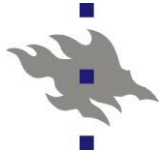
    n Piece and peer selection algorithms

    n Measurement work

n Avalanche from Microsoft

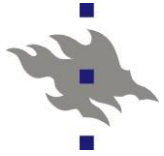    n Same as BitTorrent, but uses network coding

n P2P Streaming
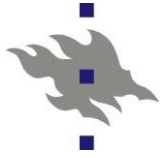
    n Lots of proposals

# Network Coding: Motivation

n Main question in BitTorrent:

   *Which piece to download next?*

n BitTorrent's answer: Rarest first!

n Rare blocks might disappear from system?

n How about advanced coding techniques?

   n Erasure codes at the source

   n Network coding by the peers

# Erasure Codes at the Source

n Original source creates redundant blocks

    n Reed-Solomon, Tornado, Digital Fountain codes, …

n Original file has $N$ blocks, source creates $K$ blocks

n Generally applies:

Any $N$ out of $N+K$ are sufficient to reconstruct file

    n Some coding schemes require a bit more than $N$

n Problem: Redundant data generated by one entity

n Problem: $K$ is fixed

    n Some codes allow for "infinite $K$"

# Network Coding

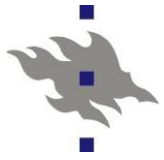n  How about letting peers do the coding?
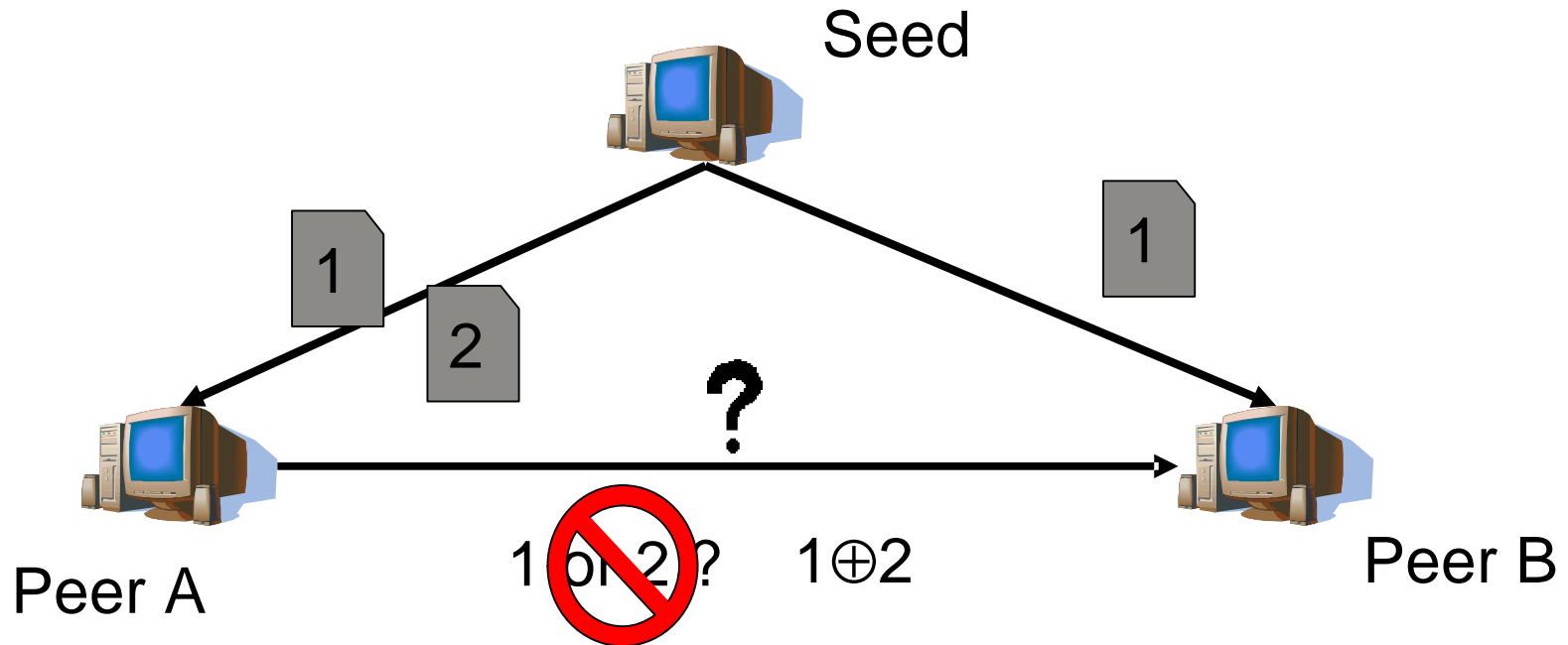
n  Solution: Network coding

Idea:

n  Peers create linear combinations of blocks they have

n  Send block with coefficients to others

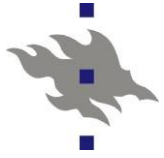n  Sufficient number of linearly independent blocks allows for reconstruction of file
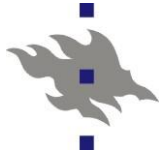
# Network Coding



Seed

Peer A

Peer B

1 or 2?    1⊕2

- Peer A sends a linear combination of blocks 1 and 2
- Peer B already has block 1, can compute block 2

# Network Coding: Math

n File has *N* blocks

  n For example $F = [x_1, x_2]$ (two blocks)

n Pick two numbers $a_{i,1}$, $a_{i,2}$

n Code $E_i (a_{i,1}, a_{i,2}) = a_{i,1} * x_1 + a_{i,2} * x_2$

  n Can create an infinite number of $E_i$'s

n When we have two linearly independent $E_i (a_{i,1}, a_{i,2})$'s, we can re-create original blocks $x_1$ and $x_2$

  n Same as solving a set of linear equations

n Note: All math in some finite field, e.g., GF(216)

# Practical Considerations

n In principle, *information* is spread out better

n Need for linear independence

    n May need to download more data than file has
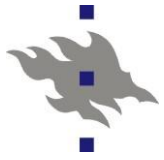
    n No analysis of relevance of this done yet

n Performance has been found to be good

    n Download times shorter and more *predictable*

n System very robust against departures of seeds

n Does not need altruistic users to stay logged on

# Peer-to-Peer Multicast and Streaming

n Streaming content delivery gaining importance

    n Mainly for video-on-demand systems

n Many systems for P2P streaming or multicast

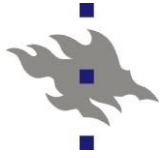|  | One recipient | Multiple recipients |
|---|---|---|
| Bulk data | Traditional unicast | Multicast |
| Media data | Unicast streaming | Multicast streaming |

# Why Multicast?

n   Multicast is sending data to many receivers at the same time

n   Compare: Broadcast is sending data to everybody

n   Two-kinds of multicasting schemes:

1.   *1-to-m multicast*

   n   One sender, many receivers, e.g., video streaming

2.   *n-to-m multicast*

   n   Many senders, many receivers, e.g., video conferencing

n   For content distribution 1-to-m multicast is more relevant (and more researched)

n   Multicast has two main components:

   n   Group management
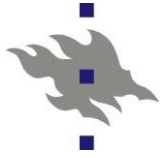
   n   Routing

# Multicast Group Management and Routing

n Group is another name for the receivers
n Group management is about how and who manages the membership in the group
  - n Centralized, de-centralized, …

n Multicast routing is how to route the packets
n Multicast property: Ideally, each packet will go over any given network link at most once
n Achieved with a multicast tree, rooted at the sender
n Requires support from the routers in the network
n IPv4 has multicast address block (block D, 224.x.x.x)
  - n Support from routers not guaranteed
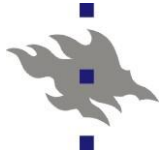n IPv6 should have native support for multicast

# Why ALM?

n IP-level multicast is efficient, but requires support from all the routers in the network

n Most routers do not have multicast support turned on

n Result: No Internet-wide multicast possible

    n mbone overlay network offers some support

n Application level multicast performs multicasting on the application level

n Benefit: No support from routers needed

n Disadvantage: Not as efficient as native multicast

n ALM builds an overlay network

n Overlay typically consists of the receivers and senders

    n Some overlays assume fixed nodes in network

# ALM Details

n ALM can be classified according to:

1. How they build distribution tree?

   n First mesh, then tree, or directly as spanning tree?

2. How to build overlay?

   n Use a structured network, e.g., DHT as basis?

   n Build directly on top of end-user machines?

n ALM metrics:

   n Relative Delay Penalty: How much does the one-way delay increase because of overlay (ratio of overlay/direct path)

   n Throughput

   n Stress: How many times a packet traverses a given link

   n Control overhead
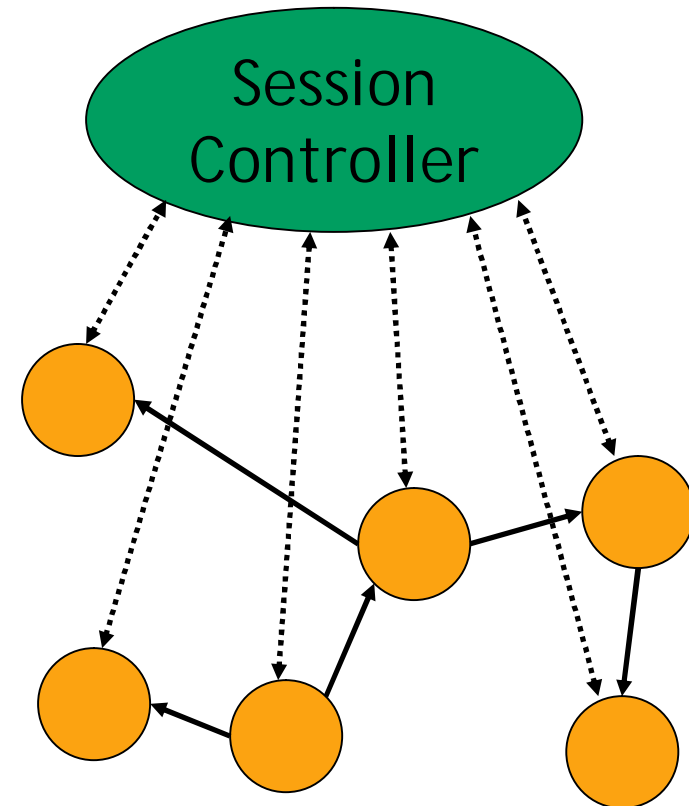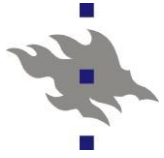
# ALM Architectures

n Architectures based on unstructured overlays

  n Centralized architecture

  n De-centralized architecture

n Architectures based on structured overlays

  n Flooding-based multicast

  n Tree-based multicast

n Structured overlays are DHTs in this case
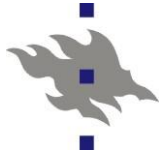
# Application Level Multicast Infrastructure

n Also known as ALMI

n ALMI is a centralized, unstructured overlay system

n One session controller controls the distribution

    n Not part of actual data transfer, just control

n Session controller calculates optimal multicast tree

**Session Controller**

# ALMI Details

n When a node wants to join the tree, it contacts the session controller

n Node is assumed to know address of session controller

    n For example, well-known URL, etc.

n Session controller assigns ID to node and gives information about where the node is in the tree

n Joining node contacts its parent in the tree

    n Parent adds new node as child
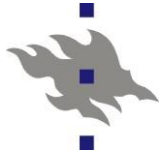
n When node leaves, it informs the session controller

# ALMI Pros and Cons

Advantages

n High control over overlay

n Easy to implement

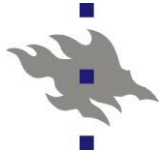n Easy to detect malicious nodes

    n Because of centralized nature

Disadvantages

n Poor scalability

n Session controller is a single point of failure

    n Backup controllers monitor main controller

    n When main controller fails, one of the backups takes over

# End System Multicast (ESM)
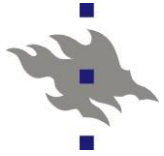
n Fully distributed ALM scheme based on end nodes

n Group management handled by Narada protocol

n Idea of Narada:

1. Overlay nodes organize themselves in a mesh

2. Source-specific multicast trees are built on top of mesh

n Hence, quality of multicast depends on quality of mesh

n Goal of Narada is to match overlay requirements with underlay quality

   n Meaning: Overlay link should have same properties as the corresponding underlay link between the same nodes

n Limited out-degree of nodes to limit load

# Narada Group Management

- Every node participates in group management
- Does not scale well, but Narada's target is small to medium sized multicast networks
- Joining node contacts some existing members
  - Again, addresses assumed to be somehow known
- Every member sends periodic heartbeats messages
- Heartbeats announce e.g., arrival of new node
- Eventually, all nodes know about the new node
- Departures detected by missing heartbeats
- Can also repair network partitions:
  - If heartbeats missing, then send probe to that node
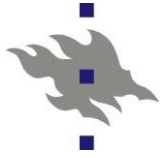  - If node answers, then reconnect overlay
  - Otherwise assume node is gone

# Narada Performance

Relative Delay Penalty

n For longer underlay delays, RDP between 2 and 4

n For short underlay delays, RDP up to 14

n Explanation: For short underlay links, even a slightly
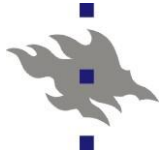suboptimal overlay will lead to high RDP

Stress

n Narada limits maximum stress to 9
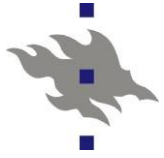
n Sending everything over unicast has maximum stress 128

# Flooding-Based Replication

n Use a structured overlay to construct multicast trees

n This case: Use Content Addressable Network, CAN

n Multicast in CAN works as follows

1. Multicast group has a well-known identifier which is mapped to some point in CAN

2. When node wants to join group, it contacts the responsible node for the group ID

3. Responsible node redirects new node to a mini-CAN

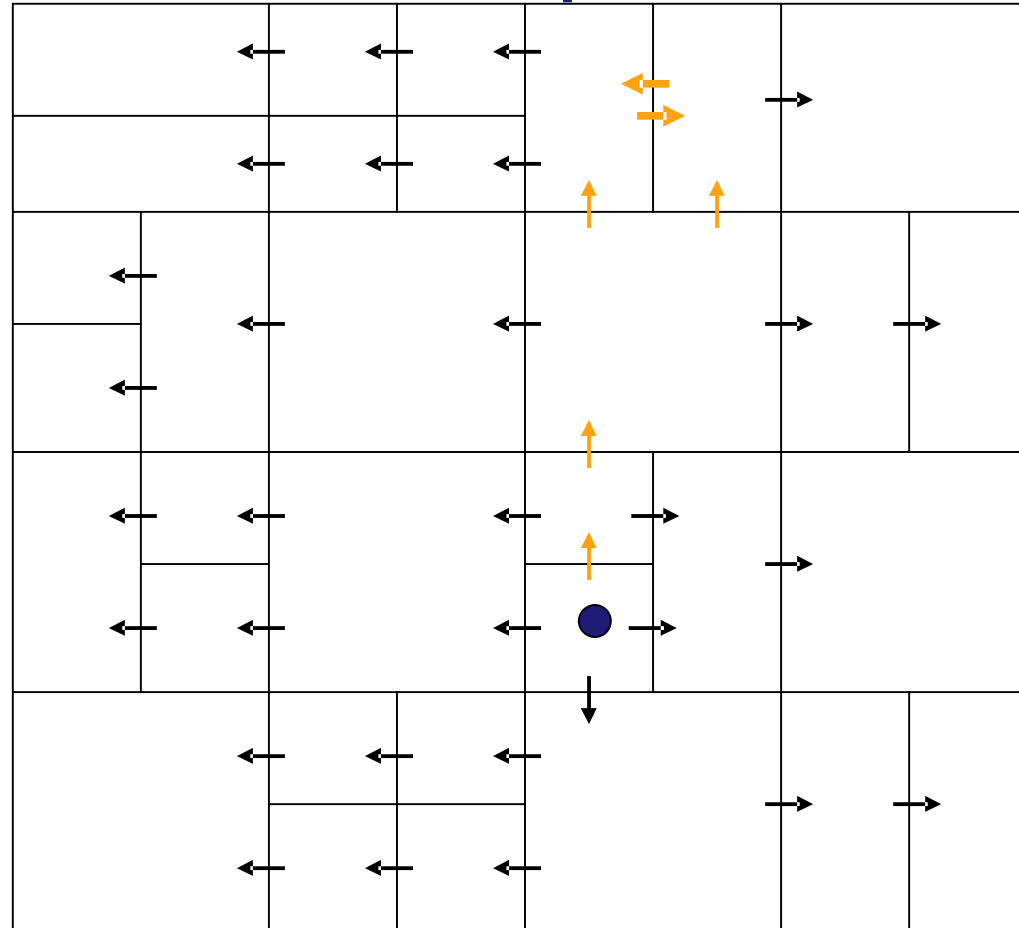   n Mini-CAN has only nodes who are members of multicast group

# CAN Multicast Rules

1. *Source forwards message to all its neighbors*
2. *When node receives message, it forwards it to all its neighbors, except the ones in the direction from which the message came*
3. *If message has traversed more than half of the distance across any dimension, message is not forwarded*
4. *Nodes cache sequence numbers of messages and discard any duplicates they see*

- Rules 1 and 2 ensure messages reach everywhere
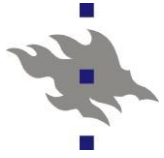- Rule 3 prevents routing loops
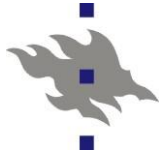- Rule 4 is just performance enhancement

# CAN Multicast Example



n Orange arrows lead to duplicates

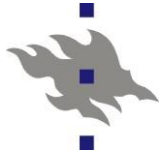n Note: Not all arrows shown in figure

# Tree-Based Example

n Tree-based multicast with structured overlays

n Example: Scribe

   n Based on Pastry DHT

n Idea: Multicast group ID maps to a node in DHT

n That node is root of the multicast tree

   n Also called rendezvous point (RP)

n When a node joins, it sends message to RP

n Message routed in DHT

n Every node on the path checks if it knows about that multicast group

   n If yes, then add new node as child

   n If not, add new node, forward message towards RP

# Scribe: Sending Data

n When a node wants to send data to a multicast group, it sends it to the rendezvous point

n Join messages have created the multicast tree rooted at the rendezvous point

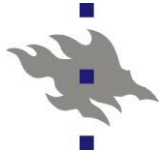n RP forwards the message to its children, who forward it to their children, etc.

# Comparison of CAN and Scribe

**Relative Delay Penalty**

n CAN has higher RDP, typically between 6 and 8

n Scribe has RDP about 2

**Stress**

n Average stress a bit lower in CAN

n Maximum stress in CAN much lower

n Scribe better for delay-sensitive applications

    n For example, video conference

n CAN better for non delay-critical applications

    n For example, TV broadcast

# Chapter Summary

n Content distribution overview

n Why P2P content distribution?

n Network coding

n Peer-to-peer multicast