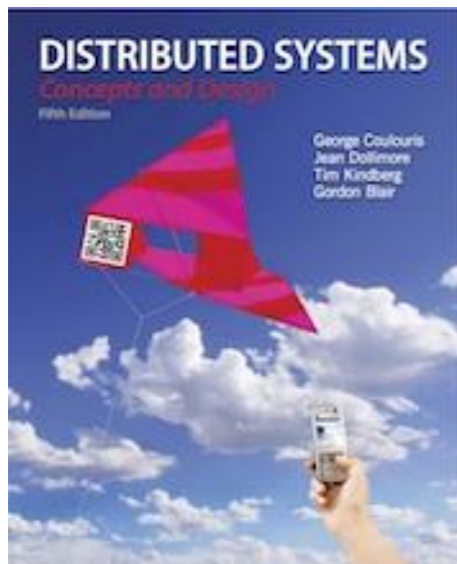


Slides for Chapter 10: Peer-to-Peer Systems



From **Coulouris, Dollimore, Kindberg and Blair**
**Distributed Systems:
Concepts and Design**

Edition 5, © Addison-Wesley 2012

Overview of Chapter

- Introduction
- Napster and its legacy
- Peer-to-peer middleware
- Routing overlays
- Overlay case studies: Pastry, Tapestry
- Application case studies: Squirrel, OceanStore, Ivy

Introduction

Goal:

- Provide fully decentralized and self-organizing, dynamically balancing the storage and processing loads as nodes join and leave

Characteristics of peer-to-peer systems:

- Each user contributes resources (files, computing cycles, etc.)
- Each node has similar functional capability
- No centrally administered systems
- Provide anonymity to providers and users of resources
- Require algorithms for data placement, and for workload balances so that nodes do not suffer undue overhead

Introduction (cont.)

Three generations of peer-to-peer systems:

- Napster (music exchange)
- Freenet, Gnutella, Kazaa, BitTorrent (file sharing)
- Pastry, Tapestry, CAN, Chord, Kademlia (peer-to-peer middleware)
- Resources identified by GUIDs (Globally Unique Identifiers) using secure hashing
- Suitable for storing immutable objects (music, video)
- Overlay (application-level) routing used instead of IP routing

Figure 10.1: Distinctions between IP and overlay routing for peer-to-peer applications

	<i>IP</i>	<i>Application-level routing overlay</i>
<i>Scale</i>	IP v4 is limited to 2 ³² addressable nodes. The IP v6 name space is much more generous (2 ¹²⁸), but addresses in both versions are hierarchically structured and much of the space is pre-allocated according to administrative requirements.	Peer-to-peer systems can address more objects. The GUID name space is very large and flat (>2 ¹²⁸), allowing it to be much more fully occupied.
<i>Load balancing</i>	Loads on routers are determined by network topology and associated traffic patterns.	Object locations can be randomized and hence traffic patterns are divorced from the network topology.
<i>Network dynamics (addition/deletion of objects/nodes)</i>	IP routing tables are updated asynchronously on a best-efforts basis with time constants on the order of 1 hour.	Routing tables can be updated synchronously or asynchronously with fractions of a second delays.
<i>Fault tolerance</i>	Redundancy is designed into the IP network by its managers, ensuring tolerance of a single router or network connectivity failure. <i>n</i> -fold replication is costly.	Routes and object references can be replicated <i>n</i> -fold, ensuring tolerance of <i>n</i> failures of nodes or connections.
<i>Target identification</i>	Each IP address maps to exactly one target node.	Messages can be routed to the nearest replica of a target object.
<i>Security and anonymity</i>	Addressing is only secure when all nodes are trusted. Anonymity for the owners of addresses is not achievable.	Security can be achieved even in environments with limited trust. A limited degree of anonymity can be provided.

Introduction (cont.)

- Can also be used for distributed computation (e.g. SETI)

Napster and its legacy

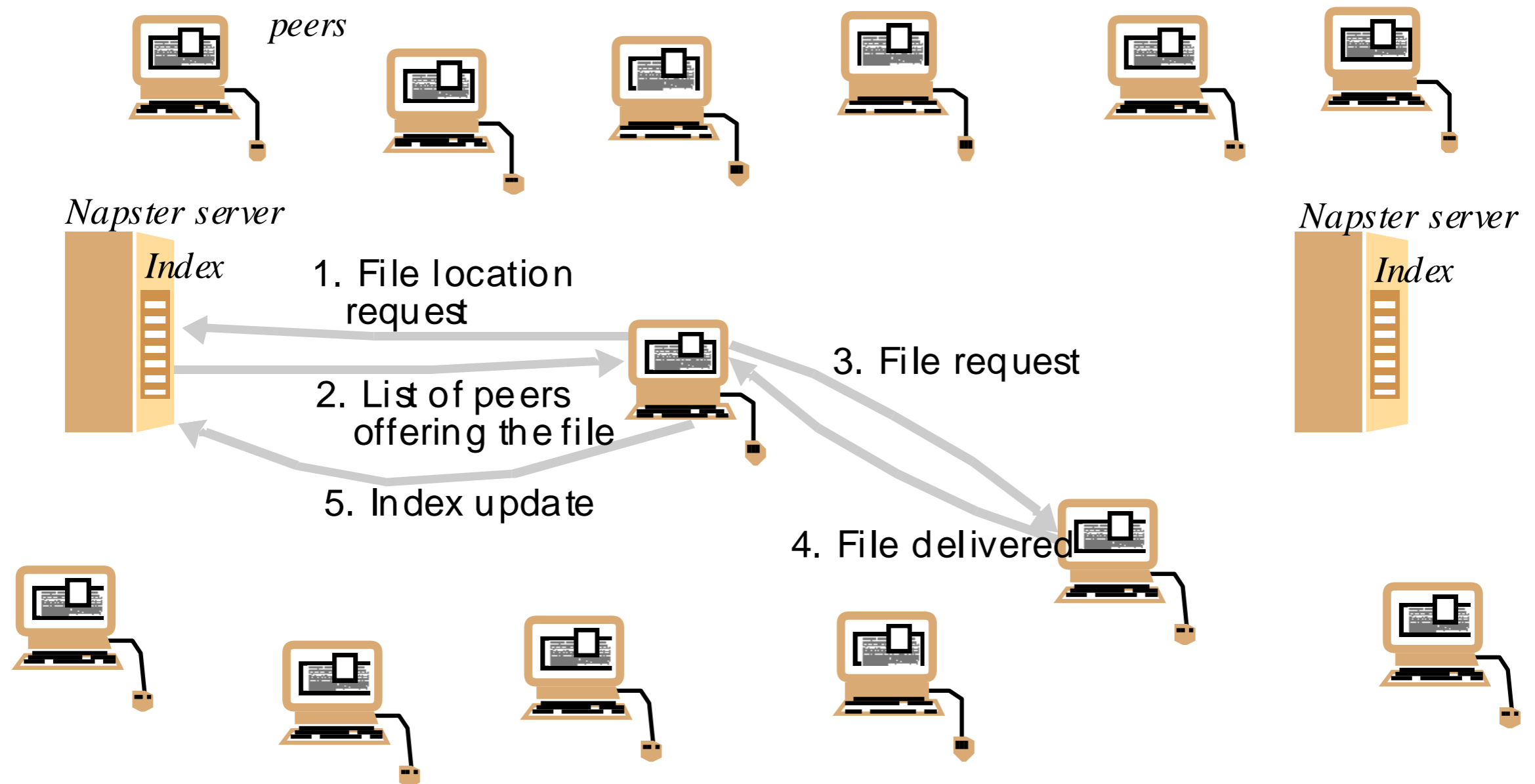
Goal:

- Distributed file sharing system
- Users supply files for sharing (stored on their own computers)
- Launched in 1999, had several million users
- Closed for music copyright infringement violations

System architecture:

- Centralized indexing system for locating files
- Users share file by linking to indexing service
- Load distribution mechanism to locate closest file copy to a requester
- Assumes files are static (do not change)
- Does not worry about consistency of replicas (of same song)

Figure 10.2: Napster: peer-to-peer file sharing with a centralized, replicated index



Peer-to-peer middleware

Functional requirements:

- Simplify the construction of services that are implemented across many hosts
- Enable client to locate any resource
- Add and remove resources dynamically
- Add and remove hosts (computers)

Non-functional requirements:

- Global scalability
- Load balancing
- Optimizing local interactions among neighboring peers
- Highly dynamic host availability
- Security, anonymity, deniability, resistance to censorship

Routing overlays

Requirements:

- Distributed algorithm responsible for locating nodes and objects
- Routing algorithm in the *application layer*, different than network layer (IP) routing
- Objects replicated and placed on nodes and can be relocated without client involvement
- Routing can locate 'nearest' copy of desired object
- GUID an example of a 'pure name' that does not reveal object location

Main tasks of routing overlay:

- Routes request to access object via its GUID to a replica
- Adding a new object requires computing a new GUID and announces it
- Deleting an object makes all copies unavailable
- Adding and removing new nodes

Routing overlays

GUID:

- Computed from part of object state (e.g. its name) using hashing
- Each GUID must be unique
- Sometimes called DHT (Distributed Hash Tables)

Distributed object location and routing (DOLR) layer:

- Maintains mapping between GUIDs and nodes where object is stored
- DOLR approach separated locating object from other routing functions
- DOLR may or may not be used

Figure 10.3: Distribution of information in a routing overlay

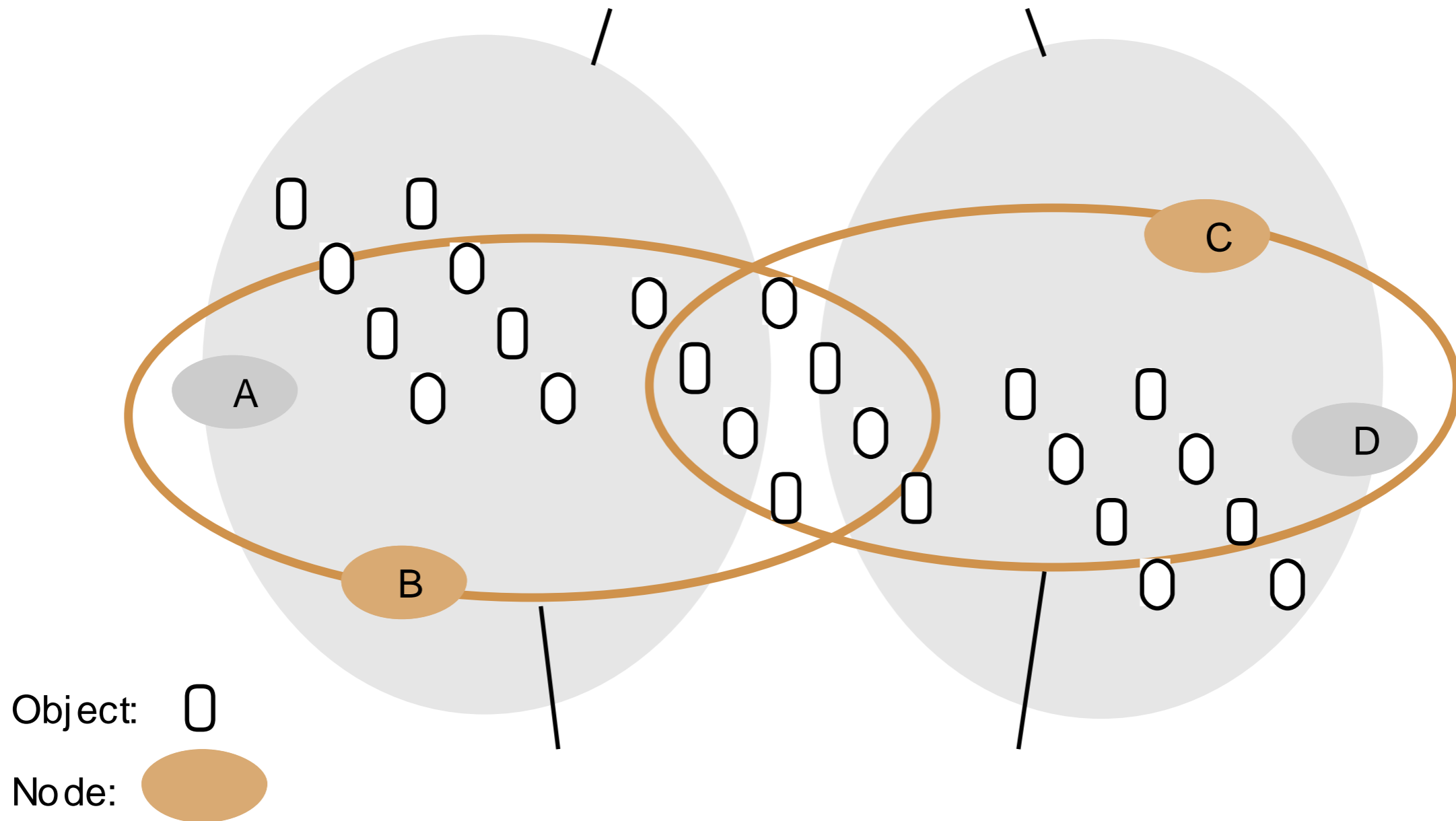


Figure 10.4: Basic programming interface for a distributed hash table (DHT) as implemented by the PAST API over Pastry

put(*GUID*, *data*)

The *data* is stored in replicas at all nodes responsible for the object identified by *GUID*.

remove(*GUID*)

Deletes all references to *GUID* and the associated data.

value = *get*(*GUID*)

The data associated with *GUID* is retrieved from one of the nodes responsible it.

Figure 10.5: Basic programming interface for distributed object location and routing (DOLR) as implemented by Tapestry

publish(*GUID*)

GUID can be computed from the object (or some part of it, e.g. its name). This function makes the node performing a *publish* operation the host for the object corresponding to *GUID*.

unpublish(*GUID*)

Makes the object corresponding to *GUID* inaccessible.

sendToObj(*msg*, *GUID*, [*n*])

Following the object-oriented paradigm, an invocation message is sent to an object in order to access it. This might be a request to open a TCP connection for data transfer or to return a message containing all or part of the object's state. The final optional parameter [*n*], if present, requests the delivery of the same message to *n* replicas of the object.

Pastry

- Routing overlay with 128-bit GUIDs
- GUID computed by a secure hash function such as SHA-1 (see Chapter 11) using public key of node
- Typically hash function is applied to the object name (or another known part of the object state)
- $0 \leq \text{GUID} \leq 2^{128} - 1$

Routing performance:

- Order of $\log N$ steps when N nodes participate in system
- Routing overlay built over UDP
- Network distance between nodes based on hop counts and round trip latency – used to set up routing tables at each node

Pastry

Routing algorithm stage 1:

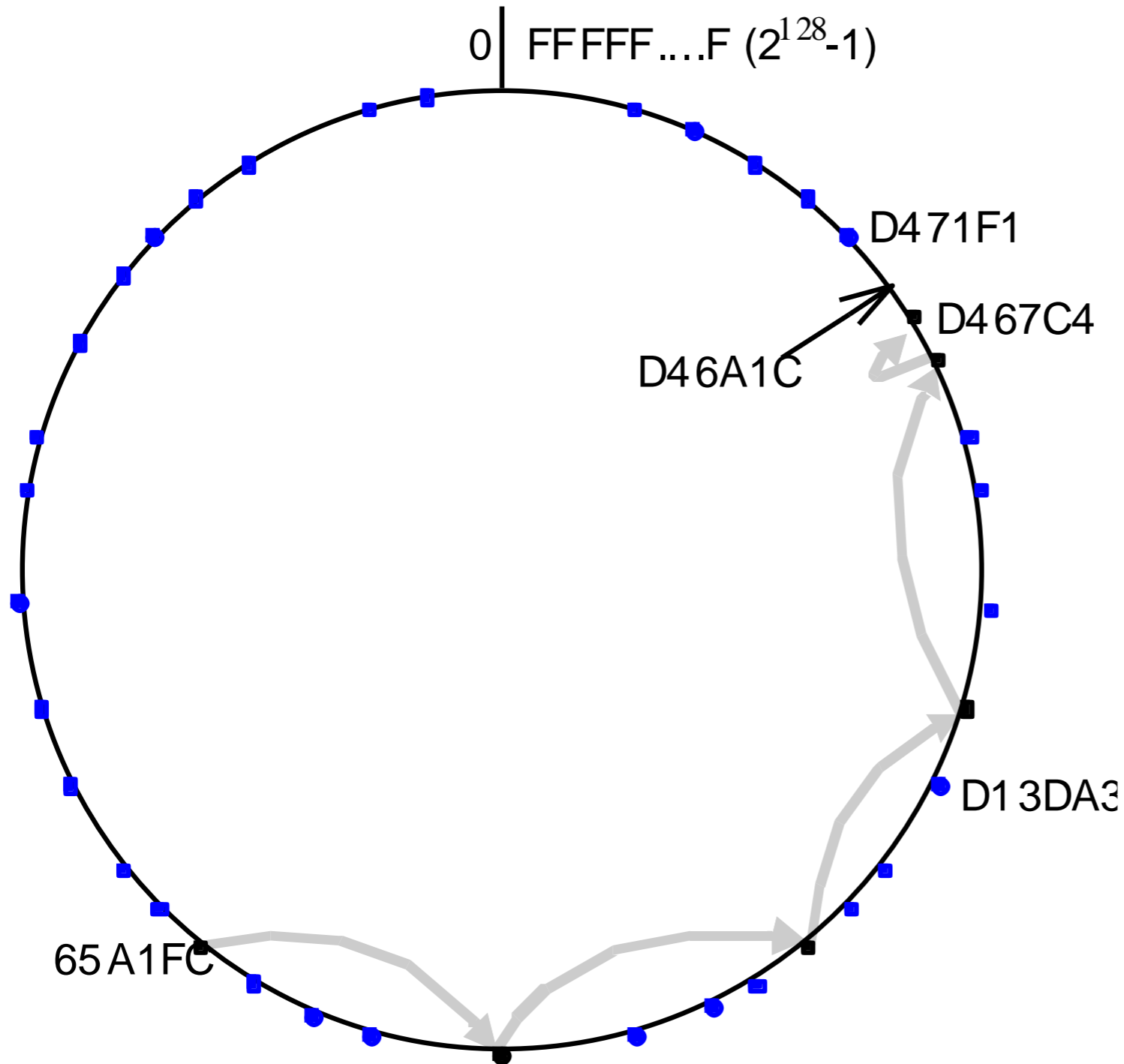
- Each node keeps a *leaf set* contains nodes closest to current node
- GUID space can be visualized as a circle
- If node A receives a routing request for node D, it finds closest node in its leaf set to D to forward the message

Routing algorithm stage 2:

- Each node keeps a *tree-structured* routing table
- Entries include (GUID, IP address) of a set of nodes

Figure 10.6: Circular routing alone is correct but inefficient

Based on Rowstron and Druschel [2001]



The dots depict live nodes. The space is considered as circular: node 0 is adjacent to node $(2^{128}-1)$. The diagram illustrates the routing of a message from node 65A1FC to D46A1C using leaf set information alone, assuming leaf sets of size 8 ($l = 4$). This is a degenerate type of routing that would scale very poorly; it is not used in practice.

Figure 10.7: First four rows of a Pastry routing table

$p =$	<i>GUID prefixes and corresponding nodehandles n</i>															
0	<i>0</i>	<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>	<i>6</i>	<i>7</i>	<i>8</i>	<i>9</i>	<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>	<i>F</i>
	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>		<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>
1	<i>60</i>	<i>61</i>	<i>62</i>	<i>63</i>	<i>64</i>	<i>65</i>	<i>66</i>	<i>67</i>	<i>68</i>	<i>69</i>	<i>6A</i>	<i>6B</i>	<i>6C</i>	<i>6F</i>	<i>6E</i>	<i>6F</i>
	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>		<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>
2	<i>650</i>	<i>651</i>	<i>652</i>	<i>653</i>	<i>654</i>	<i>655</i>	<i>656</i>	<i>657</i>	<i>658</i>	<i>659</i>	<i>65A</i>	<i>65B</i>	<i>65C</i>	<i>65D</i>	<i>65E</i>	<i>65F</i>
	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>		<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>
3	<i>65A0</i>	<i>65A1</i>	<i>65A2</i>	<i>65A3</i>	<i>65A4</i>	<i>65A5</i>	<i>65A6</i>	<i>65A7</i>	<i>65A8</i>	<i>65A9</i>	<i>65AA</i>	<i>65AB</i>	<i>65AC</i>	<i>65AD</i>	<i>65AE</i>	<i>65AF</i>
	<i>n</i>		<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>	<i>n</i>

The routing table is located at a node whose GUID begins 65A1. Digits are in hexadecimal. The n 's represent [GUID, IP address] pairs specifying the next hop to be taken by messages addressed to GUIDs that match each given prefix. Grey-shaded entries indicate that the prefix matches the current GUID up to the given value of p : the next row down or the leaf set should be examined to find a route. Although there are a maximum of 128 rows in the table, only $\log_{16} N$ rows will be populated on average in a network with N active nodes.

Figure 10.8: Pastry routing example Based on Rowstron and Druschel [2001]

Routing a message from node 65A1FC to D46A1C. With the aid of a well-populated routing table the message can be delivered in $\sim \log_6(N)$ hops.

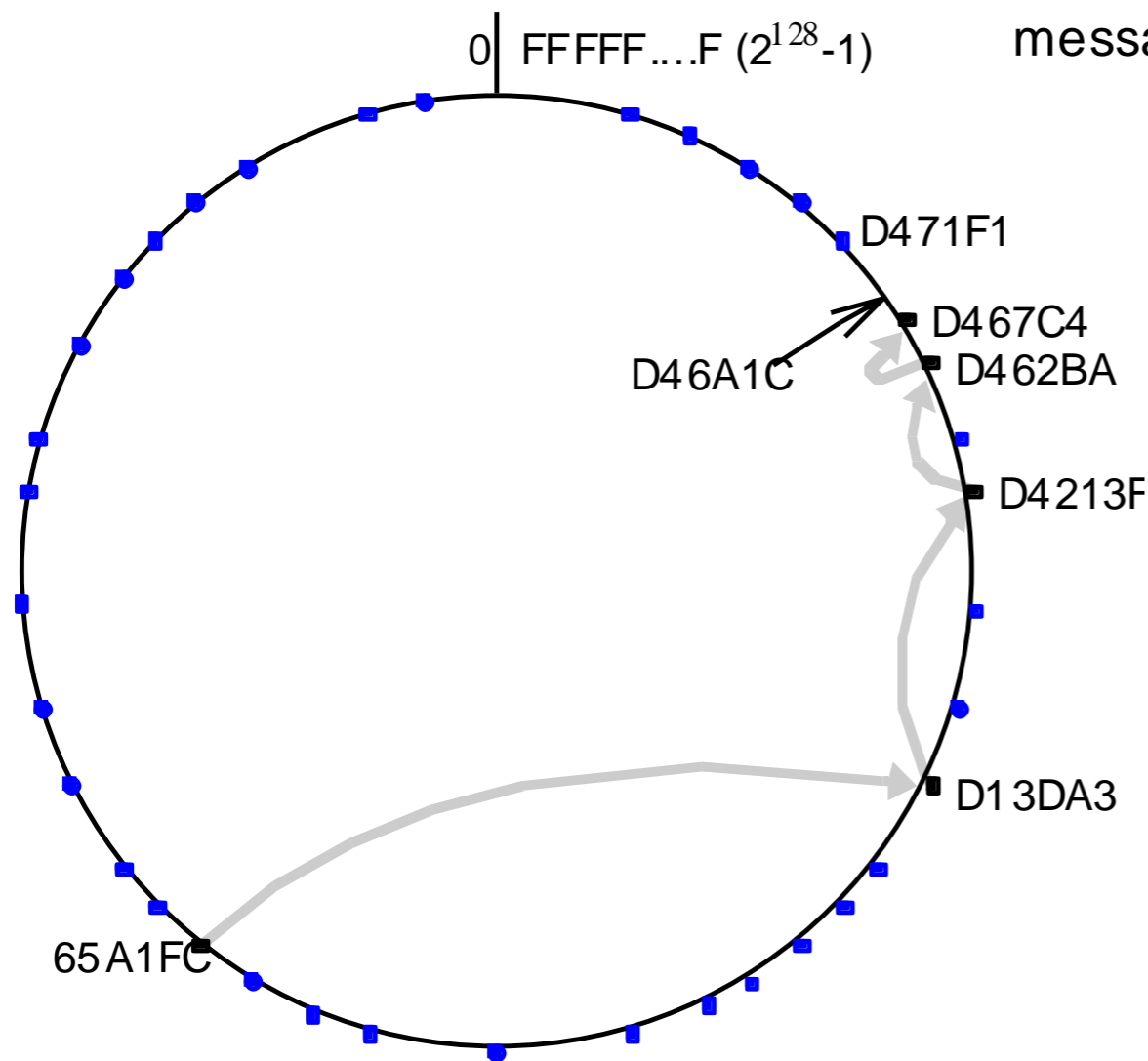


Figure 10.9: Pastry's routing algorithm

To handle a message M addressed to a node D (where $R[p, i]$ is the element at column i , row p of the routing table):

1. If $(L_{-1} < D < L_l)$ { // the destination is within the leaf set or is the current node.
2. Forward M to the element L_i of the leaf set with GUID closest to D or the current node A .
3. } else { // use the routing table to dispatch M to a node with a closer GUID
4. find p , the length of the longest common prefix of D and A . and i , the $(p+1)^{\text{th}}$ hexadecimal digit of D .
5. If $(R[p, i] \neq null)$ forward M to $R[p, i]$ // route M to a node with a longer common prefix.
6. else { // there is no entry in the routing table
7. Forward M to any node in L or R with a common prefix of length i , but a GUID that is numerically closer.
- }
- }
- }

Tapestry

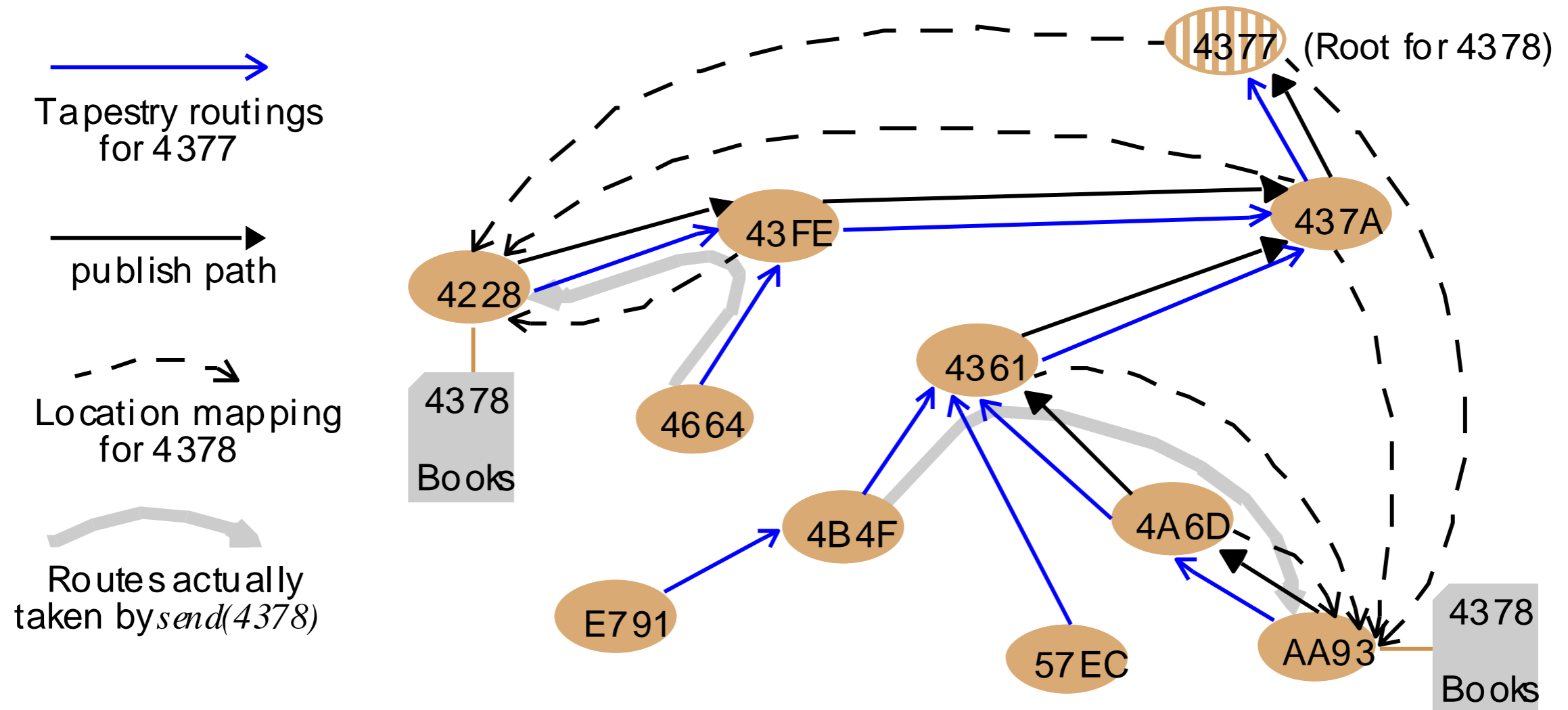
- Uses DOLR to hide the distributed hash table from applications
- Use 160-bit identifiers
- $0 \leq \text{GUID} \leq 2^{160} - 1$
- Identifiers refer both to objects (GUID) and to nodes (NodeId)

Structured versus unstructured peer-to-peer:

- Pastry, tapestry known as *structured* peer-to-peer because they maintain hash table
- In highly dynamic systems, overhead of maintaining hash tables
- Unstructured peer-to-peer tries to reduce overhead of central control
- Joining node establishes connections with local neighbors –

Figure 10.10: Tapestry routing

From [Zhao et al. 2004]



Replicas of the file *Phil's Books* ($G=4378$) are hosted at nodes 4228 and AA93. Node 4377 is the root node for object 4378. The Tapestry routings shown are some of the entries in routing tables. The publish paths show routes followed by the publish messages laying down cached location mappings for object 4378. The location mappings are subsequently used to route messages sent to 4378.

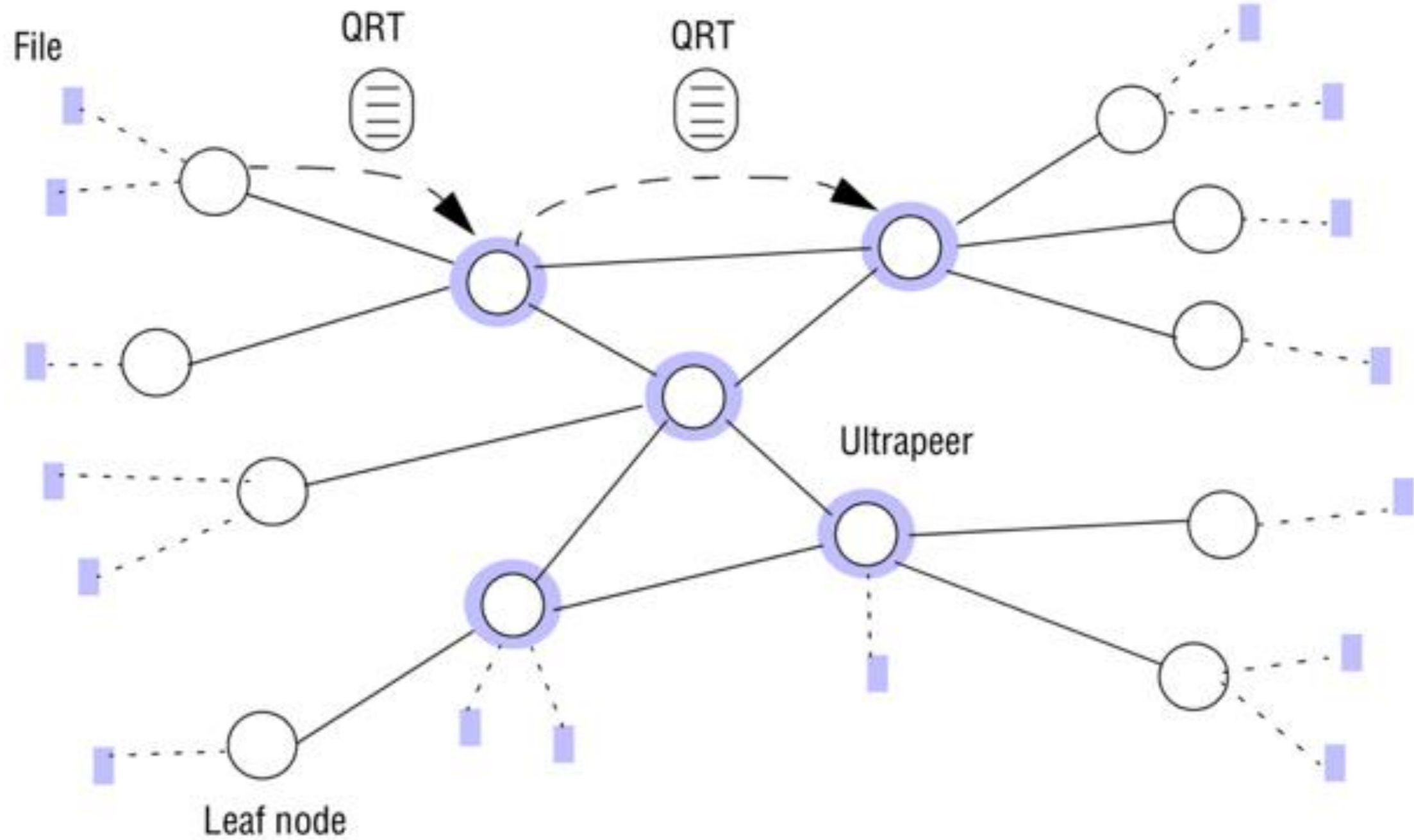
Figure 10.11: Structured versus unstructured peer-to-peer systems

	<i>Structured peer-to-peer</i>	<i>Unstructured peer-to-peer</i>
<i>Advantages</i>	Guaranteed to locate objects (assuming they exist) and can offer time and complexity bounds on this operation; relatively low message overhead.	Self-organizing and naturally resilient to node failure.
<i>Disadvantages</i>	Need to maintain often complex overlay structures, which can be difficult and costly to achieve, especially in highly dynamic environments.	Probabilistic and hence cannot offer absolute guarantees on locating objects; prone to excessive messaging overhead which can affect scalability.

Gnutella

- Successful file sharing system
- Divides node types into regular peers (leaves) and ultrapeers
- *Ultrapeers* form the heart of the network
- Leaves connect to ultrapeers, which do the routing
- Ultrapeers heavily connected to each other to reduce number of hops

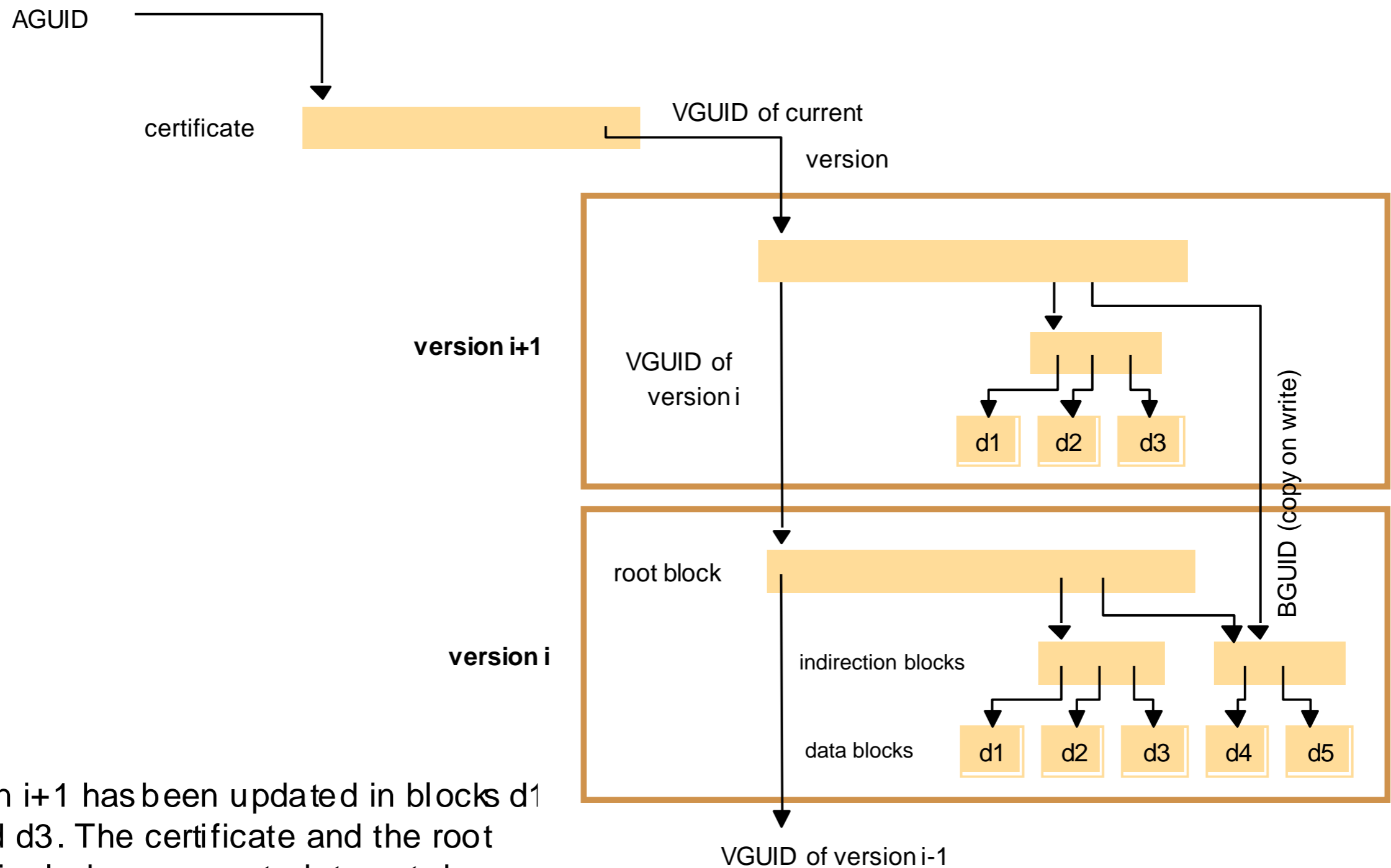
Figure 10.12: Key elements in the Gnutella protocol



OceanStore

- Peer-to-peer store for *mutable* files
- Applications may include distributed file system, email hosting, etc.
- Both mutable and immutable objects can be replicated
- Prototype called Pond

Figure 10.13: Storage organization of OceanStore objects



Version i+1 has been updated in blocks d1 d2 and d3. The certificate and the root blocks include some metadata not shown. All unlabelled arrows are BGUIDs

Figure 10.14: Types of identifier used in OceanStore

<i>Name</i>	<i>Meaning</i>	<i>Description</i>
BGUID	block GUID	Secure hash of a data block
VGUID	version GUID	BGUID of the root block of a version
AGUID	active GUID	Uniquely identifies all the versions of an object

Figure 10.15: Performance evaluation of the Pond prototype emulating NFS

<i>Phase</i>	<i>LAN</i>		<i>WAN</i>		<i>Predominant operations in benchmark</i>
	<i>Linux NFS</i>	<i>Pond</i>	<i>Linux NFS</i>	<i>Pond</i>	
1	0.0	1.9	0.9	2.8	Read and write
2	0.3	11.0	9.4	16.8	Read and write
3	1.1	1.8	8.3	1.8	Read
4	0.5	1.5	6.9	1.5	Read
5	2.6	21.0	21.5	32.0	Read and write
Total	4.5	37.2	47.0	54.9	

Ivy

- Emulates SUN NFS file system
- Supports multiple readers and writers over an overlay routing layer
- Distributed hash-addressed data store
- Stores logs of file updates and can reconstruct current state of file from log
- Log records stored in Dhash distributed hash addressed storage service

Figure 10.16: Ivy system architecture

