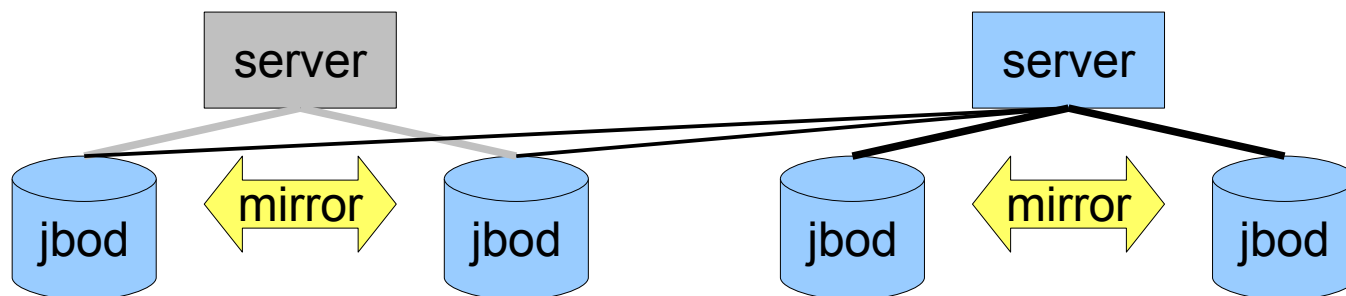


# iSCSI at CERN

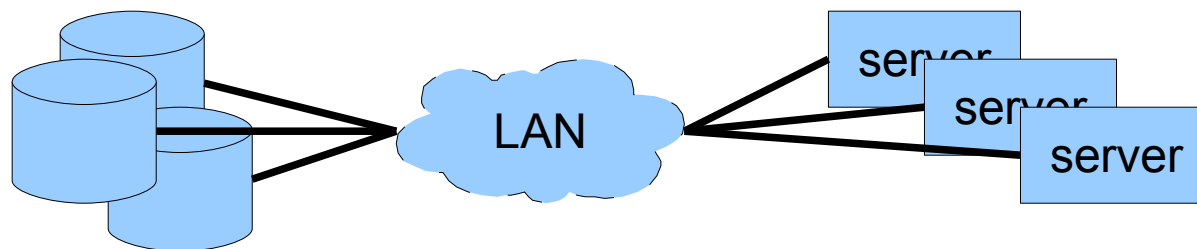
[Andras.Horvath@cern.ch](mailto:Andras.Horvath@cern.ch), 2009

- “Good Enough (tm)” performance...
- Large-scale deployment
  - Reliability and robustness
  - Automatisation
  - Monitoring
- Cost
  - hardware
  - operational

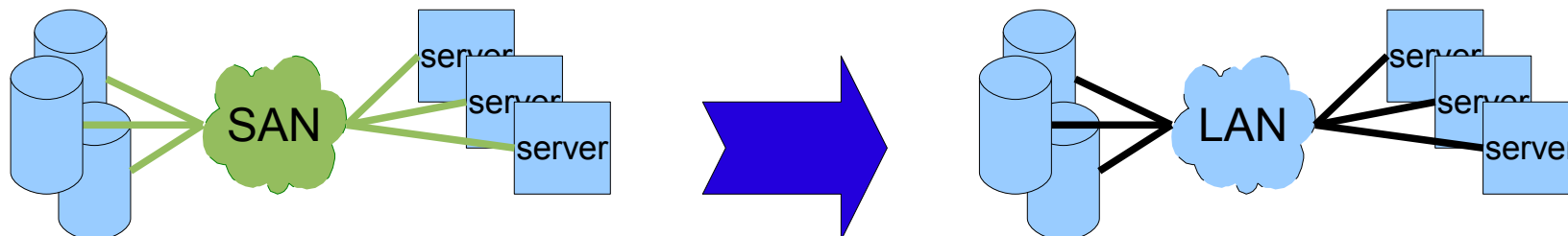
increase data availability



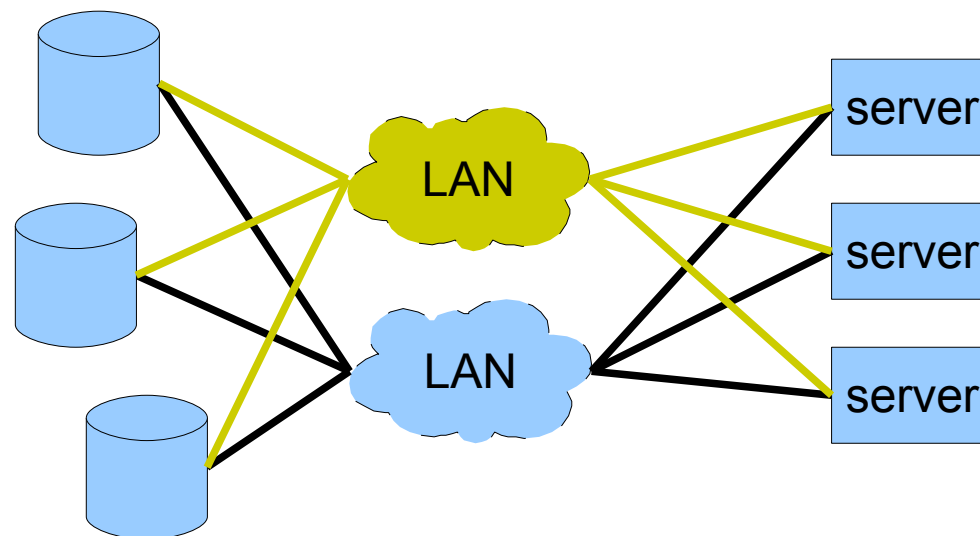
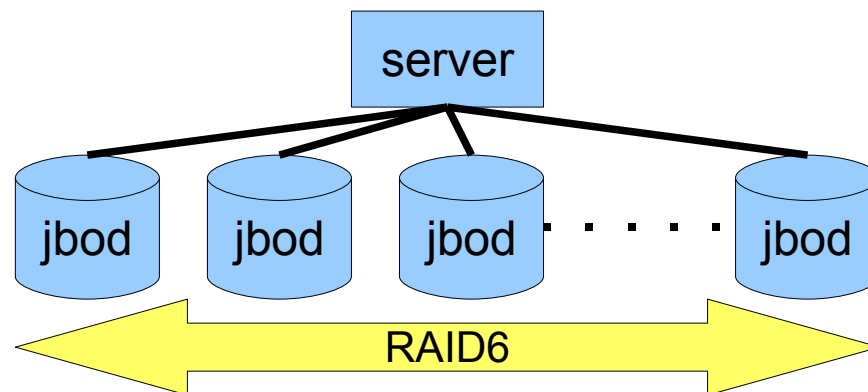
more flexibility in scaling



consolidate network infrastructure



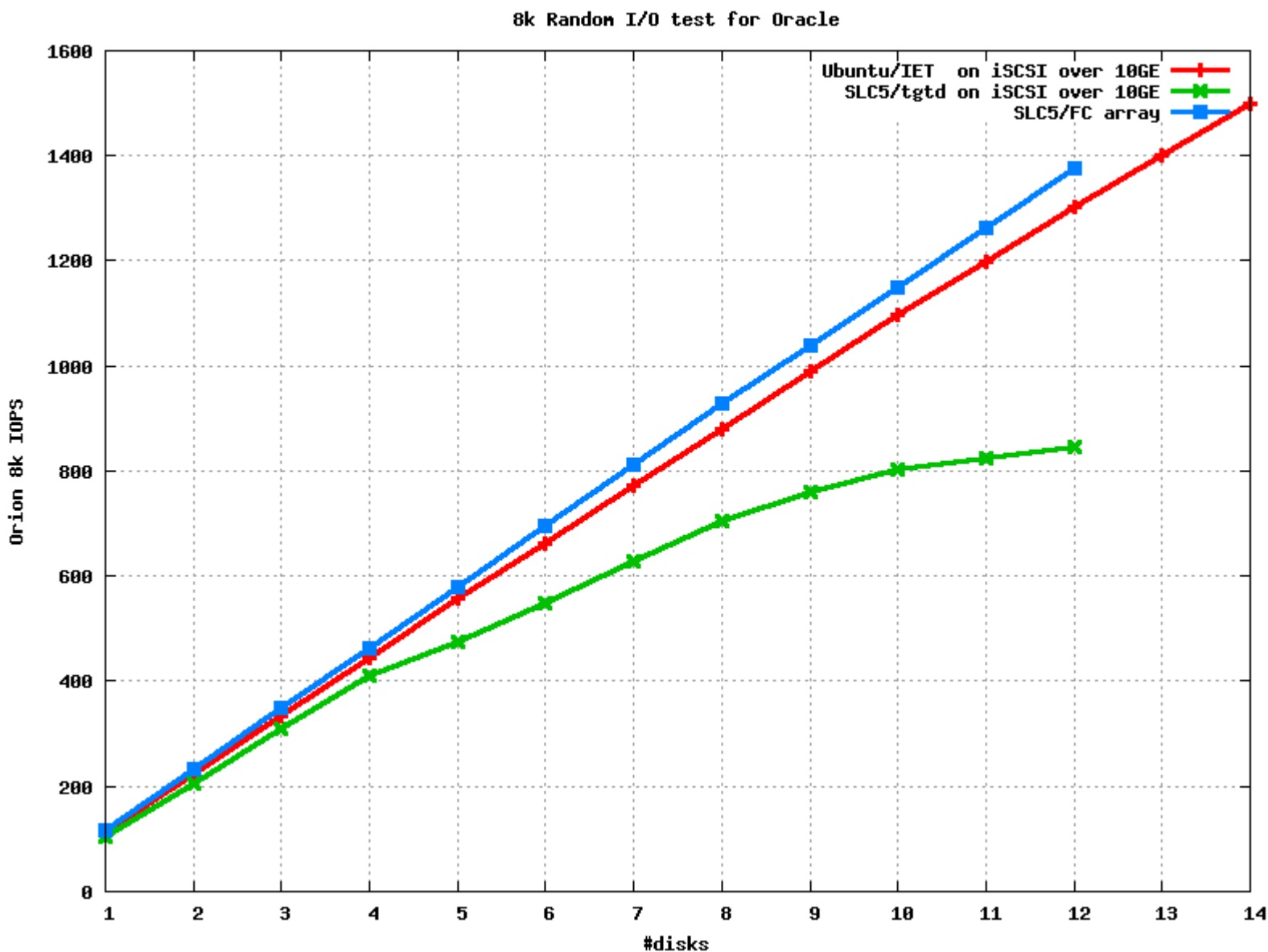
- Bulk storage
- Databases (RAC)
- Virtualisation
- Backup space
- Lustre OST/MDT
- ....

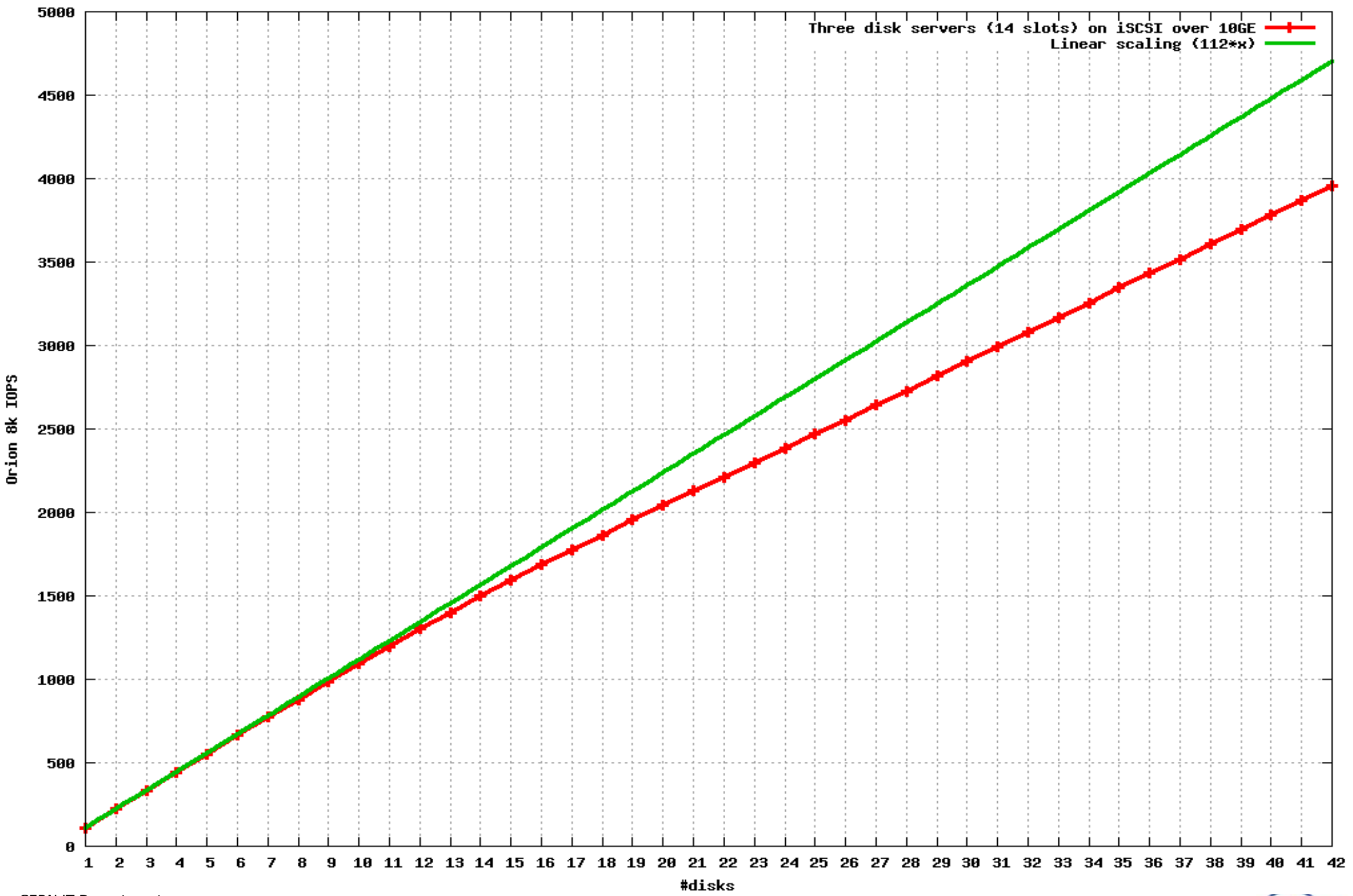


- Linux initiators: OK in SL4 and SL5
- Linux targets:
  - default RHEL5 stack (tgt): single-threaded
  - IET target stack: multithreaded, faster
- Robustness
  - client can survive target reboots (!)
  - zero fsprobe errors so far
  - watch out for persistent connections



- Infortrend low-end box
  - inexpensive, 1GE only
- Dell EqualLogic “medium” box
  - 3x1GE, redundant everything, load balancing
- Storage-in-a-box server
  - GE or 10GE
  - Ubuntu Jaunty (IET stack) or SLC5 (tgttd)
  - 3ware RAID card in JBOD mode
  - 2x Clovertown CPUs, 8G RAM, 14+2 disks







- Test and learn corner cases (MD RAID..)
- Quattor NCM components
  - initiator- and target-side components
  - ncm-filesystems/blockdevices will be reusable
  - Linux-HA integration
- Monitoring
  - MD monitoring is available
- Recovery procedures and tools
- Resource management

- 1) Proprietary targets, pilot production
- 2) Linux target, simplest (mirrored) case
  - w/o failover
- 3) Oracle test setup
  - JBODs, Oracle manages data
- 4) Mirror with automated failover
  - Linux-HA integration
- 5) Lustre OST
  - Test for Lustre, production for iSCSI

- Improve MD rebuild speeds
  - write-intent bitmap (on SSD)?
  - parallel rebuild? (~RAID5E / 5EE / 6E)
- Data integrity
  - T10 DIF (data integrity)?
  - read/check parity on read?