# WalB: Block-level WAL for Efficient Incremental Backup

Dec 2, 2010

Takashi HOSHINO

Cybozu Labs, Inc.

# Contents

- Motivation
- WalB
  - Architecture
  - Main Algorithm
  - Data Format
  - Pros and Cons
- Current Progress
- Summary

# Motivation

- There is no good backup solution
  - Online
  - Small performance overhead
  - Supports various applications
  - Cost-effective

- We need
  - Consistent full backup
  - Incremental backup
  - Block-level backup
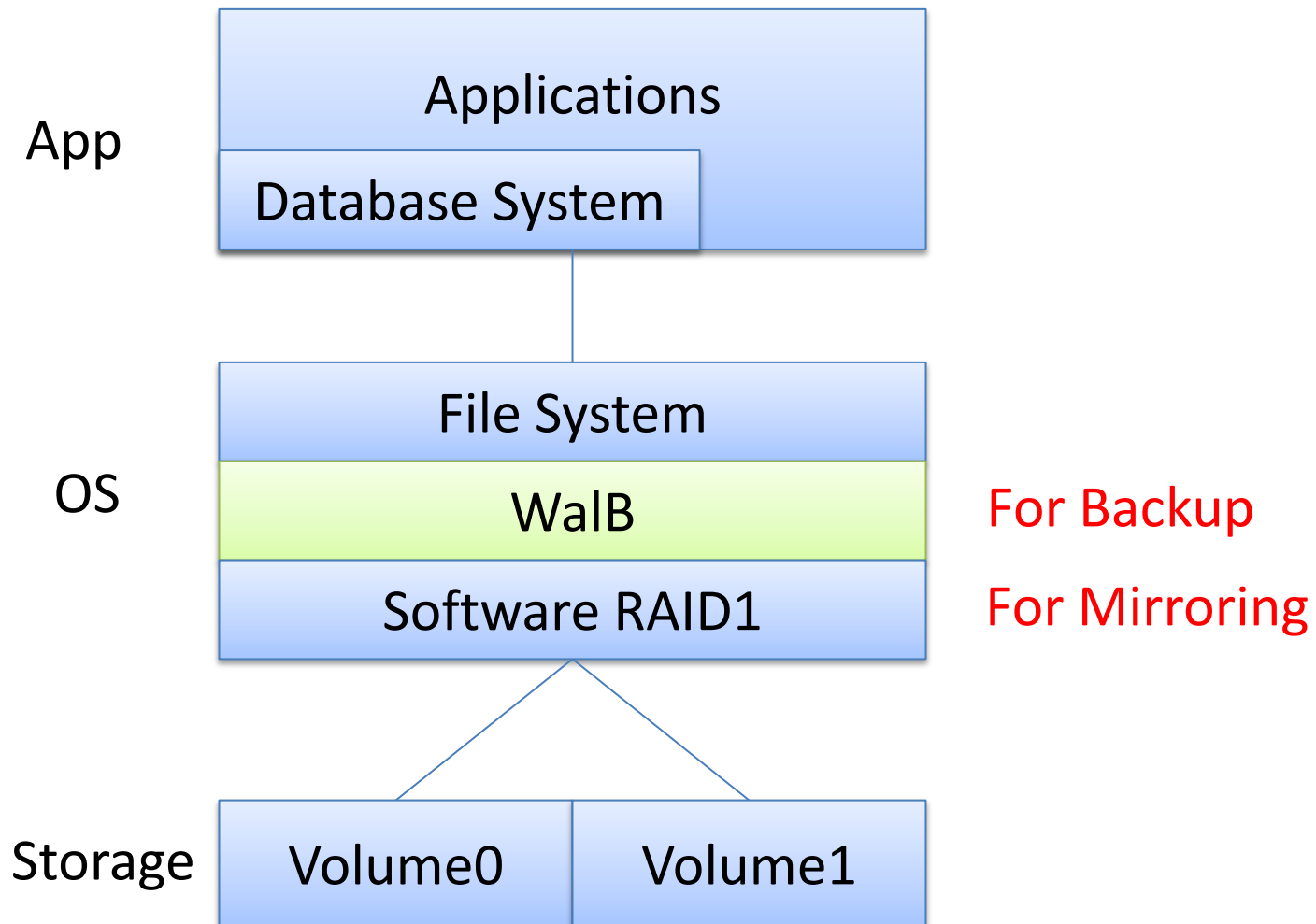  - To use commodity hardware and free software only

# Requirements inside Cybozu

- Guarantee backup interval
  - within 5-10min
- Keep multiple backup archives
  - also in remote site
- Keep multiple snapshots
  - per 1day for 1week

# WalB

- A wrapper block device driver
  - Data device to store data
  - Log device to store WAL (write-ahead log)
- Related user-land tools
  - Device controller
  - Log extractor

- Target OS and architecture
  - Latest Linux kernel (2.6.36)
  - x86_64 host
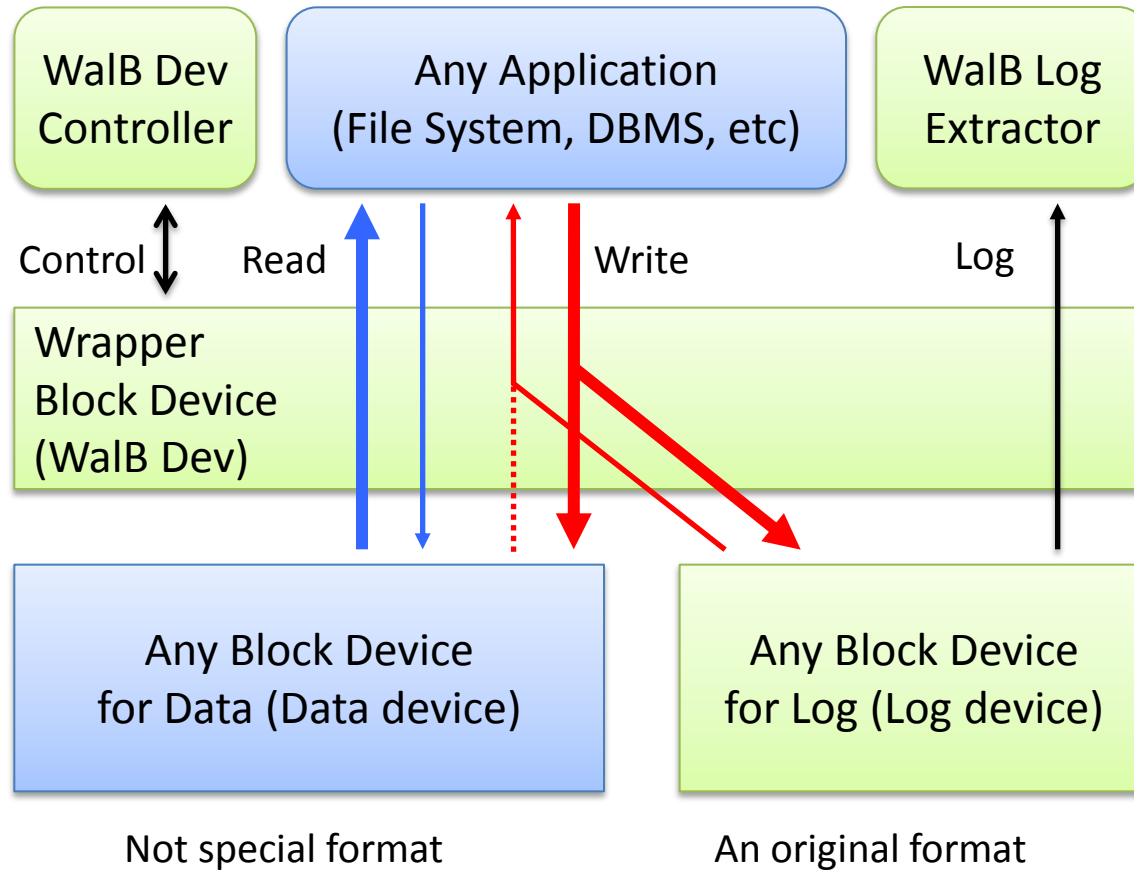
# System Architecture with WalB

App

Applications

Database System

OS

File System

WalB — For Backup

Software RAID1 — For Mirroring

Storage

Volume0

Volume1

# Backup vs Mirroring

| | Backup | Mirroring |
|---|---|---|
| Typical methods | Making snapshot, Logging writes | RAID1, Replication |
| Recoverable failures | Operation miss, Application bug | Failure of facilities |
| Can keep latest data? | No | Yes |

We need both functionalities to save data from lost.

# WalB Architecture



WalB Dev Controller

Any Application (File System, DBMS, etc)

WalB Log Extractor

Control

Read

Write

Log

Wrapper Block Device (WalB Dev)

Any Block Device for Data (Data device)

Any Block Device for Log (Log device)

Not special format

An original format
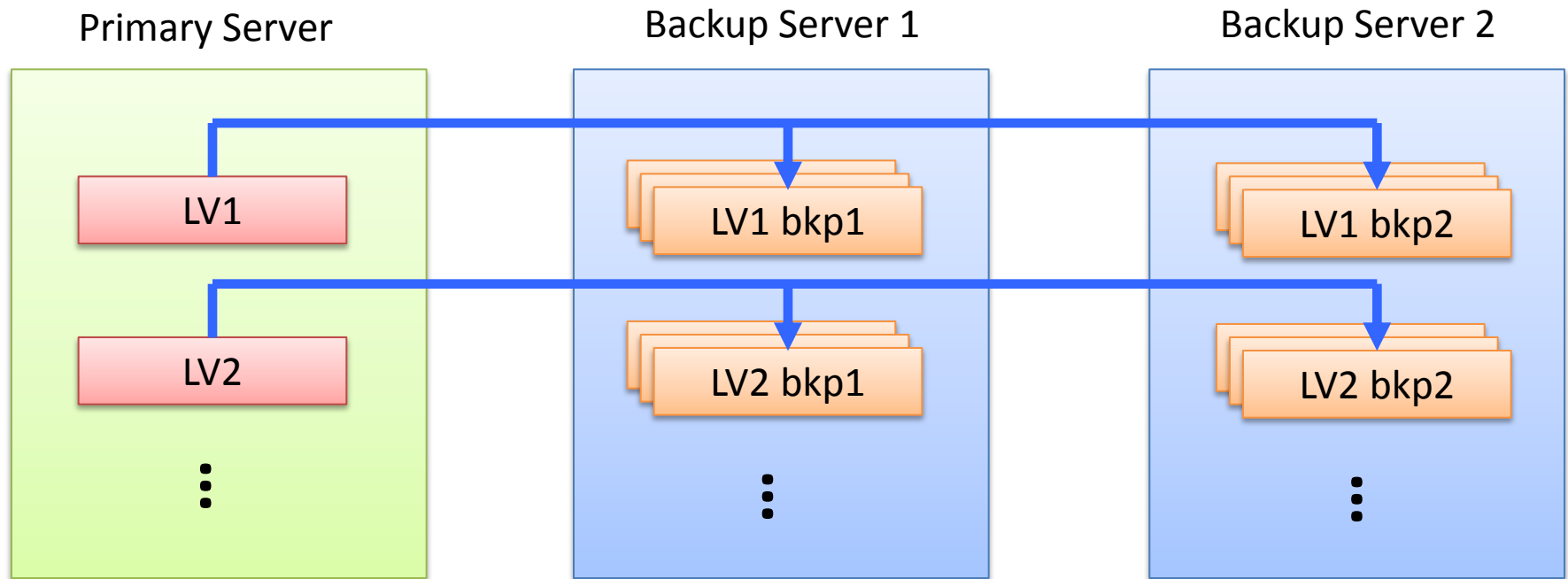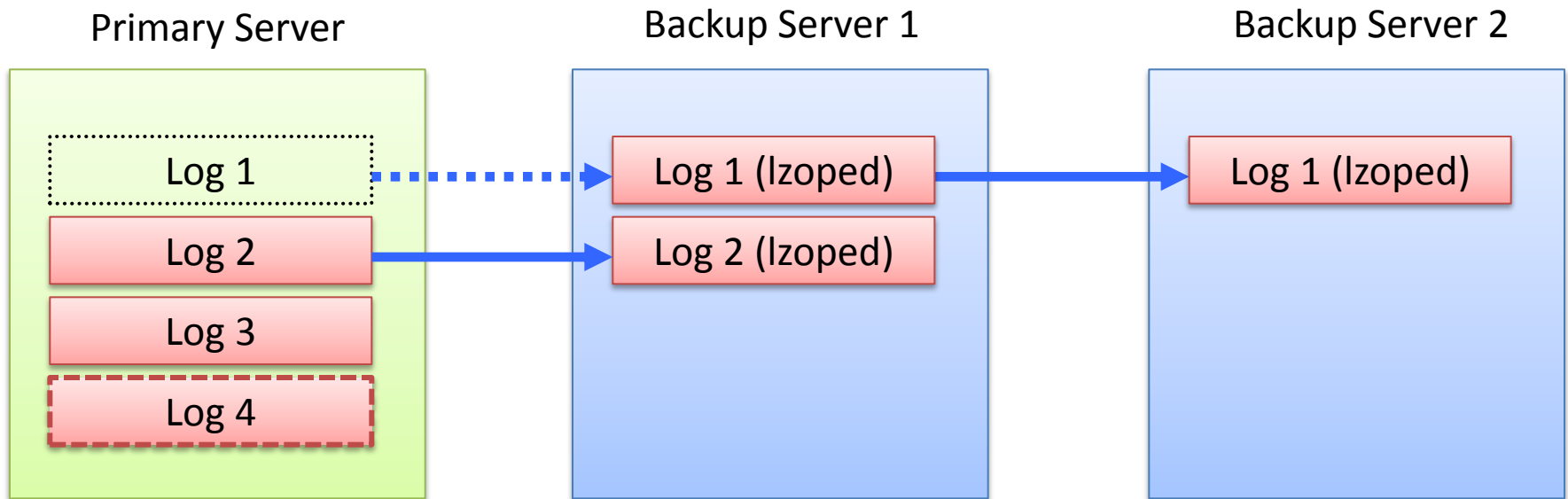
# WalB Functionalities

- Online incremental backup

- Online consistent full backup


- Snapshot creation/deletion
  - not accessible due to no index

- Volume resize
  - not resize of log capacity
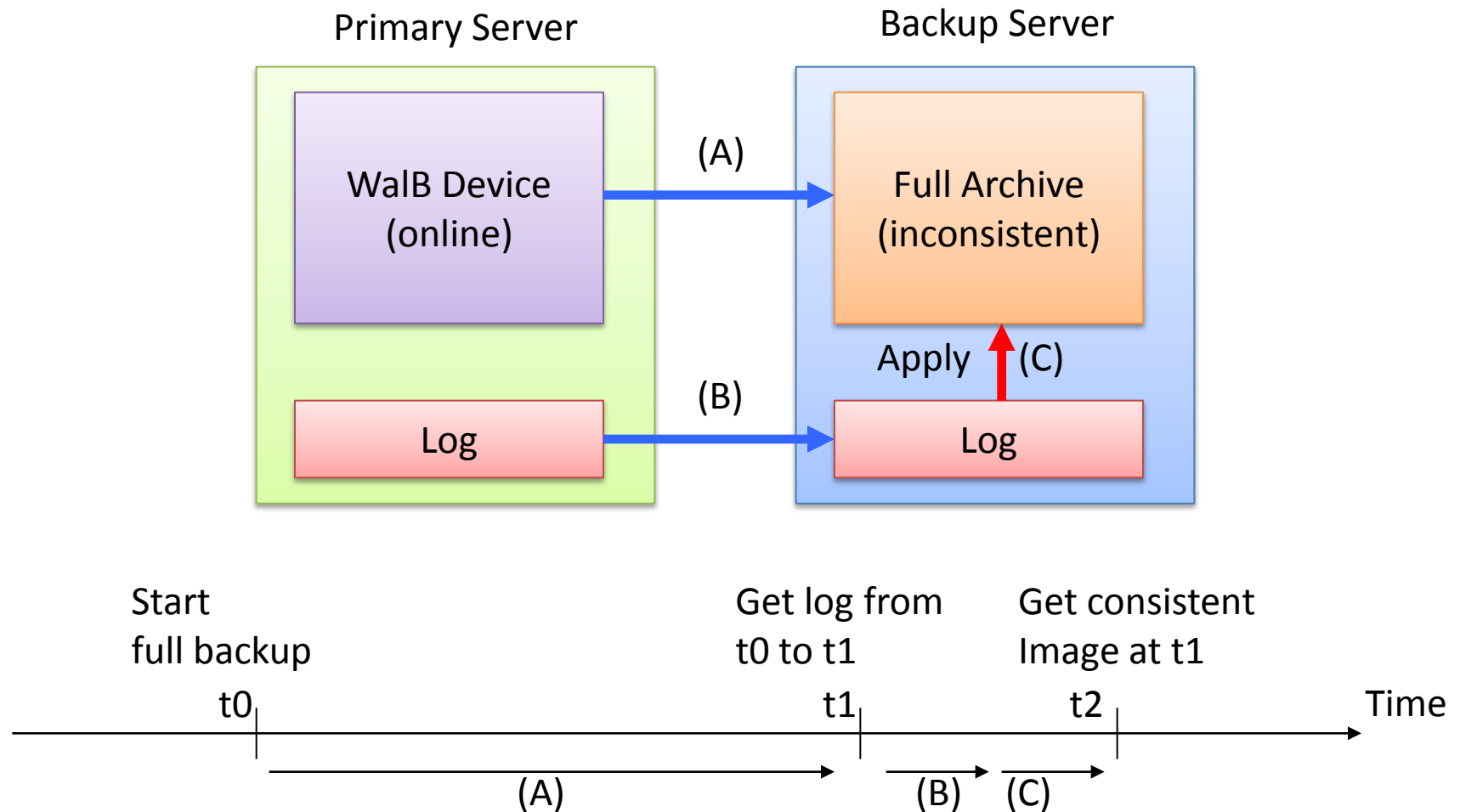
# Incremental Backup with WalB

# Log Transfer



- Each log is compressed with lzop and transferred via ssh/rsh.
- Proxy may be useful for pipelined transfer to multiple sites.
- Delete original log in primary server after all backup servers have a replica.

# Consistent Full Backup with WalB

# Read/Write Algorithm (1)

Request
Queue

Write

Read

• Generate logpack header
• Generate request for log device
• Issue the request

Wait completion

• Generate and issue to the data device

Wait completion

• Send completion for upper layer

Simple Algorithm

• Generate request for data device
• Issue the request

Wait completion

• Send completion for upper layer
• Update written_lsid
• Free logpack header buffer

# Read/Write Algorithm (2)

Request
Queue

Write

Read

• Generate logpack header
• Allocate data buffer and copy data for data device write
• Generate request for log device
• Issue the request

Wait completion

• Search PendingTree
  If all data are in the tree
  Then copy data and send completion for upper layer
  Else generate request for lack data and issue the request to the data device

Wait completion

• Insert to PendingTree
• Send completion for upper layer
• Generate request for data device
• Issue the request

Wait completion

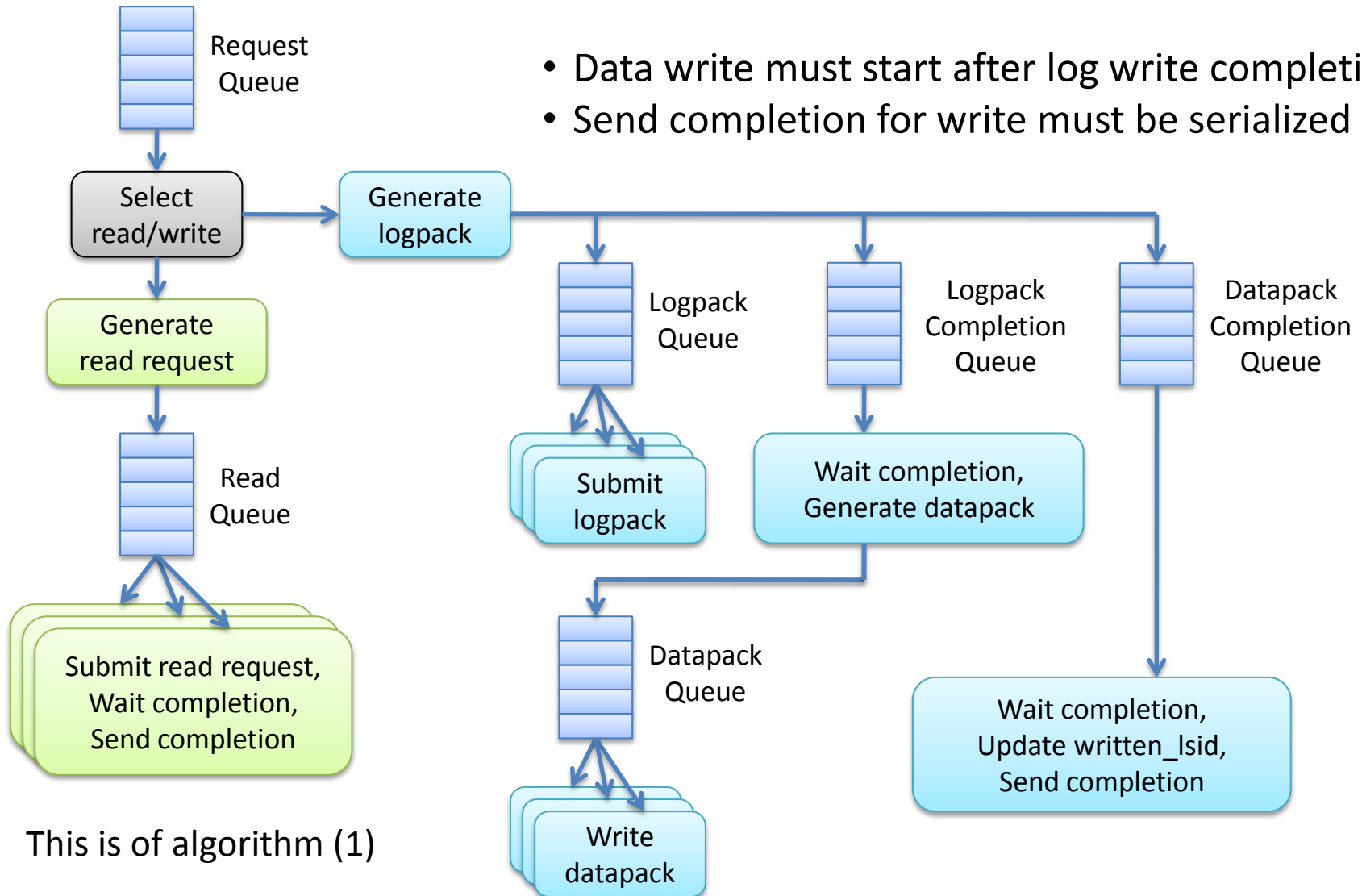• If all requests have finished
Then send completion for upper layer

Faster but more complex than (1)

• Delete from PendingTree
• Free data buffer for data device write
• Update written_lsid
• Free logpack header buffer

# Parallel Task Processing

- Data write must start after log write completion
- Send completion for write must be serialized

Request Queue

Select read/write

Generate logpack

Generate read request

Logpack Queue

Logpack Completion Queue

Datapack Completion Queue

Read Queue

Submit logpack

Wait completion, Generate datapack

Submit read request, Wait completion, Send completion

Datapack Queue

Wait completion, Update written_lsid, Send completion

This is of algorithm (1)

Write datapack
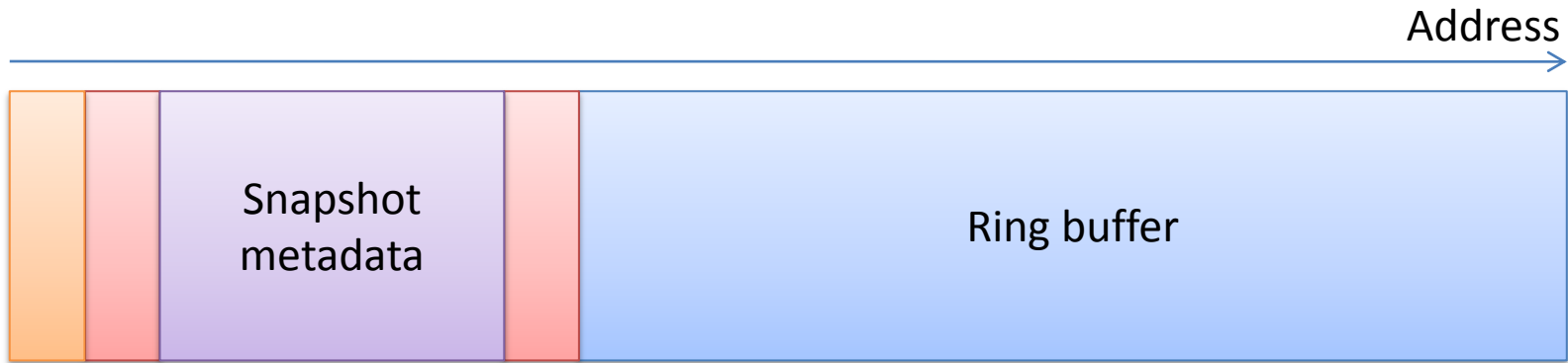
# WalB Data Format

- Data device
  - The same image as wrapping block device
- Log device
  - Overview
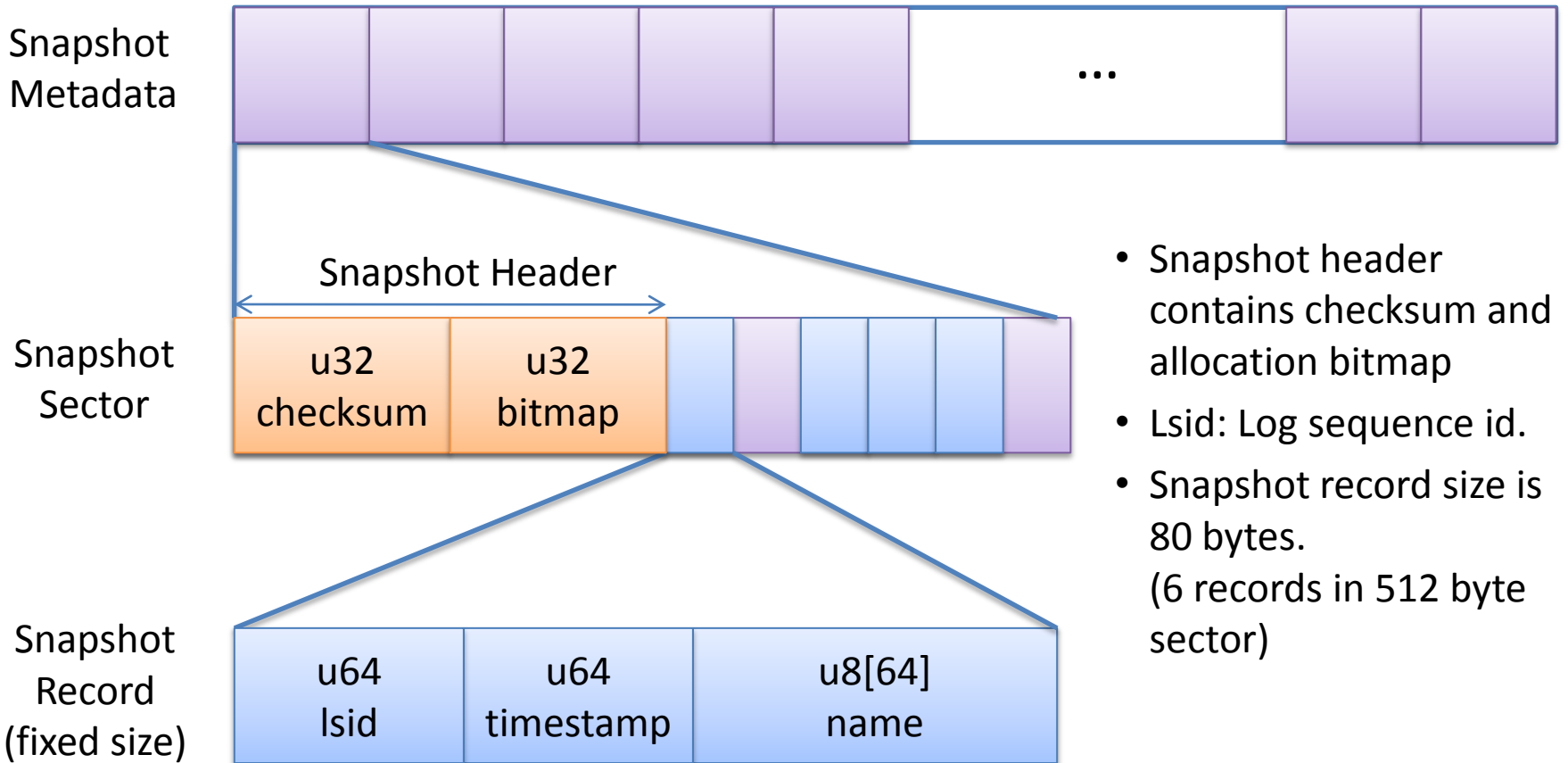  - Snapshot metadata
  - Ring buffer
  - Logpack

# Log Device Format

Address



- Snapshot metadata
  - The size is determined at device creation.
- Ring buffer
  - Stores write-ahead log.
  - The size is determined at device creation.
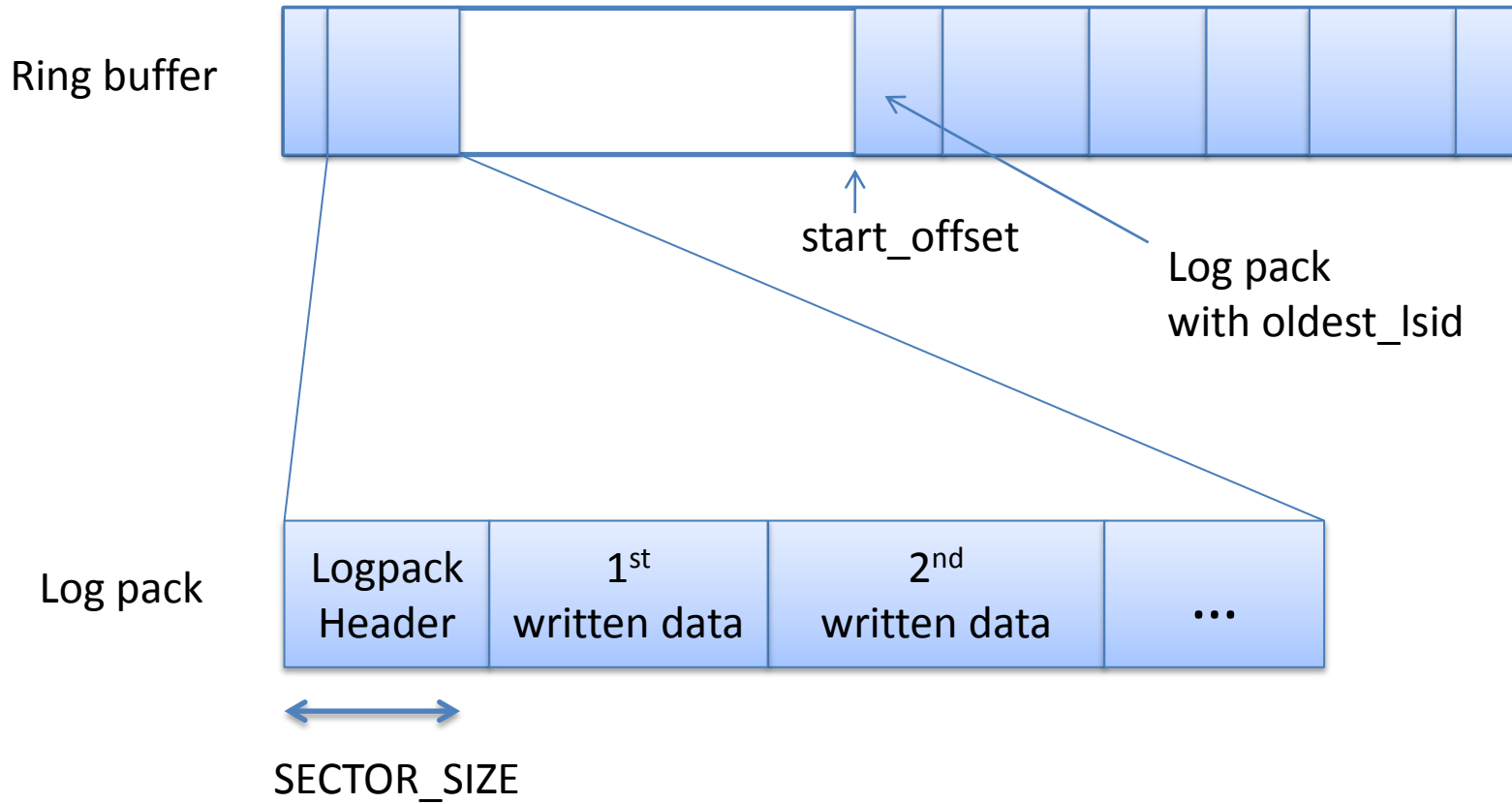  - The size can not be changed.

Superblock
(SECTOR_SIZE)

Reserved (PAGE_SIZE)

PAGE_SIZE = 4096 bytes
SECTOR_SIZE = 512 or 4096 bytes

# Snapshot Metadata

**Snapshot Metadata**



**Snapshot Sector**

Snapshot Header

| u32 checksum | u32 bitmap |
| --- | --- |

**Snapshot Record (fixed size)**

| u64 lsid | u64 timestamp | u8[64] name |
| --- | --- | --- |

- Snapshot header contains checksum and allocation bitmap
- Lsid: Log sequence id.
- Snapshot record size is 80 bytes.
  (6 records in 512 byte sector)

# Ring Buffer



Ring buffer

start_offset

Log pack
with oldest_lsid

Log pack

| Logpack Header | 1st written data | 2nd written data | ... |

SECTOR_SIZE

# Logpack Header

**Logpack Header**

| u32 checksum | u16 # of IO | u16 Total IO size | 1st | 2nd | 3rd | ... | |
|---|---|---|---|---|---|---|---|

Log Record

**Log Record**

```
u64 lsid;        /* Log sequence id */
u64 offset;      /* IO offset by the sector. */
u16 lsid_local;  /* local sequence id
                    as the data offset in the log record. */
u16 size;        /* IO size by the sector. */
u16 is_exist;    /* 0 if this record does not exist. */
u16 reserved1;
```

# Pros and Cons

| | WalB (redo log) | Snapshot with redo log | Snapshot with undo log |
|---|---|---|---|
| Read | No overhead | Index search | Bitmap search |
| Write | Write twice (1) / No overhead (2) | Index modification | Bitmap search/modification and old data copy |
| Typical Soft. | --- | ZFS, BtrFS | LVM |

WalB + Index ~= Block device with snapshot management

# Current Progress

- Survey and study linux kernel programming
- Design of rough architecture
- <span style="color:red">Prototype implementation</span>
- Basic evaluation and redesign if required
- Implementation of full functionalities and test
- Operation inside Cybozu
- Publication as GPLv2

- Merging to device-mapper if required
- Merging to main repository (hopefully)

# Summary

- Motivaion
  - No good backup solution
    covering various applications

- WalB
  - Is a block device driver with WAL
  - Provides efficient incremental backup