# MARS Light: Replicating  Block Devices
## over Long Distances

**LinuxTag 2014 Presentation by Thomas Schöbel-Theuer**

# Agenda

■ **Use Cases DRBD/proxy vs MARS Light**

■ **Working Principle**

■ **Behaviour at Network Bottlenecks**

■ **Multinode Metadata Propagation (Lamport Clock)**

■ **Example Scenario with 4 Nodes**

■ **Current Status / Future Plans**

# Use Cases DRBD vs MARS Light

**1&1**

## 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 >= 100GB in `/mars/` for transaction logfiles
  - dedicated spindle(s) recommended
  - RAID with BBU recommended

# 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

# MARS Working Principle

Multiversion Asynchronous Replicated Storage

Datacenter A
(primary)

Datacenter B
(secondary)

`/dev/mars/mydata`

mars.ko

Similar to MySQL replication

mars.ko

`/dev/lv-x/mydata`

`/mars/trans-logfile`

`/mars/trans-logfile`

`/dev/lv-x/mydata`

# Network Bottlenecks (1) DRBD



network throughput

additional throughput
needed for re-sync, not possible

DRBD throughput

wanted application throughput, not possible

(potential) incident ->

automatic disconnect

automatic re-connect

decreasing throughput limit

**Permanently inconsistent!**

mirror inconsistency ...

time

# Network Bottlenecks (2) MARS



network throughput

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



flaky throughput limit

network throughput

MARS application throughput

Best possible throughput behaviour

MARS network throughput

corresponding DRBD inconsistency

time

**1&1**

**Problems for ≥ 3 nodes:**
- simultaneous updates
- races

Host B
(secondary)

Host A
(primary)

Host C
(secondary)

**Solution: symlink tree + Lamport Clock => next slides**
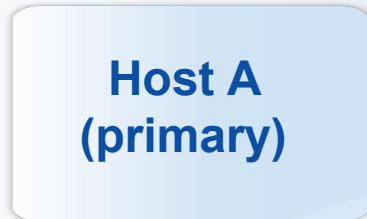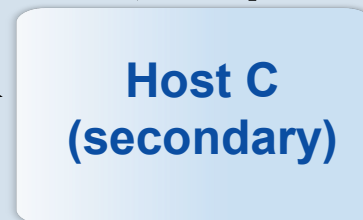
# Metadata Propagation (2)

**1&1**

**Symlink tree = key->value store**

**Originator context encoded in key**

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

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

**Host B
(secondary)**

**Host A
(primary)**

**Host C
(secondary)**

**Anyone knows anything about others**

**But later**

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

# Metadata Propagation (3)

**1&1**

**Lamport Clock = virtual timestamp**

**Propagation goes never backwards!**

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

**Host A (primary)**

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

**Host B (secondary)**

**Host C (secondary)**

**Races are compensated**

**Propagation paths play no role**

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

# Productive Scenario since 03/2014 (1&1 eShop / ePages)   1&1

**Datacenter A**

**Datacenter B**

← georedundancy (BGP) →

**AppCluster A1 (primary)**

**AppCluster B1 (secondary)**

room-to-room

room-to-room

**AppCluster A2 (secondary)**

**AppCluster B2 (secondary)**

→ potential data flow

→ actual data flow (in this scenario)

# Current Status / Future Plans

■ Source / docs at

    github.com/schoebel/mars

■ Productive on customer data since 03/2014

■ Database support / near-synchronous modes planned for end of 2014

■ Further challenges:
  – community revision at LKML planned
  – split into 3 parts:
    • Generic `brick` framework
    • `XIO`/`AIO` personality (1st citizen)
    • MARS Light (1st application)
  – hopefully attractive for other developers!

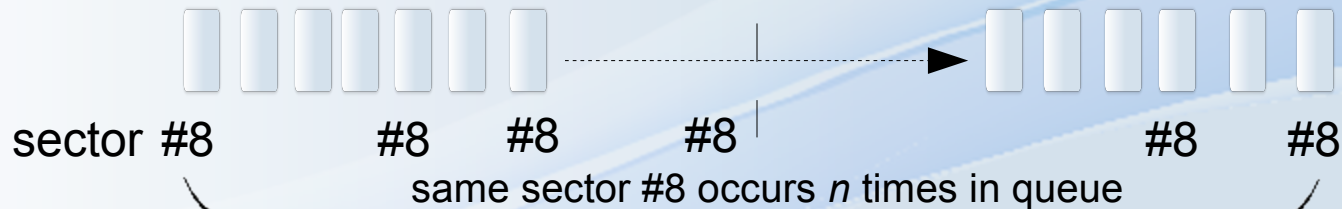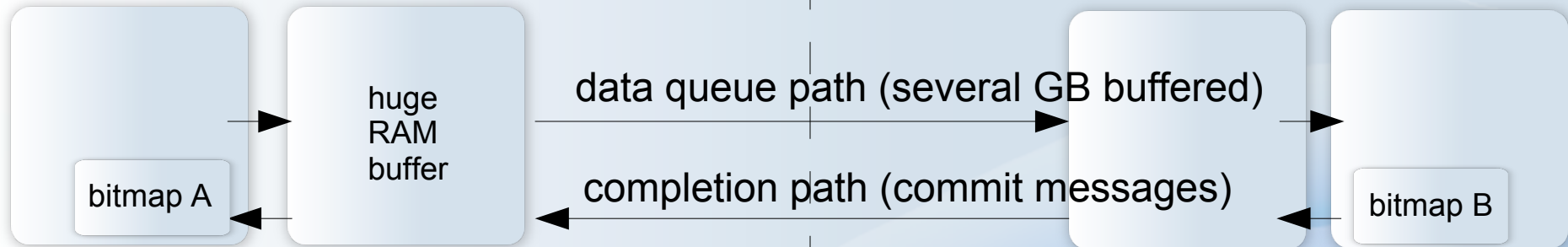# Appendix

# DRBD+proxy Architectural Challenge

**DRBD Host A
(primary)**

**Proxy A'**        A != A' possible

**Proxy B'
(essentially
unused)**

**DRBD Host B
(secondary)**

bitmap A

huge
RAM
buffer

data queue path (several GB buffered)

completion path (commit messages)

bitmap B

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

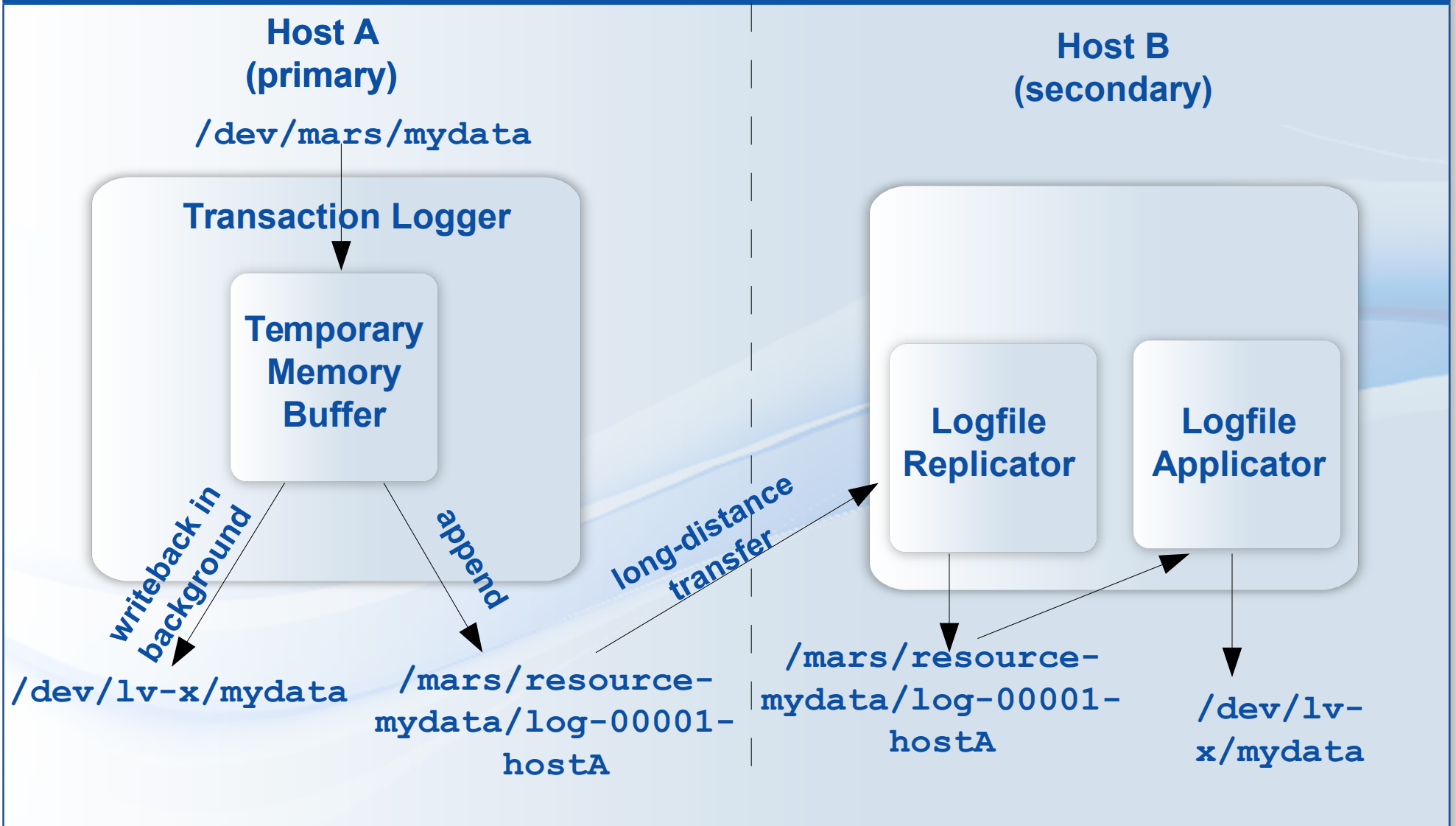same sector #8 occurs *n* times in queue

*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 Light 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
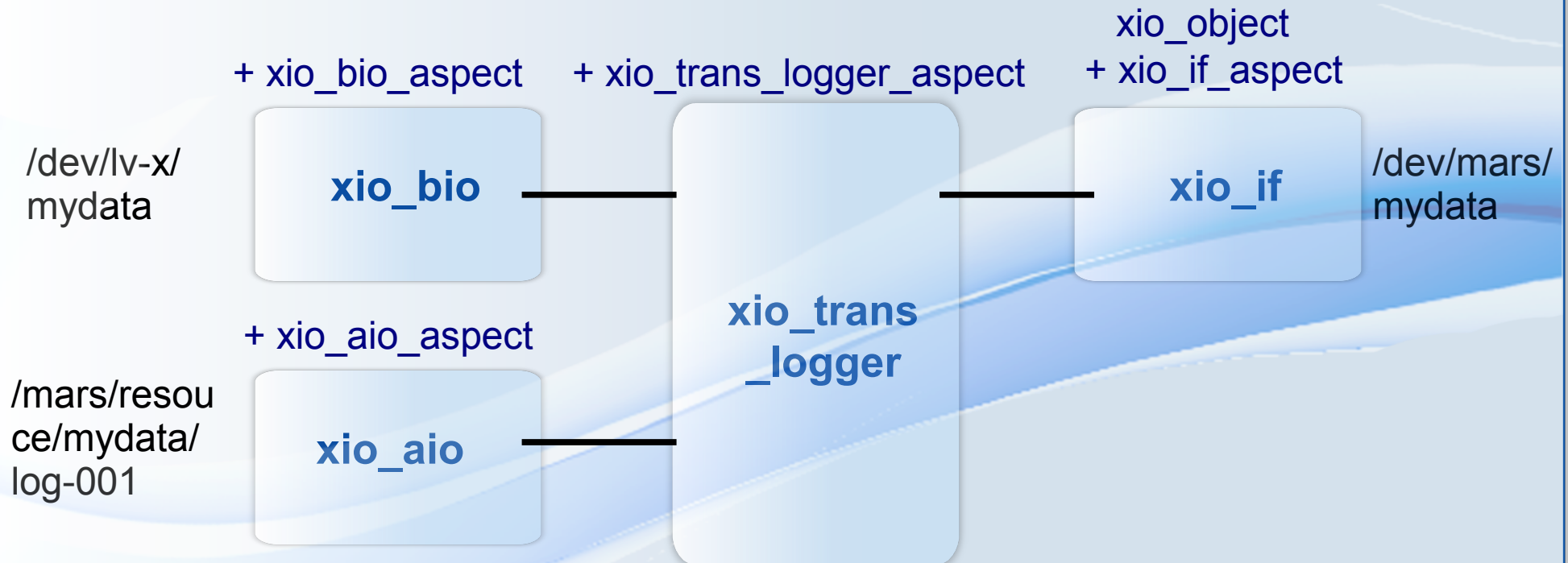+ AOP = Aspect Oriented Programming
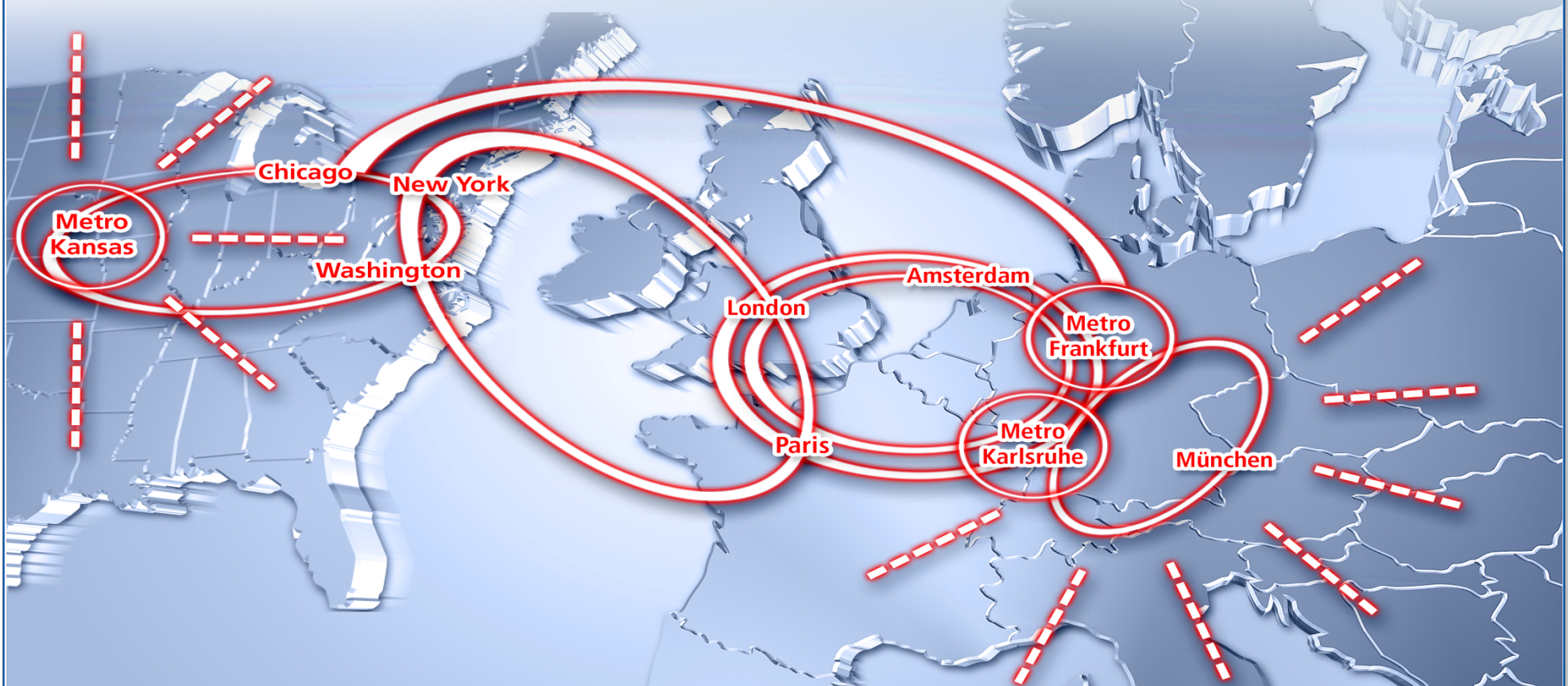
Generic Bricks

Generic Objects

Generic Aspects s

# Bricks, Objects + Aspects (Example)

**1&1**

+ xio_bio_aspect

+ xio_trans_logger_aspect

xio_object
+ xio_if_aspect

/dev/lv-x/
mydata

**xio_bio**

**xio_trans
_logger**

**xio_if**

/dev/mars/
mydata

+ xio_aio_aspect

/mars/resou
ce/mydata/
log-001

**xio_aio**

## Aspects are automatically attached on the fly

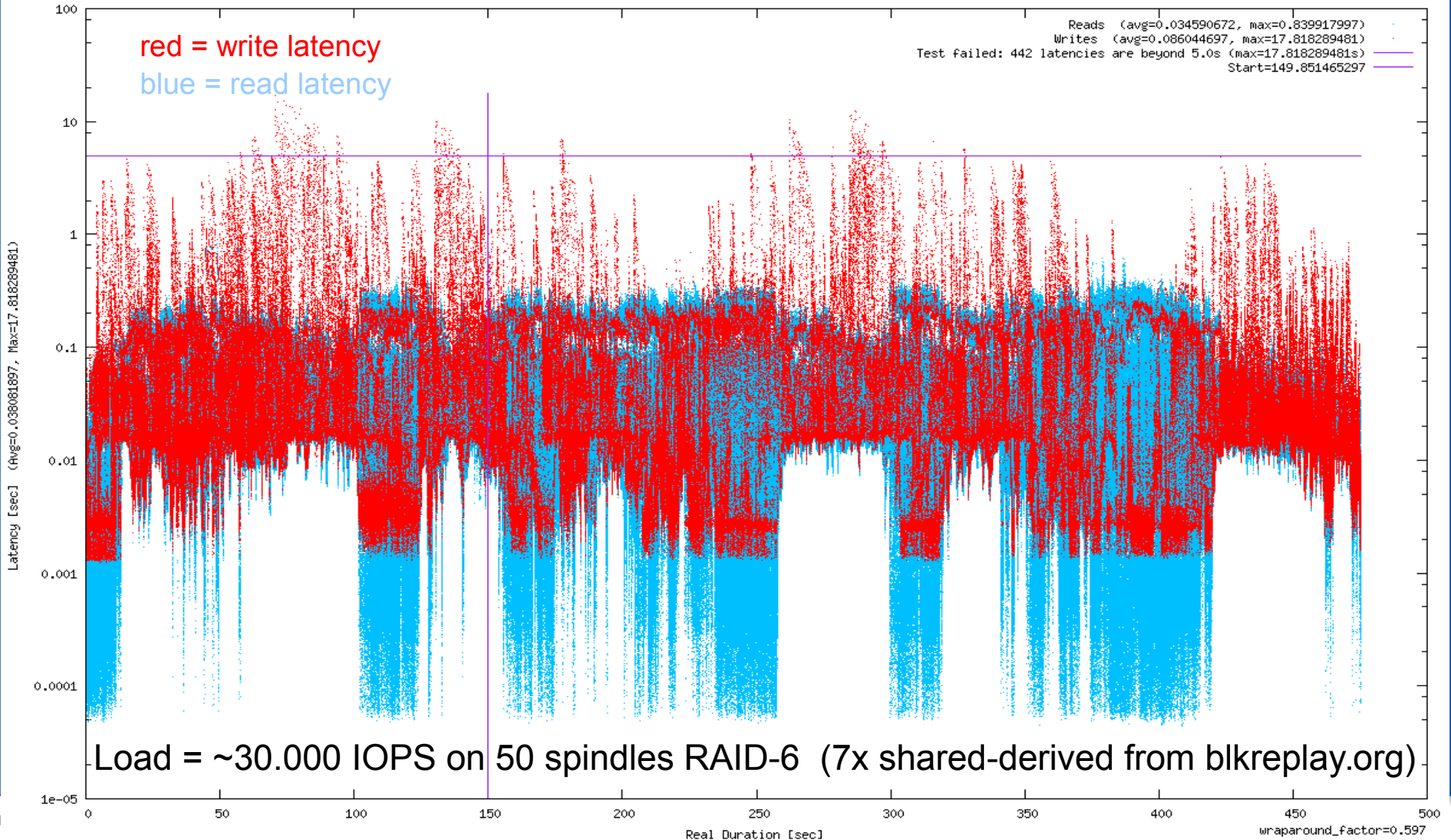# Appendix: 1&1 Wide Area Network Infrastructure

- Global external bandwidth > 285 GBit/s
- Peering with biggest internet exchanges on the world
- Own metro networks (DWDM) at the 1&1 datacenter locations
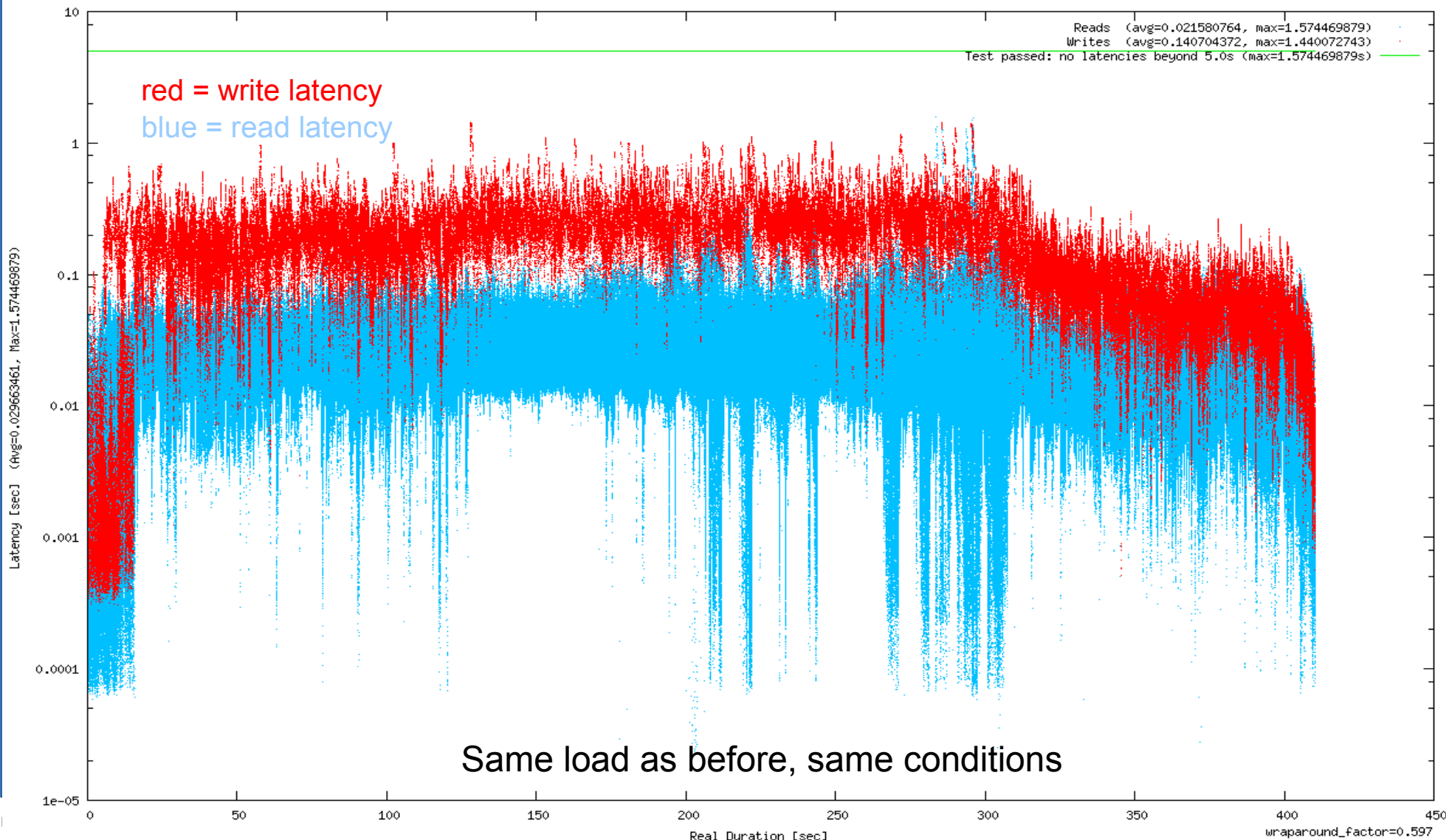
# IO Latencies over loaded Metro Network (1) DRBD



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

red = write latency

blue = read latency

Reads  (avg=0.034590672, max=0.839917997)
Writes  (avg=0.086044697, max=17.818289481)
Test failed: 442 latencies are beyond 5.0s (max=17.818289481s)
Start=149.851465297

Load = ~30.000 IOPS on 50 spindles RAID-6  (7x shared-derived from blkreplay.org)

wraparound_factor=0.597

# 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

Reads  (avg=0.021580764, max=1.574469879)
Writes (avg=0.140704372, max=1.440072743)
Test passed: no latencies beyond 5.0s (max=1.574469879s)

red = write latency
blue = read latency

Same load as before, same conditions