

Arpeggio:

Metadata Searching and Content Sharing with Chord

Austin T. Clements, Dan R. K. Ports, David R. Karger

{aclements, drkp, karger}@mit.edu

MIT Computer Science and Artificial Intelligence Laboratory



Outline

- **Overview**
- Searching
- Index Gateways
- Content Distribution
- Conclusion

The Content-Sharing Problem

• Goals

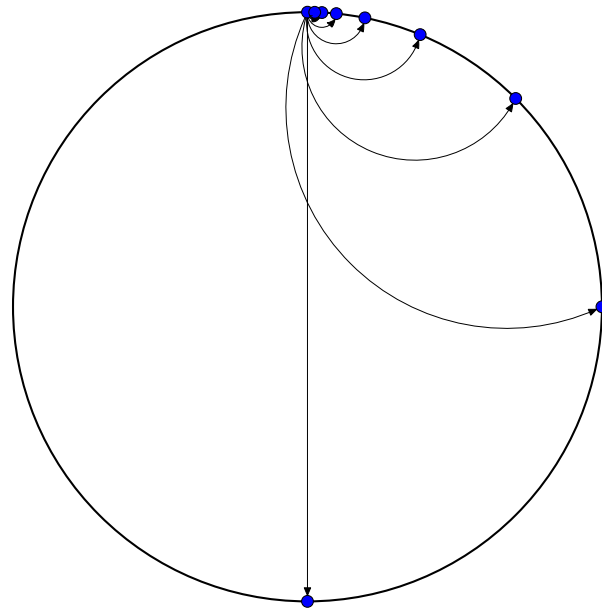
- Find files matching a search query
- Identify sources for a file
- Want full decentralization

• Assumptions

- Only searching *metadata*
- Metadata is small (compared to actual data)
- Highly dynamic and unstable network topology, content, and sources

DHTs (Almost) to the Rescue

- Great for lookup-by-name
- Insufficient for efficient search-by-content
- Powerful underlying LOOKUP abstraction



Outline

- Overview
- **Searching**
- Index Gateways
- Content Distribution
- Conclusion

Index Entries

(debian, disk1, iso)



name Debian Disk1.iso


file ID cdb79ca3db1f39b1940ed5...

size 586MB


type application/x-iso9660-image

⋮

Distributed Inverted Indexing


(debian, disk1, iso)	
name	Debian Disk1.iso
file ID	cdb79ca3db1f3...


Distributed Inverted Indexing

(debian, disk1, iso) 	
name	Debian Disk1.iso
file ID	cdb79ca3db1f3...




Distributed Inverted Indexing


(debian, disk1, iso)	
name	Debian Disk1.iso
file ID	cdb79ca3db1f3...


(debian, disk2, iso)	
name	Debian Disk2.iso
file ID	5ccbf54d7e502...

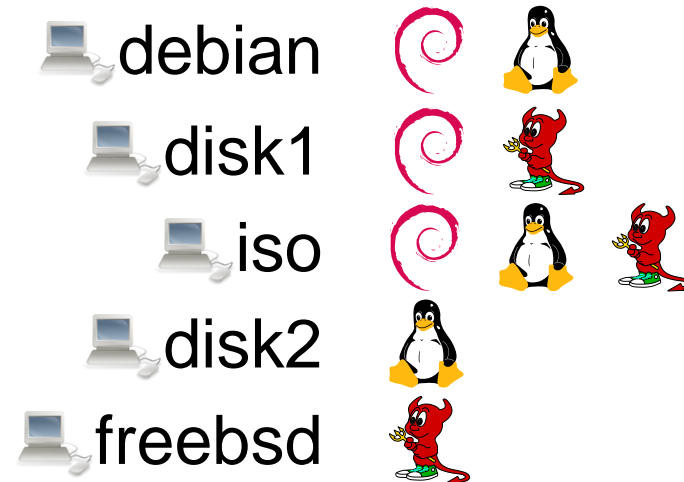


Distributed Inverted Indexing


(debian, disk1, iso)	
name	Debian Disk1.iso
file ID	cdb79ca3db1f3...


(debian, disk2, iso)	
name	Debian Disk2.iso
file ID	5ccbf54d7e502...


(disk1, freebsd, iso)	
name	Freebsd Disk1.iso
file ID	fbcbfdff31f27de...

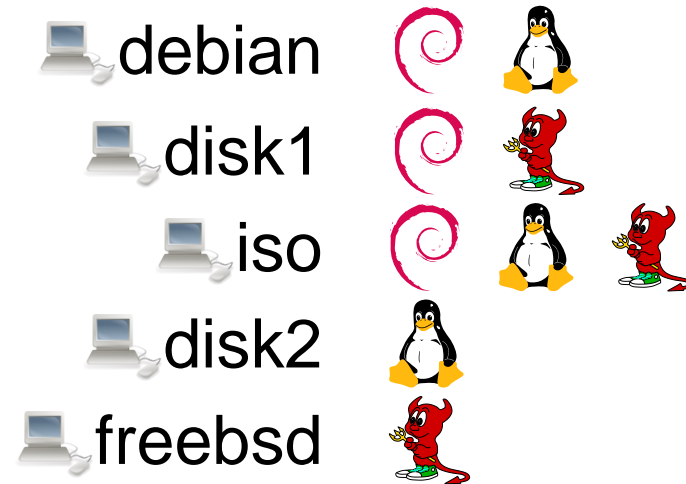


Distributed Inverted Indexing

(debian, disk1, iso)	
name	Debian Disk1.iso
file ID	cdb79ca3db1f3...


(debian, disk2, iso)	
name	Debian Disk2.iso
file ID	5ccbf54d7e502...


(disk1, freebsd, iso)	
name	Freebsd Disk1.iso
file ID	fbcbfdff31f27de...




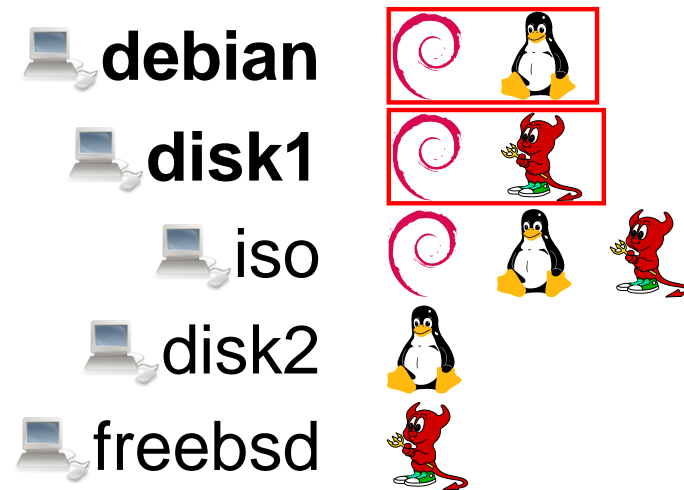
 “debian disk1?”

Distributed Inverted Indexing


(debian, disk1, iso) 
name Debian Disk1.iso
file ID cdb79ca3db1f3...


(debian, disk2, iso) 
name Debian Disk2.iso
file ID 5ccbf54d7e502...


(disk1, freebsd, iso) 
name FreeBSD Disk1.iso
file ID fbcbfdf31f27de...



Distributed Inverted Indexing

(debian, disk1, iso)	
name	Debian Disk1.iso
file ID	cdb79ca3db1f3...

(debian, disk2, iso)	
name	Debian Disk2.iso
file ID	5ccbf54d7e502...

(disk1, freebsd, iso)	
name	Freebsd Disk1.iso
file ID	fbcbfdff31f27de...





Problem: Network hoarding


Index-Side Filtering

- Keywords are small, so store keywords in index
- Pick one index node
- Send full query
- Index node performs filtering and returns only relevant results
- Can also include other *filterable metadata*, e.g. file size, MP3 bitrate, etc.

Index-Side Filtering

(debian, disk1, iso)	
name	Debian Disk1.iso
file ID	cdb79ca3db1f3...

(debian, disk2, iso)	
name	Debian Disk2.iso
file ID	5ccbf54d7e502...

(disk1, freebsd, iso)	
name	Freebsd Disk1.iso
file ID	fbcbfdff31f27de...

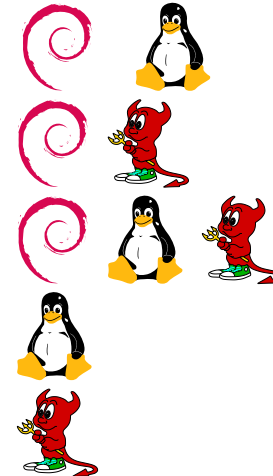
 debian

 disk1


 iso


 disk2


 freebsd



Index-Side Filtering

(debian, disk1, iso)	
name	Debian Disk1.iso
file ID	cdb79ca3db1f3...


(debian, disk2, iso)	
name	Debian Disk2.iso
file ID	5ccbf54d7e502...


(disk1, freebsd, iso)	
name	Freebsd Disk1.iso
file ID	fbcbfdff31f27de...




 “debian disk1?”

Index-Side Filtering


(debian, disk1, iso)	
name	Debian Disk1.iso
file ID	cdb79ca3db1f3...


(debian, disk2, iso)	
name	Debian Disk2.iso
file ID	5ccbf54d7e502...


(disk1, freebsd, iso)	
name	Freebsd Disk1.iso
file ID	fbcbfdff31f27de...



Index-Side Filtering

(debian, disk1, iso)	
name	Debian Disk1.iso
file ID	cdb79ca3db1f3...

(debian, disk2, iso)	
name	Debian Disk2.iso
file ID	5ccbf54d7e502...

(disk1, freebsd, iso)	
name	Freebsd Disk1.iso
file ID	fbcbfdff31f27de...



 “debian disk1?”


{  }

Problem: Poor query load-balancing

Keyword-Set Indexing

- Build index on *keyword sets* rather than keywords
- Store subsets of size $\leq K$
- More keyword-set indexes, but each is shorter
- Single-keyword indexes are less important, so can be truncated
 - $< 29\%$ of web searches have only 1 keyword.
[Reynolds & Vahdat 2003]
- To search: send filtered query to any K -size subset index

Keyword-Set Indexing

(debian, disk1, iso) 	
name	Debian Disk1.iso
file ID	cdb79ca3db1f3...

($K = 2$)

 debian

 disk1

 iso


 debian disk1


 debian iso

 disk1 iso



Keyword-Set Indexing

(debian, disk1, iso)	
name	Debian Disk1.iso
file ID	cdb79ca3db1f3...

(debian, disk2, iso)	
name	Debian Disk2.iso
file ID	5ccbf54d7e502...

($K = 2$)

 debian

 disk1

 iso

 debian disk1

 debian iso

 disk1 iso


 disk2


 debian disk2


 disk2 iso



Keyword-Set Indexing

(debian, disk1, iso) 	
name	Debian Disk1.iso
file ID	cdb79ca3db1f3...

(debian, disk2, iso) 	
name	Debian Disk2.iso
file ID	5ccbf54d7e502...

(disk1, freebsd, iso) 	
name	Freebsd Disk1.iso
file ID	fbcbfdff31f27de...

($K = 2$)

 debian

 disk1

 iso

 debian disk1

 debian iso

 disk1 iso

 disk2


 debian disk2


 disk2 iso


⋮



Keyword-Set Indexing

(debian, disk1, iso) 	
name	Debian Disk1.iso
file ID	cdb79ca3db1f3...

(debian, disk2, iso) 	
name	Debian Disk2.iso
file ID	5ccbf54d7e502...

(disk1, freebsd, iso) 	
name	Freebsd Disk1.iso
file ID	fbcbfdff31f27de...

($K = 2$)

 debian

 disk1

 iso

 debian disk1

 debian iso

 disk1 iso

 disk2

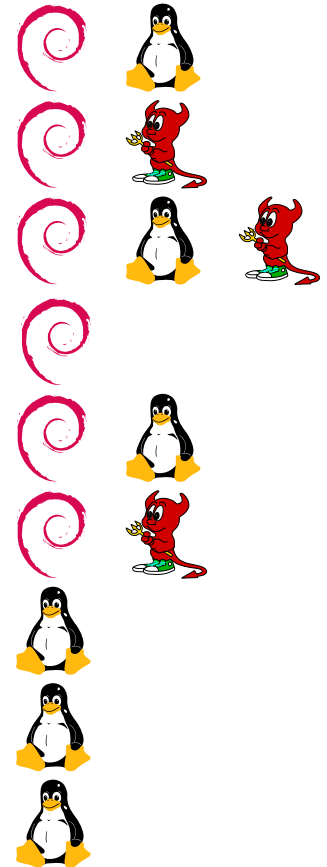
 debian disk2

 disk2 iso


⋮





“debian disk1?”



Keyword-Set Indexing

(debian, disk1, iso) 	
name	Debian Disk1.iso
file ID	cdb79ca3db1f3...

(debian, disk2, iso) 	
name	Debian Disk2.iso
file ID	5ccbf54d7e502...

(disk1, freebsd, iso) 	
name	Freebsd Disk1.iso
file ID	fbcbfdff31f27de...

($K = 2$)

 debian

 disk1

 iso

 **debian disk1**

 debian iso

 disk1 iso

 disk2

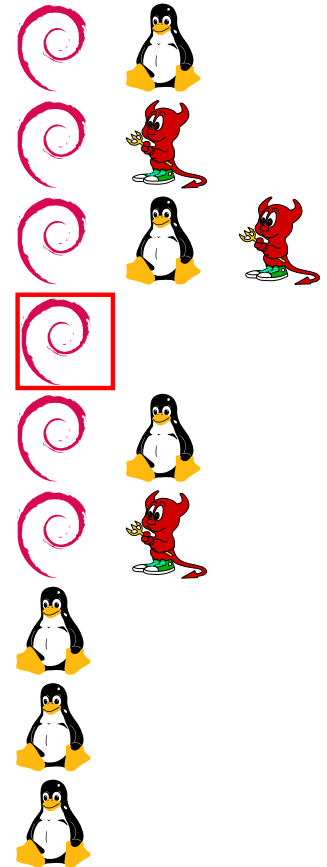
 debian disk2

 disk2 iso

⋮



“debian disk1?”



Indexing Cost

m = metadata keywords

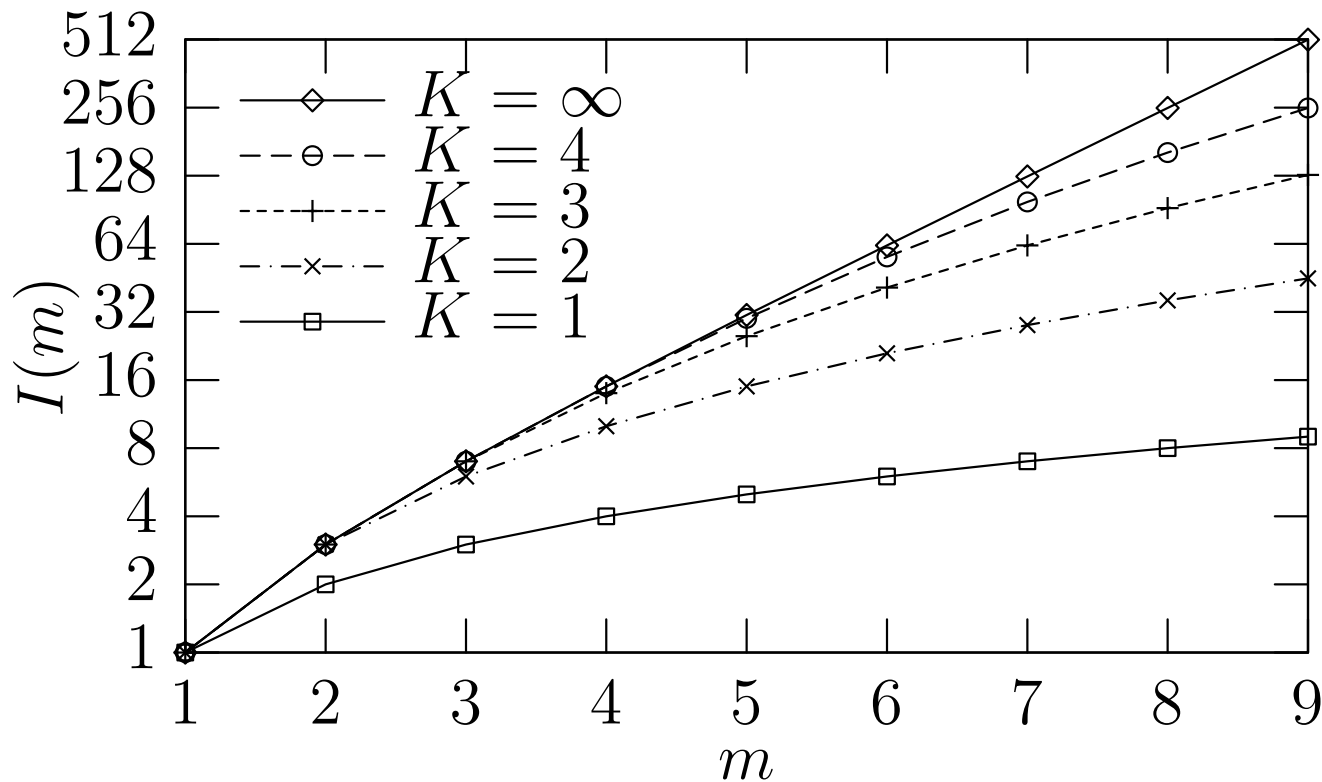
K = maximum subset size parameter

$I(m)$ = index entries

$$= \sum_{i=1}^K \binom{m}{i} = \begin{cases} 2^m - 1 & \text{if } m \leq K \\ O(m^K) & \text{if } m > K \end{cases}$$

Indexing Cost

$$I(m) = \sum_{i=1}^K \binom{m}{i} = \begin{cases} 2^m - 1 & \text{if } m \leq K \\ O(m^K) & \text{if } m > K \end{cases}$$



Indexing Cost

$$I(m) = \sum_{i=1}^K \binom{m}{i} = \begin{cases} 2^m - 1 & \text{if } m \leq K \\ O(m^K) & \text{if } m > K \end{cases}$$

For files with many metadata keywords, $I(m)$ is polynomial in m .

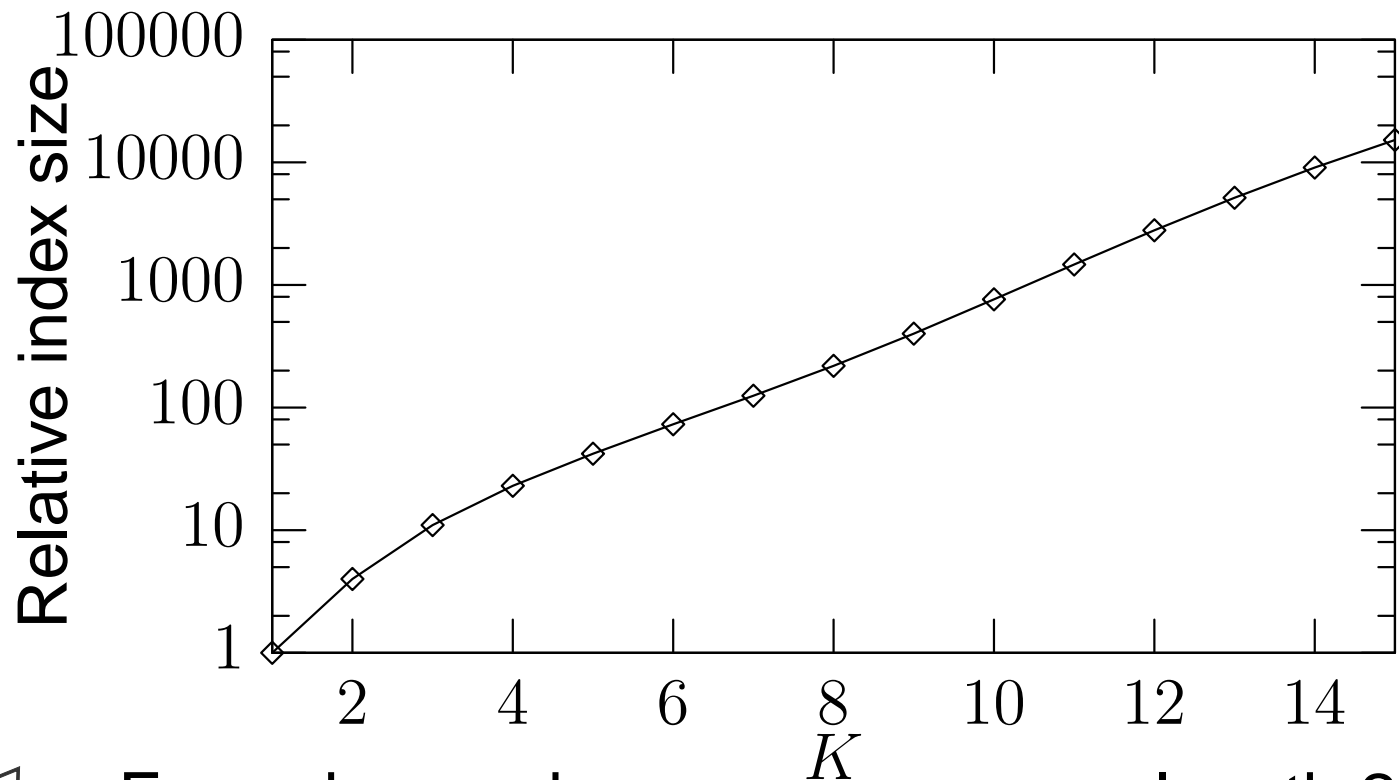
Storage Costs (FreeDB)

Number of songs	21,195,244
Total index entries ($K = 1$)	134,403,379
Index entries per song ($K = 1$)	6.274406
<hr/>	
Total index entries ($K = 3$)	1,494,688,373
Index entries per song ($K = 3$)	66.078093

⇒ Total storage cost only an order of magnitude more than required for $K = 1$ inverted index.

Choosing K

- Larger K improves query load distribution, increases indexing costs

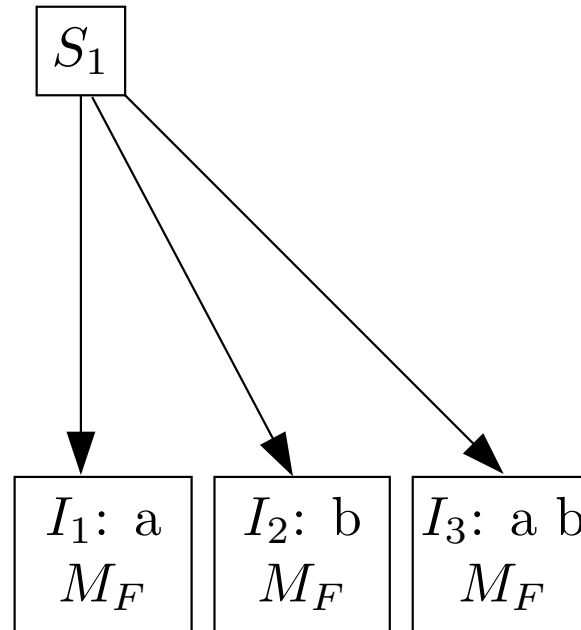


For web searches: average query length 2.53

Outline

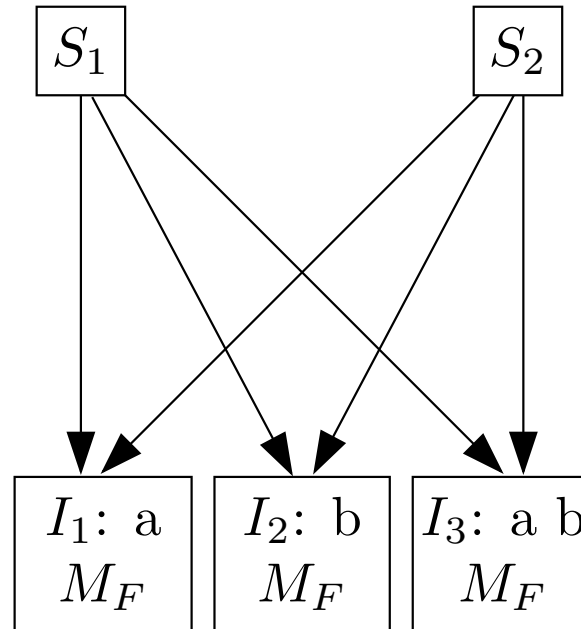
- Overview
- Searching
- **Index Gateways**
- Content Distribution
- Conclusion

Index Gateways



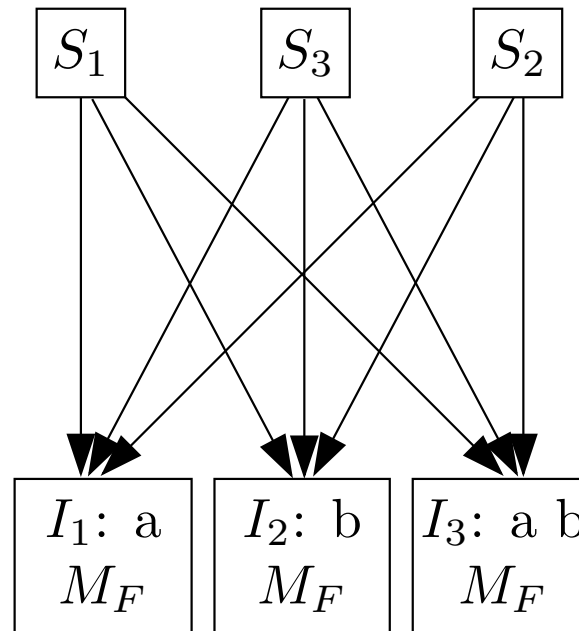
- Each file has one metadata block, stored in $I(m)$ indexes.

Index Gateways



- Each peer sharing the file will insert the *same* metadata block into each index.

Index Gateways

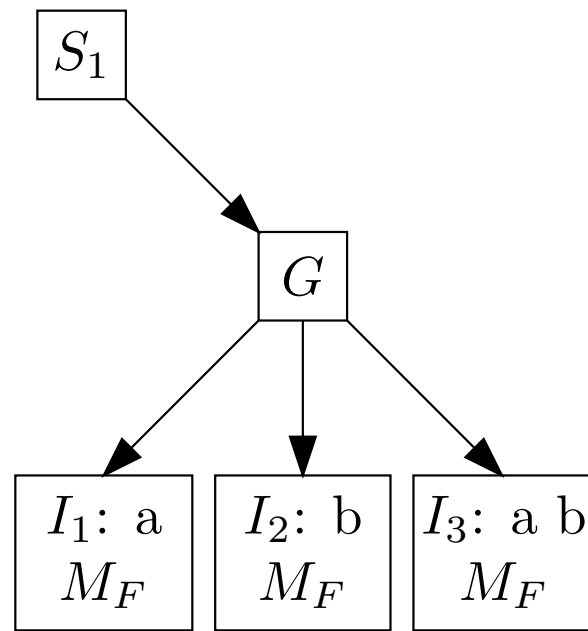


- Total insertion cost for m metadata keywords and s source peers: $sI(m)$ messages.

Problem: expensive and redundant

Index Gateways

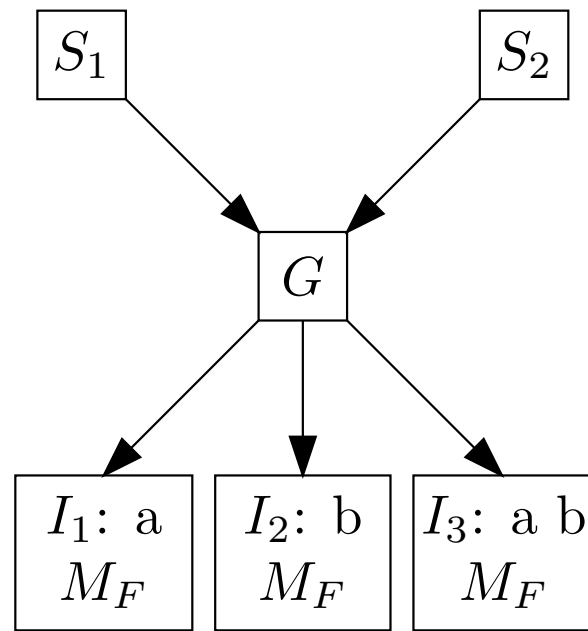
- Solution: aggregate updates at an *index gateway*
- Receives metadata blocks from sources and sends to indexes only when necessary



Insertion cost is now $s + I(m)$ (vs. $sI(m)$)!

Index Gateways

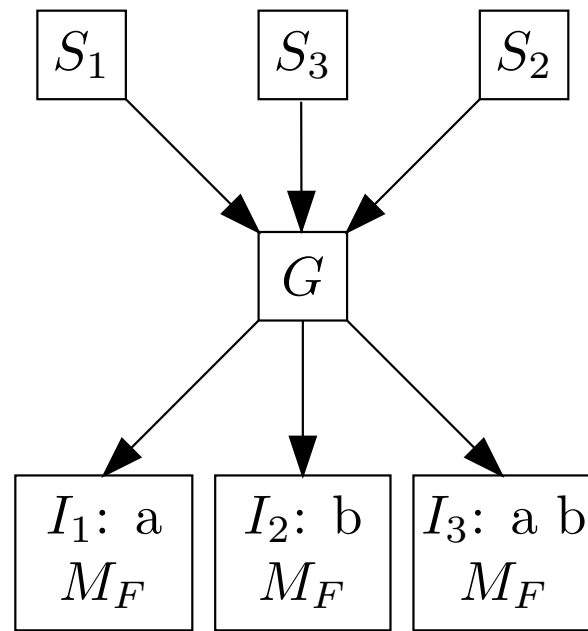
- Solution: aggregate updates at an *index gateway*
- Receives metadata blocks from sources and sends to indexes only when necessary



Insertion cost is now $s + I(m)$ (vs. $sI(m)$)!

Index Gateways

- Solution: aggregate updates at an *index gateway*
- Receives metadata blocks from sources and sends to indexes only when necessary



Insertion cost is now $s + I(m)$ (vs. $sI(m)$)!

Outline

- Overview
- Searching
- Index Gateways
- **Content Distribution**
- Conclusion

Direct Storage?

- Content is large
- Network has churn
 - Kazaa median session length 2.4 minutes [Gummadi et al. 2003]

Problem: DHT storage of content is impractical

Indirect Storage

- Add indirection
- Store only small pointers in the network

keywords

↓ keyword-set index search

file IDs

↓

sources

Content-Sharing Subrings

- How to identify sources for a file?
- Simple solution: “tracker node”
- Instead, file sources form subrings
 - To find source, LOOKUP random ID in file’s subring
 - Search and maintenance costs same as with tracker, but distributed over network
 - No single point of failure

Outline

- Overview
- Searching
- Index Gateways
- Content Distribution
- **Conclusion**

Conclusion

- Supports search with distributed keyword-set index
- Extends standard DHT interface with network-side processing
 - Index-side filtering
 - Index gateways
- Content distribution with indirect storage
 - Indexing via subrings



Bonus Section



Bonus 1 – Postfetching

- Indirect storage leads to unavailable content
- Long metadata expiration leads to visible, unavailable content
- Insert request blocks into the network
- When source node rejoins and locates request blocks actively push requested content into the caches of randomly-selected nodes.

Bonus 2 – Metadata Expiration

- File availability constantly changes as peers join/leave
- Expiration rather than polling
- Peers must be able to handle references to unavailable files
- Long expiration times to track access attempts for unavailable files

Bonus 3 – Index Replication

- Replication instead of erasure coding
- Only weak consistency required
- Updates propagated periodically
- Expirations performed independently

Bonus 4 – Segmentation

- At this level, content is atomic
- Can't share partially downloaded content

Problem: Doesn't utilize upload bandwidth

Bonus 4 – Segmentation

- Split content into chunks and share on the chunk level
- Typical file sharing networks contain many similar files

Problem: Underutilization of content sources

Bonus 4 – Segmentation

- Place chunk boundaries based on content data
- Identify chunks by hash of their contents
- Files are now just a list of their constituent chunk IDs

