

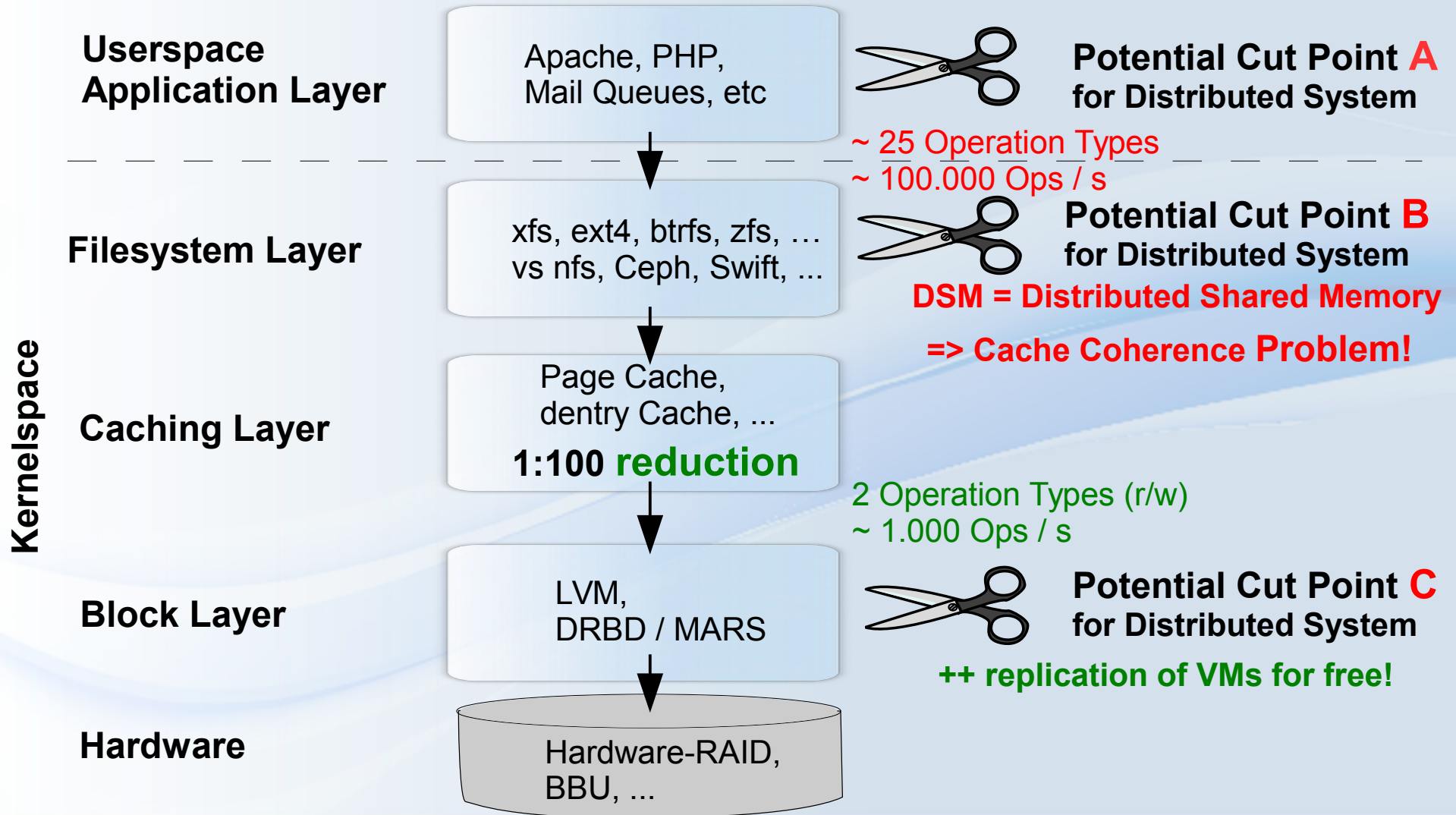
MARS: Replicating Petabytes over Long Distances



GUUG 2016 Presentation by Thomas Schöbel-Theuer

- **Long Distances: Block Level vs FS Level**
- **Long Distances: Big Cluster vs Sharding**
- **Use Cases DRBD vs MARS Light**
- **MARS Working Principle**
- **Behaviour at Network Bottlenecks**
- **Multinode Metadata Propagation (Lamport Clock)**
- **Example Scenario with 4 Nodes**
- **Current Status / Future Plans**

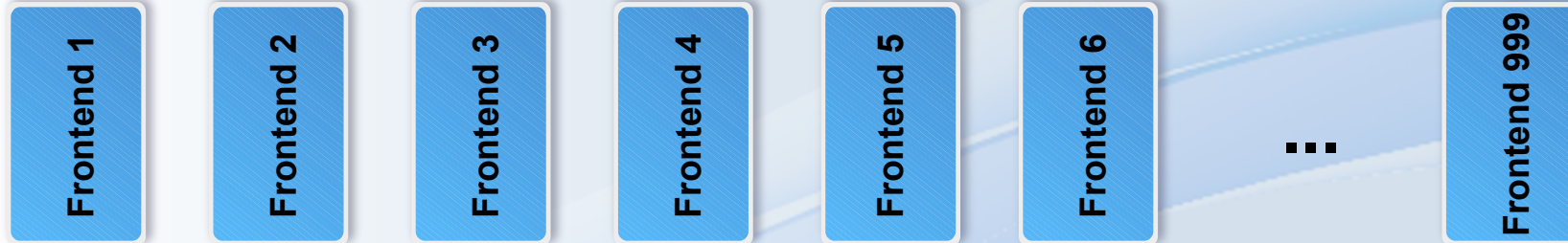
Replication at Block Level vs FS Level



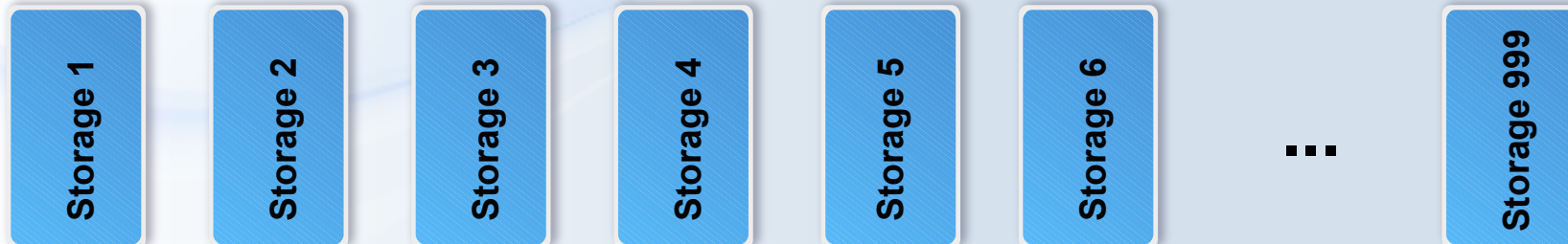
Scaling Architectures (1): **Big Cluster** vs Sharding

User 1
User 2
User 3
User 4
User 5
User 6
User 7
User 8
User 9
User 10
User 11
User 12
User 13
User 14
⋮
User 999999

Internet $O(n \cdot k)$



Internal Storage (or FS) Network $O(n^2)$



x 2 for geo-redundancy

Scaling Architectures (2): Big Cluster vs **Sharding**

User 1
User 2
User 3
User 4
User 5
User 6
User 7
User 8
User 9
User 10
User 11
User 12
User 13
User 14
⋮
User 999999

Internet $O(n*k)$ ✓

Storage + Frontend 1

Storage + Frontend 2

Storage + Frontend 3

Storage + Frontend 4

Storage + Frontend 5

Storage + Frontend 6

⋮

+++ big scale out +++

Storage + Frontend 999

++ local scalability: spare RAID slots, ...

=> method scales to petabytes

✓ X 2 for geo-redundancy

DRBD (GPL)

Application area:

- Distances: **short** (<50 km)
- Synchronously
- Needs **reliable** network
 - “**RAID-1 over network**”
 - best with crossover cables
- Short inconsistencies during re-sync
- Under pressure: long or even permanent inconsistencies possible
- Low space overhead

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 $\geq 100\text{GB}$ in `/mars/` for transaction logfiles
 - dedicated spindle(s) recommended
 - RAID with BBU recommended

MARS Working Principle

Multiversion Asynchronous Replicated Storage

Datacenter A
(primary)



`/dev/mars/mydata`

`mars.ko`

`/dev/lv-
x/mydata`

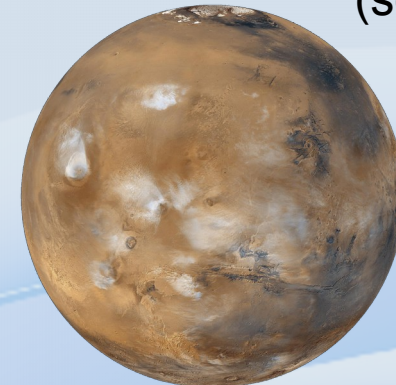
`/mars/trans-
logfile`

Similar to MySQL replication

`/mars/trans-
logfile`

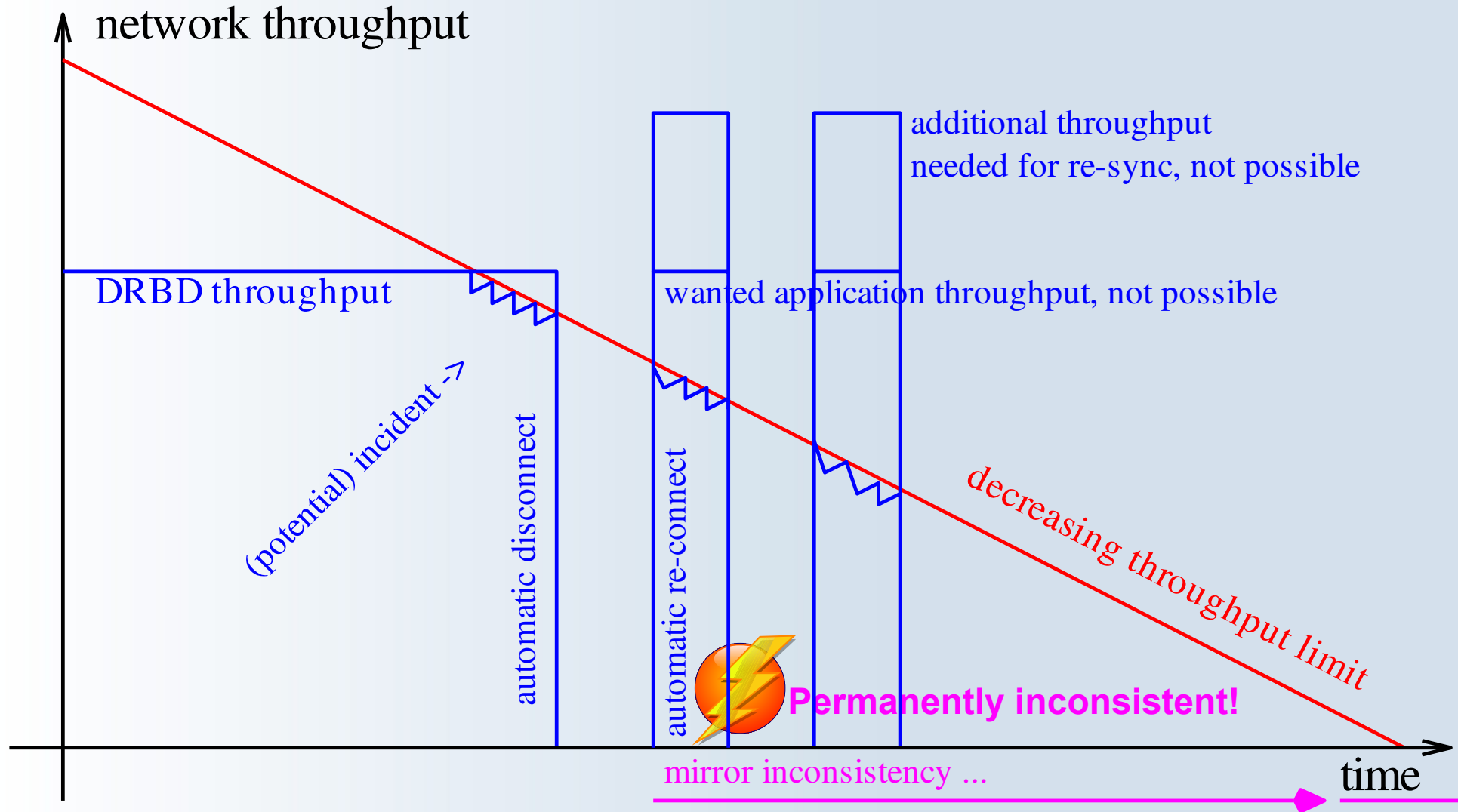
`/dev/lv-
x/mydata`

Datacenter B
(secondary)

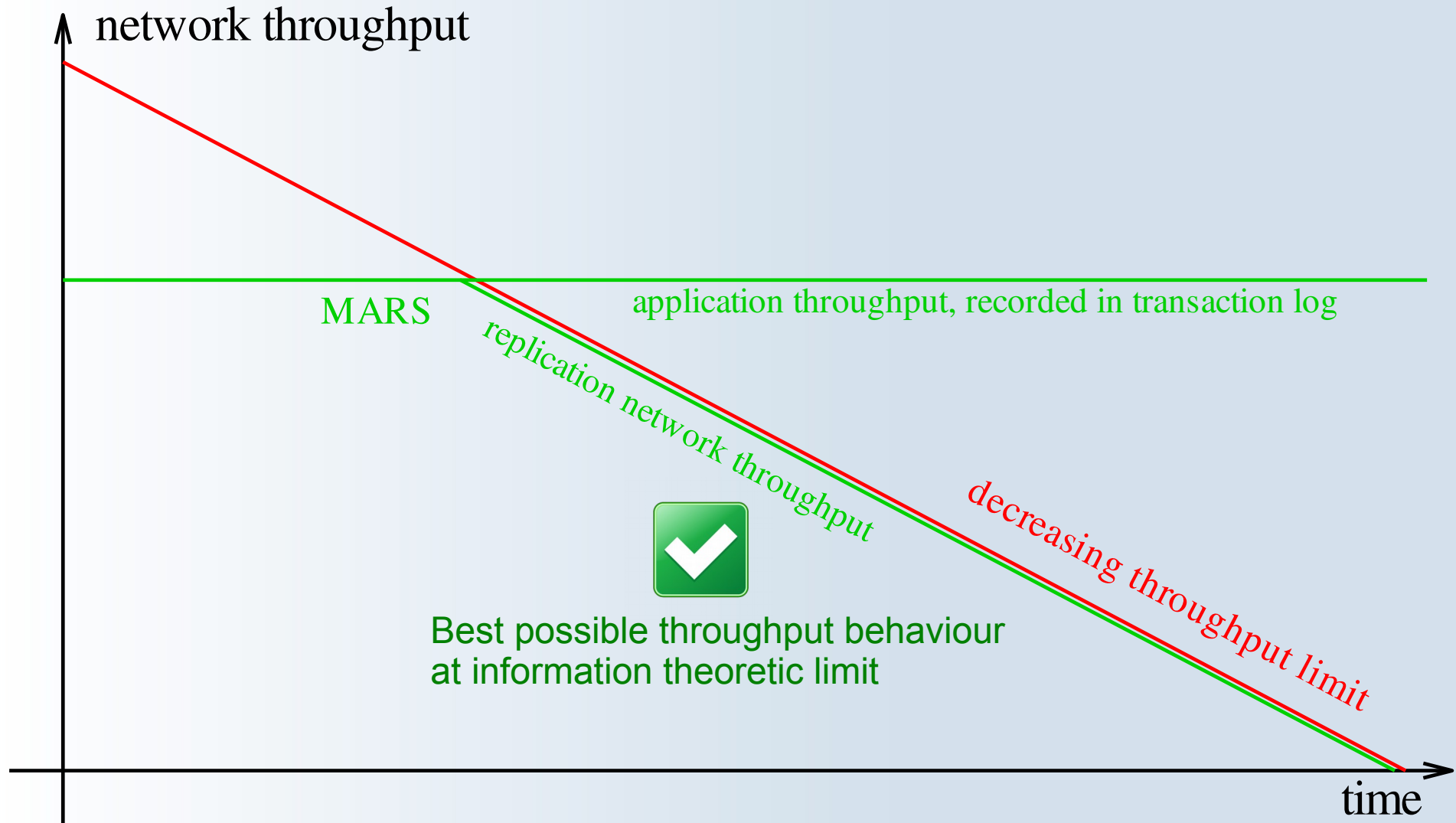


`mars.ko`

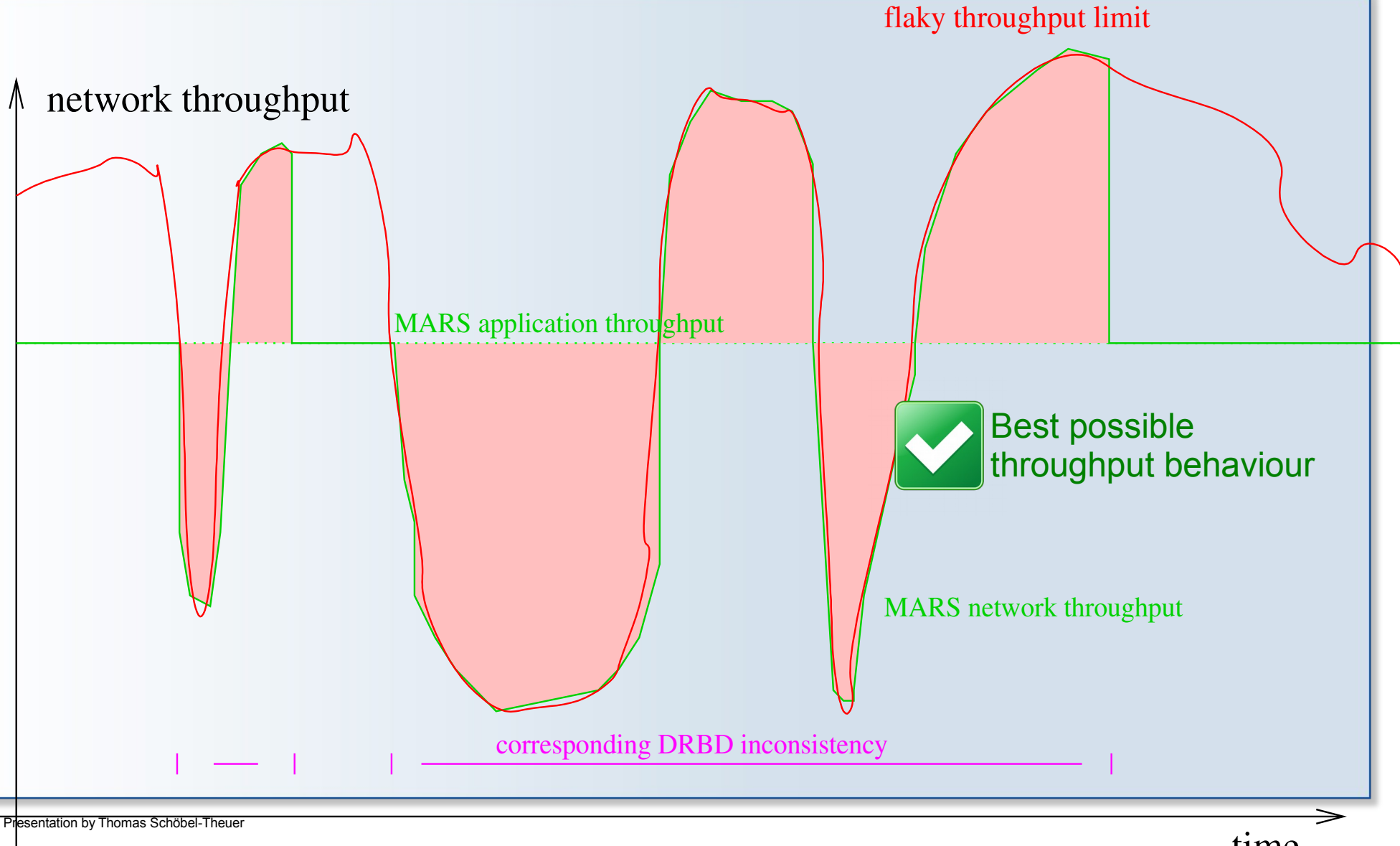
Network Bottlenecks (1) DRBD



Network Bottlenecks (2) MARS



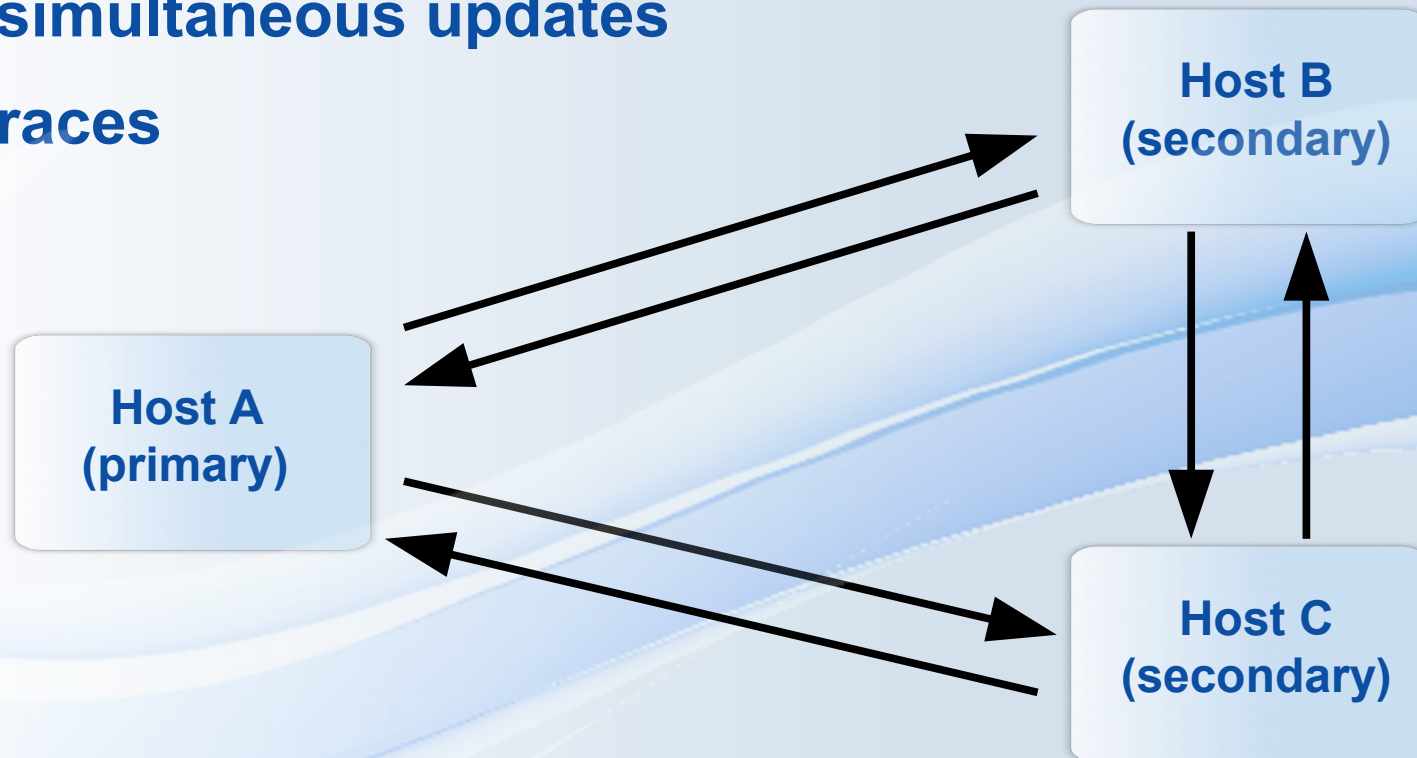
Network Bottlenecks (3) MARS



Metadata Propagation (1)

Problems for ≥ 3 nodes:

- simultaneous updates
- races



Solution: symlink tree + Lamport Clock => next slides

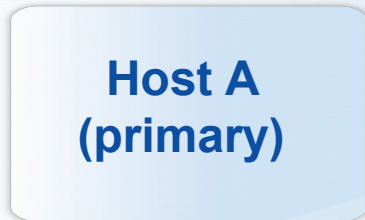
Metadata Propagation (2)



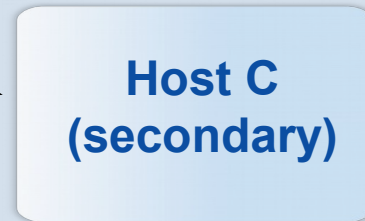
symlink tree = key->value store

Originator context encoded in key

`/mars/resource-mydata/size-hostA -> 1000`



`/mars/resource-mydata/size-hostA -> oldvalue`



Anyone knows anything about others

But later

`/mars/resource-mydata/size-hostA -> oldvalue`

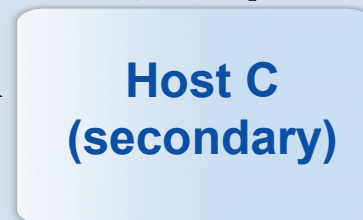
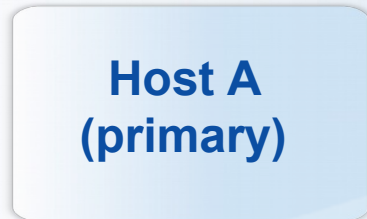
Metadata Propagation (3)

Lamport Clock = virtual timestamp

Propagation goes never backwards!

`/mars/resource-mydata/size-hostA -> veryveryoldvalue`

`/mars/resource-mydata/size-hostA -> 1000`

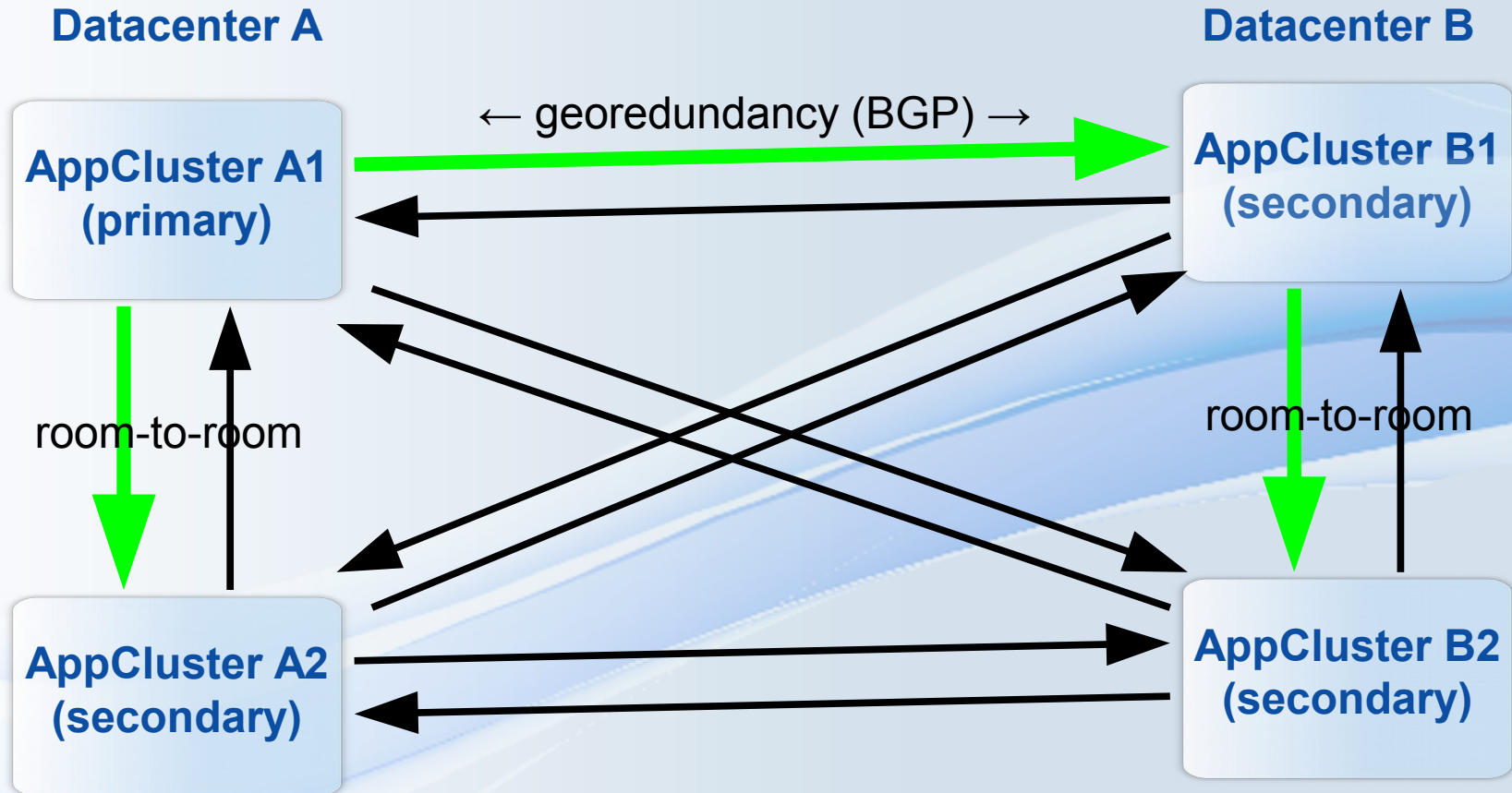


Races are compensated

Propagation paths play no role

`/mars/resource-mydata/size-hostA -> 1000`

Productive Scenario since 02/2014 (1&1 eShop / ePages)



→ potential data flow
→ actual data flow (in this scenario)

Current Status

- Source / docs at
 - `github.com/schoebel/mars`
 - `mars-manual.pdf` ~ 100 pages
- light0.1stable productive on customer data since 02/2014
- MARS status Feb 2016:
 - > 1700 servers (shared hosting + databases)
 - > 2x8 Petabyte total
 - ~ 10 billions of inodes in > 3000 xfs instances
 - > 8 millions of operating hours
- Socket Bundling (light0.2beta)
 - Up to 8 parallel TCP connections per resource
 - easily saturates 1Gbit uplink between Karlsruhe/Europe and Lenexa/USA
- WIP-remote-device
 - `/dev/mars/mydata` can appear anywhere
- WIP-compatibility:
 - no kernel prepatch needed anymore
 - currently tested with vanilla kernels 3.2 ... 4.4



- md5 checksums on underlying disks
- Mass-scale clustering
- Database support / near-synchronous modes

- Further challenges:
 - community revision at LKML planned
 - replace symlink tree with better representation
 - split into 3 parts:
 - Generic `brick` framework
 - `XIO` / `AIO` personality (1st citizen)
 - MARS Light (1st application)
 - hopefully attractive for other developers!



Appendix



DRBD+proxy (proprietary)

Application area:

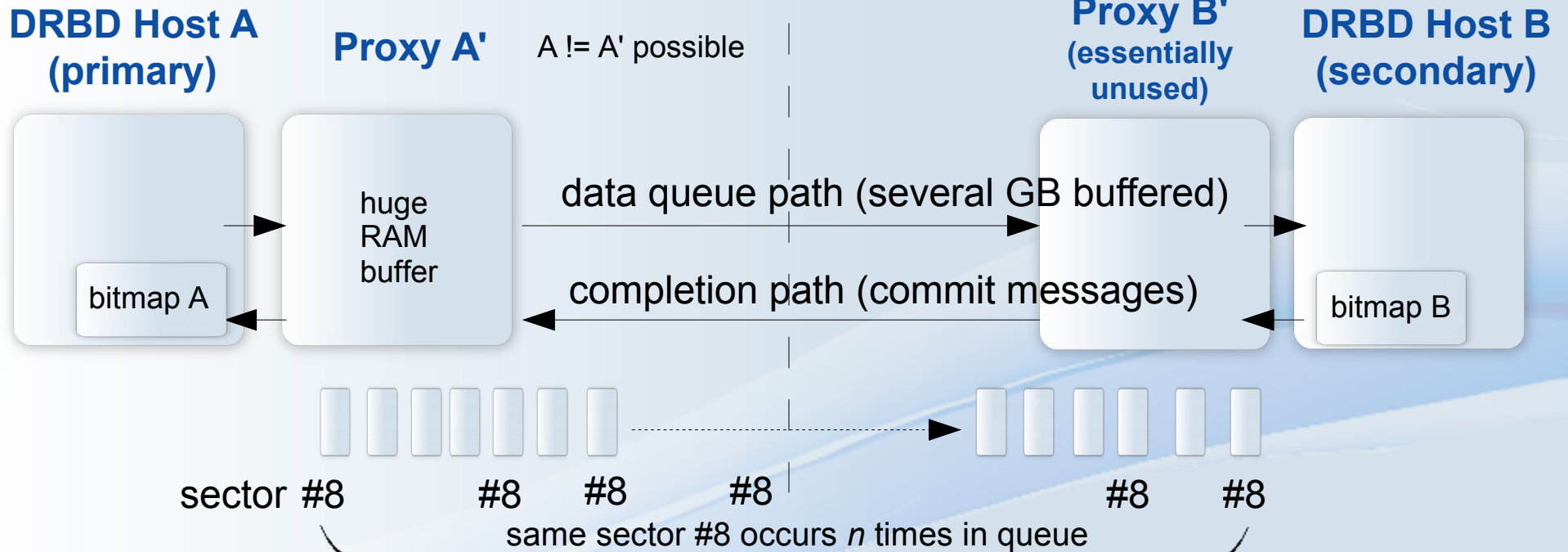
- Distances: any
- Asynchronously
 - **Buffering in RAM**
- Unreliable network leads to **frequent re-syncs**
 - RAM buffer gets lost
 - at cost of actuality
- **Long** inconsistencies during re-sync
- Under pressure: **permanent** inconsistency possible
- High memory overhead
- Difficult scaling to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

MARS Light (GPL)

Application area:

- Distances: **any** ($\gg 50$ 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 ≥ 100 GB in `/mars/` for transaction logfiles
 - dedicated spindle(s) recommended
 - RAID with BBU recommended
- Easy scaling to $k > 2$ nodes

DRBD+proxy Architectural Challenge



n times

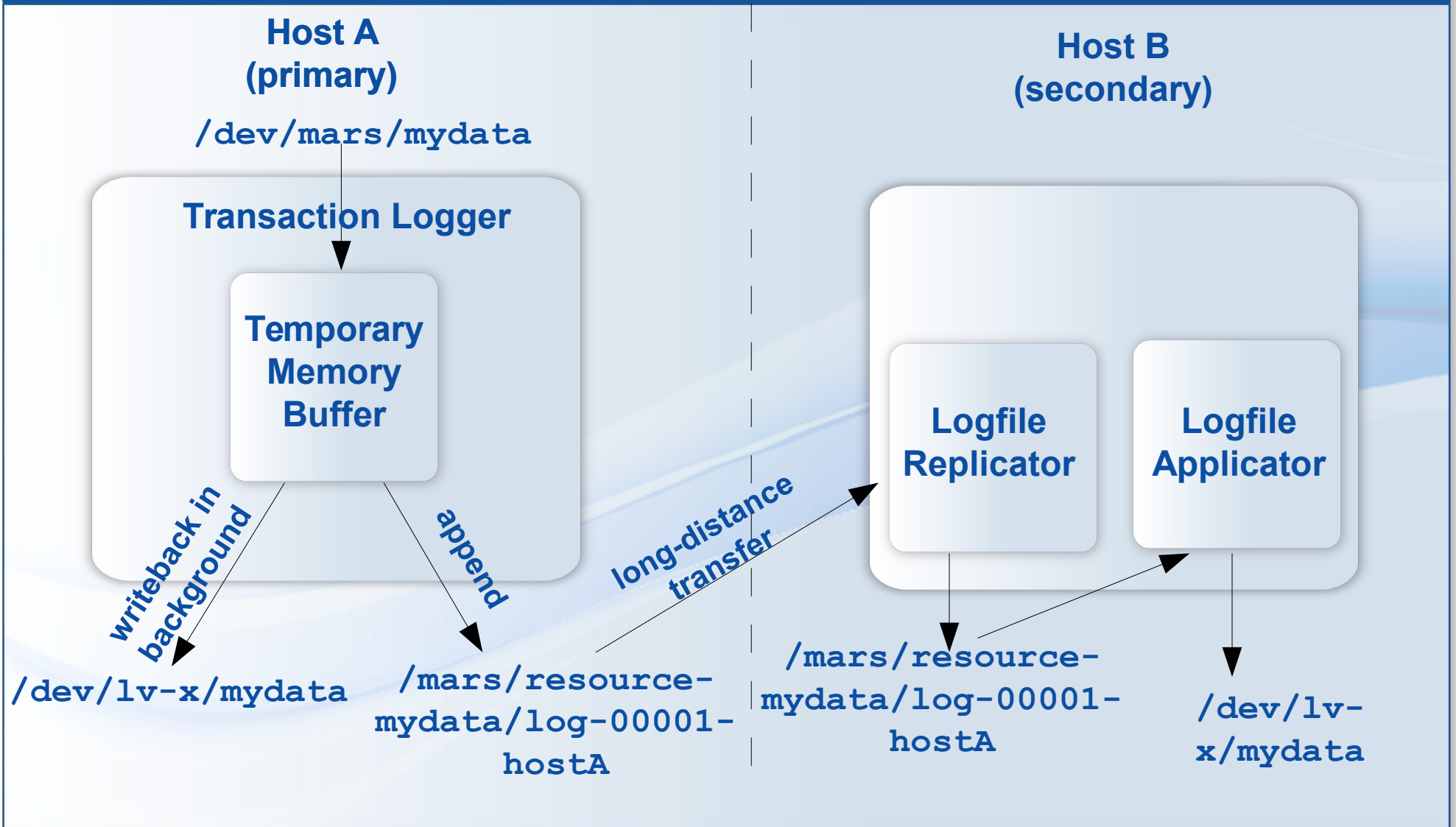
=> need $\log(n)$ bits for counter

=> but DRBD bitmap has only 1 bit/sector

=> workarounds exist, but complicated

(e.g. additional dynamic memory)

MARS Light Data Flow Principle



Framework Architecture

for MARS + future projects



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 \approx AIO

**XIO
bricks**

**future
Strategy
bricks**

**other future
Personalities
and their bricks**

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

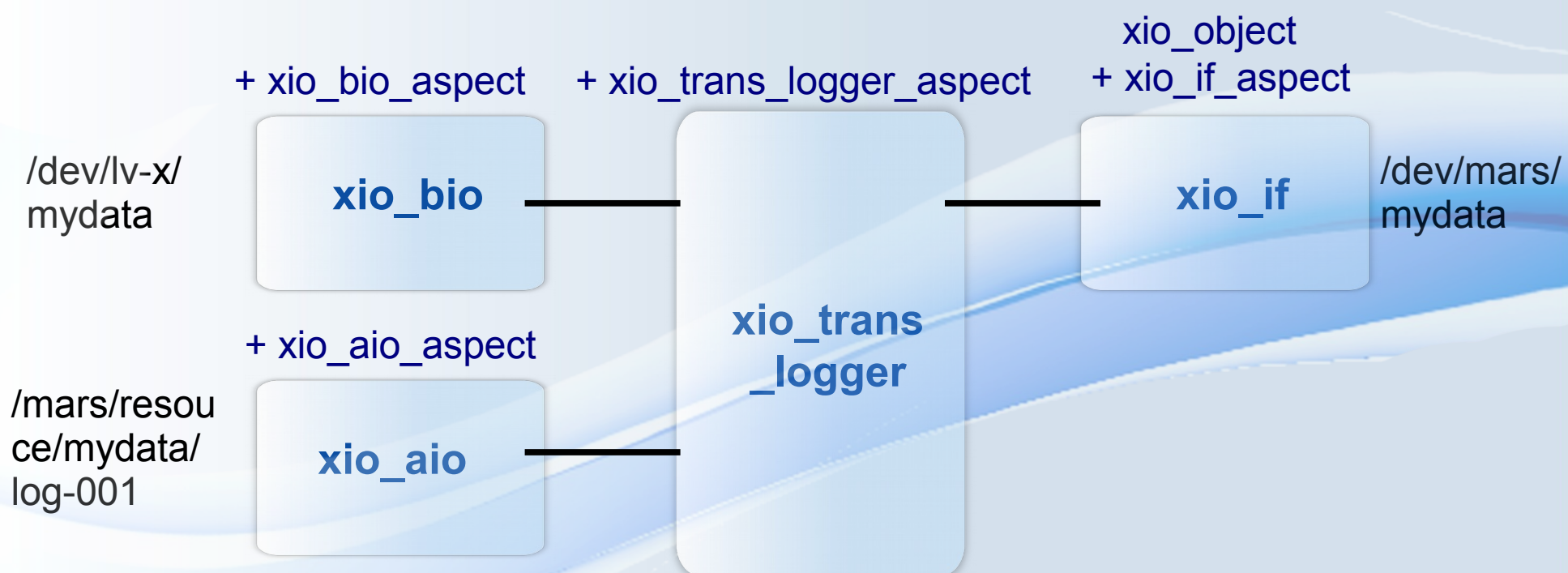
Generic Bricks

Generic Objects

Generic Aspects

S

Bricks, Objects + Aspects (Example)



Aspects are automatically attached on the fly

Appendix: 1&1 Wide Area Network Infrastructure

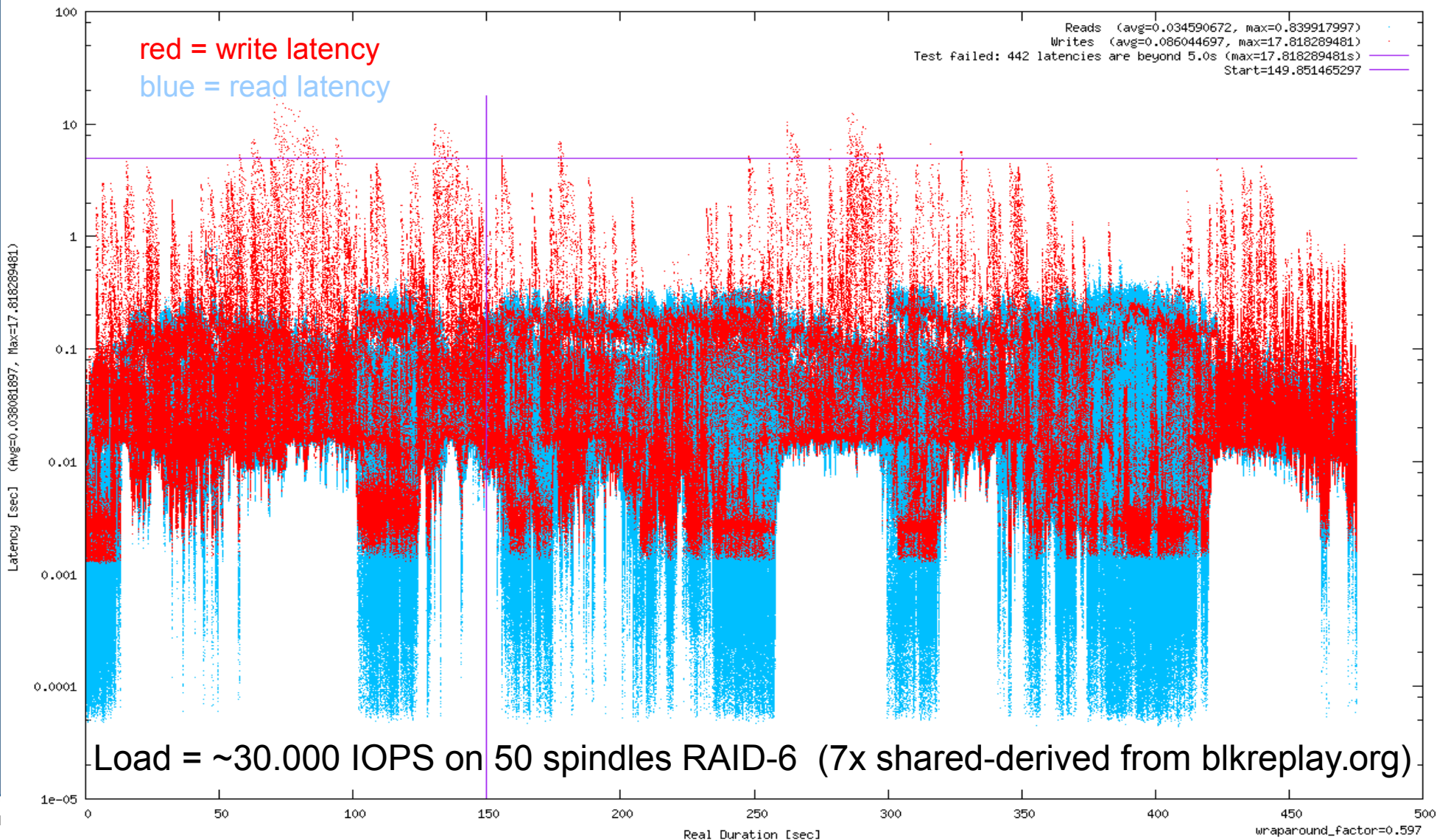
- Global external bandwidth > 285 GBit/s
- Peering with biggest internet exchanges on the world
- Own metro networks (DWDM) at the 1&1 datacenter locations



IO Latencies over loaded Metro Network (1) DRBD



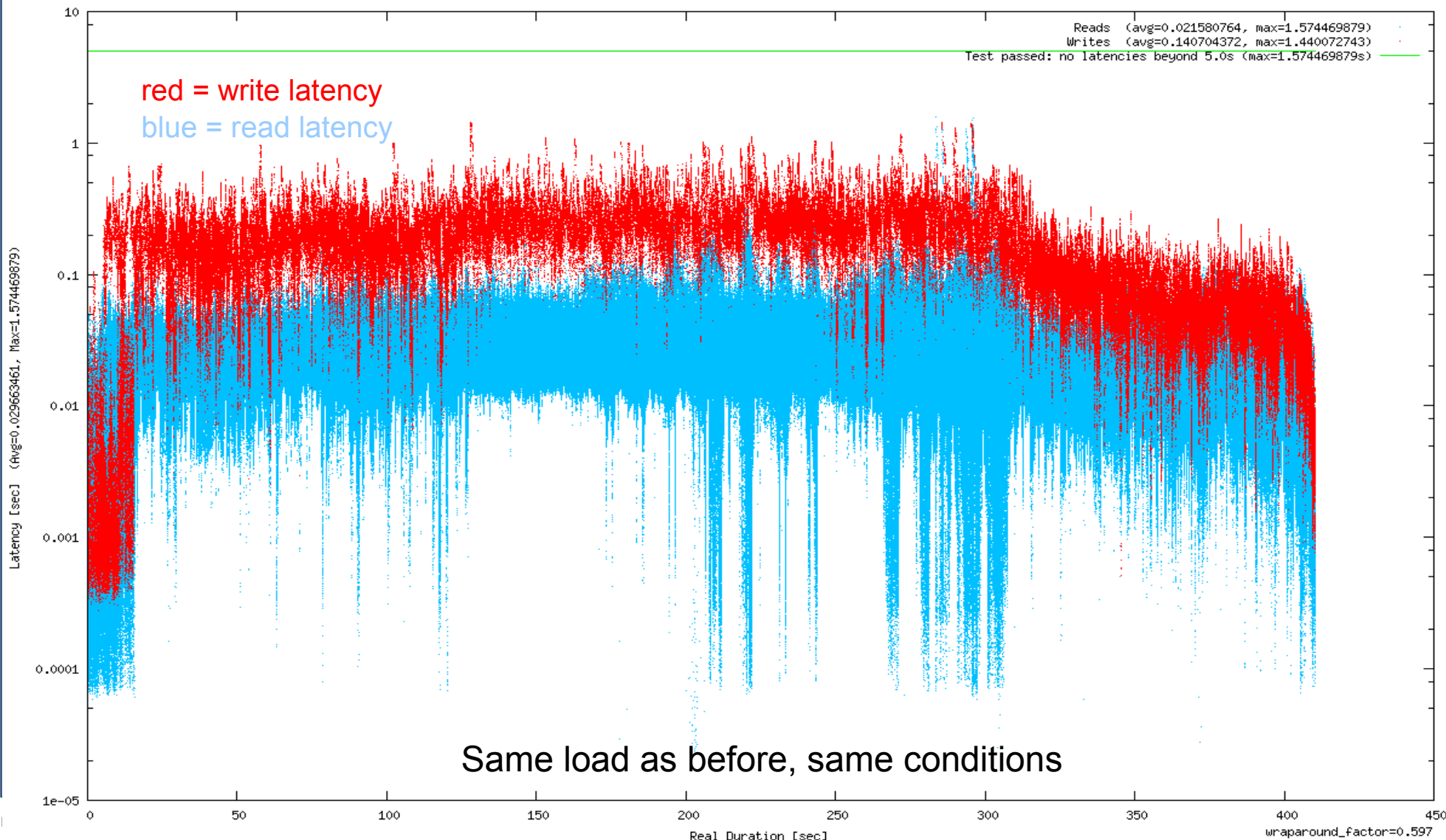
MARS-DRBD-COMPARISON.shared-derived.drbd-8.3.13.g01.latency.realtime Wed Sep 4 16:19:16 2013



IO Latencies over loaded Metro Network (2) MARS



MARS-DRBD-COMPARISON.shared-derived.mars-lvm.mars.g01.latency.realtime Wed Sep 4 17:12:41 2013



Same load as before, same conditions

Performance of Socket Bundling Europe↔USA

