# Petabytes Data Migration
# and Load Balancing
## with Football + MARS on Enterprise-Critical Data

**1&1**

## LCA 2019 Presentation by Thomas Schöbel-Theuer

# Container Football: Agenda
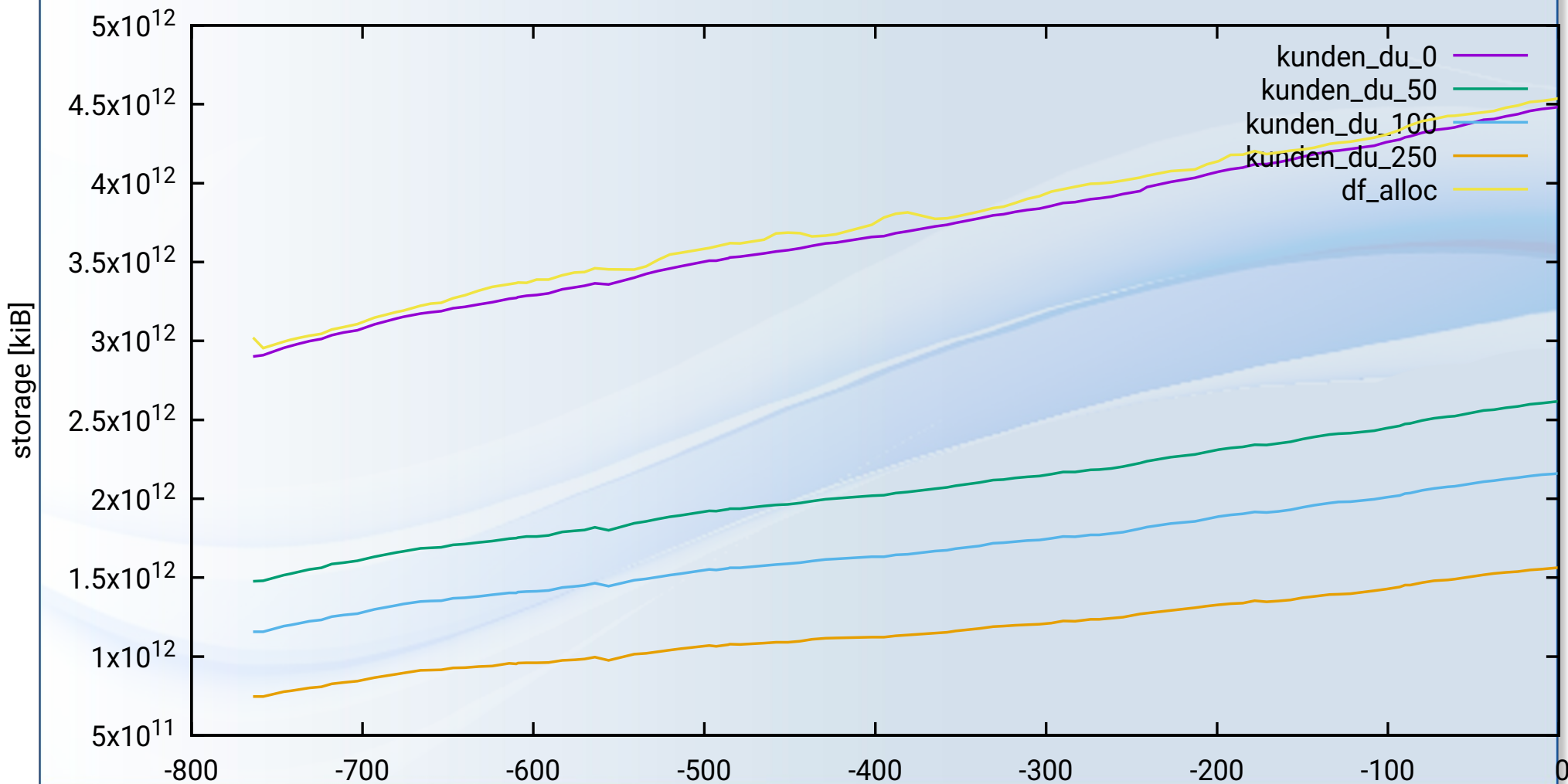
**1&1**

> **New method for load balancing**

- **Motivation: data growth > 20% / year**

- **HOWTO Container Football = Background Migration of LVs**

  **e.g. for load balancing, HW lifecycle, etc**

- **The Football Automation Project**

- **Current Status / Future Plans**

# Motivation: Growth at 1&1 ShaHoLin = Shared Hosting Linux

1&1

- **~ 9 millions of customer home directories**

- **~ 10 billions of inodes**

- **> 4.5 petabytes *allocated* in ~ 2700 xfs instances, LVM ~ 8 PB x 2 for geo-redundancy via MARS**

- **Growth rate ~ 21 % / year**

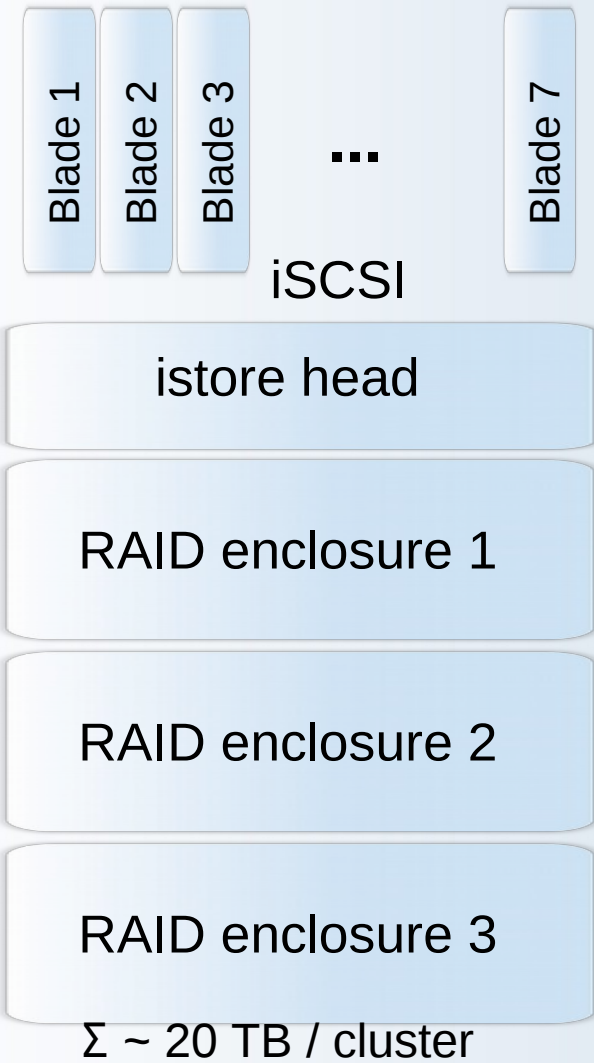# Motivation: Growth at 1&1 ShaHoLin = Shared Hosting Linux

1&1

storage-icpus+istores Tue Jan  1 00:00:00 CET 2019

Tue Jan  1 00:00:00 CET 2019 -- Days in the Past [d]

# Hardware Lifecycle

**1&1**

Blade 1  Blade 2  Blade 3  ...  Blade 7

iSCSI

istore head

RAID enclosure 1

RAID enclosure 2

RAID enclosure 3

Σ ~ 20 TB / cluster

**x 2** for geo-redundancy

"Efficiency project"
x 1800 old blades
in 4 datacenters

~ 7 LXC containers

OLD          NEW

==>

1 U hypervisor

64 CPU threads, local RAID

* better power consumption
* better customer performance
* better TCO

**x 2** for geo-redundancy

# Football: Basic Operations

**1&1**

- **migrate: move pairs of LVs to new pairs of hypervisors**
  ```
  ./football.sh migrate infong4711 cluster1234
  ```

- **shrink: use local rsync for in-place resizing (downtime)**
  ```
  ./football.sh shrink infong4711 75%
  ```

- **expand: online lvresize + marsadm resize + xfs_growfs**
  ```
  ./football.sh expand infong4712 75%
  ```

- **several combined operations, e.g. migrate+shrink (less network traffic)**

# HOWTO Container Football = Background Migration of LVs

**1&1**

## HOST A (old) VM is running → HOST B (new) has spare space

- lvdisplay /dev/vg/$mydata

- 

- 

- (meanwhile VM is altering data)

- $vmmanager stop /dev/mars/$mydata

- 

- 

- marsadm leave-resource $mydata

- lvremove /dev/vg/$mydata
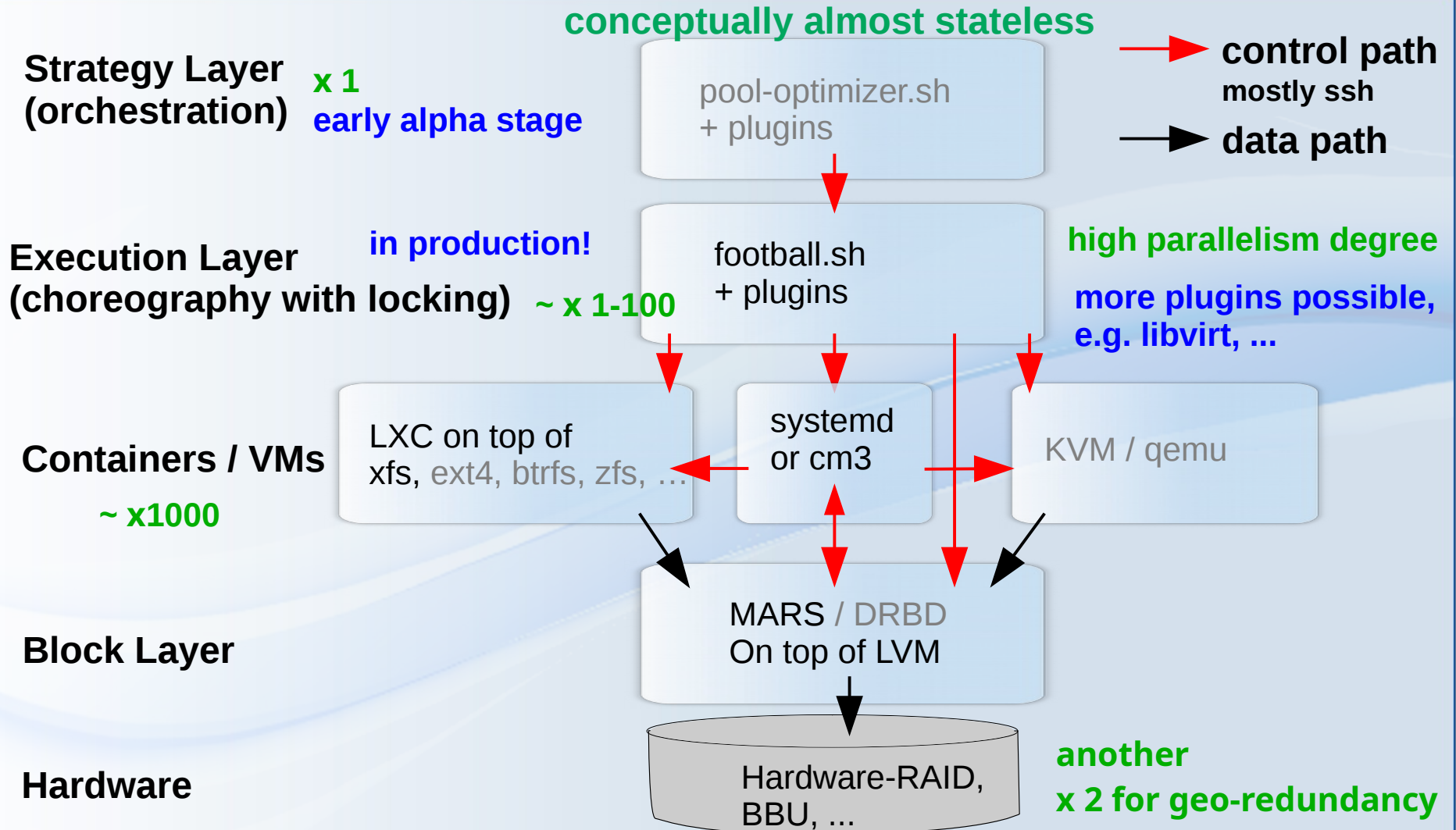
- 

- lvcreate -L $size -n $mydata vg

- marsadm join-resource $mydata \
  /dev/vg/$mydata

- marsadm view: wait for UpToDate

- 

- marsadm primary $mydata

- $vmmanger start /dev/mars/$mydata

- 

- 

*start kick*

*done kick*

*cleanup kick*

### => also works with 2 old replicas → 2 new replicas

**Toolset: football.sh in github.com/schoebel/football**

# Football Architecture (grey = not yet implemented)

1&1

**Strategy Layer
(orchestration)**   x 1
**early alpha stage**

**conceptually almost stateless**

pool-optimizer.sh
+ plugins

control path
**mostly ssh**
**data path**

**Execution Layer
(choreography with locking)**   ~ x 1-100

**in production!**

football.sh
+ plugins

**high parallelism degree**

**more plugins possible,
e.g. libvirt, ...**

**Containers / VMs**

~ x1000

LXC on top of
xfs, ext4, btrfs, zfs, ...

systemd
or cm3

KVM / qemu

**Block Layer**

MARS / DRBD
On top of LVM

**Hardware**

Hardware-RAID,
BBU, ...

**another
x 2 for geo-redundancy**

# Football Current Status

- GPL with lots of plugins, some generic, some 1&1-specific
  - about 2/3 of code is generic
  - plugins/football-basic.sh uses `systemd`
  - https://github.com/schoebel/football
  - https://github.com/schoebel/mars
- Multiple operations:
  - migrate $vm $target_cluster
    - low downtime (seconds to few minutes)
  - shrink $vm $target_percent
    - uses local incremental rsync, more downtime
  - expand $vm $target_percent
    - online, no downtime
- In production at internal Efficiency project
  - get rid of old hardware
  - Concentrate ~ 7 LXC containers on 1 hypervisor

  - currently >50 „kicks" per week
    - limited by hardware deployment speed
    - Proprietary Planner (for HW lifecycle)
  - Almost finished: ~70% of ~1800 blades already migrated (mid of January 2019) and mostly shrunk

# Sponsoring (MARS + Football)

- Best for > 1 PiB of enterprise-critical data
  - Example: ShaHoLin (slide3)
  - More plugins in future, e.g. for KVM, ...
- Future pool-optimizer will deliver similar functionality than **Kubernetes**
  - but on stateful storage + containers instead of stateless Docker containers
  - State is in the storage and in the machines, but not in orchestration
- Long-term perspective
  - MARS is largely complementary to DRBD
  - Geo-redundancy with OpenSource components
  - distances > 50km possible, tolerates flaky replication networks
  - **Price / performance** better than anything else (see mars-manual.pdf)
  - Architectural **reliability** better than BigCluster with cheaper hw + network!
- ask me: decades of experience with enterprise-critical data and long-distance replication

# Appendix

# MARS Current Status

■ MARS source under GPL + docs:

> `github.com/schoebel/mars`
> `mars-manual.pdf` ~ **100 pages**

■ mars0.1stable productive since 02/2014

■ Backbone of the 1&1 geo-redundancy feature

■ MARS status January 2018:

> \> 5800 servers (shared hosting + databases)
>
> \> 2x12 petabyte total
>
> ~ 10 billions of inodes in > 2500 xfs instances, biggest ~ 40 TB
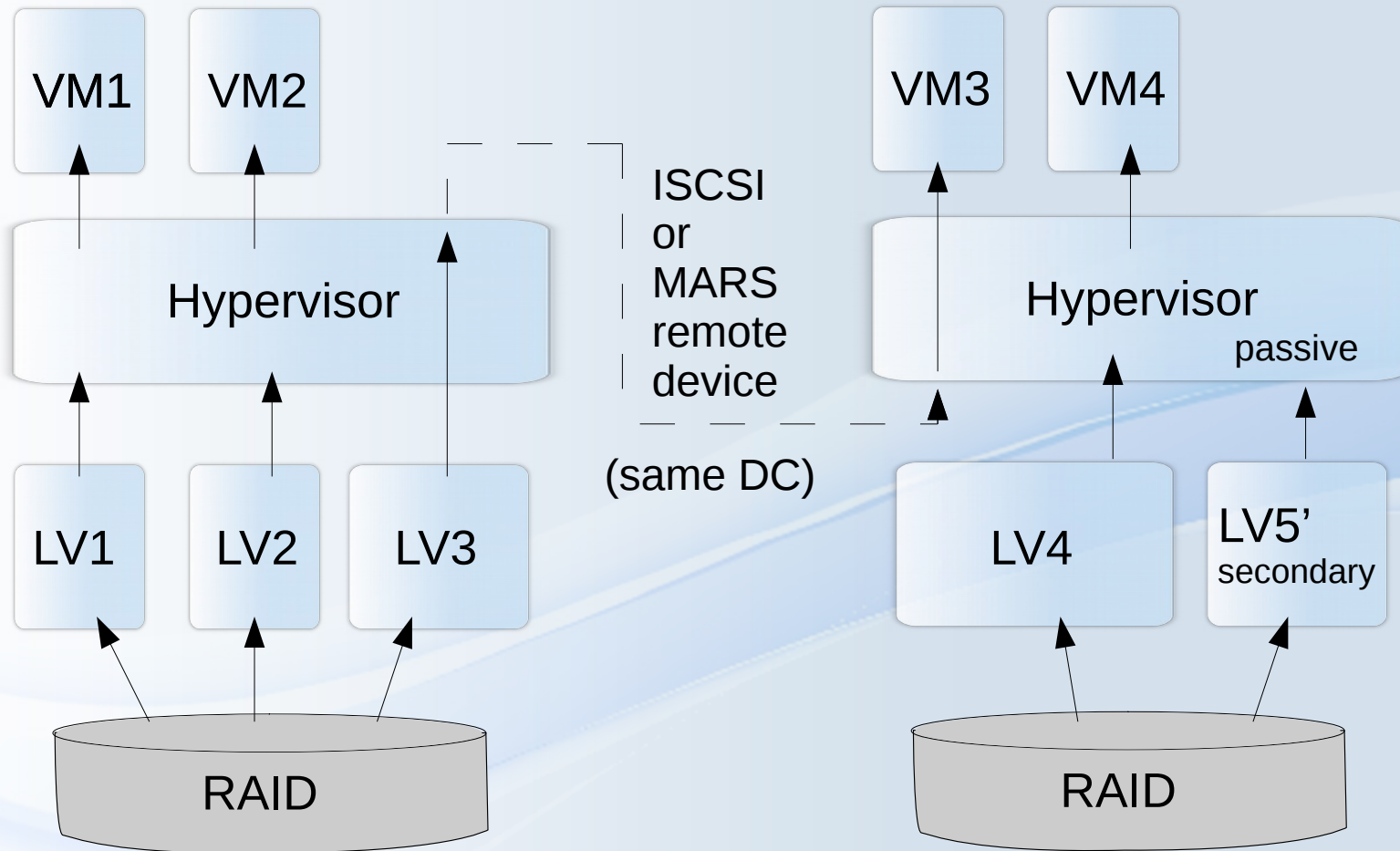>
> <= 10 LXC Containers on 1 Hypervisor

■ New internal Efficiency project
  – Concentrate more LXC containers on 1 hypervisor

  – New public branch mars0.1b with many new features, e.g. mass-scale clustering, socket bundling, remote device, etc
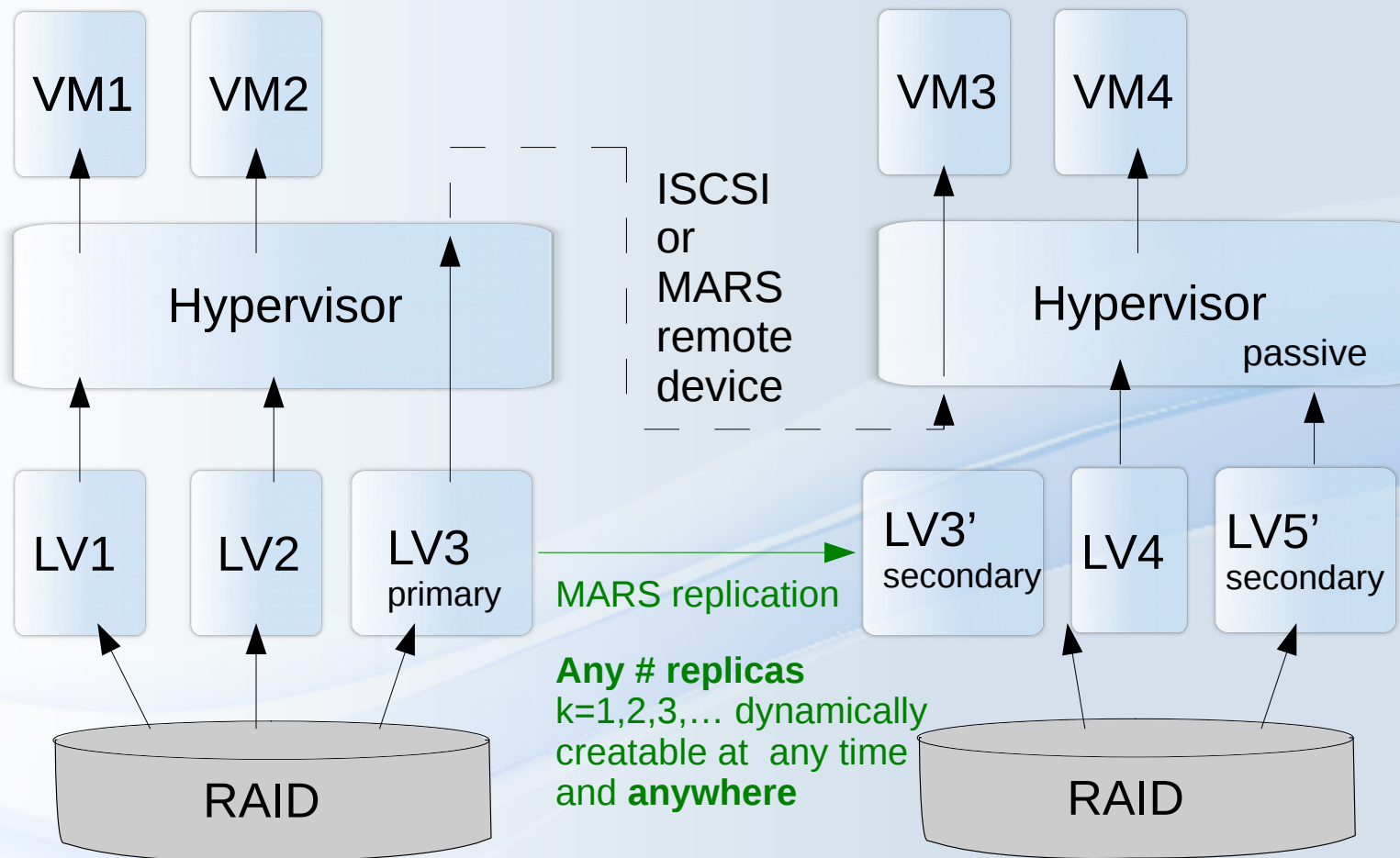
  – mars0.1b currently in ALPHA stage

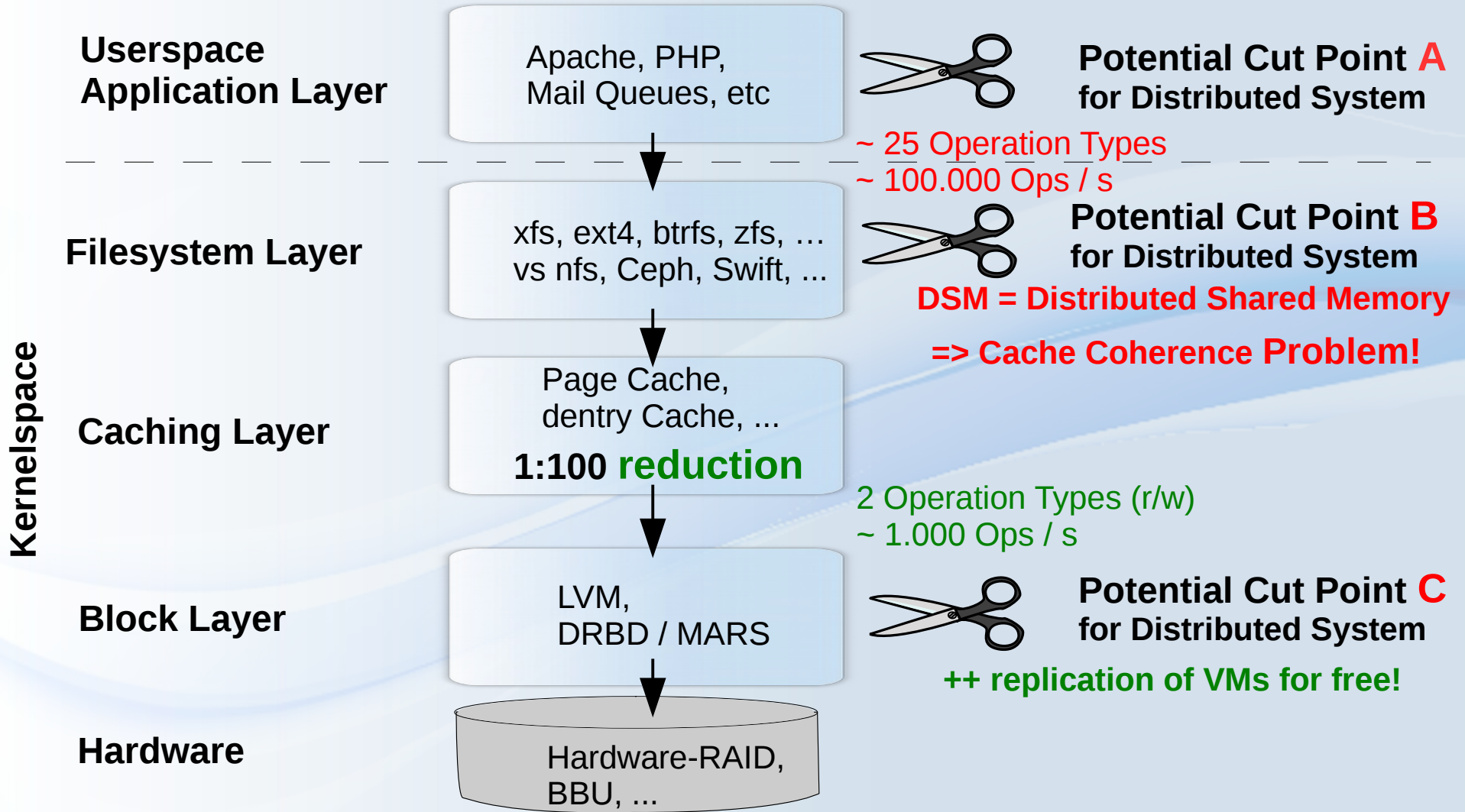# Flexible MARS Sharding + Cluster-on-Demand



any hypervisor works in client and/or server role
and preferably **locally** at the same time

# Flexible MARS Background Migration



=> any hypervisor may be source or destination of some LV replicas at the same time

# Replication at Block Level vs FS Level

**1&1**

**Kernelspace**

| Userspace Application Layer | Apache, PHP, Mail Queues, etc | ✂ **Potential Cut Point A** for Distributed System |

~ 25 Operation Types
~ 100.000 Ops / s

| Filesystem Layer | xfs, ext4, btrfs, zfs, … vs nfs, Ceph, Swift, ... | ✂ **Potential Cut Point B** for Distributed System |

**DSM = Distributed Shared Memory**

**=> Cache Coherence Problem!**

| Caching Layer | Page Cache, dentry Cache, ... **1:100 reduction** |

2 Operation Types (r/w)
~ 1.000 Ops / s

| Block Layer | LVM, DRBD / MARS | ✂ **Potential Cut Point C** for Distributed System |

**++ replication of VMs for free!**

| Hardware | Hardware-RAID, BBU, ... |

# Use Cases DRBD+proxy vs MARS Light

1&1

## DRBD+proxy
**(proprietary)**

**Application area:**
- Distances: any
- Aynchronously
  - **Buffering in RAM**
- Unreliable network leads to **frequent re-syncs**
  - RAM buffer gets lost
  - at cost of actuality
- **Long** inconsistencies during re-sync
- Under pressure: **permanent** inconsistency possible
- High memory overhead
- Difficult scaling to k>2 nodes

## MARS Light
**(GPL)**

**Application area:**
- Distances: **any** ( >>50 km )
- Asynchronously
  - near-synchronous modes in preparation
- Tolerates **unreliable network**
- Anytime consistency
  - no re-sync
- Under pressure: no inconsistency
  - possibly at cost of actuality
- Needs >= 100GB in `/mars/` for transaction logfiles
  - dedicated spindle(s) recommended
  - RAID with BBU recommended
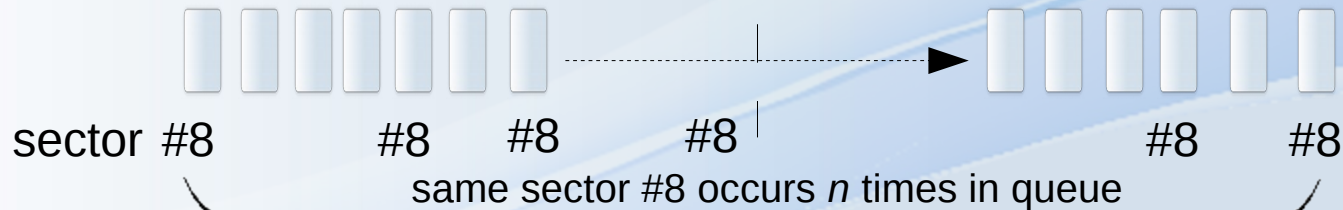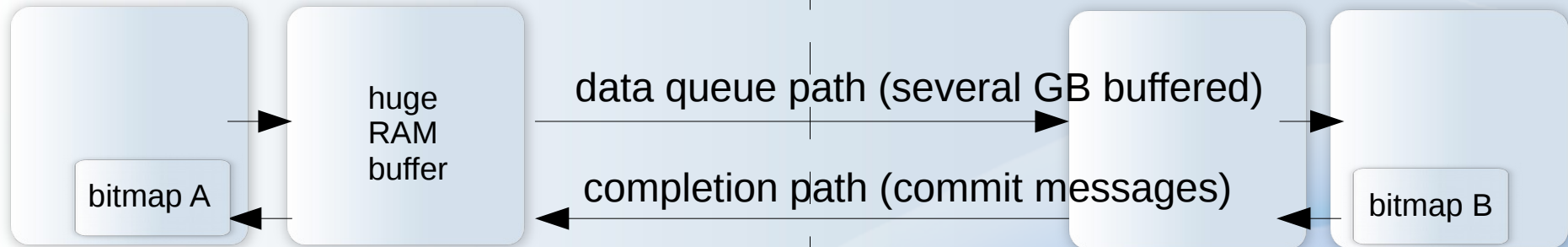- Easy scaling to k>2 nodes

# DRBD+proxy Architectural Challenge

**1&1**

**DRBD Host A (primary)**     **Proxy A'**     A != A' possible     **Proxy B' (essentially unused)**     **DRBD Host B (secondary)**

huge RAM buffer

data queue path (several GB buffered)

completion path (commit messages)

bitmap A

bitmap B

sector #8          #8     #8          #8                         #8          #8

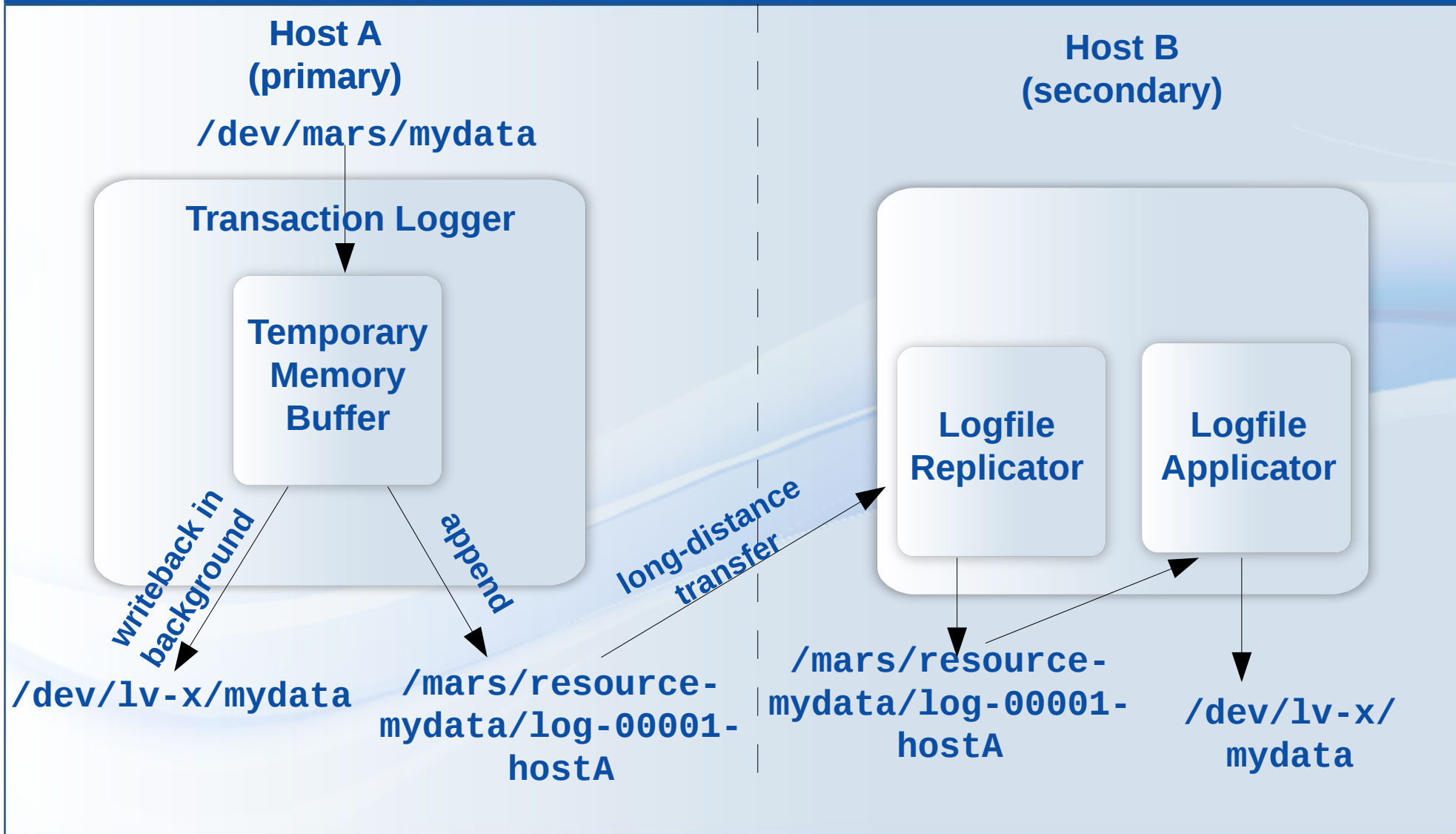same sector #8 occurs *n* times in queue

*n* times
=> need *log*(n) bits for counter
=> but DRBD bitmap has only 1 bit/sector
=> workarounds exist, but complicated
(e.g. additional dynamic memory)

# MARS Data Flow Principle

**1&1**

**Host A
(primary)**

`/dev/mars/mydata`

**Transaction Logger**

**Temporary
Memory
Buffer**

*writeback in
background*

*append*

`/dev/lv-x/mydata`

`/mars/resource-
mydata/log-00001-
hostA`

*long-distance
transfer*

**Host B
(secondary)**

**Logfile
Replicator**

**Logfile
Applicator**

`/mars/resource-
mydata/log-00001-
hostA`

`/dev/lv-x/
mydata`

# Framework Architecture for MARS + future projects

1&1

**External Software, Cluster Managers, etc**

**Userspace Interface `marsadm`**

**Framework Application Layer**
MARS Light, MARS Full, etc

| MARS Light | MARS Full | ... |

**Framework Personalities**
XIO = eXtended IO ≈ AIO

| XIO bricks | future Strategy bricks | other future Personalities and their bricks |

**Generic Brick Layer**
IOP = Instance Oriented Programming
+ AOP = Aspect Oriented Programming

Generic Bricks

Generic Objects

Generic Aspects
s