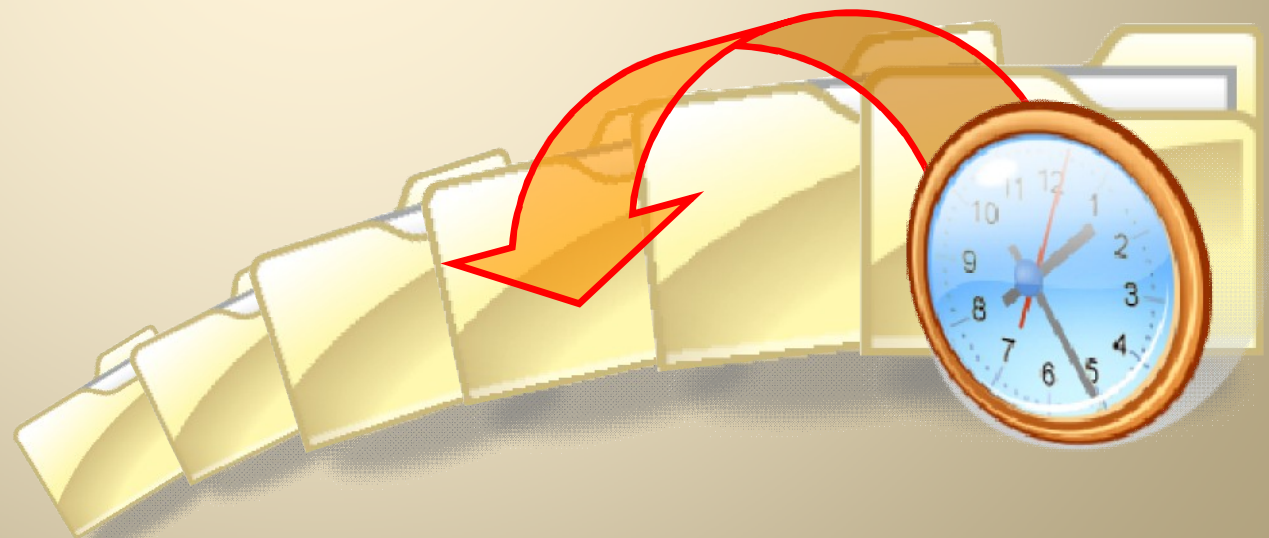# The NILFS2 Filesystem: Review and Challenges

Ryusuke KONISHI

NTT Cyberspace Laboratories

NTT Corporation

# Agenda

- NILFS Introduction
- File System Design
- Development Status
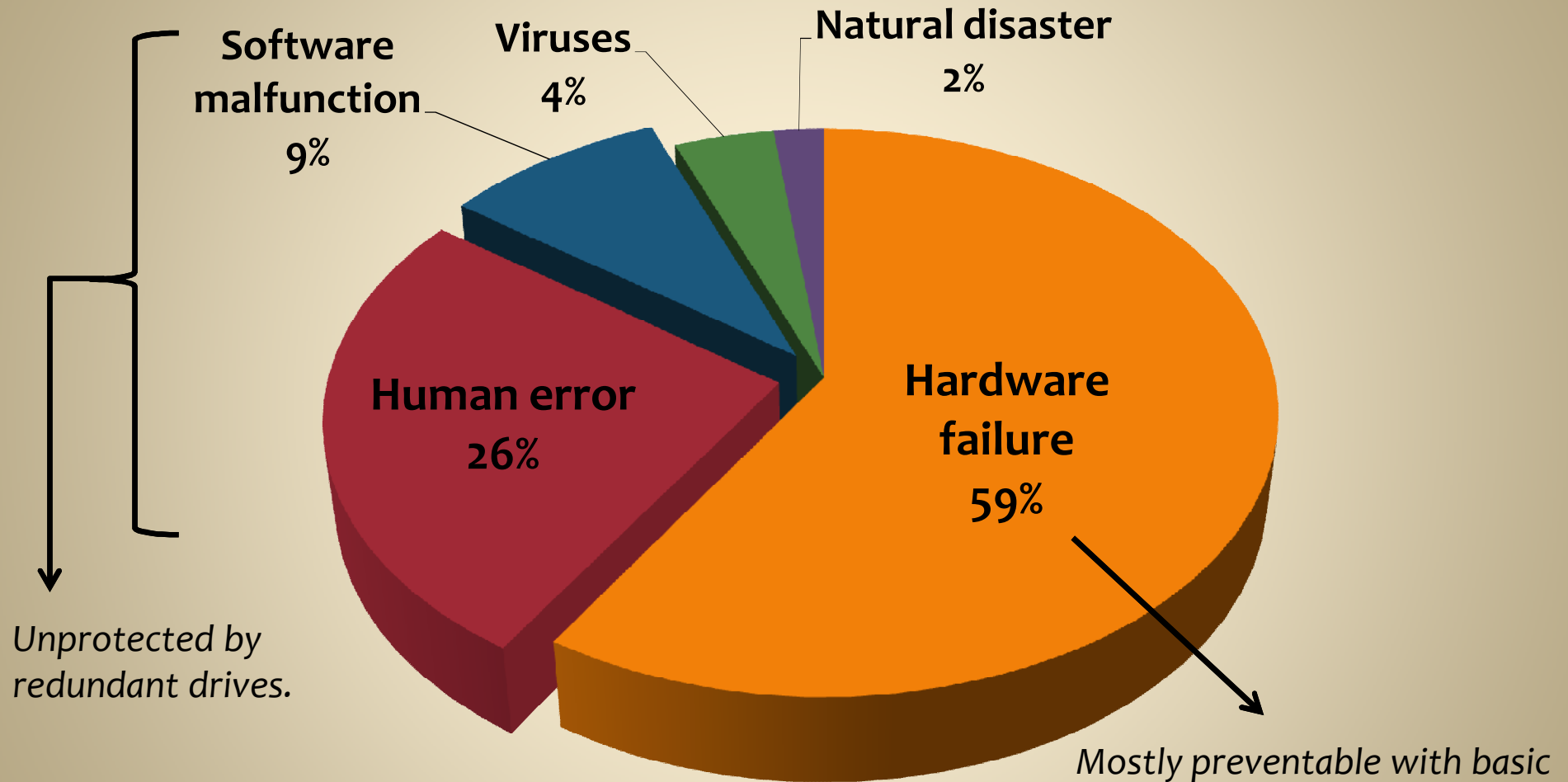- Wished features & Challenges

# What's NILFS

- **NILFS** is the Linux file system supporting "continuous snapshotting"
  - ✓ Provides versioning capability of entire file system
  - ✓ Can retrieve previous states before operation mistake
    - ▪ even restores files mistakenly overwritten or destroyed just a few seconds ago.
  - ✓ Merged into the mainline kernel 2.6.30

# Why NILFS ? (1)

## CAUSE OF DATA LOSS



Software malfunction
9%

Viruses
4%

Natural disaster
2%

Human error
26%

Hardware failure
59%

Unprotected by redundant drives.

Mostly preventable with basic high-integrity system (Redundant configuration)
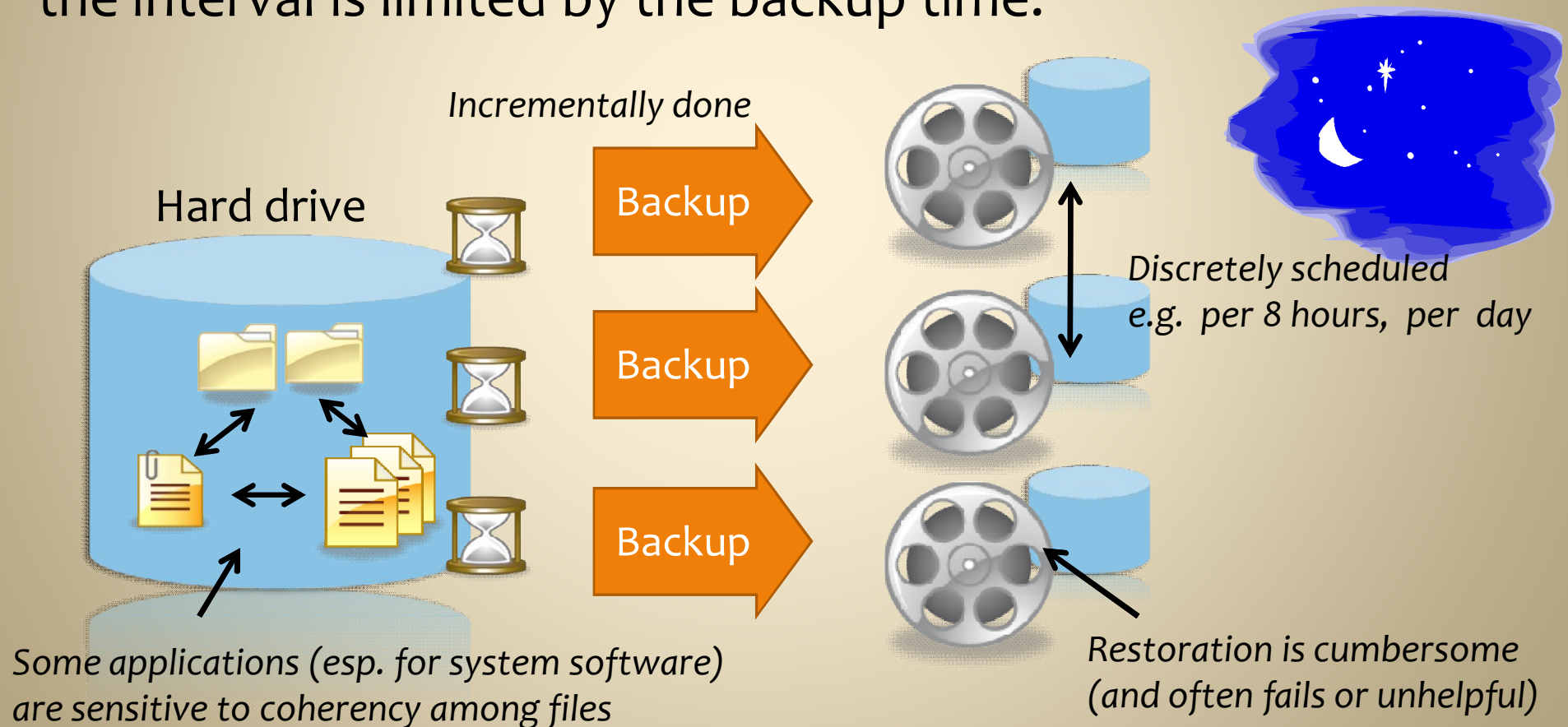
Source:  Ontrack Data Recovery, Inc.
Including office PC. The data is based on actual data recoveries performed by Ontrack.

# Why NILFS ? (2)

**ISSUES IN BACKUP**

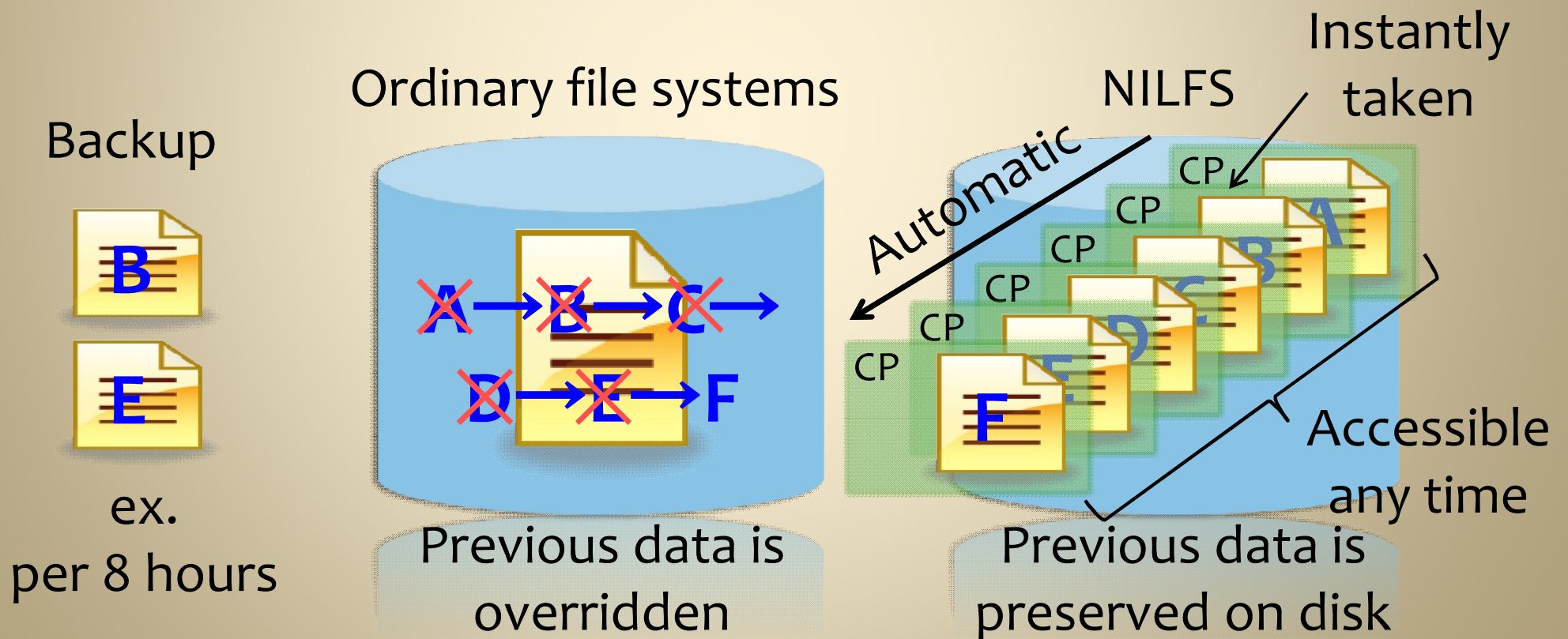- Changes after the last backup are not safe
- But, frequent backups place burden on the system as well as the interval is limited by the backup time.

*Incrementally done*

Hard drive

Backup

Backup

Backup

*Discretely scheduled e.g. per 8 hours, per day*

*Some applications (esp. for system software) are sensitive to coherency among files*

*Restoration is cumbersome (and often fails or unhelpful)*

# NILFS Solution

- Adopt Log-structured File System approach to continually save data on disk.
- Checkpoints are created every time user makes a change, and each checkpoint is mutable to snapshots later on.



Backup

Ordinary file systems

NILFS

Instantly taken

Automatic

Accessible any time

ex. per 8 hours

Previous data is overridden

Previous data is preserved on disk

# Comparison of Snapshots

| File System (solution) | Maximum Number of Snapshots | Instant Snapshotting | Writable Snapshots | Retroactive Snapshots | Incremental Backup |
|---|---|---|---|---|---|
| NTFS Volume Shadow Copy | 64 | | | | Optional (Third party product) |
| ZFS | Unlimited*1 | √ | √ | | √ |
| Btrfs | Unlimited*1 | √ | √ | | Planned |
| NILFS2 | Unlimited*1 | √ | | √ | Requested |
| Apple Time Machine | Thinned out automatically | | | – | √ |
| CDP | Unlimited*1 | – | – | – | √ |

Backup Solutions { Apple Time Machine, CDP

*1: No practical limits (bounded by disk capacity)

# NILFS Disk Write

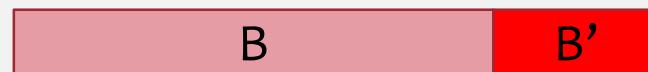- Only modified blocks are incrementally written to disk
  - ✓ This write scheme is applied even to metadata and intermediate blocks

**Application view:**

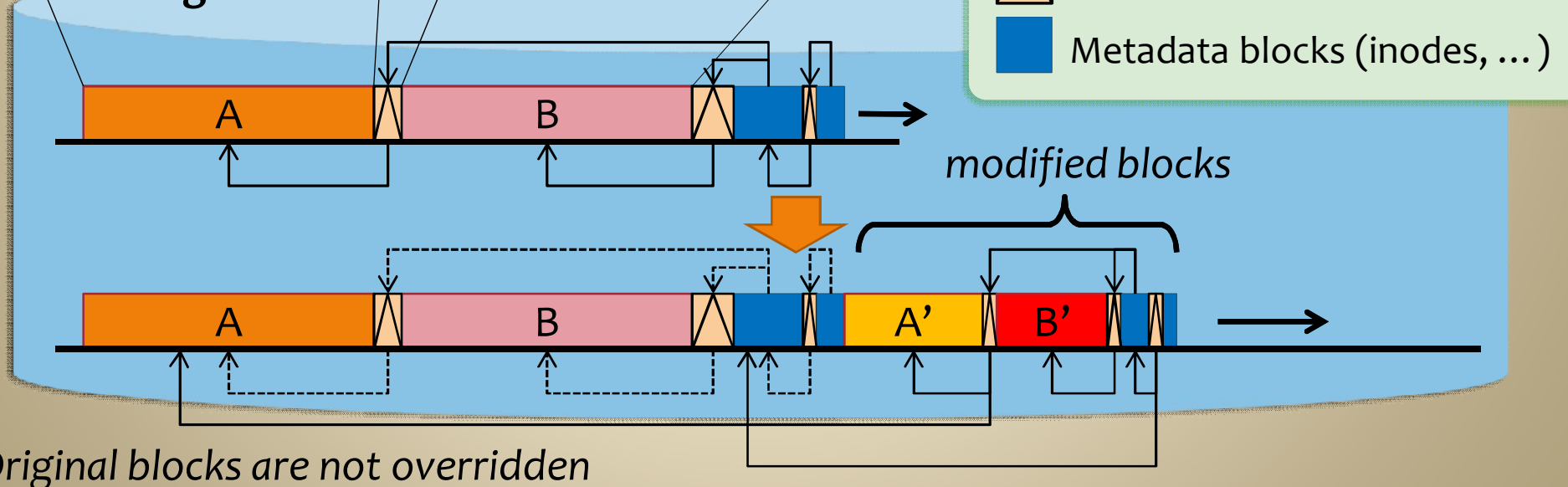File A (modified)

| A | A' | A |

File B (appended)

| B | B' |

**On disk images:**

B-Tree intermediate blocks

Metadata blocks (inodes, …)

| A | | B | | | |

*modified blocks*

| A | | B | | | A' | B' | |

*Original blocks are not overridden*

# NILFS Metadata Hierarchy

Super blocks

Super root block

## CPFILE
checkpoint information

## SUFILE
Segment usage

## DAT
Disk block address translation

## Block mapping (B-tree)

Inode

Node blocks
(Intermediate blocks *1)

x255

x255

Data blocks

*1 B-tree depth varies with the number of data blocks

## IFILE
contains Inodes

checkpoint versions

cno=100
cno=101
cno=102
cno=103
cno=104
cno=105

Files, Directories, Symbolic links     ......

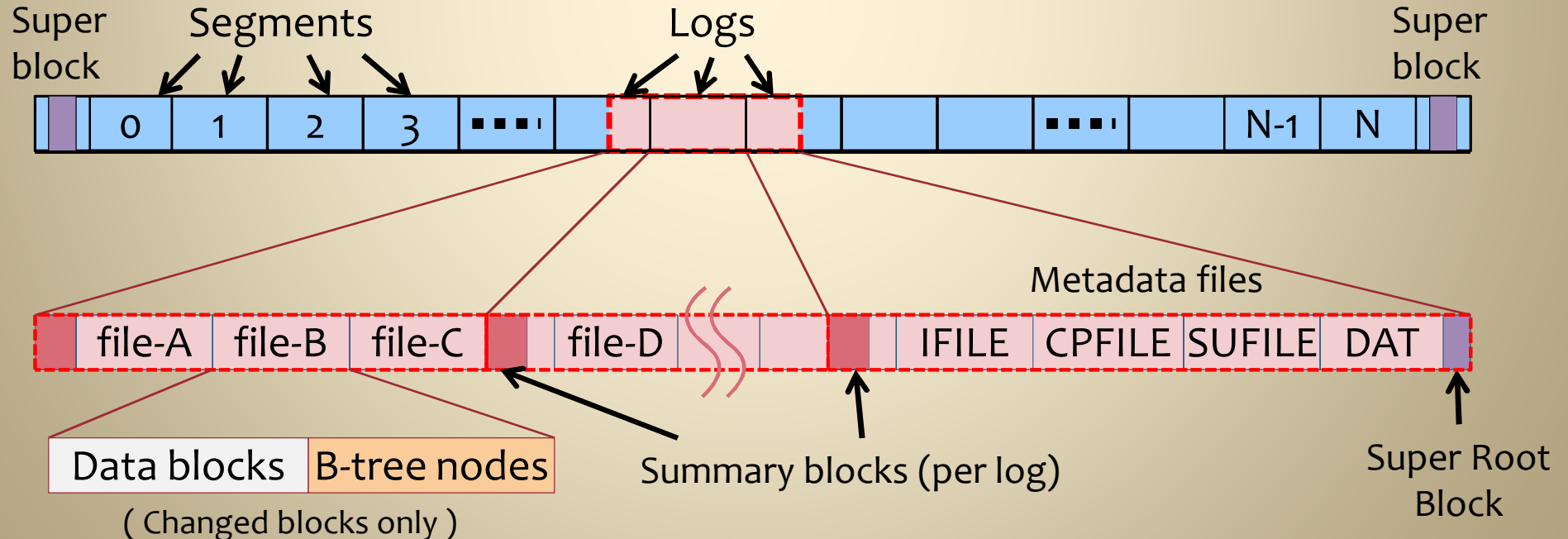ino=123     ino=124     ino=125     ino=126

# Disk Layout Summary

- Segments
  - ✓ Disk space is allocated or freed per segment
  - ✓ Each segment is filled with logs

- Logs
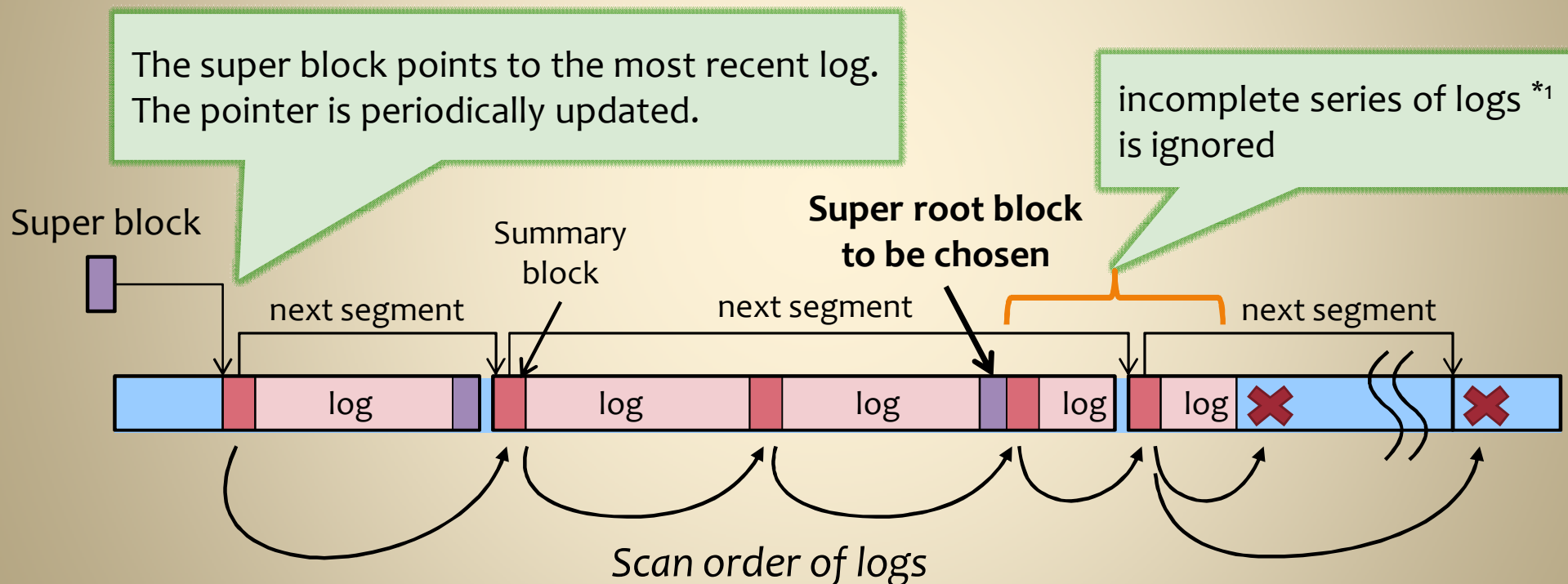  - ✓ Organize delta of data and metadata per file
  - ✓ Compose a new version of metadata hierarchy every checkpoint

Super block | Segments | Logs | Super block

| 0 | 1 | 2 | 3 | ▪▪▪▪ | | | | | | ▪▪▪▪ | N-1 | N |

Metadata files

| file-A | file-B | file-C | | file-D | | IFILE | CPFILE | SUFILE | DAT |

Data blocks | B-tree nodes
( Changed blocks only )

Summary blocks (per log)

Super Root Block

# Mount-time Recovery

- How does NILFS recover from unclean status?
  - ✓ Finds the last log which has a super root block, and done!
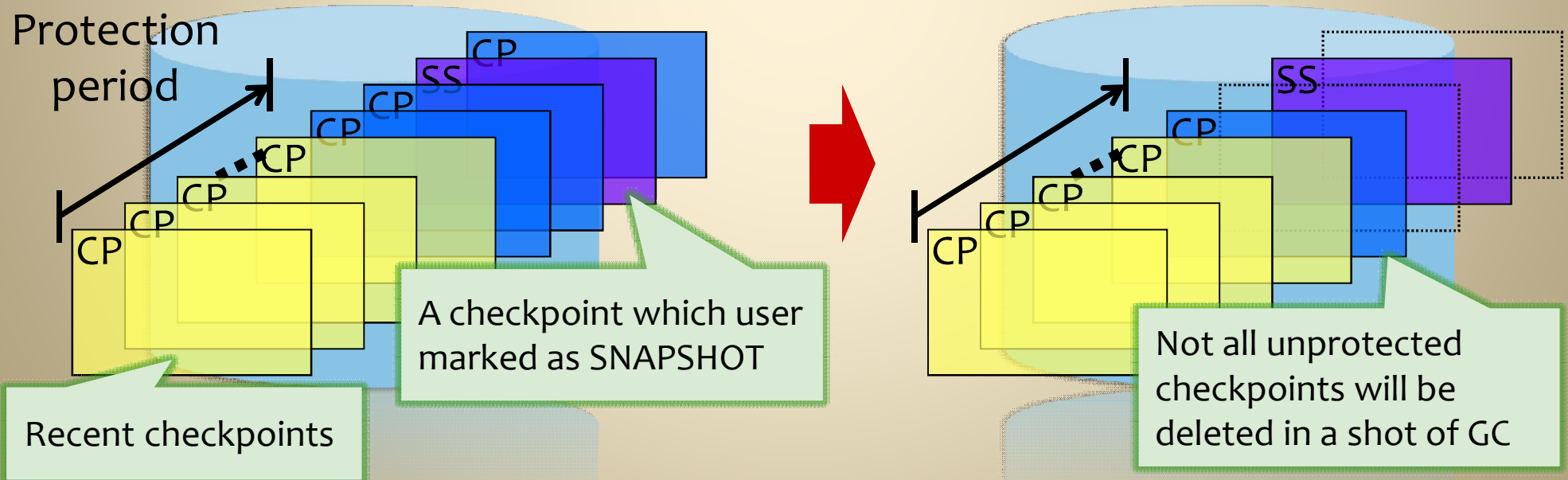  - ✓ Each log is validated with checksums

The super block points to the most recent log.
The pointer is periodically updated.

incomplete series of logs [1]
is ignored

Super block

Summary
block

**Super root block
to be chosen**

next segment                    next segment                    next segment

Super block    log        log              log          log      log

*Scan order of logs*

[1]  Series of logs may not have the super root block.  This type of variant is allowed for optimizations
     to make synchronous write operation faster.
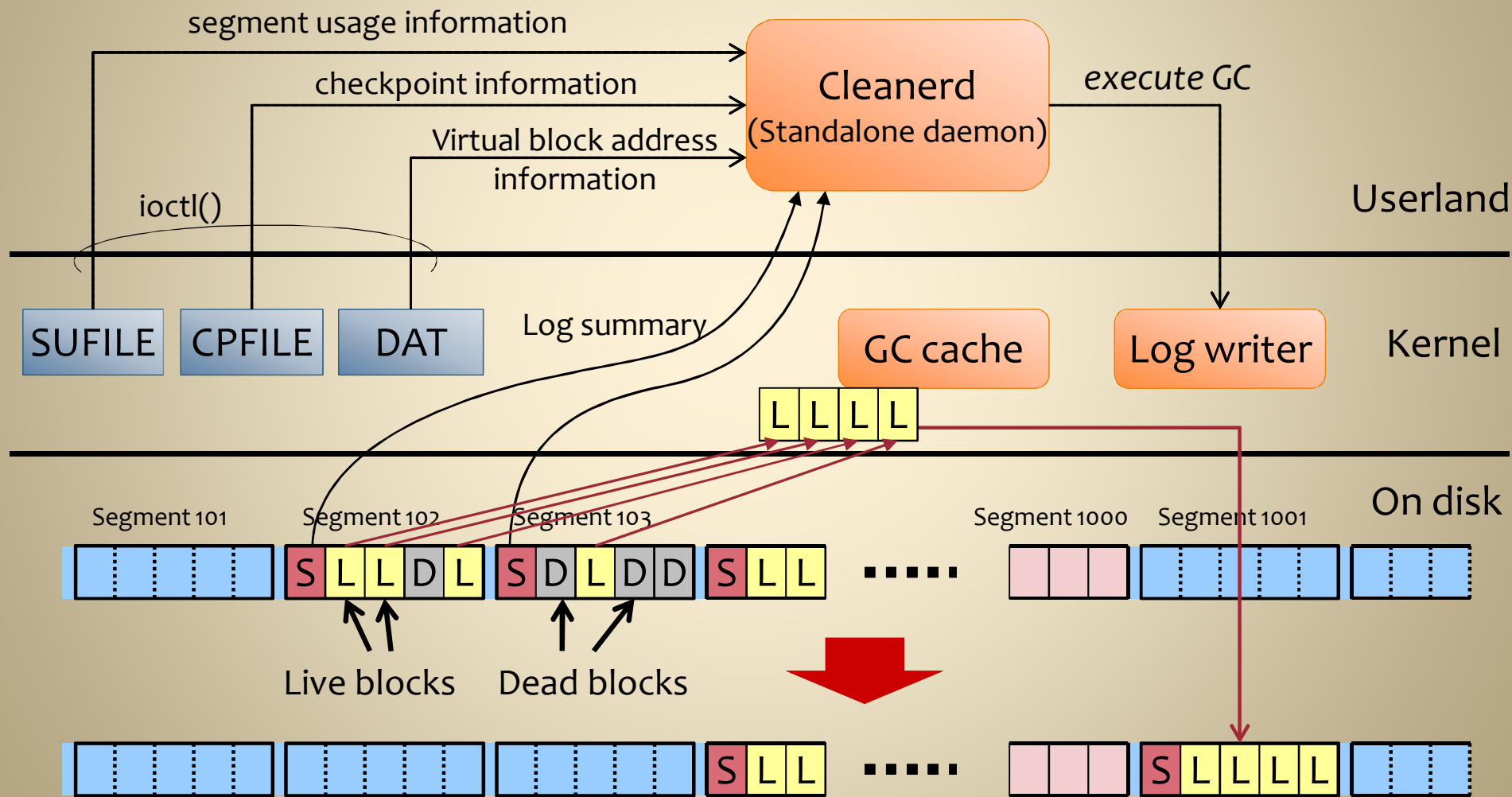
# Garbage Collection (1)

- Creates new disk space to continue writing logs
  - ✓ Essential function of Log structured File Systems
- A disk block is in-use if it belongs to a snapshot or recent checkpoints; unused blocks are freed with their checkpoints

Preserving checkpoints as snapshots

Protection period

CP CP CP CP CP CP SS CP

A checkpoint which user marked as SNAPSHOT

Recent checkpoints

CP CP CP CP SS CP

Not all unprotected checkpoints will be deleted in a shot of GC

# Garbage Collection (2)

## Overall view



segment usage information

checkpoint information

Virtual block address information

Cleanerd
(Standalone daemon)

*execute GC*

ioctl()

Userland

SUFILE   CPFILE   DAT

Log summary

GC cache

Log writer

Kernel

L L L L

Segment 101   Segment 102   Segment 103   Segment 1000   Segment 1001   On disk

S L L D L   S D L D D   S L L   ......

S L L   ......   S L L L L
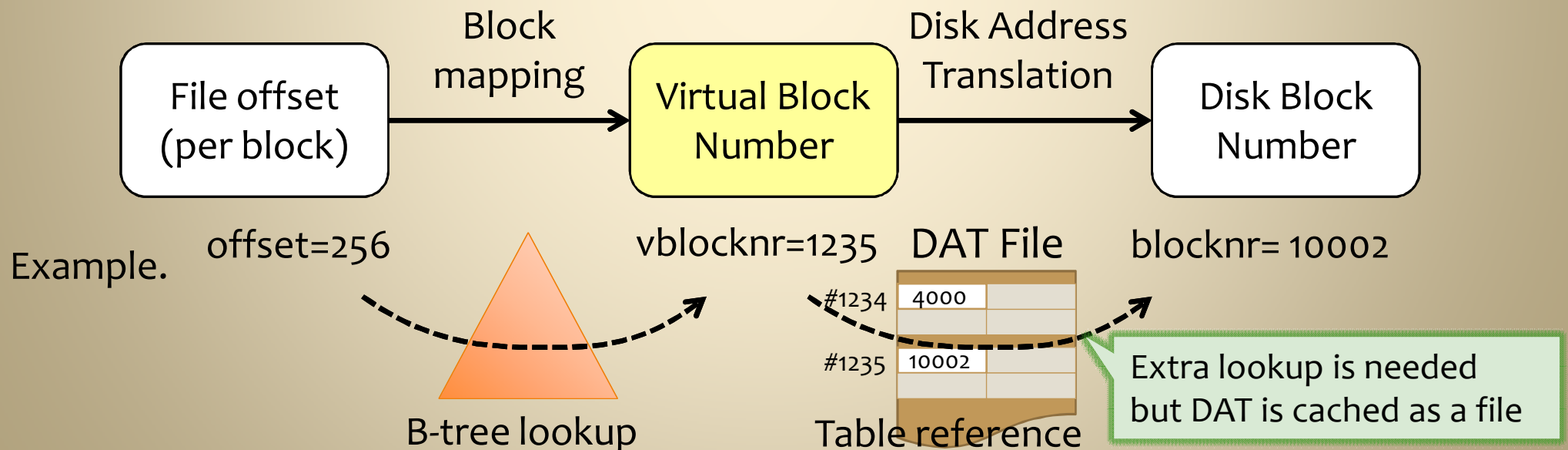
Live blocks   Dead blocks

# Garbage Collection (3)
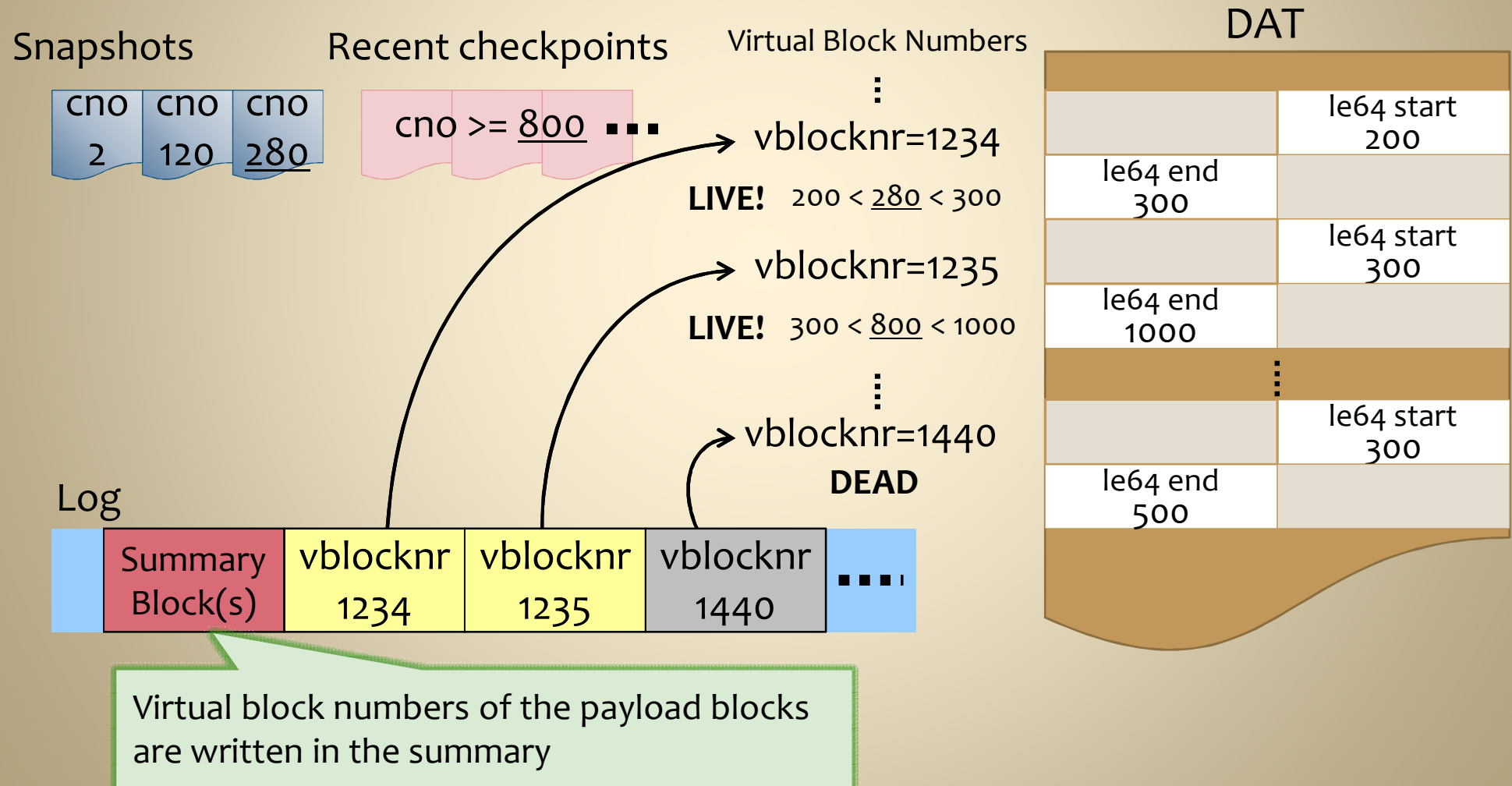
**block addressing**

- Issue for moving disk blocks
  - ✓ Must rewrite b-tree node blocks and inodes having a pointer to moved blocks
  - ✓ Disk blocks are pointed from many parent blocks because NILFS makes numerous versions
- Solution
  - ✓ Use virtual (i.e. indirect) block numbers instead of real disk block numbers



Block mapping → Disk Address Translation

| File offset (per block) | Virtual Block Number | Disk Block Number |

Example.

offset=256    vblocknr=1235    DAT File    blocknr= 10002

B-tree lookup

#1234  4000
#1235  10002

Table reference

Extra lookup is needed but DAT is cached as a file

# Garbage Collection (4)

## Live or dead determination

- Cleanerd determines if each disk block is **LIVE** or **DEAD** from DAT

**Snapshots**

| cno 2 | cno 120 | cno 280 |
|---|---|---|

**Recent checkpoints**

cno >= 800 ···

**Virtual Block Numbers**

⋮

vblocknr=1234

**LIVE!** $200 < 280 < 300$

vblocknr=1235

**LIVE!** $300 < 800 < 1000$

⋮

vblocknr=1440

**DEAD**

**DAT**

| | le64 start 200 |
|---|---|
| le64 end 300 | |
| | le64 start 300 |
| le64 end 1000 | |
| ⋮ | |
| | le64 start 300 |
| le64 end 500 | |

**Log**

| Summary Block(s) | vblocknr 1234 | vblocknr 1235 | vblocknr 1440 | ···· |
|---|---|---|---|---|

Virtual block numbers of the payload blocks are written in the summary

# Current Development Status (1)

- Achievements
  - ✓ Snapshots
    - ▪ Automatically and continuously taken
    - ▪ Mountable as read-only file systems
    - ▪ Mountable concurrently with the writable mount (convenient for online backup)
    - ▪ Quick listing
    - ▪ Easy administration

  - ✓ Online disk space reclamation
    - ▪ Can maintain multiple snapshots

# Current Development Status (2)

- Achievements
  - ✓ Other Features
    - Quick recovery on-mount after system crash
    - B-tree based file and meta data management
    - 64-bit data structures; support many files, large files and disks
    - Block sizes smaller than page size (e.g. 1KB or 2KB)
    - Redundant super blocks (automatic switch)
    - 64-bit on-disk timestamps which are free of the year 2038 problem
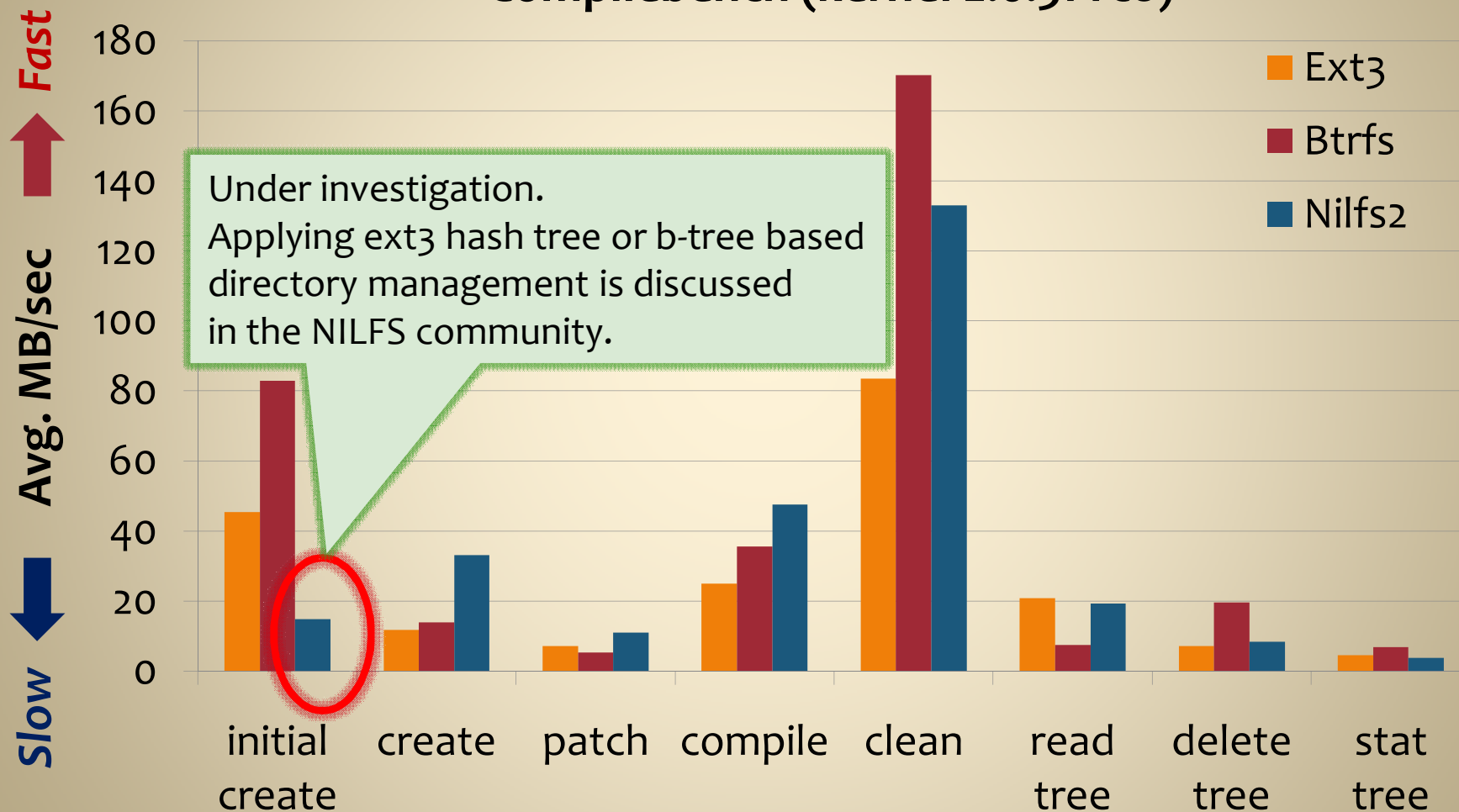    - Nano second timestamps

# Current Development Status (3)

- Todo
  - ✓ On-disk atime
  - ✓ Extended attributes (work in progress)
  - ✓ POSIX ACLs
  - ✓ O_DIRECT write
    - ▪ Currently fallback to buffered write
  - ✓ Fsck
  - ✓ Resize
  - ✓ Quotas

- Performance issues
  - ✓ Directory operations
  - ✓ Write performance
  - ✓ Optimization for silicon disks (esp. for SSD)

# Performance issues (1)

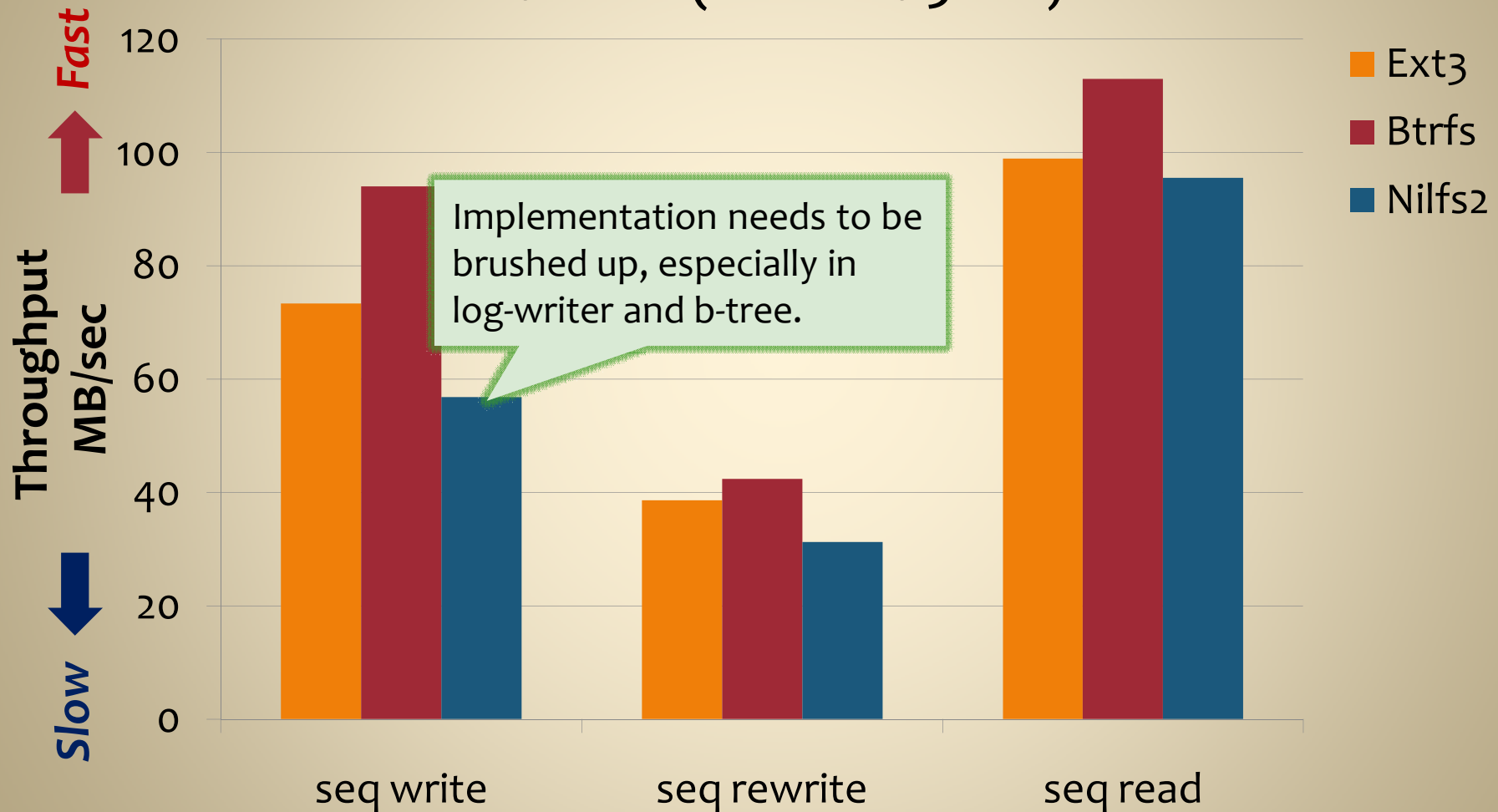## Compilebench (kernel 2.6.31-rc8)



Under investigation.
Applying ext3 hash tree or b-tree based
directory management is discussed
in the NILFS community.

Legend: Ext3, Btrfs, Nilfs2

Y-axis: Avg. MB/sec — Fast / Slow — 0, 20, 40, 60, 80, 100, 120, 140, 160, 180

X-axis: initial create, create, patch, compile, clean, read tree, delete tree, stat tree

Hardware specs:
Processor: Pentium Dual-Core CPU E5200 @ 2.49GHz, Chipset: Intel 4 Series Chipset + ICH10R,  Memory: 2989MB, Disk: ST3500620AS

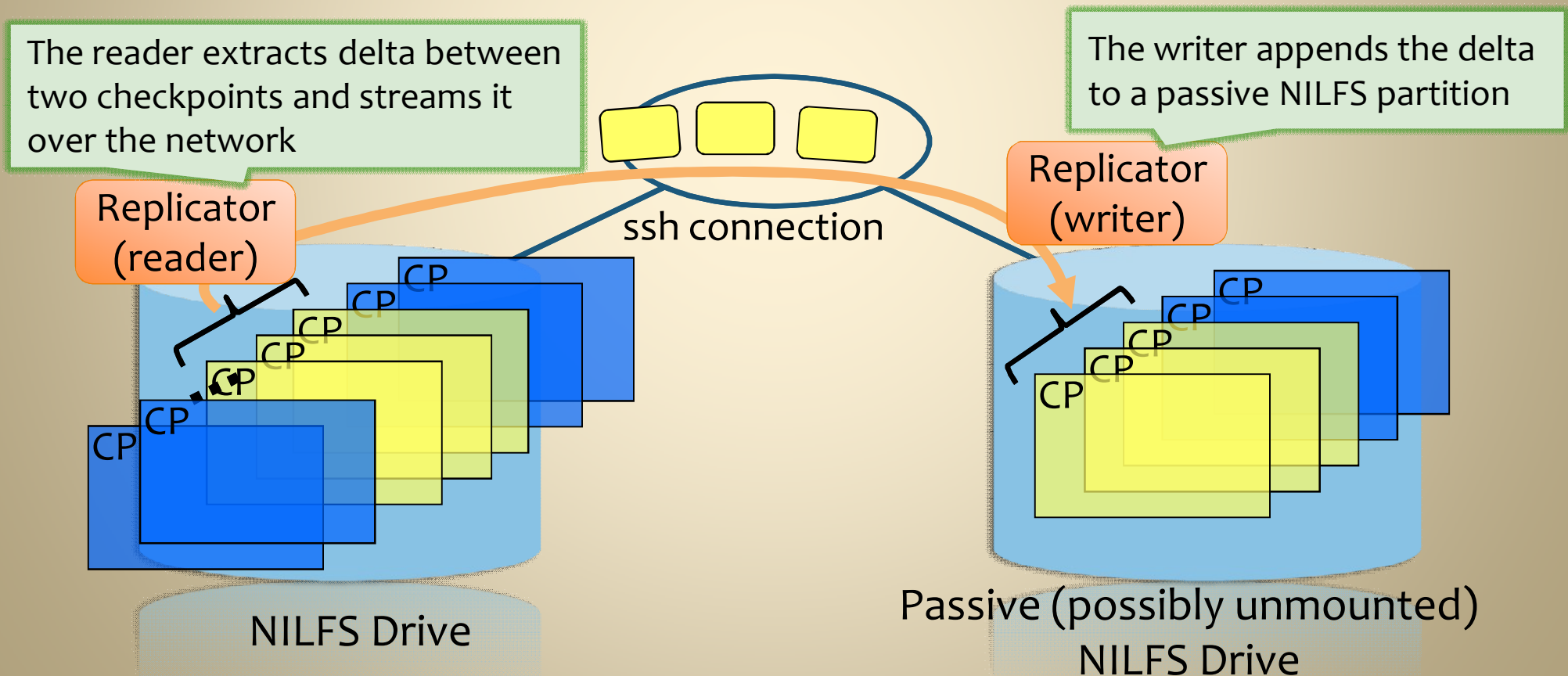# Performance issues (2)

**Bonnie++ (kernel 2.6.31-rc8)**



Implementation needs to be brushed up, especially in log-writer and b-tree.

Hardware specs:
Processor: Pentium Dual-Core CPU E5200 @ 2.49GHz, Chipset: Intel 4 Series Chipset + ICH10R, Memory: 2989MB, Disk: ST3500620AS

# Checkpoint-based Replication (1)

- Faster and robust online backup like ZFS
  - ✓ Back up checkpoints instead of usual files
  - ✓ Similar features are planned for btrfs, TUX3, and the Device Mapper (dm replication)

The reader extracts delta between two checkpoints and streams it over the network

The writer appends the delta to a passive NILFS partition

Replicator (reader)

Replicator (writer)

ssh connection

CP

NILFS Drive

Passive (possibly unmounted) NILFS Drive

# Checkpoint-based Replication (2)

**KEY CHALLENGES DISCUSSED IN THE NILFS COMMUNITY**

- How to extract delta between two checkpoints ?
  - ✓ Two approaches
    - ▪ Scan logs in creation order (just gets delta from logs)
    - ▪ Scan DAT to gather blocks changed during given period
  - ✓ Have pros and cons
    - ▪ The former seems to be efficient, but has a limit due to GC.
    - ▪ Replicator may use either or both of these methods

- Rollback on the destination file system
  - ✓ Needed before starting replication especially to thin out the backups with GC
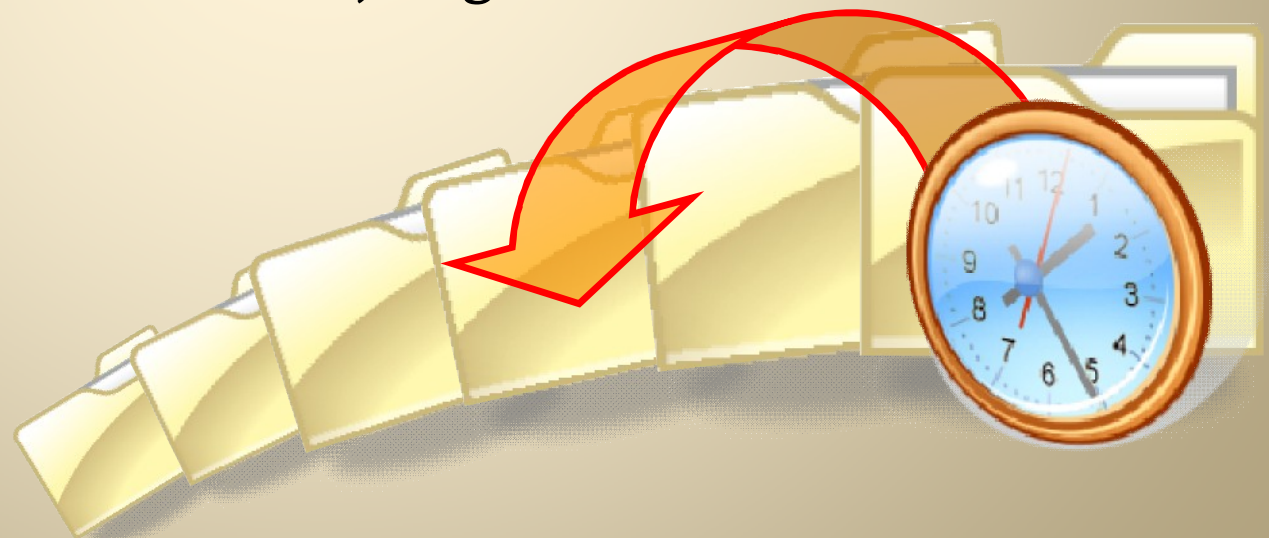
# Refining GC

- Better garbage collector is much needed
  - ✓ Better data retention policy to prevent disk full
  - ✓ Self-regulating speed
  - ✓ Smarter selection algorithm of target segments to reduce I/O

- Further chance of optimization and enhancement
  - ✓ Background data verification
  - ✓ Defragmentation
  - ✓ De-duplication
  - ✓ Background disk format upgrade

# Conclusion

- NILFS is in the mainline kernel
  - ✓ You can go back in time just before you scream "**Ohhh Nooo...!!**"
  - ✓ Instant failure recovery. Simple administration.
  - ✓ Potential for innovative application
  - ✓ … and most importantly, WORKING STABLY :)

- Contribution is welcome
  - ✓ Various topics in GC, snapshot tools, and time-oriented tools.
  - ✓ Let's drop the (EXPERIMENTAL) flag!

# Questions?

- Project page
  - ✓ http://www.nilfs.org/

- Mailing-list
  - ✓ users (at) nilfs.org
  - ✓ users-ja (at) nilfs.org

- Contact Information
  - ✓ Ryusuke KONISHI <ryusuke (at) osrg.net>

# Thank you for listening !