# hashFS

## Applying Hashing to Optimize File Systems for Small File Reads

Paul Lensing, Dirk Meister, André Brinkmann
Paderborn Center for Parallel Computing
University of Paderborn

# Contents

- Motivation and Problem

- Design Idea

- Implementation

- Evaluation

- Conclusion

# Contents

- **Motivation and Problem**
- Design Idea
- Implementation
- Evaluation
- Conclusion

UNIVERSITÄT PADERBORN
*Die Universität der Informationsgesellschaft*

- ## Web traffic, e.g. profile images, web site archive
  - Internet Archive
    - > 2500 nodes, > 6000 hard disks
  - Similar: Facebook profile images
    - > 6.5 billion images, most images 5-20 KB

- ## Challenges:
  - Access to small files
  - Multiple IOs per read request
  - ➜ High Overhead
  - ➜ Limits reads/s per disk

**UNIVERSITÄT PADERBORN**
*Die Universität der Informationsgesellschaft*

- Small files (4 KB – 20 KB)
- Accesses are almost exclusively reads

- Filenames are randomly generated or calculated
- High number of concurrent users
- ➔ No name or directory locality
- ➔ Accesses are randomly distributed

- RAM/disk ratio low

- Some POSIX features are not important
  - Directory permission
  - Last file access time (atime)

Example: Read file /var/image/20

Read dentries /

Read dentries /var

Read dentries /var/image

Read file /var/image/20

UNIVERSITÄT PADERBORN
Die Universität der Informationsgesellschaft

# Partial solution

- Store small file data in inode
    - Eliminates 1 IO operation
    - But that operation is not seeking anyway

- Directory lookups are a major problem

# Goal:

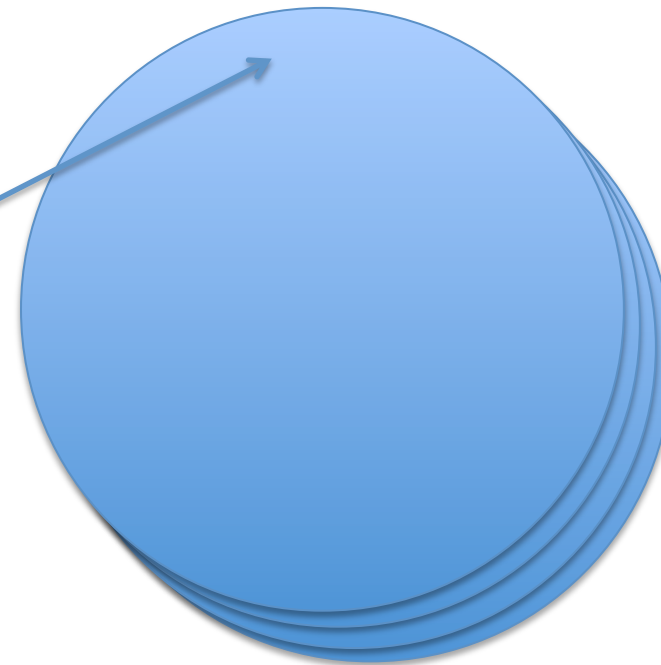1 IO operation per read request
(in web scenario)

# Contents

- Motivation and Problem
- **Design Idea**
- Implementation
- Evaluation
- Conclusion

- Compute location of file
  - No (recursive) lookup
- hash file path to position on disk

hash(/var/images/2010/05/03/10)

- Exact position is not a good approach
  - Collisions

- Hash pathname to "track"
  - Read complete track
    - Overhead, but not prohibitive in random workload
  - Needs geometry data of disk
  - Also possible to use each data region as "track"

# Contents

UNIVERSITÄT PADERBORN
*Die Universität der Informationsgesellschaft*

- Prototype filesystem

- Based on ext2

- Transparent to application and 3rd party tools
  - No library, etc.
  - Except directory permission checks

- Metadata
  - **Track inode** per hashed file
  - Stored in **track inode block** per track

**UNIVERSITÄT PADERBORN**
*Die Universität der Informationsgesellschaft*

- **Every file has additional track inode**
  - Only vital information
    - Inode number, size, security,
    - Identity hash, collision bit
    - n direct pointers
- **Track inode block**
  - Stored all track inodes of track
  - First file system block of track
  - Default: 113 track inodes per track inode block

**UNIVERSITÄT PADERBORN**
*Die Universität der Informationsgesellschaft*

```
Hash pathname to track

Hash pathname to identity hash

Write normal inode

Read track inode block

Search track inode via identity hash

If track inode exists: # Collision
    Set collision bit

Else:

    Create track inode
```

Hash pathname to track

Hash pathname to identity hash

Read track inode block

Search for track inode via identity
  hash

If not exists or if collision bit is
  set:

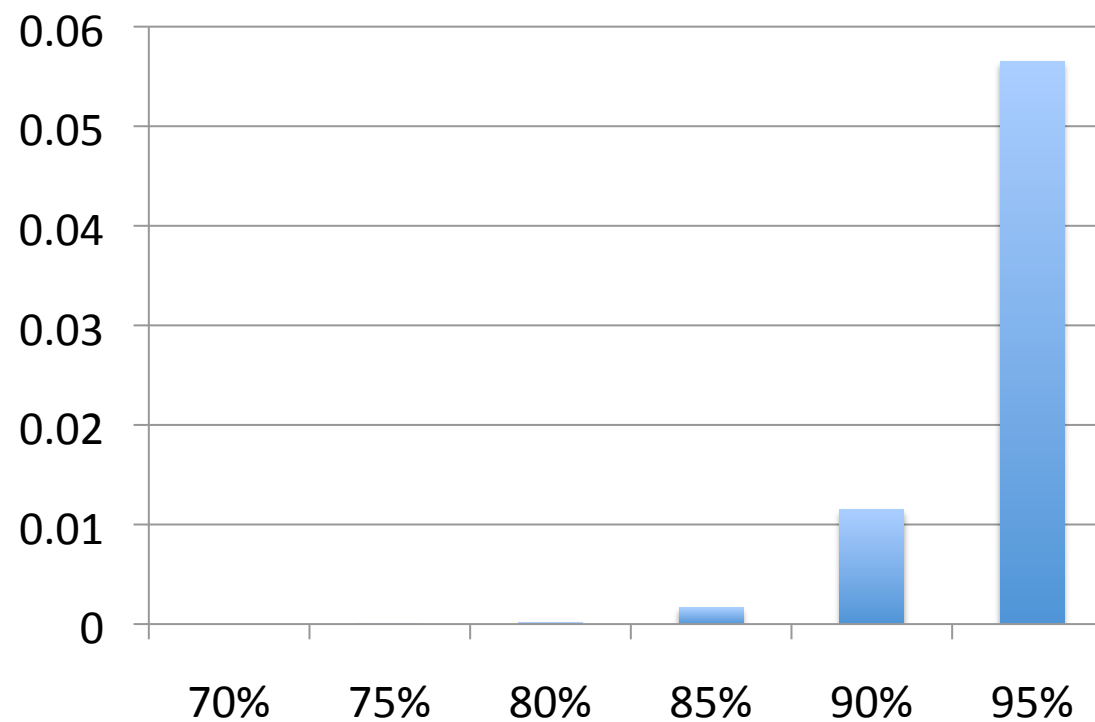  Fail # Use normal lookup

Else:

  Success

# Block Allocation

- Allocate block for „hashed file"

- Possible issue:
  - No free blocks on hashed track

- Basic approach:
  - Keep ratio of track free for hashed files
  - Depending on file system utilization
  - Use hashing approach only for small files

- Hashed track completely used by fs meta data
  - Prevented by eliminating track from geometry data

- Duplicate identity hashes

- No free track inodes in track inode block

- ➔ Alternative lookup fails
  - Normal lookup necessary as backup

Simulation using file set distribution as on central AFS filesystem at UPB

- Modification of VFS layer

  – Bypass recursive directory lookup by pathname method

  – Additional function pointer

    - Transparent to other file systems

  – If pathname lookup has no result

    - Use recursive directory lookup

- hashFS as additional filesystem module

  – geometry data via module options

# Contents

- Motivation and Problem

- Design Idea

- Implementation

- **Evaluation**

- Conclusion

- Random file access distribution
  - Uniform distribution
  - Zipf-like distribution
- Web traffic is assumed to be zipf-like
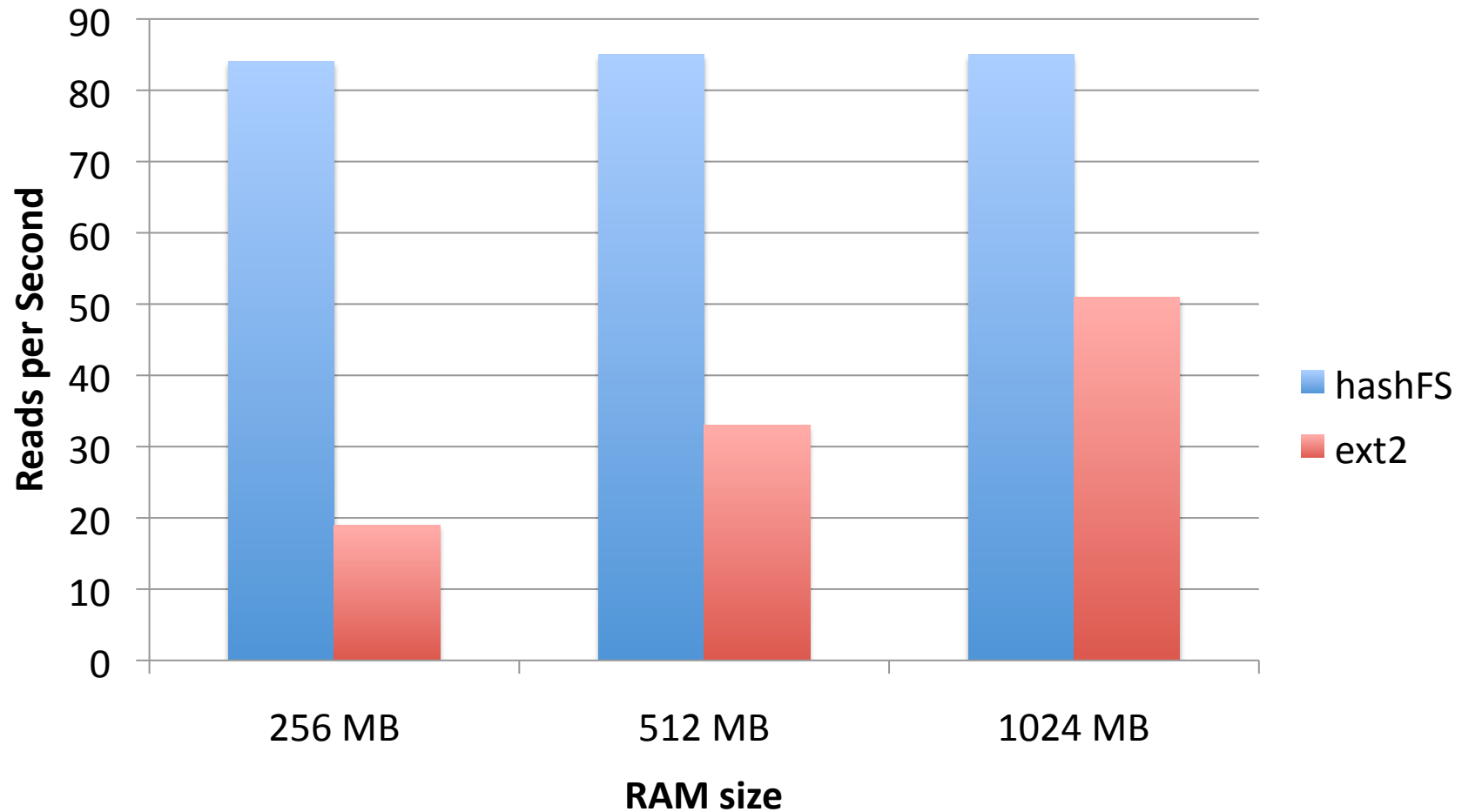  - CDNs, memcache or other caching layers flatten the skewness

- File size: 4 KB

- File count: 18 Million

- Track size: 465 sectors per track
  - Heuristic: tracks have different sizes on disk

- Partition: 90 GB
  - Small partition to limit evaluation time
  - Adjust cache size accordingly
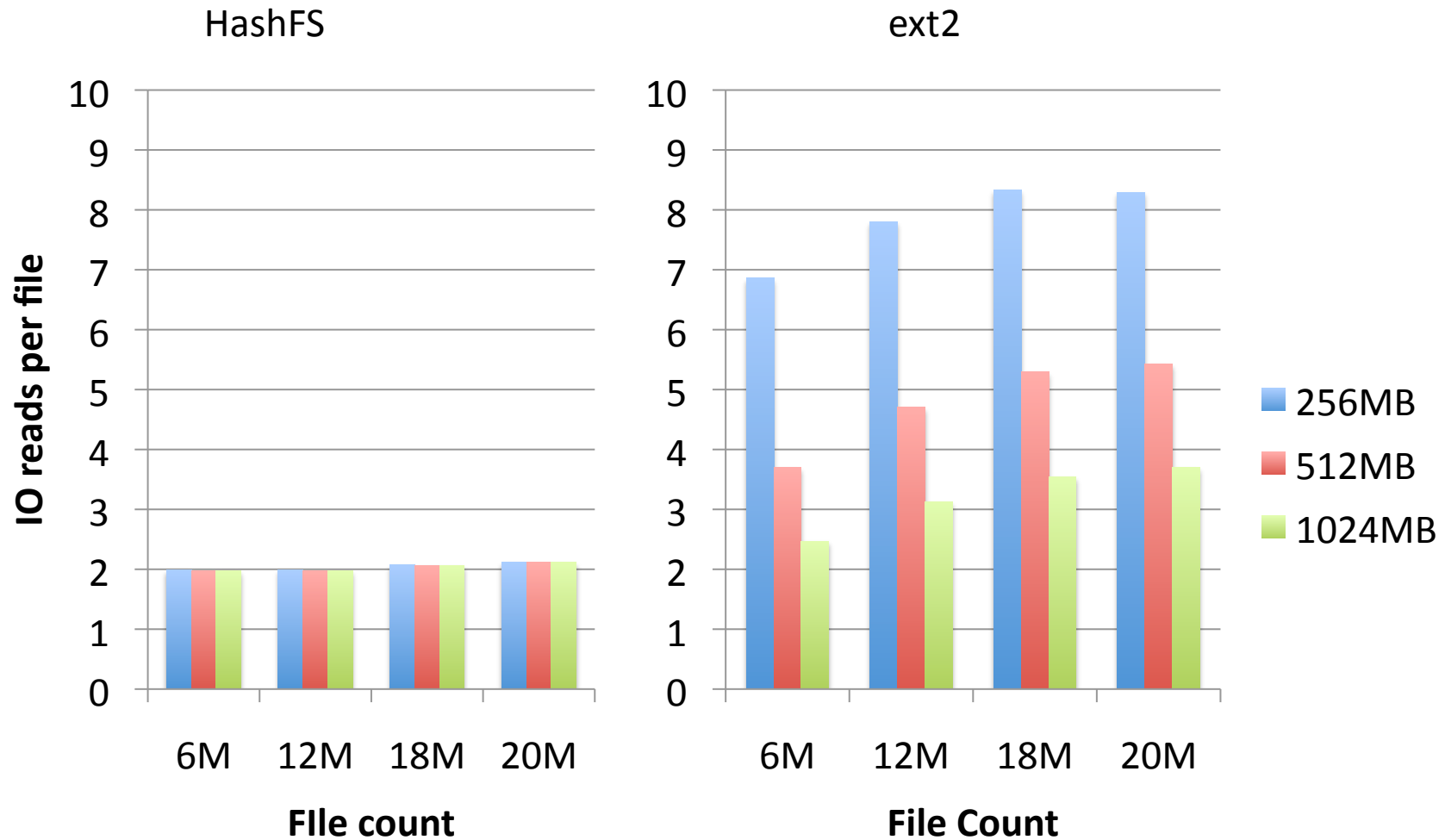
- RAM size: 256 MB, 512 MB, 1024 MB

- **Flat File Set:**
  - Files per directory: 10,000
  - Max depth of directory: 2
- **Deep File Set:**
  - Files per directory: 100
  - Max depth of directory: 6
- **Deep File Set better for ext2**
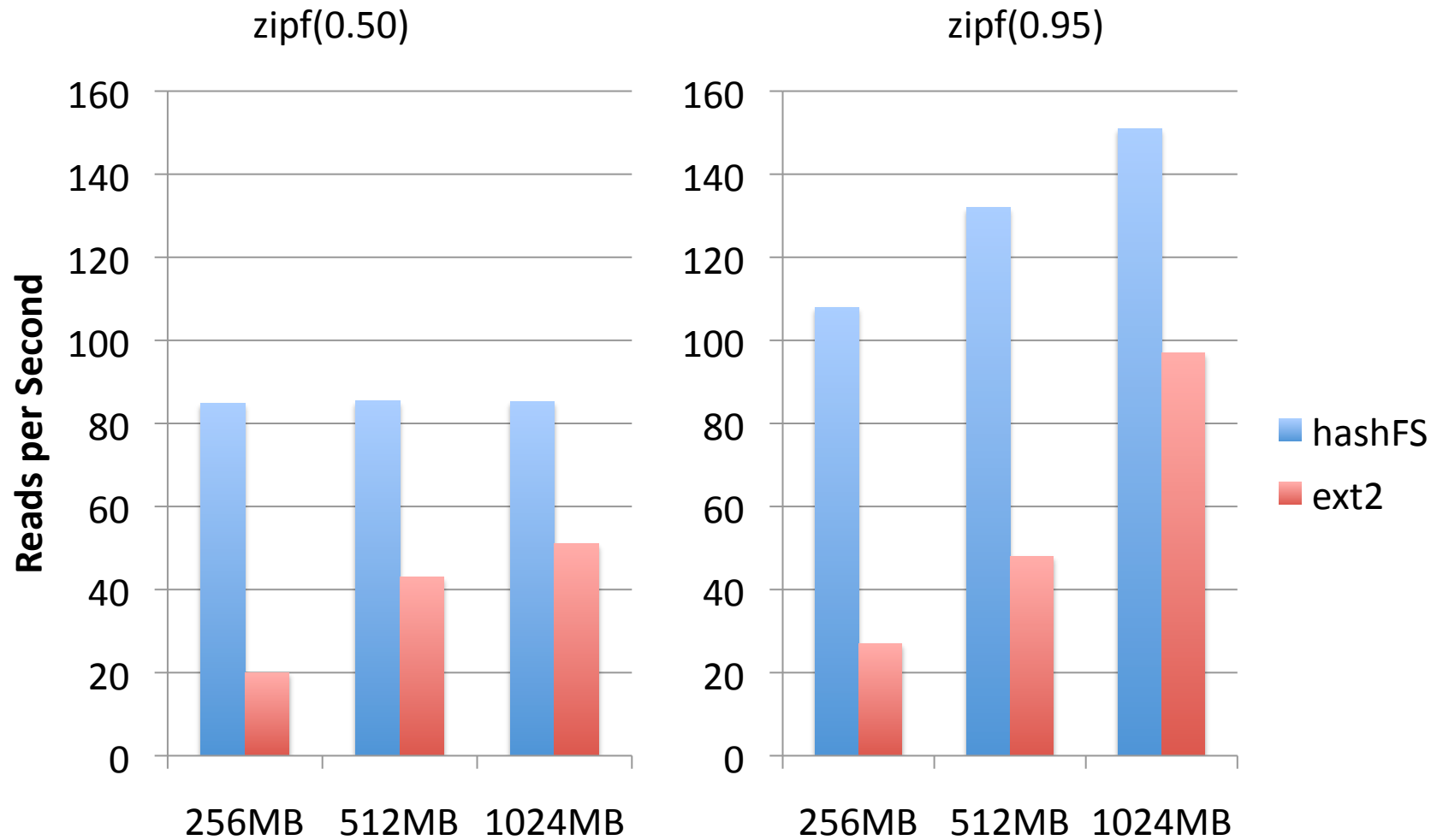  - Except with large caches
- **Here: Only results for Deep File Set**

HashFS

ext2

# Zipf distribution

- ## hashFS
  - Nearly constant access time
  - Only influenced by allocation problems due to utilization

- ## ext2
  - Very different performance values
  - Depends on
    - Cache size and cache state
    - Number and depth of directories
    - Number of files

UNIVERSITÄT PADERBORN
*Die Universität der Informationsgesellschaft*

# Contents

- Motivation and Problem

- Design Idea

- Implementation

- Evaluation

- **Conclusion**

- No full POSIX conformity
  - Directory access permissions are not checked
  - noatime implicit

- Negative performance impact for
  - Large file reads
  - Writes (update to track metadata)

- Current state
  - Prototype implementation in Linux 2.6.18
  - Writes are implemented poorly
  - Evaluation in more complex settings

UNIVERSITÄT PADERBORN
*Die Universität der Informationsgesellschaft*

# Conclusion

- **Hashing pathnames is viable approach**
  - limited to certain situations

- **Approach not limited to ext2**
  - All general-purpose file systems use recursive directory lookup approach
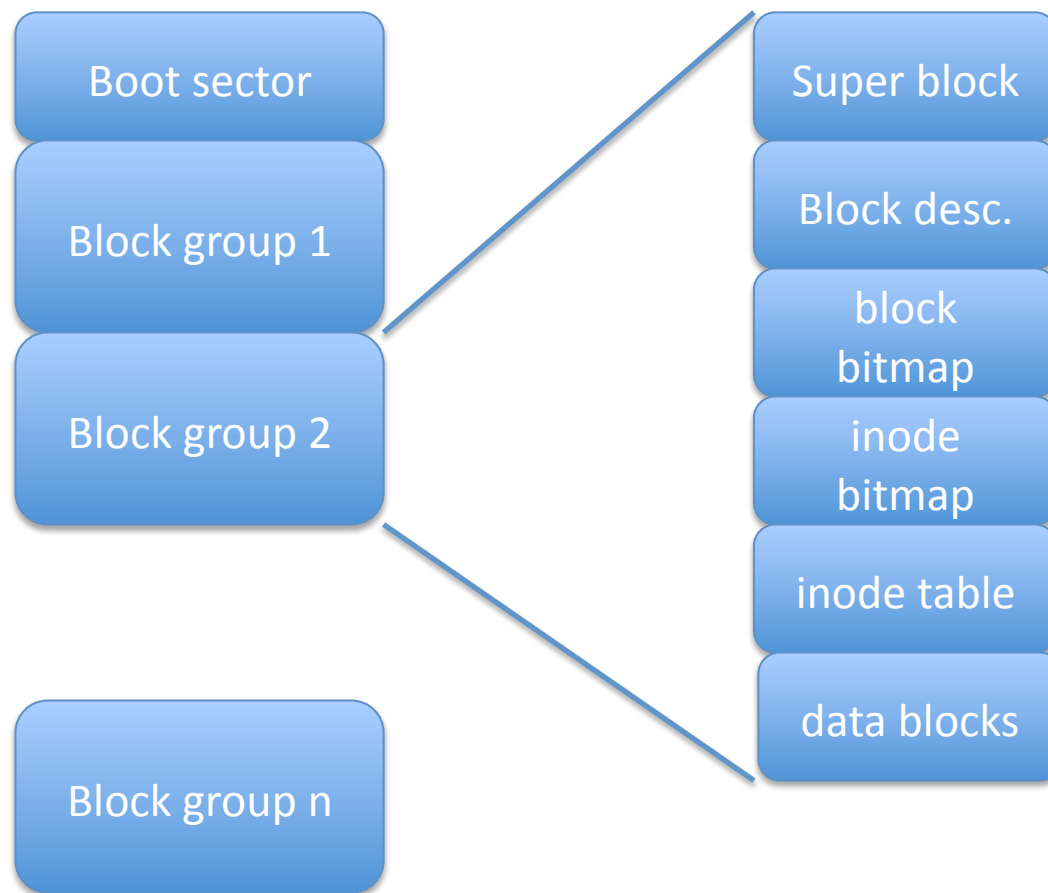  - IO operations for lookup at cache miss

**UNIVERSITÄT PADERBORN**
*Die Universität der Informationsgesellschaft*

# Questions

Boot sector

Block group 1

Block group 2

Block group n

Super block

Block desc.

block bitmap

inode bitmap

inode table

data blocks

PADERBORN
CENTER FOR
PARALLEL
COMPUTING

mode

owner

size

time stamps

...

block 1

...

block 12

12 direct pointers

1 indirect pointers

1 double indirect pointers

1 triple indirect pointers

block 13

block 14

data block

...

1036

UNIVERSITÄT PADERBORN
Die Universität der Informationsgesellschaft

Warm-up Phase

# Track inode issues

| number of files | Disk Utilization | Track inodes issues (average) | Ratio (average) |
|---|---|---|---|
| 1-13 million | <69.6% | 0 | 0% |
| 14 million | 74.7% | 1.9 | < 0.0001% |
| 15 million | 79.8% | 25.7 | 0.0002% |
| 16 million | 84.9% | 277.9 | 0.0017% |
| 17 million | 90.0% | 1946.6 | 0.0115% |
| 18 million | 95.1% | 10,170.0 | 0.0565% |

Simulation using file set distribution as on central AFS filesystem at UPB

UNIVERSITÄT PADERBORN
Die Universität der Informationsgesellschaft