

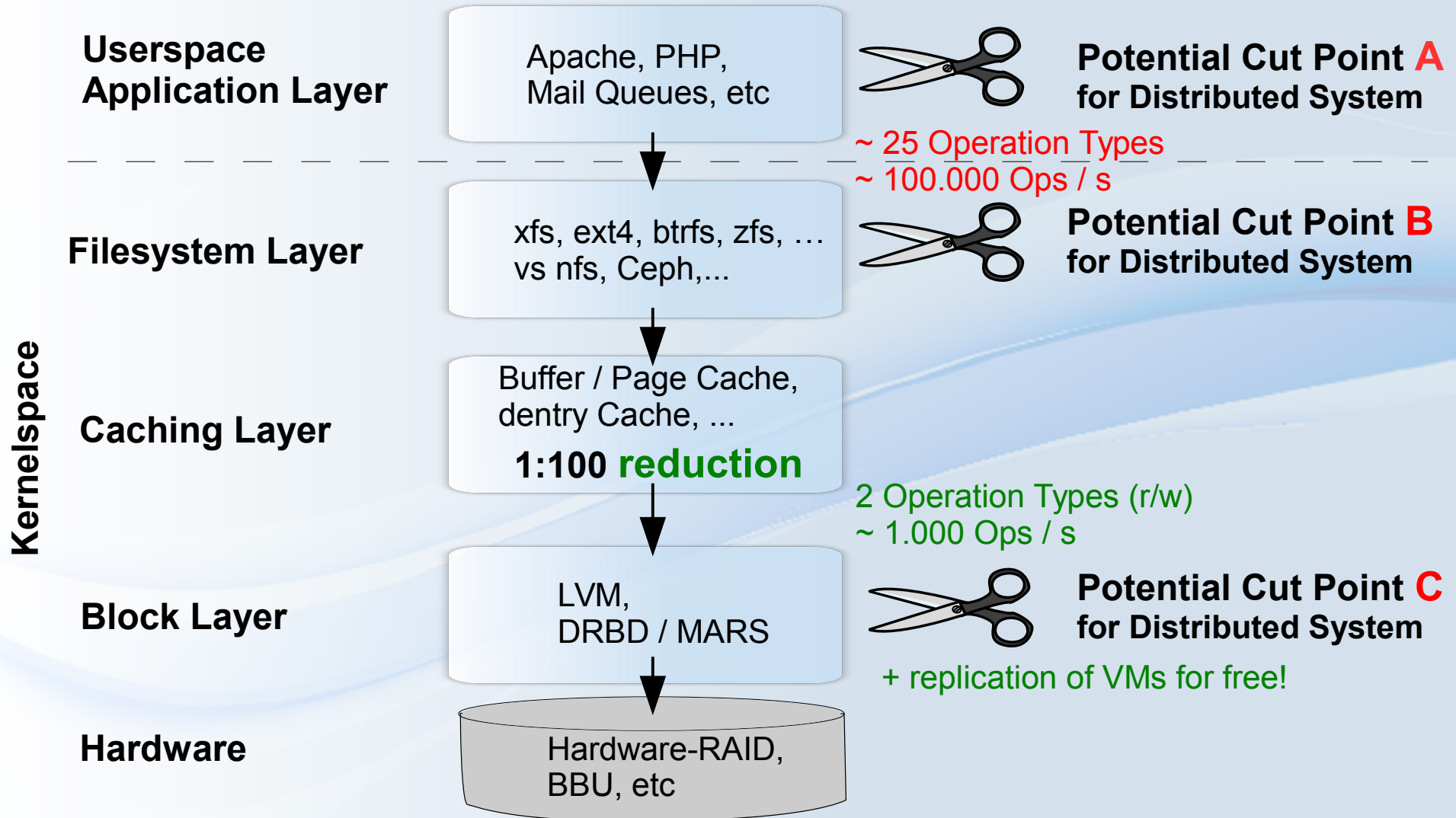
MARS: Replicating Petabytes over Long Distances



FROSCON 2015 Presentation by Thomas Schöbel-Theuer

- **Long Distances: Block Level vs FS Level**
- **Use Cases DRBD/proxy 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



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

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

MARS Working Principle

Multiversion Asynchronous Replicated Storage

Datacenter A
(primary)



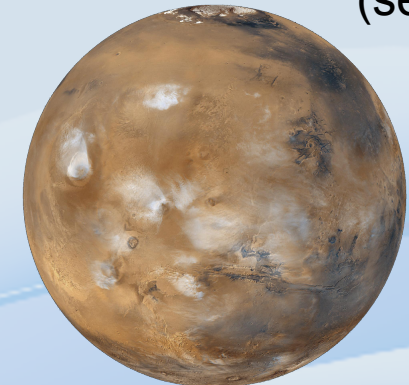
`/dev/mars/mydata`

`mars.ko`

`/dev/lv-x/mydata`

`/mars/trans-
logfile`

Similar to MySQL replication



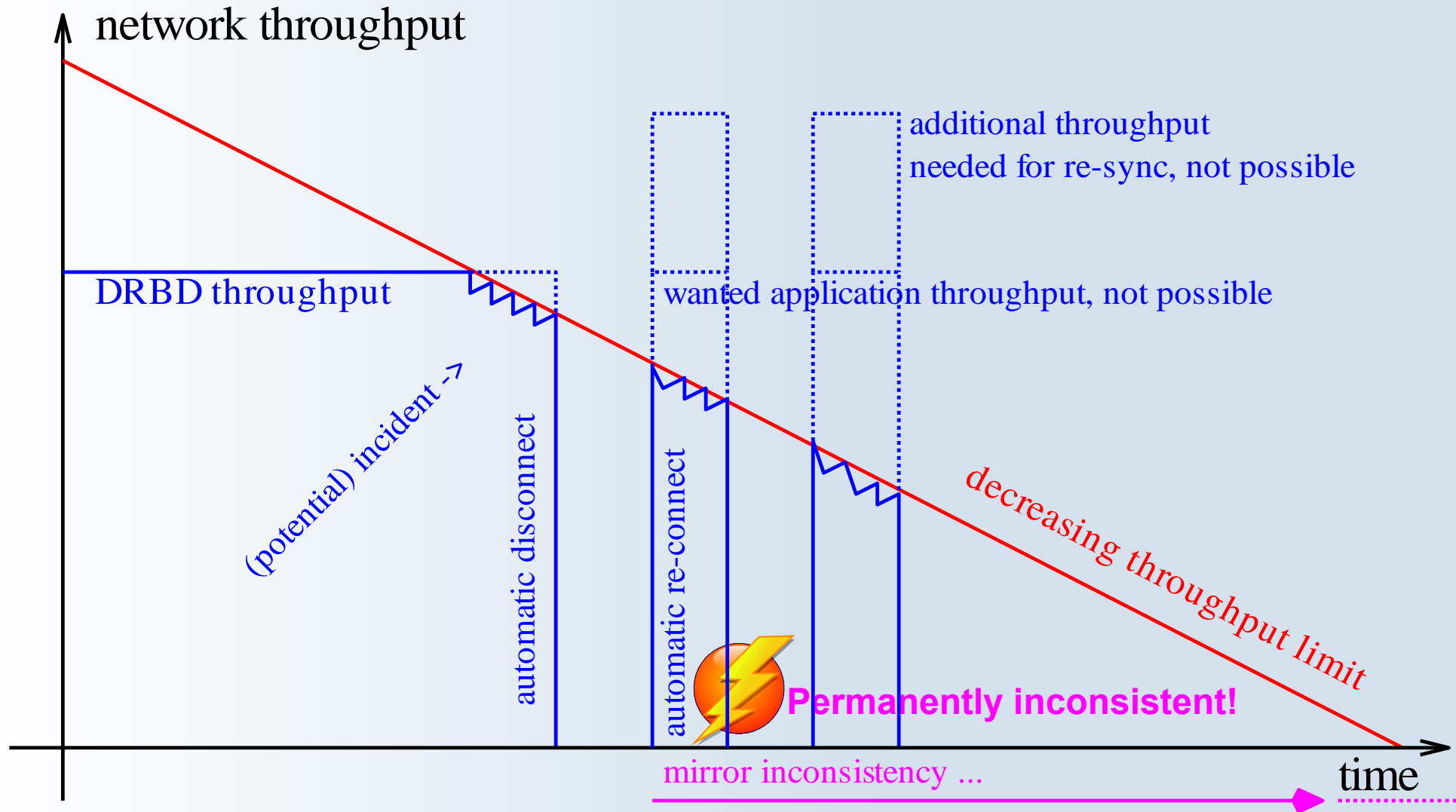
Datacenter B
(secondary)

`mars.ko`

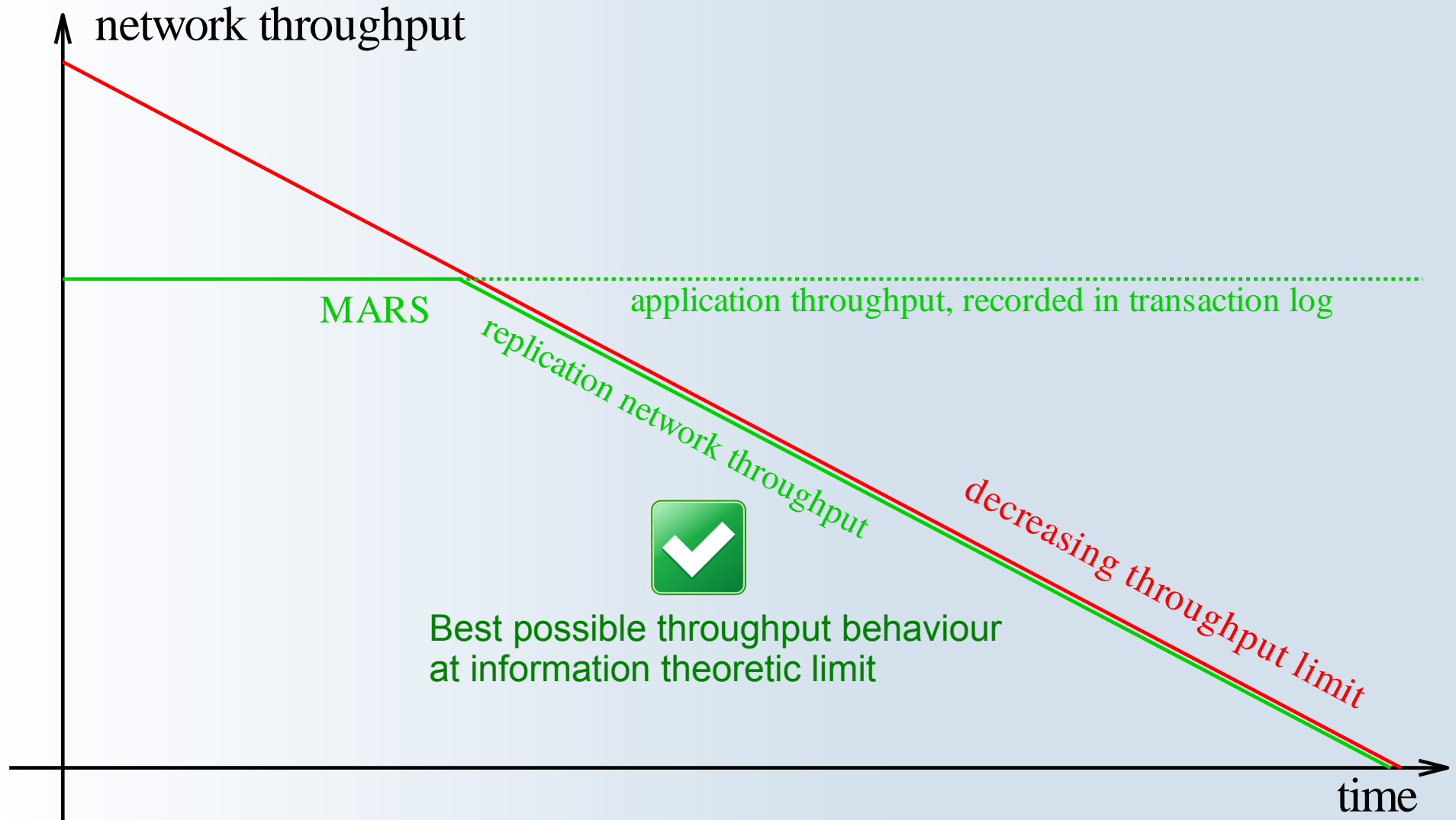
`/mars/trans-
logfile`

`/dev/lv-x/mydata`

Network Bottlenecks (1) DRBD



Network Bottlenecks (2) MARS



MARS

application throughput, recorded in transaction log

replication network throughput

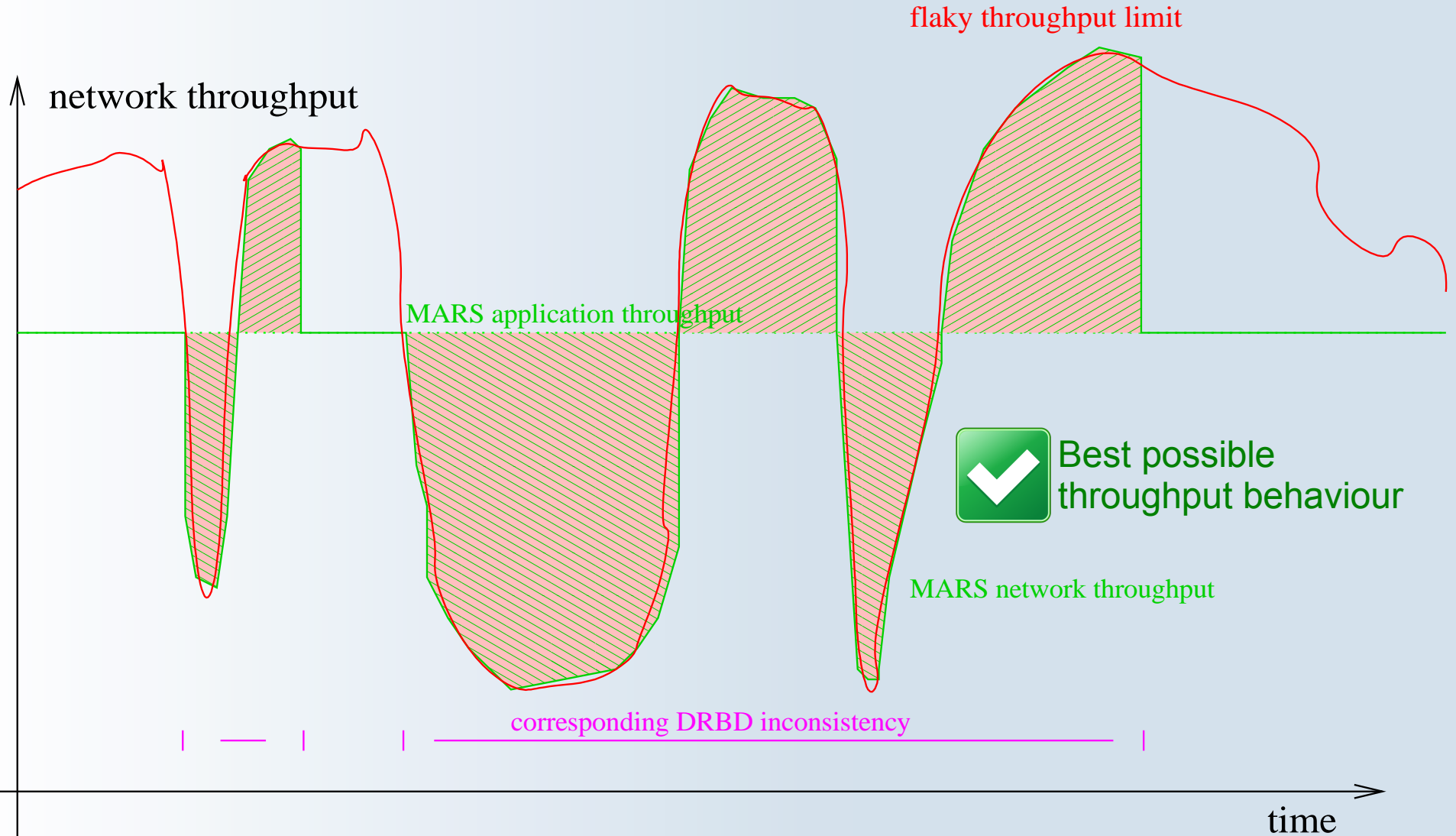
decreasing throughput limit



Best possible throughput behaviour at information theoretic limit

time

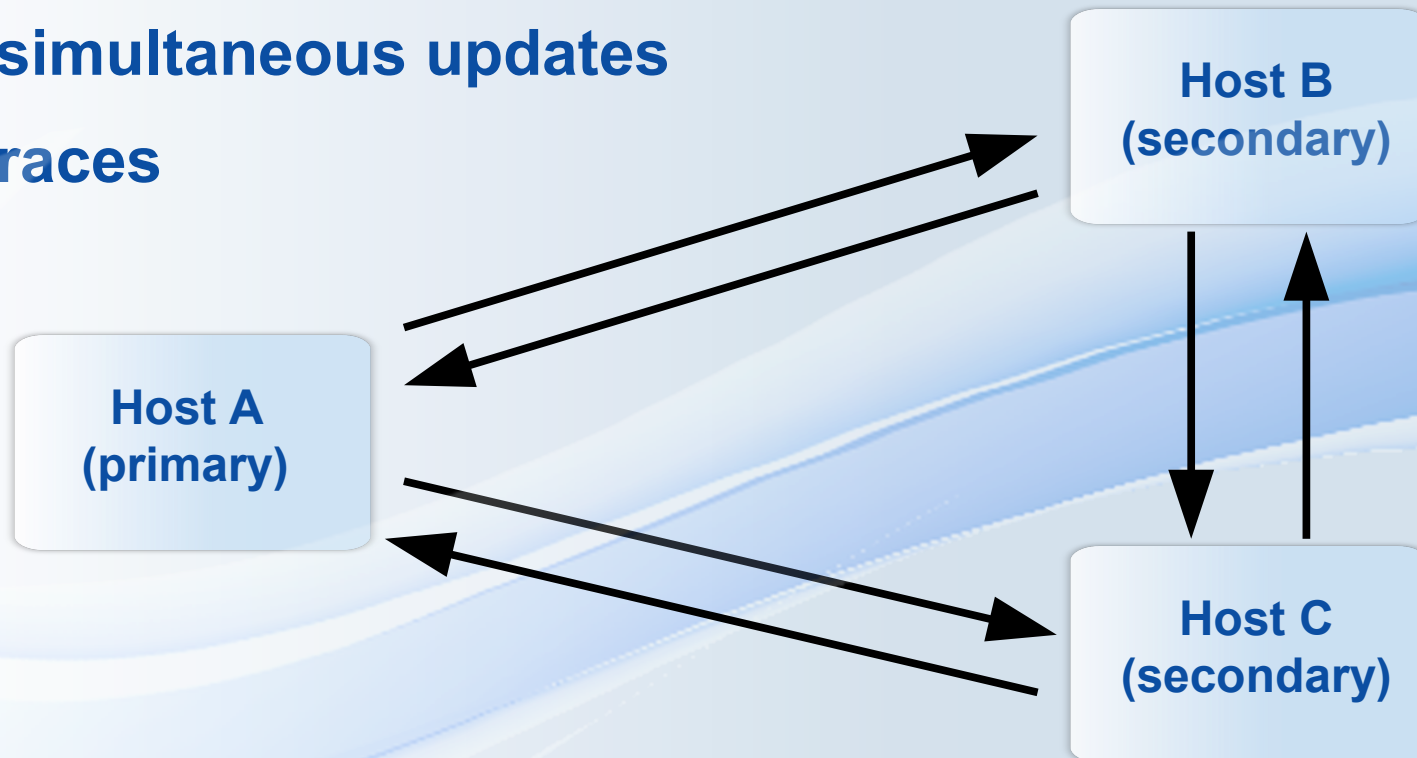
Network Bottlenecks (3) MARS



Metadata Propagation (1)

Problems for ≥ 3 nodes:

- simultaneous updates
- races



Solution: symlink tree + Lamport Clock => next slides

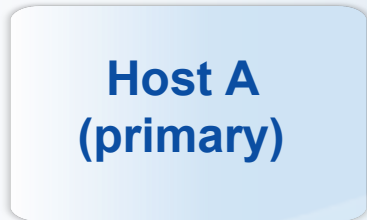
Metadata Propagation (2)



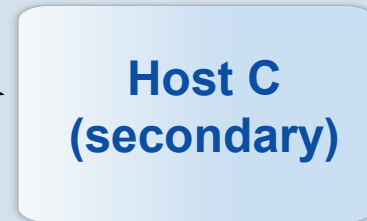
symlink tree = key->value store

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`

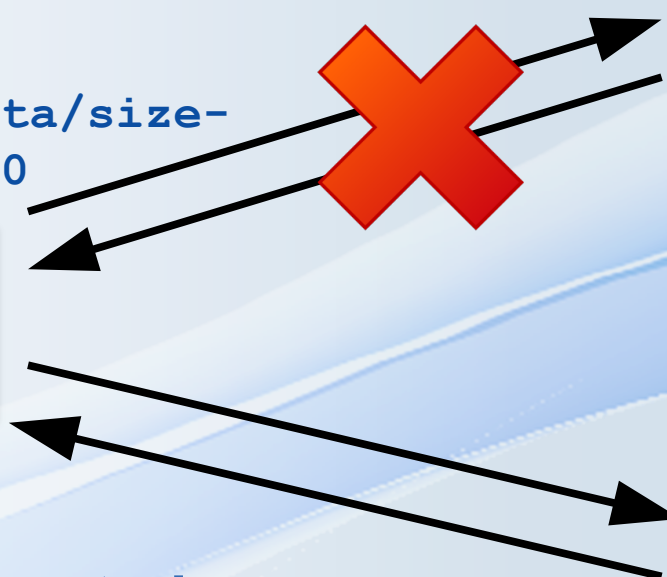
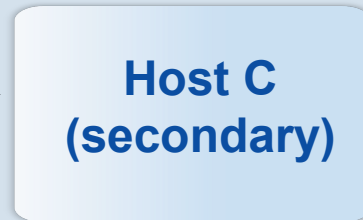
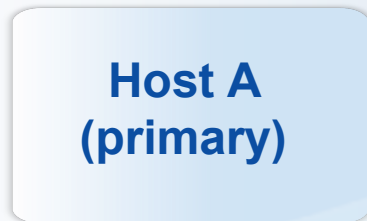
Metadata Propagation (3)

Lamport Clock = virtual timestamp

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

Propagation goes never backwards!

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

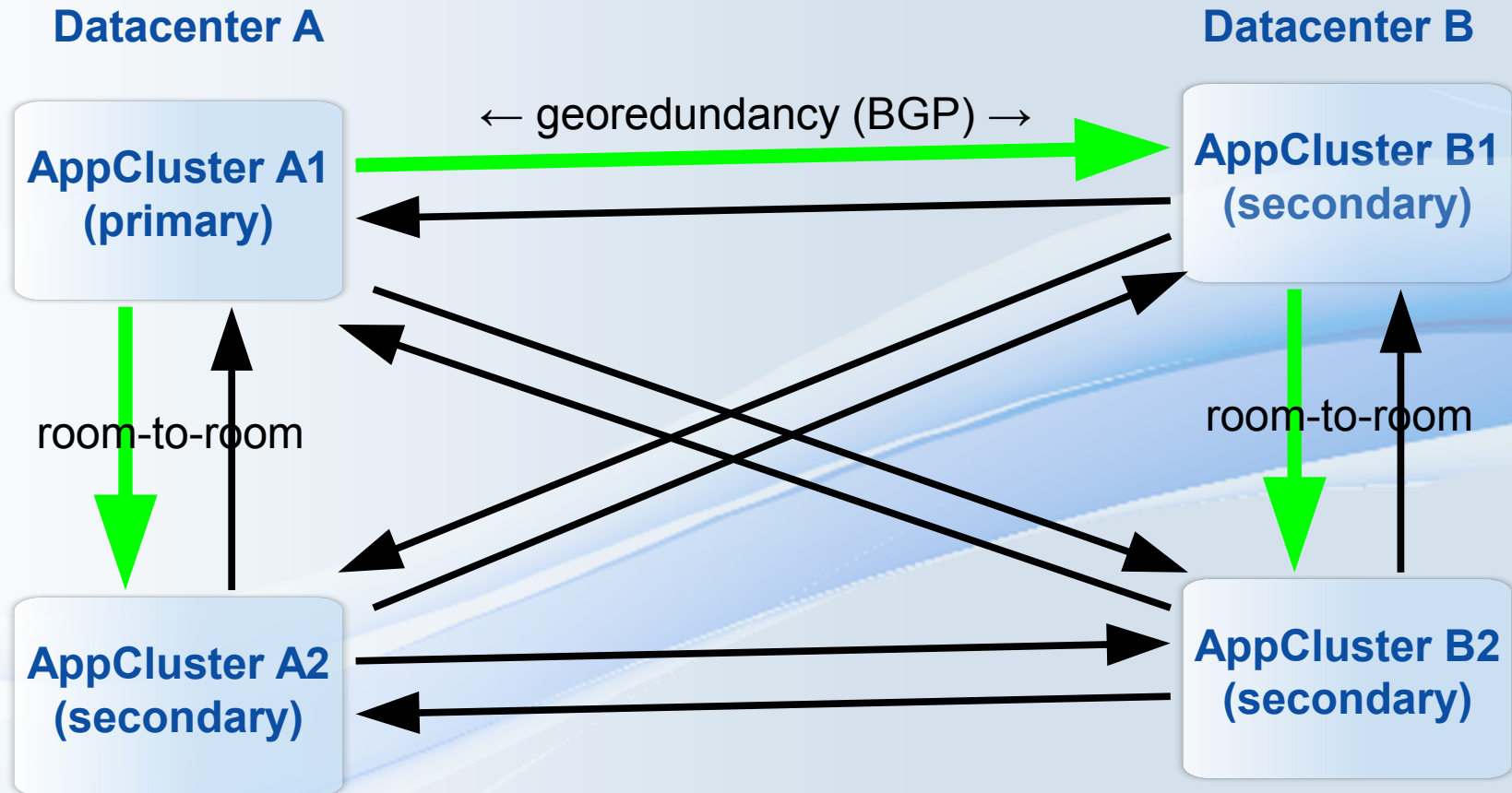


Races are compensated

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 August 2015:

- > 700 central storage servers
- > 2x6 Petabyte total
- ~ 10 billions of inodes in > 3000 xfs instances
- > 3 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-compatibility:

- no kernel prepatch needed anymore
- currently tested with vanilla kernels 3.2 ... 4.2



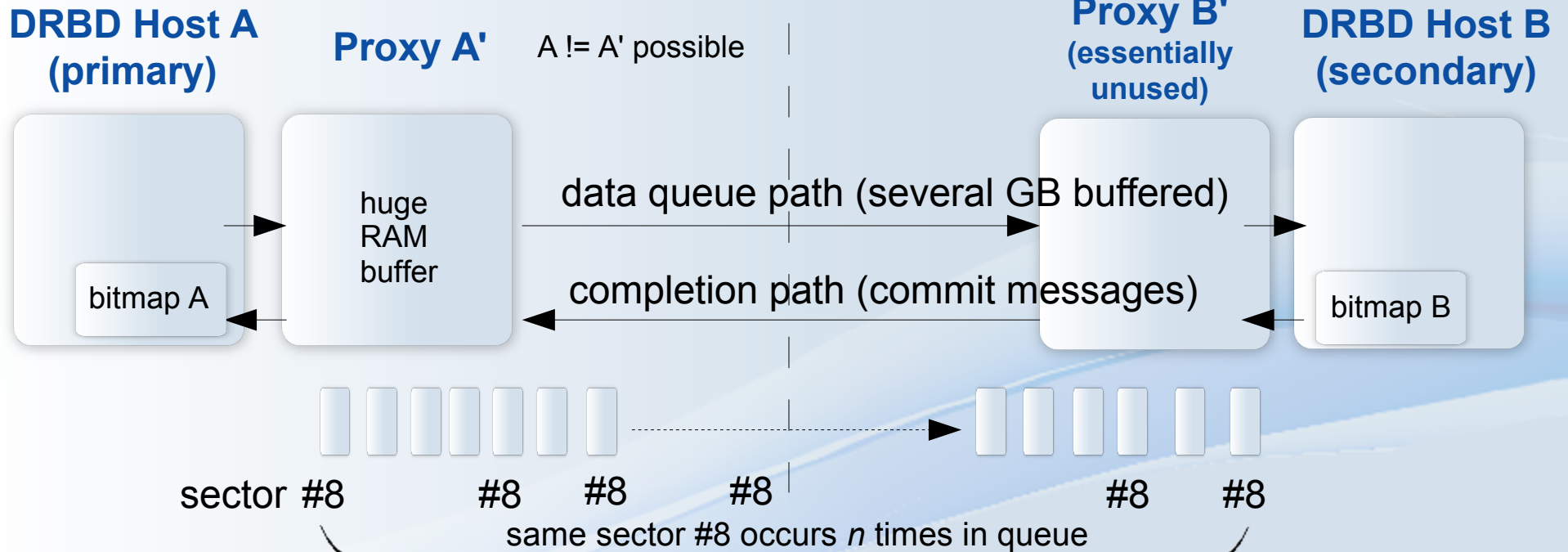
- Remote Device, substitute iSCSI
 - Mass-scale clustering
 - Database support / near-synchronous modes
-
- Further challenges:
 - community revision at LKML planned
 - replace symlink tree with better representation
 - split into 3 parts:
 - Generic `brick` framework
 - `XIO` / `AIO` personality (1st citizen)
 - MARS Light (1st application)
 - hopefully attractive for other developers!



Appendix



DRBD+proxy Architectural Challenge



n times

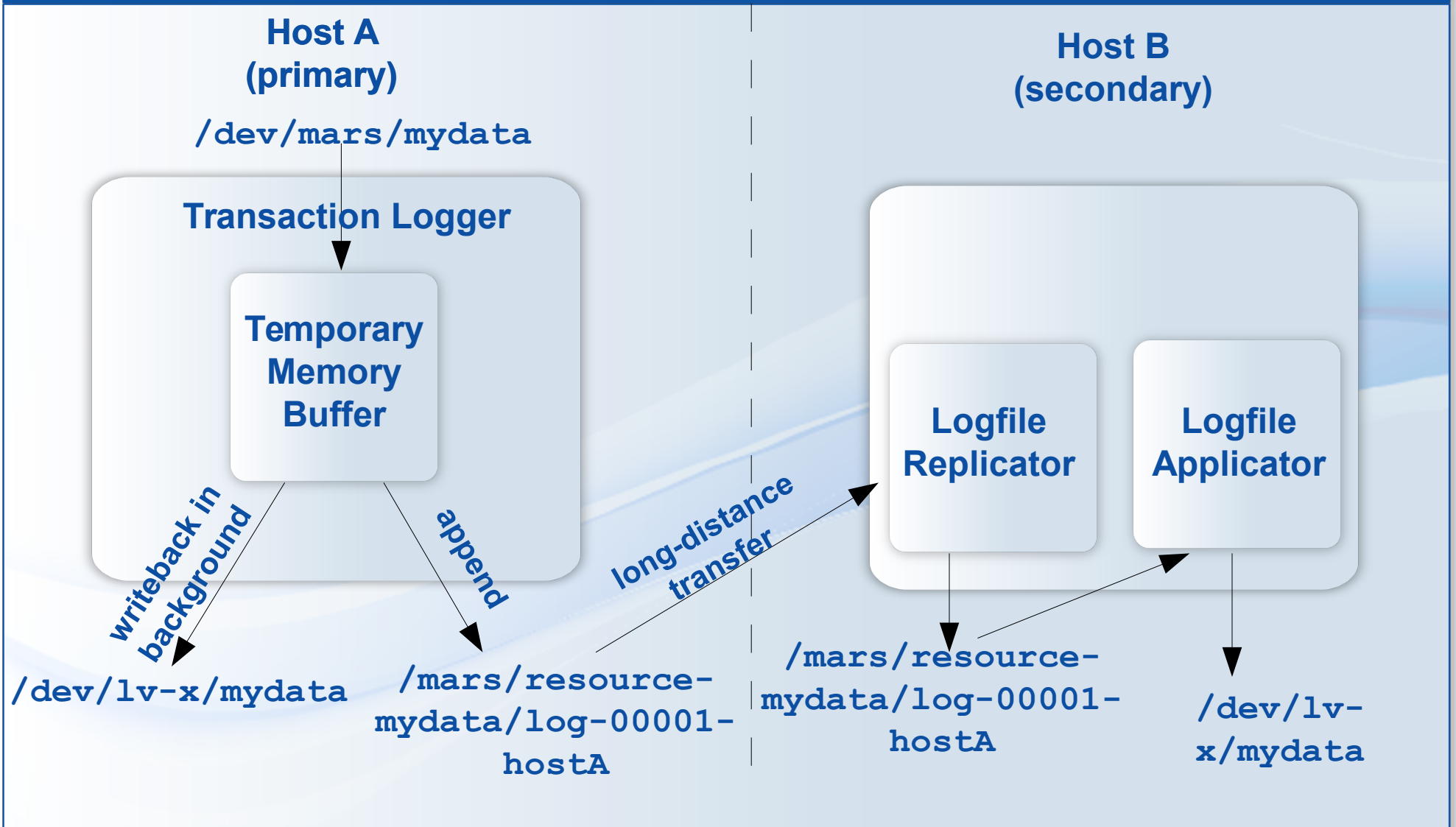
=> need $\log(n)$ bits for counter

=> 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

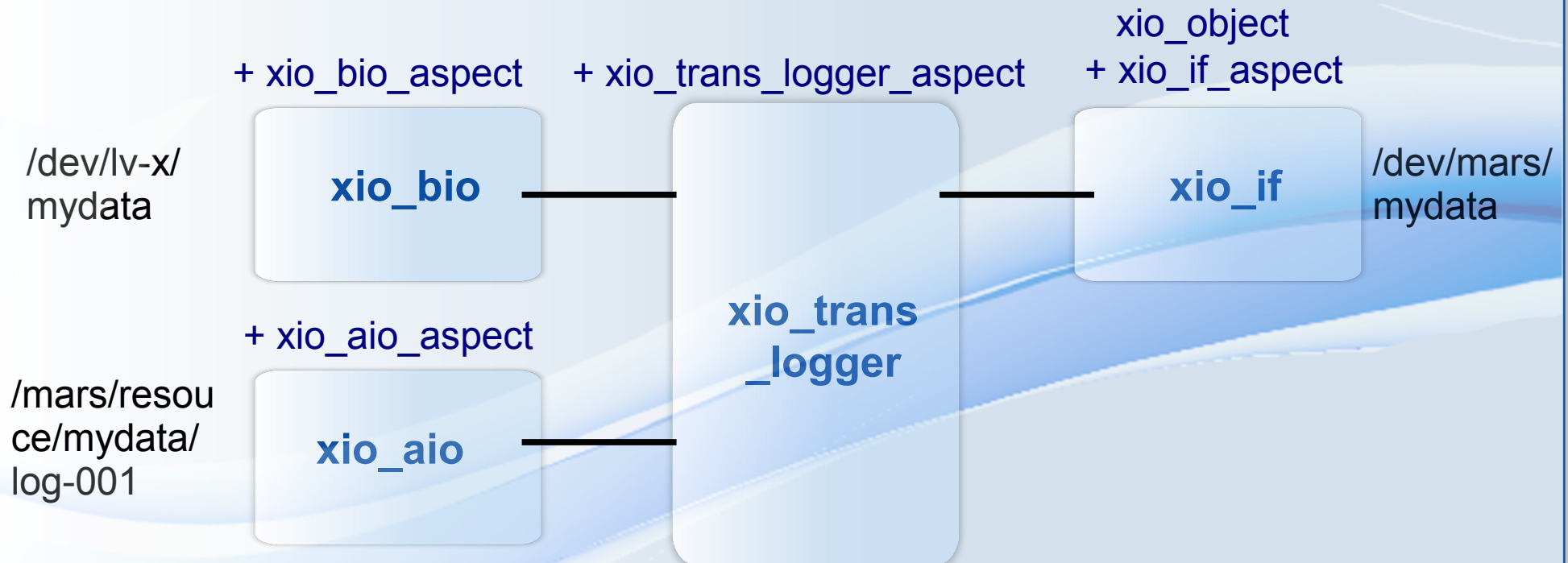
Generic Bricks

Generic Objects

Generic Aspects

S

Bricks, Objects + Aspects (Example)



Aspects are automatically attached on the fly

Appendix: 1&1 Wide Area Network Infrastructure

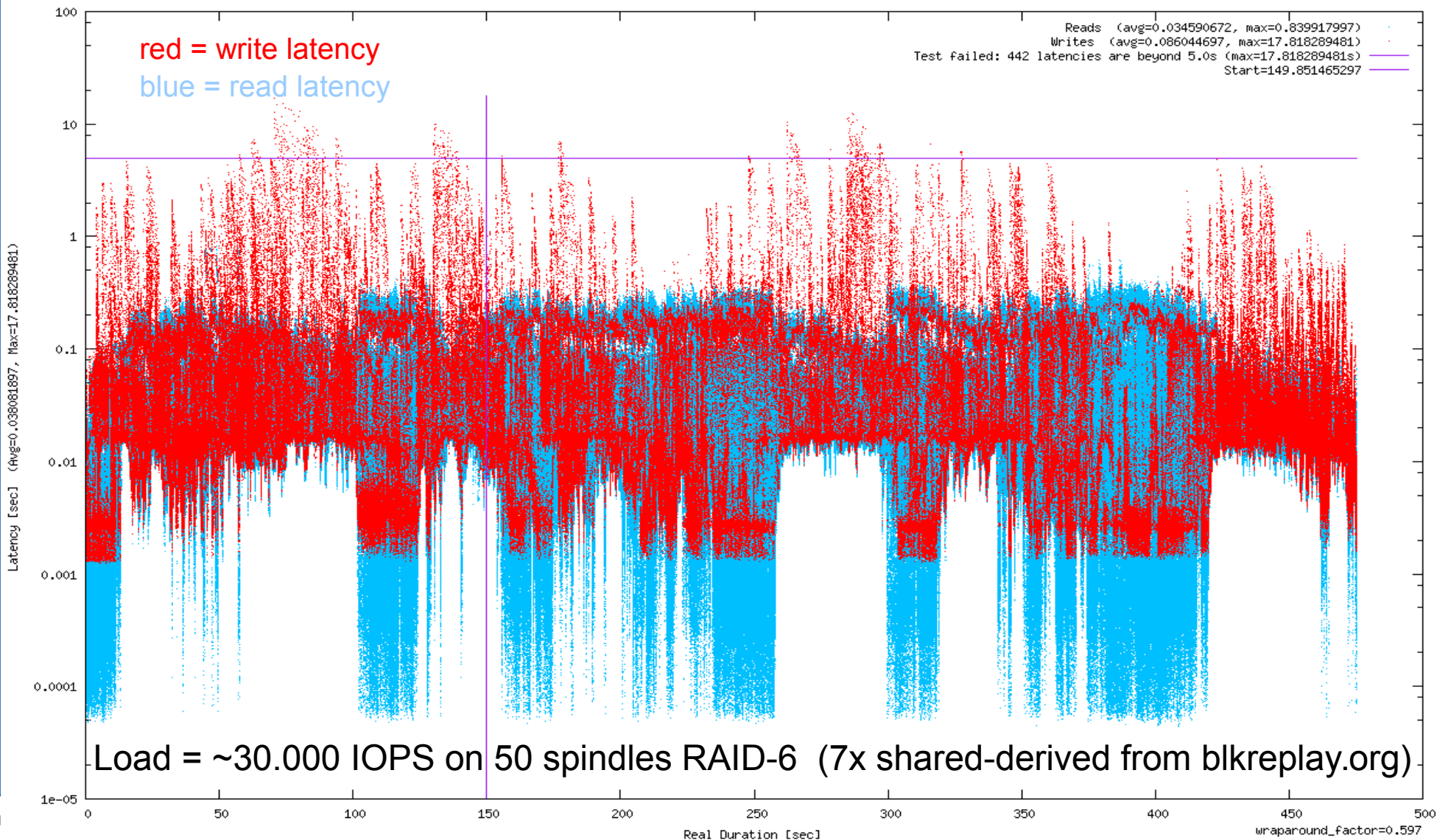
- Global external bandwidth > 285 GBit/s
- Peering with biggest internet exchanges on the world
- Own metro networks (DWDM) at the 1&1 datacenter locations



IO Latencies over loaded Metro Network (1) DRBD



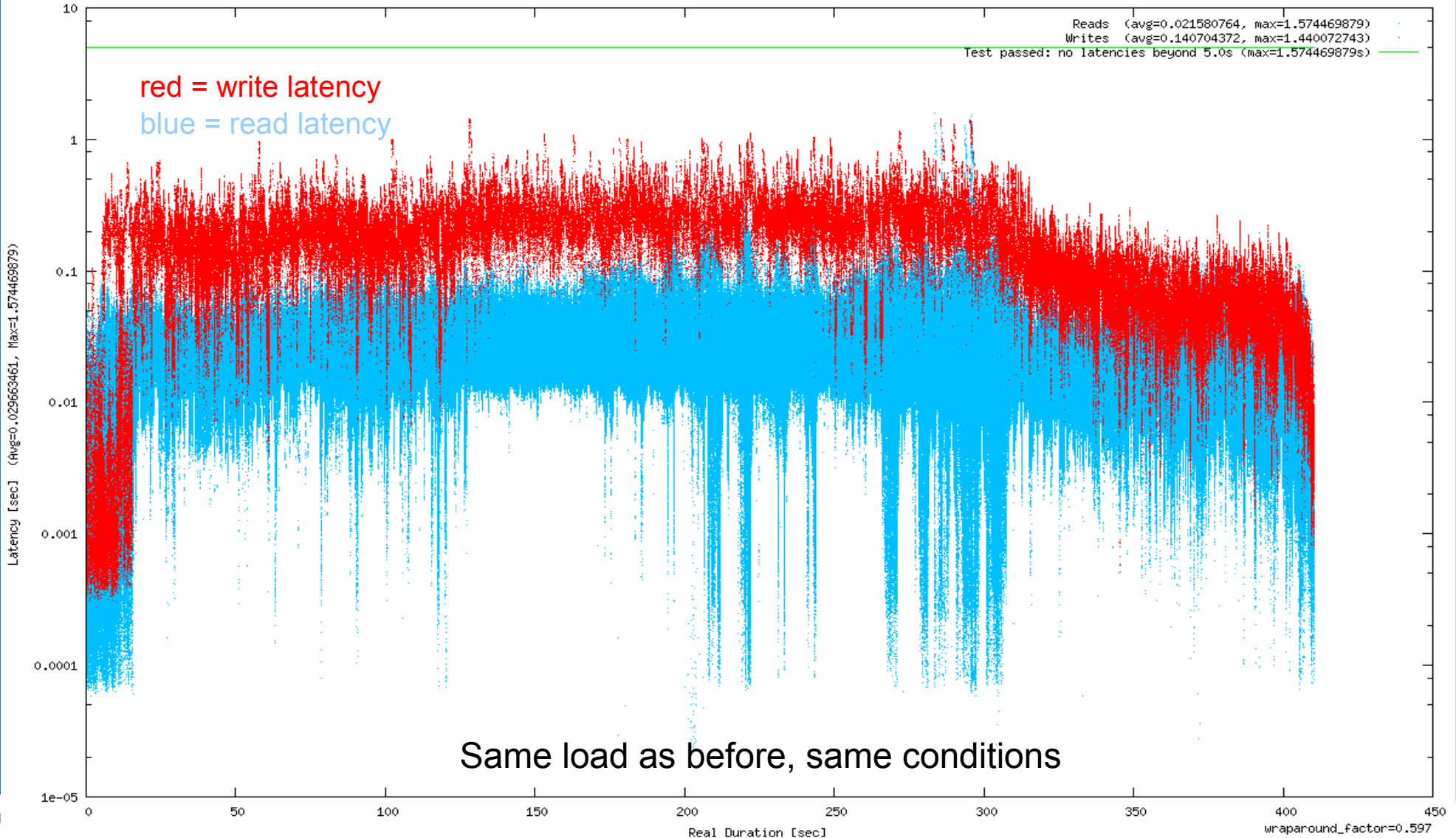
MARS-DRBD-COMPARISON.shared-derived.drbd-8.3.13.g01.latency.realtime Wed Sep 4 16:19:16 2013



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

