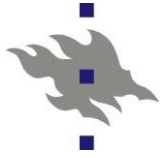**HELSINGIN YLIOPISTO**
**HELSINGFORS UNIVERSITET**
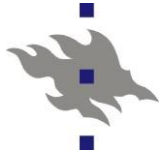**UNIVERSITY OF HELSINKI**

# Peer-to-Peer and Grid Computing

Chapter 4: Peer-to-Peer Storage
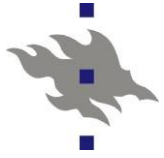
# Chapter Outline

n Using DHTs to build more complex systems

    n How DHT can help?

    n What problems DHTs solve?

    n What problems are left unsolved?

n P2P storage basics, with examples

    n Splitting into blocks (CFS)

    n Wide-scale replication (OceanStore)

    n Modifiable filesystem with logs (Ivy)

n Future of P2P filesystems

# How to Use a DHT?

n Recall: DHT maps keys to values

n Applications based on DHTs must need this functionality

  n Or: Must be designed in this way!

  n Possible to design an application in several ways

n Keys and values are application specific

  n For filesystem: Value = file

  n For email: Value = email message

  n For distributed DB: Value = contents of entry, etc.

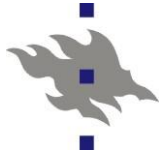n Application stores values in DHT and uses them

  n Simple, but a powerful tool

# Problems Solved by DHT

- DHT solves the problem of mapping keys to values in the distributed hash table
- Efficient storage and retrieval of values
- Efficient routing
  - Robust against many failures
  - Efficient in terms of network usage
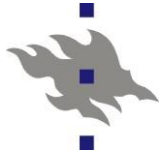
- Provides hash table-like abstraction to application

# Problems NOT Solved by DHT

n Everything else except what is on previous slide…

n In particular, the following problems

n Robustness
- n No guarantees against big failures
- n Threat models for DHTs not well-understood yet

n Availability
- n Data not guaranteed to be available
- n Only probabilistic guarantees (but possible to get high prob.)

n Consistency
- n No support for consistency
- n Data in DHT often highly replicated, consistency is a problem

n Version management
- n No support for version management
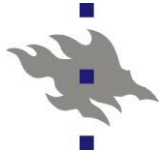- n Might be possible to support this to some degree

# P2P FS: Introduction

n P2P filesystems (FS) or P2P storage systems were the first applications of DHTs

n Fundamental principle:

*Key = filename, Value = file contents*

n Different kinds of systems

n Storage for read-only objects

n Read-write support

n Stand-alone storage systems

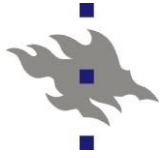n Systems with links to standard filesystems (e.g., NFS)

# P2P FS: Current State

n Only examples of P2P filesystems come from research

n Research prototypes exist for many systems

n No wide-area deployment

  n Experiments on research testbeds

  n No examples of real deployment and real usage in wide-area

n After initial work, no recent advances?

  n At least, not visible advances

n Three examples:

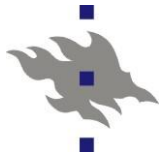  n Cooperative File System, CFS

  n OceanStore

  n Ivy

# P2P FS: Why?

n Why build P2P filesystems?

n Light-weight, reliable, wide-area storage

  n At least in principle…

n Distributed filesystems not widely deployed either…

  n Were studied already long time ago

n Gain experience with DHT and how DHTs could be used in real applications

  n DHT abstraction is powerful, but it has limitations

  n Understanding of the limitations is valuable

# P2P FS: Basic Techniques

n Three fundamental basic techniques for building distributed storage systems

1. Splitting files into blocks

2. Replicating files (or blocks!)

3. Using logs to allow modifications

n For now: Simple analysis of advantages and disadvantages and three examples

n Detailed performance analysis in Chapter 5

  n For blocks and replication

# Splitting Files into Blocks

n Why: Files are of different sizes and peers storing large files have to serve more data

Pro:

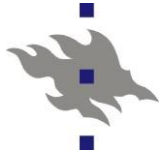n Dividing files into equal-sized blocks and storing blocks on different peers can achieve load balance

n If different files share blocks, we can save on storage

Con:

n Instead of needing one peer online, all peers with all blocks must be online (see below)

n Need metadata about blocks to be stored somewhere

n Granularity tradeoff: Small blocks -> Good load balance, but lots of overhead and vice versa

# Replication

n Why: If file (or block) is stored only on one peer and that peer is offline, data is not available

n Replicating content to multiple peers significantly increases content availability

Pro:

n High availability and reliability

  n But only probabilistic guarantees

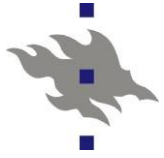Con:

n How to coordinate lots of replicas?

  n Especially important if content can change

n Unreliable network requires high degree of replication for decent availability
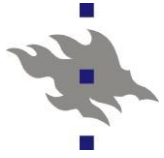
  n "Wastes" storage space

# Logs

n Why: If we want to change the stored files, we need to modify every stored replica

n Keep a log for every file (user, …) which gives information about the latest version

Pro:

n Changes concentrated in one place
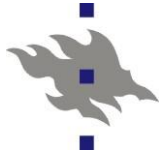
n Anyone can figure out what is the latest version

Con:

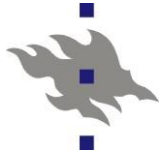n How to keep the log available?

  n By replicating it? ;-)

# P2P FS: Overview

n Three examples of P2P filesystems

n CFS (blocks and replication)

  n Basic, read-only system

  n Based on Chord

n OceanStore (replication)

  n Vision for a global storage system

  n Based on Tapestry

n Ivy (logs)

  n Read-write, provide NFS semantics

  n Based on Chord

# CFS

n CFS = Cooperative File System

n Developed at MIT, by same people as Chord

n CFS based on the Chord DHT

n Read-only system, only 1 publisher

n CFS stores blocks instead of whole files

    n Part of CFS is a generic block storage layer

n Features:

    n Load balancing

    n Performance equal to FTP

    n Efficiency and fast recovery from failures
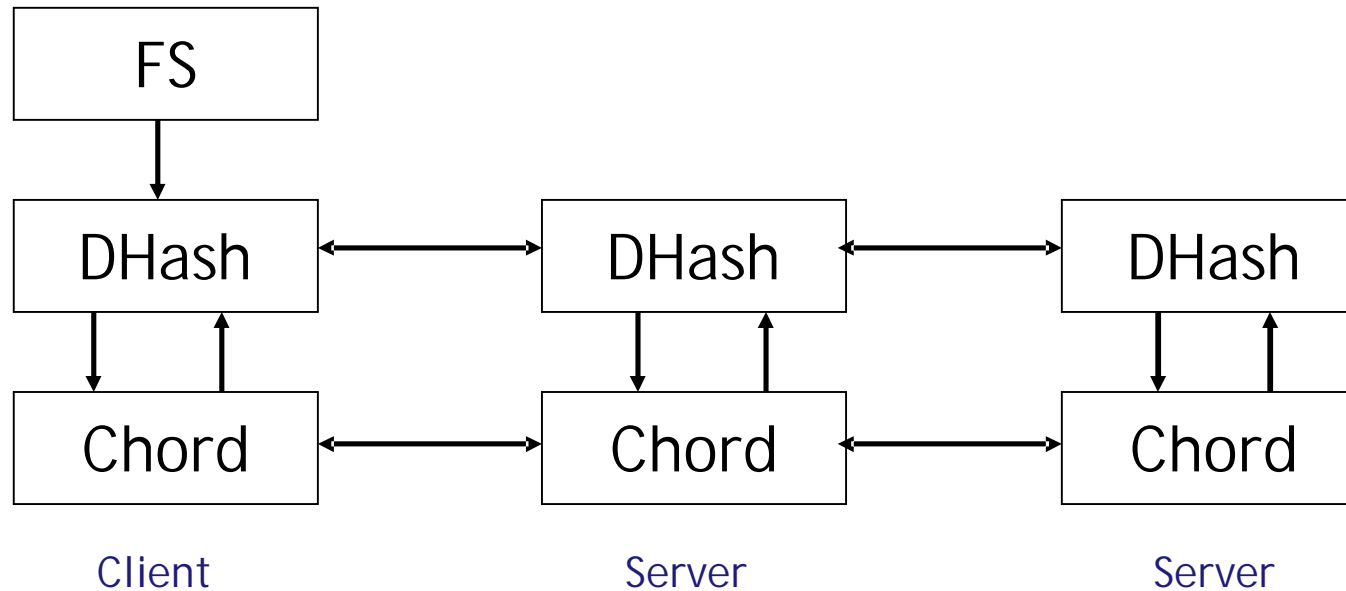
# CFS Properties

n Decentralized control

n Scalability (comes from Chord)

n Availability

    n In absence of catastrophic failures, of course…

n Load balance

    n Load balanced according to peers' capabilities

n Persistence

    n Data will be stored for as long as agreed

n Quotas

    n Possibility of per-user quotas

n Efficiency

    n As fast as common FTP in wide area
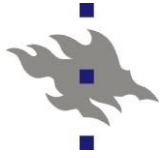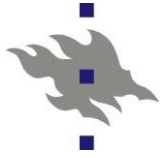
# CFS: Layers, Clients, and Servers



- n FS: provide filesystem API, interpret blocks as files
- n DHash: Provides block storage
- n Chord: DHT layer, slightly modified
- n Clients access files, servers just provide storage

# Chord Layer in CFS

n Chord layer is slightly modified from basic Chord

n Instead of 1 successor, each node keeps track of $r$ successors

n Finger tables as before

n Also, try to reduce lookup latency

  n Nodes measure latencies to other nodes

  n Report measured latencies to other nodes

# DHash Layer

- DHash stores blocks as opposed to whole files
  - Better load balancing
  - More network query traffic, but not really significant
- Each block replicated $k$ times, with $r \geq k$
- Two kinds of blocks:
  - Content block, addressed by hash of contents
  - Signed blocks (= root blocks), addressed by public key
    - Signed block is the root of the filesystem
    - One filesystem per publisher
- Blocks are cached in network
  - Most of caching near the responsible node
  - Blocks in local cache replaced according to least-recently-used
  - Consistency not a problem, blocks addressed by content
    - Root blocks different, may get old (but consistent) data

# DHash API

n DHash provides following API to clients:

| | |
|---|---|
| Put_h(block) | Store block |
| Put_s(block, pubkey) | Store or update signed block |
| Get(key) | Fetch block associated with key |

n Filesystem starts with root (= signed) block

- n Block under publisher's public key and signed with private key
- n For clients, read-only, publisher can modify filesystem by inserting a new root block
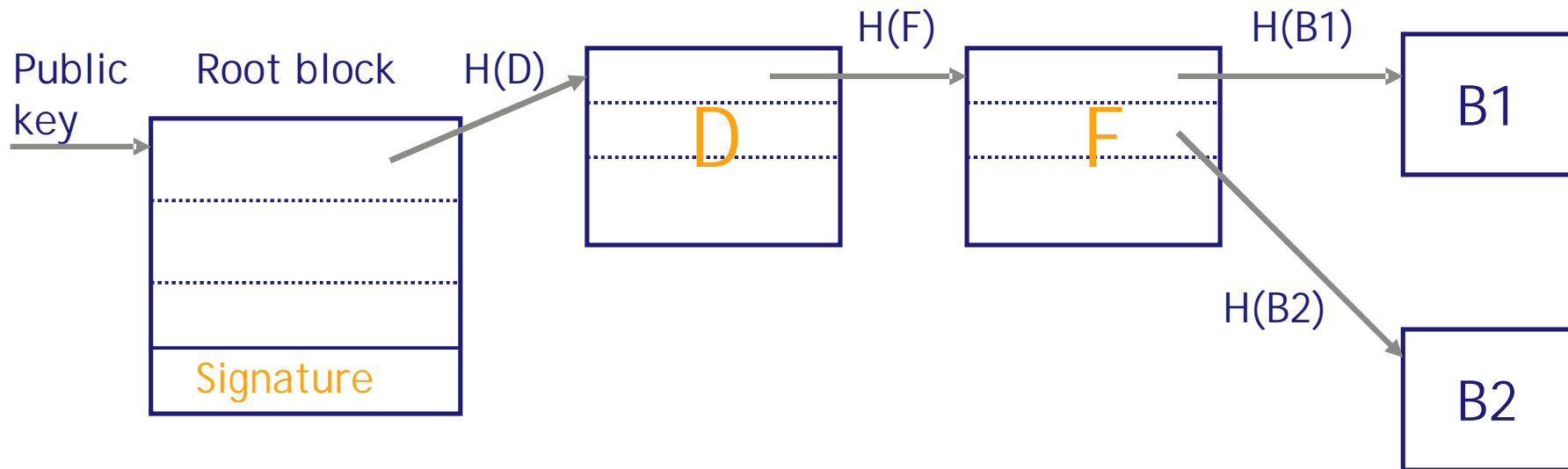- n Root block has pointers to other blocks
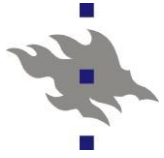  - Either piece of a file or filesystem metadata (e.g., directory)

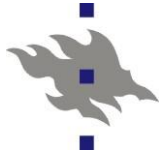n Data stored for an agreed-upon finite interval

# CFS Filesystem Structure



n Root block identified by publisher's public key

    n Each publisher has its own filesystem

    n Different publishers are on separate filesystems

n Other blocks identified based on hash of contents
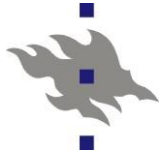
n Other blocks can be metadata or pieces of file

# Load Balancing and Quotas

n Different servers have different capabilities

n One real server can run several virtual servers

n Number of virtual servers depends on the capabilities

  n "Big" servers run more virtual servers

n CFS operates at virtual server level

  n Virtual nodes on a single real node know each other

  n Possible to use short cuts in routing


n Quotas can be used to limit storage for each node

n Quotas set on a per-IP address basis

  n Better quotas require central administration
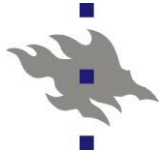
  n Some systems implement "better" quotas, e.g., PAST

# Updates in CFS

n Only publisher can change data in filesystem

n CFS will store any block under its content hash

- n Highly unlikely to find two blocks with same SHA-1 hash
- n No explicit protection for content blocks needed

n Root block is signed by publisher

- n Publisher must keep private key secret

n No explicit delete operation

- n Data stored only for agreed-upon period
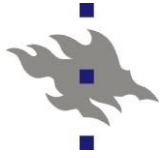- n Publisher must refresh periodically if persistence is needed

# OceanStore

n OceanStore developed at UC Berkeley

n Runs on Tapestry DHT

n Supports object modification

n Vision of ubiquitous computing:

  n Intelligent devices, transparently in the environment

  n Highly dynamic, untrusted environment

n Question: Where does persistent information reside?

n OceanStore aims to be the answer

n OceanStore's target:

  n $10^{10}$ users, each with 10000 files, i.e., $10^{14}$ files total
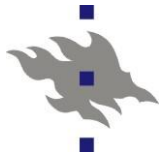
# OceanStore: Basics and Goals

- Users pay for the storage service
  - Several companies can provide services together
- Two goals:
1. Untrusted infrastructure
   - Everything is encrypted, infrastructure unreliable
   - However, assume that "most servers are working correctly most of the time"
   - One class of servers trusted to follow protocol (but not trusted with data)
2. Nomadic data
   - Anytime, anywhere
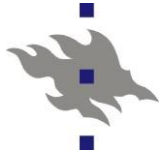   - *Introspection* used to tune system at run-time

# OceanStore Applications

n OceanStore suitable for many kinds of applications

n Storage and sharing of large amounts of data

  n Data follows users dynamically

n Groupware applications

  n Concurrent updates from many people

  n For example, calendars, contact lists, etc.

  n In particular, email and other communication applications

n Streaming applications

  n Also for sensor networks and dissemination


n Here we concentrate on storage

# System Overview

n Each object has globally unique identifier (GUID)

n Objects replicated and migrated on the fly

n Replicas located in two ways

  n Probabilistic, fast algorithm tried first

  n Slower, but deterministic algorithm used if first one fails

n Objects exist in two forms, active and archival

  n Active is the latest version with handle for updates

  n Archival is a permanent, read-only version

    - Archive versions encoded with erasure codes with lot of redundancy

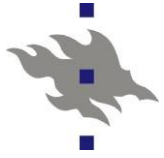    - Only a global disaster can make data disappear

# Naming and Access Control

Object naming

n Self-certifying object names

  n Hash of the object's owners key and human readable name

n Allows directories
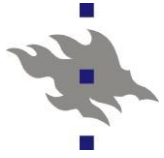
n System has no root, but user can select her own root(s)

Access control

n Restricting readers

  n All data is encrypted, can restrict readers by not giving key

  n Revoke read permission by re-encrypting everything

n Restricting writers

  n All writes must be signed and compared against ACL

# Locating Objects

n   OceanStore uses two mechanisms for locating objects

1.   Probabilistic algorithm

   n   Frequently accessed objects likely to be nearby and easily found

   n   This algorithm finds them fast

   n   Uses attenuated Bloom filters

   n   See below for more details

2.   Deterministic algorithm

   n   OceanStore based on Tapestry

   n   Deterministic routing is Tapestry's routing

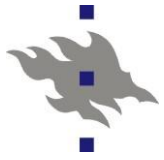   n   Guaranteed to find the object

   n   See Chapter 3 for the details

# Sidenote: Bloom Filters

n Bloom filters are "*a space-efficient probabilistic data structure that is used to test whether or not an element is a member of a set*"

  n False positives are possible

  n False negatives are NOT possible

n Bloom filter is an array of $k$ bits

  n Also need $m$ different hash functions, each maps key to a bit

n To insert, calculate all $m$ hash functions and set bits to 1

n To check, calculate all $m$ hash functions and if all bits are 1, key is "probably" in the set
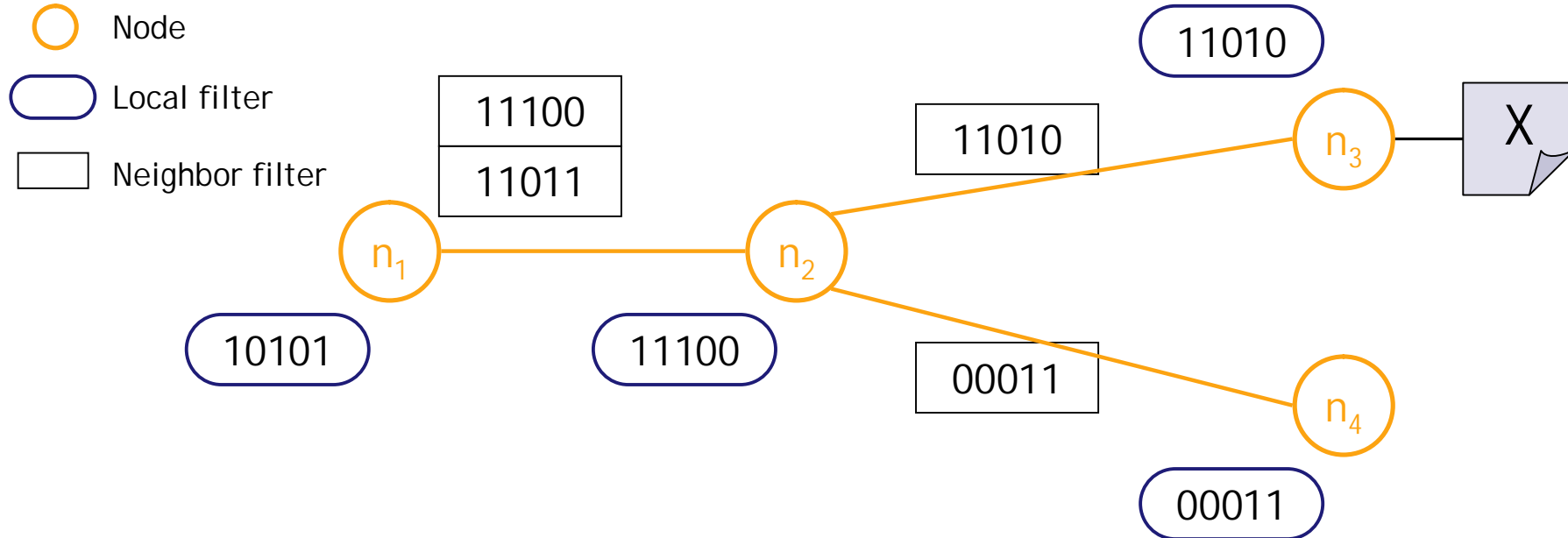
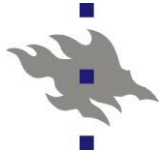  n If any bit is 0, then it is definitely not in

# Sidenote: Bloom Filters

n To insert or check for an item, Bloom filters take average $O(m)$ time

    n Fixed constant! Independent of number of entries

    n No other data structure allows for this (hash table is close)


n Attenuated Bloom filter of depth $D$ is same as an array of $D$ normal Bloom filters

    n First filter is for locally stored objects at current node

    n The $i^{th}$ Bloom filter is the union of all Bloom filters at distance $i$ through any path from current node

    n Attenuated Bloom filter for each network edge

    n Queries routed on the edge where the distance to object is shortest

# Probabilistic Query Process
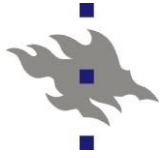
○ Node

⬭ Local filter

▭ Neighbor filter

11010

11100
11011

11010

$n_3$

X

10101

11100

$n_1$

$n_2$

00011

$n_4$

00011

n Node $n_1$ wants to find object X, X hashes to bits 0, 1, 3

n Node $n_1$ does not have 0, 1, and 3 in local filter

  n Neighbor filter for $n_2$ has them, forward query to $n_2$

n Node $n_2$ does not have them in local filter

  n Filter for neighbor $n_3$ has them, forward to $n_3$
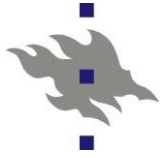
n Node $n_3$ has object

# Update Model

n Update model in OceanStore based on conflict resolution

n Update semantics

  n Each update has a list of predicates with associated actions

  n Predicates evaluated in order

  n Actions for first true predicate are atomically applied (commit)

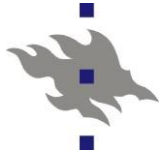  n If no predicates are true, action aborts

  n Update is logged in both cases

# Update Model: Predicates

n List of predicates is short

n Untrusted environment limits what predicates can do

n Available predicates:

  n Compare-version (metadata comparison, easy)

  n Compare-size (same as above)

  n Compare-block (easy if encryption is position-dependent block cipher)

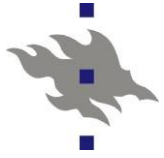  n Search (possible to search ciphertext, get boolean result)

# Available Operations

n Four operations available
   - n Replace-block
   - n Insert-block
   - n Delete-block
   - n Append

n If position-dependent cipher, Replace-block and Append are easy operations

n For Insert-block and Delete-block:
   - n Two kinds of blocks: Index and data blocks
   - n Index blocks can contain pointers to other blocks
   - n To insert a block, we replace old block with an index block which points to old block and new block
      - Actual blocks are appended to object
   - n May be susceptible to traffic analysis

# Serializing Updates

n Replicas divided into two tiers

- n Primary tier is trusted to follow protocol
- n Secondary tier is everyone else

n Primary tier cooperate in a Byzantine agreement protocol

n Secondary tier communicates with primary tier and secondary via epidemic algorithms

n Reason for two tiers:

- n Fault-tolerant protocols possible with only a small number of replicas, protocols communication-intensive
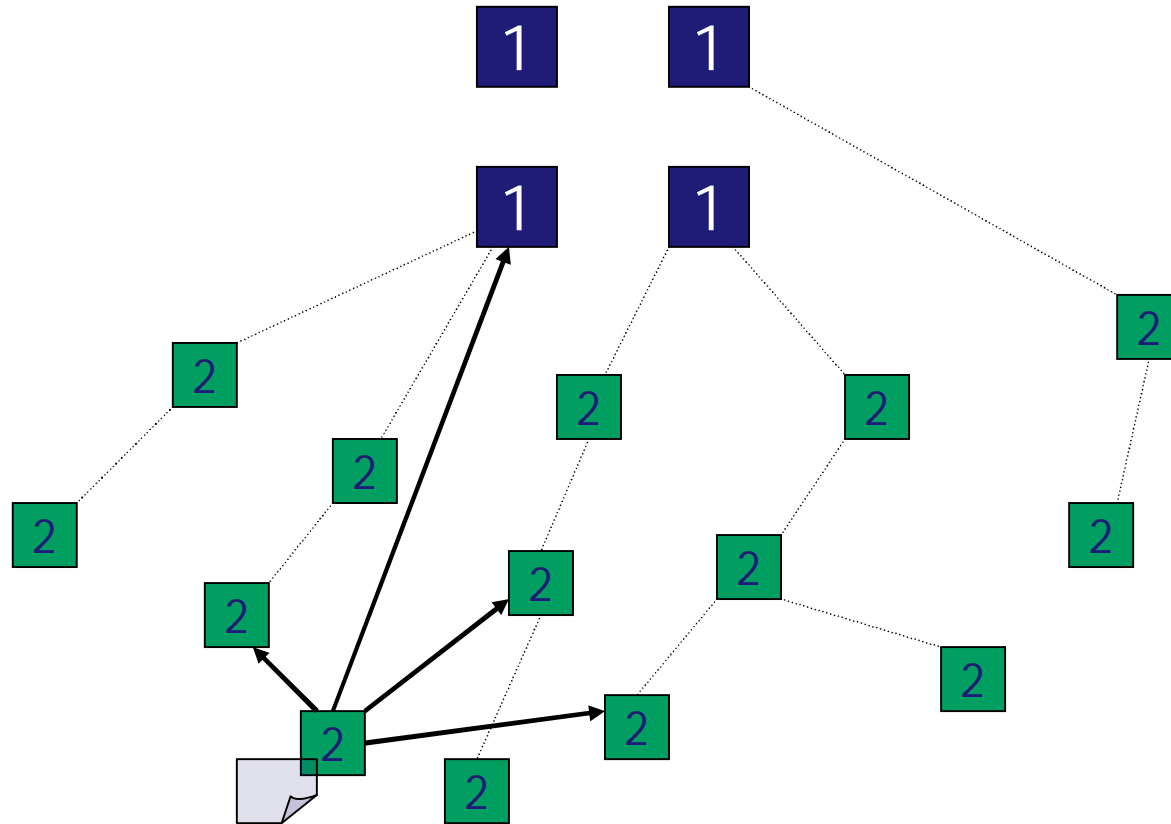- n Primary tier is well-connected and small

# Byzantine Generals Problem

n   Several divisions of the Byzantine army surround an enemy city. Each division is commanded by a general.

n   The generals communicate only through messenger

   n   Need to arrive at a common plan after observing the enemy

n   Some of the generals may be traitors

   n   Traitors can send false messages

n   Required: An algorithm to guarantee that

   1.   All loyal generals decide upon the same plan of action, irrespective of what the traitors do.

   2.   A small number of traitors cannot cause the loyal generals to adopt a bad plan.

Solution: (from L. Lamport)

n   If no more than $m$ generals out of $n = 3m + 1$ are traitors, everybody will follow the orders

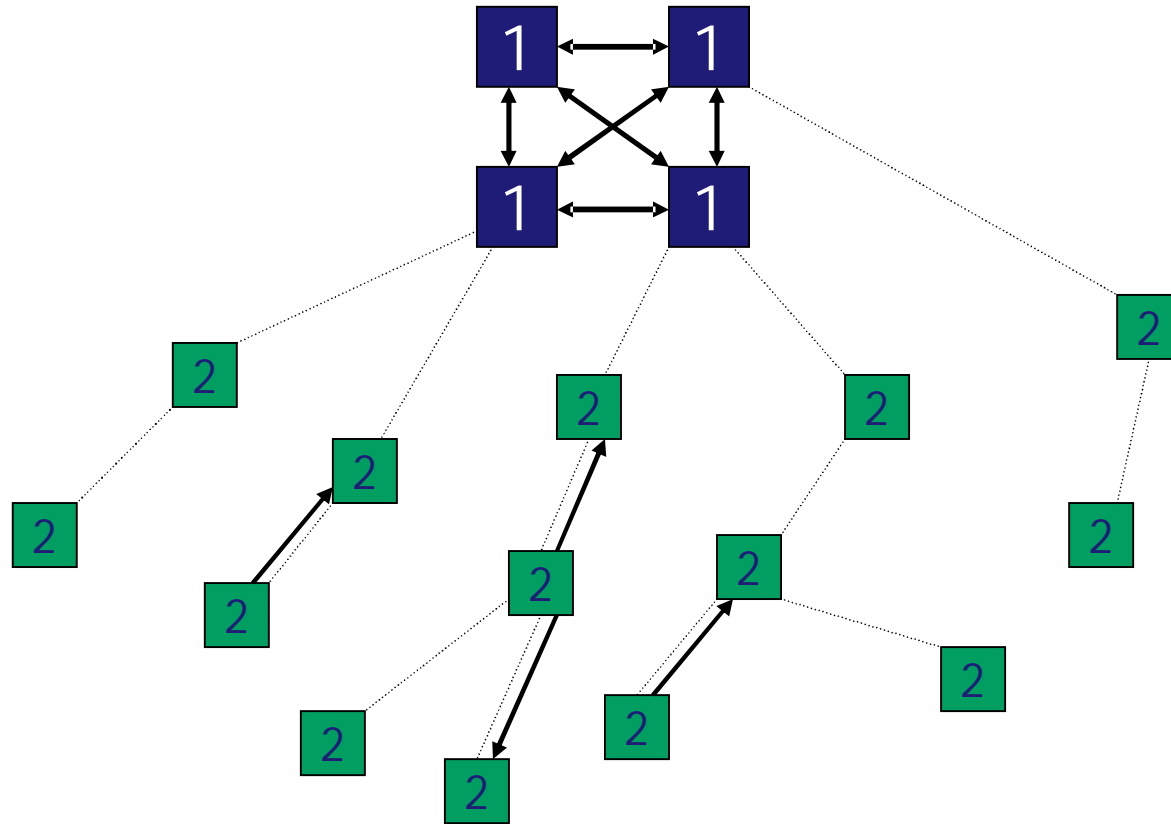# Path of an Update



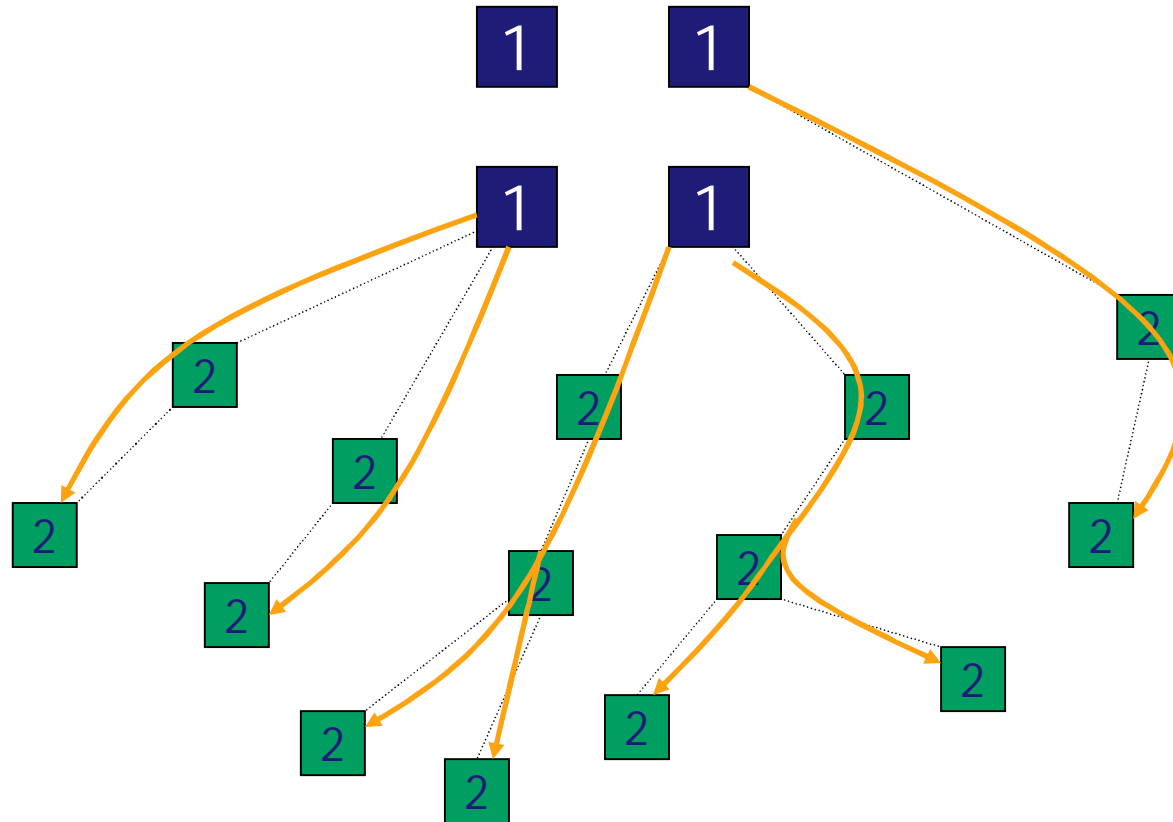n Update is sent to primary tier and to random replicas in secondary tier for that object
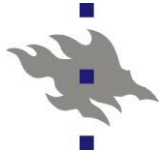
# Path of an Update



n Primary tier performs Byzantine agreement

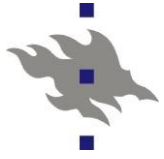n Secondary tier propagates update epidemically

# Path of an Update



n When primary tier has finished agreement protocol, the update is sent over multicast to all secondary replicas

# Deep Archival Storage
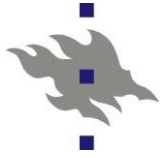
n Archival mechanism uses erasure codes

   n Reed-Solomon, Tornado, etc.

n Generate redundant data fragments

   n Created by the primary tier

n If there are enough fragments and they are spread widely, then it is likely that we can retrieve the data

n Archival copies are created when objects are changed

   n Every version is archived

   n Can be tuned to be done less frequently

# Ivy

- Ivy developed at MIT
- Based on Chord
- Provides NFS-like semantics
    - At least for fully connected networks

- Any user can modify any file
- Ivy handles everything through logs
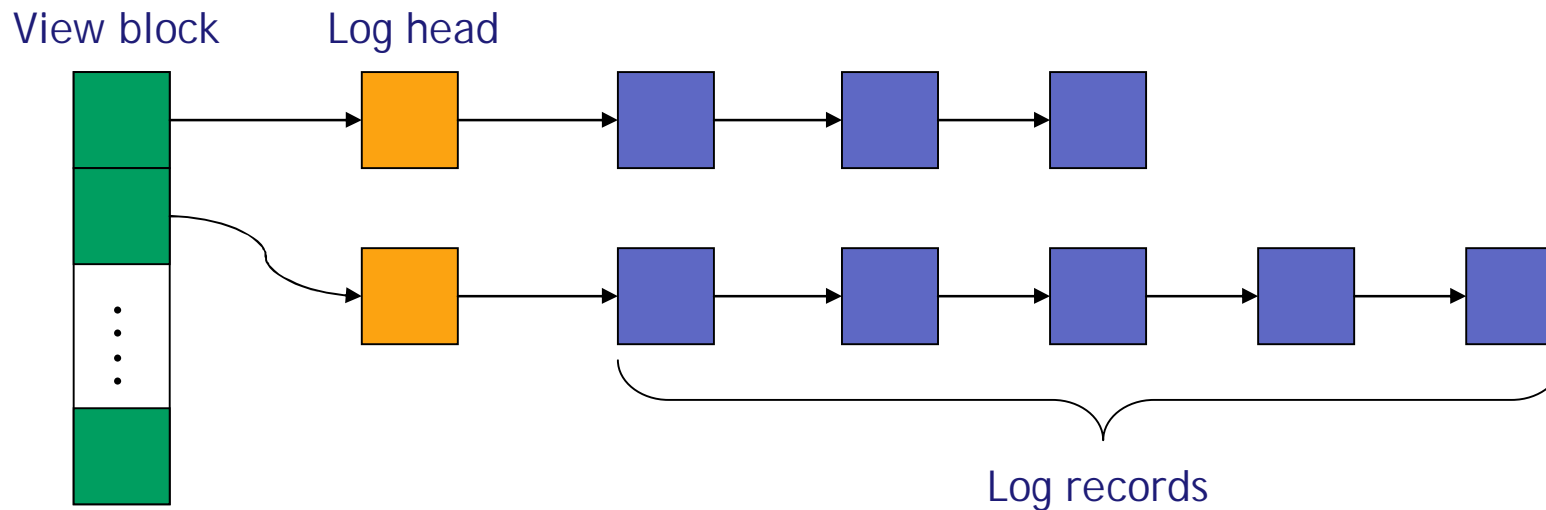
- Ivy presents a conventional filesystem interface
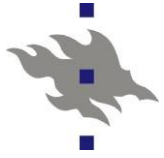
# Problems for Distributed Read/Write

1. Multiple distributed writers make it difficult to maintain consistent metadata
2. Unreliable participants make locking unattractive
   - Locking could help maintain consistency
3. Participants cannot be trusted
   - Machines may be compromised
   - Need to be able to undo
4. Distributed filesystem may become partitioned
   - System must remain operational during partitions
   - Help applications repair conflicting updates made during partitions
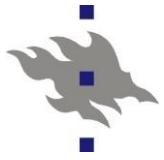
# Solution: Logs

- Each participant maintains log of its changes
  - Logs maintain filesystem data and metadata
  - Participant has private snapshot of logs of others
- Logs stored in DHash (see under CFS for details)
- Participant writes to its own log, reads all others
  - Log-head points to most recent log record
- Version vectors impose order on log records from multiple logs
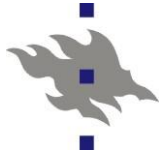
View block          Log head

Log records

# Ivy: Views

n Each user writing to a filesystem has its own log

n Participants agree on a view

    n View is a set of logs that comprise the filesystem

n View block is immutable ("root")
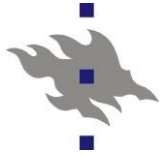
n View block has log heads for all participants

# Combining Logs

n How to determine order of modifications from logs?

   n Order should obey causality

   n All participants should agree on the order

n Each new log record has a sequence number and a version vector

   n Sequence number increasing (managed locally)

   n Version vector has sequence numbers from other logs in view

   n Version vector summarizes knowledge about log

n Log records ordered by comparing version vectors

   n Vectors $u$ and $v$ comparable if $u < v$, $v < u$, or $v = u$

   n Otherwise concurrent

n Simultaneous operations result in equal or concurrent vectors

   n Ordered by public keys of participants

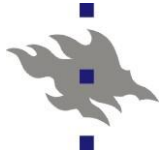   n May need special actions to return to consistency (overlapping modifications)

# Ivy: Snapshots

n Private snapshots avoid traversing whole filesystem

n Snapshot contains the entire state

n Each participant has own snapshot

- n Contents of snapshots mostly the same for all participants
- n DHash will store them only once

n To create snapshot, node:

- n Gets all logs recent than current snapshot
- n Write new snapshot

n New user must either build from scratch or take a trusted snapshot
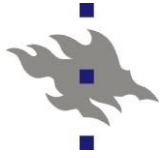
# P2P FS: General Problems

n Built on top of unreliable nodes and network

n How to achieve reliability?

>  n Replication gives reliability (see Chapter 5)
>
>  n Replication makes maintaining consistency difficult

n Nodes cannot be trusted in the general case

>  n Must encrypt all data
>
>  n Hard to do content-based conflict resolution (e.g., diff)

n Performance

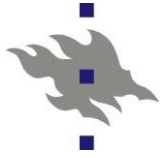>  n Distributed filesystems have much lower performance

# P2P FS: Future

n What does future hold for P2P filesystems?

n What is right area of application?

n Intranet?

    n Trusted environment, high bandwidth

    n Possibly easy to deploy?

       - Need to make a product first?

n Global Internet?

    n Lots of untrusted peers, widely distributed

    n All the problems from above

    n Can take off as a "hobby" project?
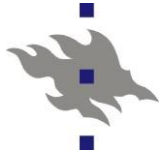
n Nowhere?

    n P2P filesystems are a total waste of time?

# P2P FS: Future

n Do we need a P2P FS to build useful applications?

n Yes, they allow efficient distributed storage

  n Working storage system is the basic building block of a useful application

n No, DHT is enough

  n Several examples of P2P applications directly on top of a DHT

n Fully reliable P2P FS would make network into one virtual computer

  n Modulo performance issues

  n Building P2P apps would be like building normal apps

# P2P FS: Real-World Examples

n Microsoft has done lot of research in this area

  n Even built prototypes

n Why no products on the market?

  n Below some wild speculation

n P2P FS would compete with traditional file servers?

  n P2P needs to be built in to the OS, would kill server market?

n Impossible to build a "good enough" P2P FS?

  n Theoretically doable, but too slow and weak in practice?

n P2P filesystem can be done, but has no advantages?

  n Possibly to build a useful system, but it costs as much as a server-based system?

# Chapter Summary

n How to build applications with DHTs

n Basic technologies of distributed storage

n As examples, 3 P2P filesystems

  n CFS

    - Read-only

  n OceanStore

    - Modifications allowed, global vision

  n Ivy

    - NFS-like semantics, traditional filesystem interface

n Discussion about the future of P2P filesystems