# The LOFAR phase II cluster

Chris Broekema

ASTRON     LOFAR     NWO
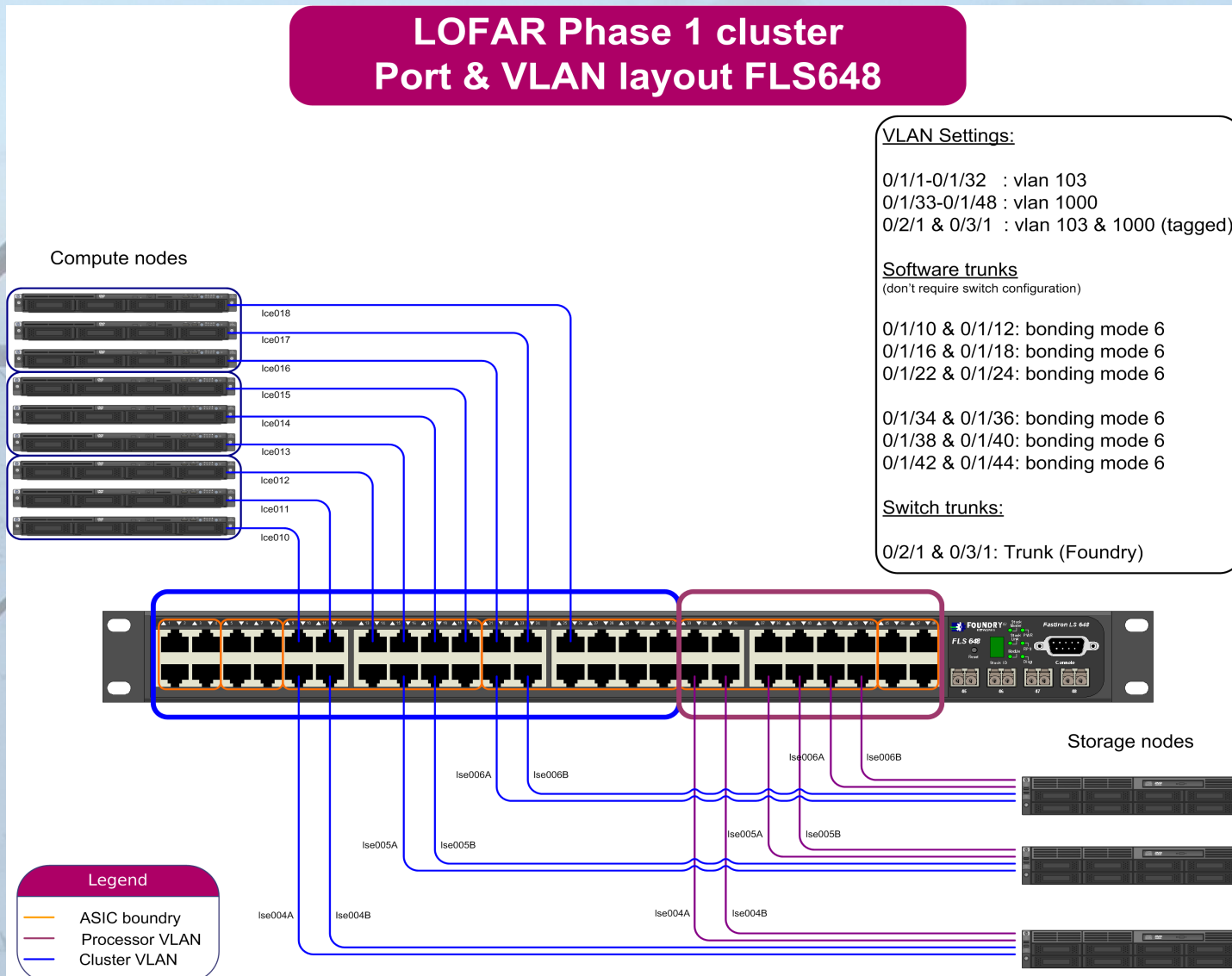
# The phase I cluster

- 500 TB storage

- Limited design bandwidth

  - 15 Gbps write (continuous)

  - 30 Gbps read (burst)

- Total compute power ~4.8 TFlops (R_peak)

  - 3.6 TFlops in compute nodes

  - 1.2 TFlops in storage nodes

- Gigabit Ethernet as interconnect

**ASTRON**     **LOFAR**     **NWO**

# The phase I cluster

- 72 Compute node
  - 2x Intel L5520
    - 4 Cores, 2.5 GHz
    - 8 Cores total, 50 Gflops
  - 16 GB main memory per node
    - 2 GB/core
  - 1 TB scratch space (2x 500 GB in RAID0)
- 24 Storage nodes
  - Identical CPUs / Memory as compute node
  - 24x 1TB disks
  - ~20 TB storage capacity available per storage node
- Configured in 8 sub-clusters
  - 3 storage nodes + 9 compute nodes = 1 sub-cluster

ASTRON    LOFAR    NWO

# The phase I cluster



**LOFAR Phase 1 cluster
Port & VLAN layout FLS648**

Compute nodes

Ice018
Ice017
Ice016
Ice015
Ice014
Ice013
Ice012
Ice011
Ice010

VLAN Settings:

0/1/1-0/1/32  : vlan 103
0/1/33-0/1/48 : vlan 1000
0/2/1 & 0/3/1  : vlan 103 & 1000 (tagged)

Software trunks:
(don't require switch configuration)

0/1/10 & 0/1/12: bonding mode 6
0/1/16 & 0/1/18: bonding mode 6
0/1/22 & 0/1/24: bonding mode 6

0/1/34 & 0/1/36: bonding mode 6
0/1/38 & 0/1/40: bonding mode 6
0/1/42 & 0/1/44: bonding mode 6

Switch trunks:

0/2/1 & 0/3/1: Trunk (Foundry)

Storage nodes

lse006A  lse006B
lse006A  lse006B
lse005A  lse005B
lse005A  lse005B
lse004A  lse004B
lse004A  lse004B

## Legend
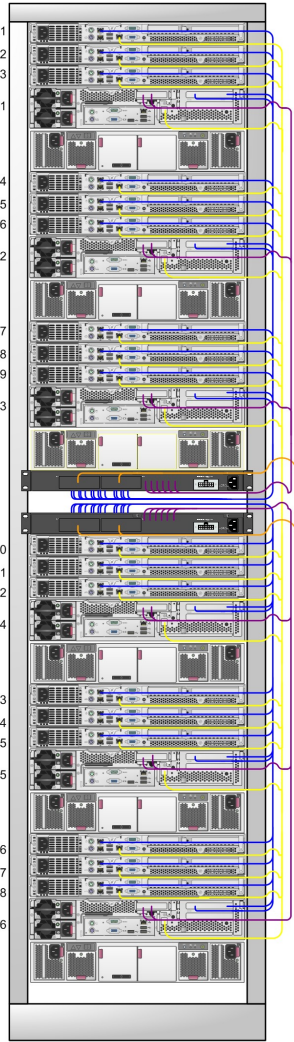| | |
|---|---|
| ASIC boundry | |
| Processor VLAN | |
| Cluster VLAN | |

ASTRON   LOFAR   NWO

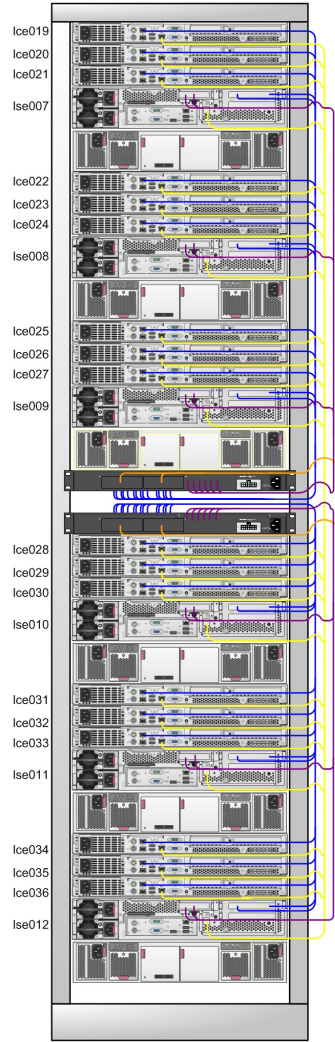# LOFAR Phase 1 cluster
# Rear view



**Legend**
- 10GbE Uplink
- 1 GbE Data line (RSP & TBB)
- 1 GbE offline data link
- Management network

**E2**

Ice001
Ice002
Ice003
Ise001
Ice004
Ice005
Ice006
Ise002
Ice007
Ice008
Ice009
Ise003
Ice010
Ice011
Ice012
Ise004
Ice013
Ice014
Ice015
Ise005
Ice016
Ice017
Ice018
Ise006

**E3**

Ice019
Ice020
Ice021
Ise007
Ice022
Ice023
Ice024
Ise008
Ice025
Ice026
Ice027
Ise009
Ice028
Ice029
Ice030
Ise010
Ice031
Ice032
Ice033
Ise011
Ice034
Ice035
Ice036
Ise012

**E4**

Ife001
Ife002
Idb001
Iexar001
Iexar002

**E5**

Ice037
Ice038
Ice039
Ise013
Ice040
Ice041
Ice042
Ise014
Ice043
Ice044
Ice045
Ise015
Ice046
Ice047
Ice048
Ise016
Ice049
Ice050
Ice051
Ise017
Ice052
Ice053
Ice054
Ise018

**E6**

Ice055
Ice056
Ice057
Ise019
Ice058
Ice059
Ice060
Ise020
Ice061
Ice062
Ice063
Ise021
Ice064
Ice065
Ice066
Ise022
Ice067
Ice068
Ice069
Ise023
Ice070
Ice071
Ice072
Ise024

ASTRON    LOFAR    NWO
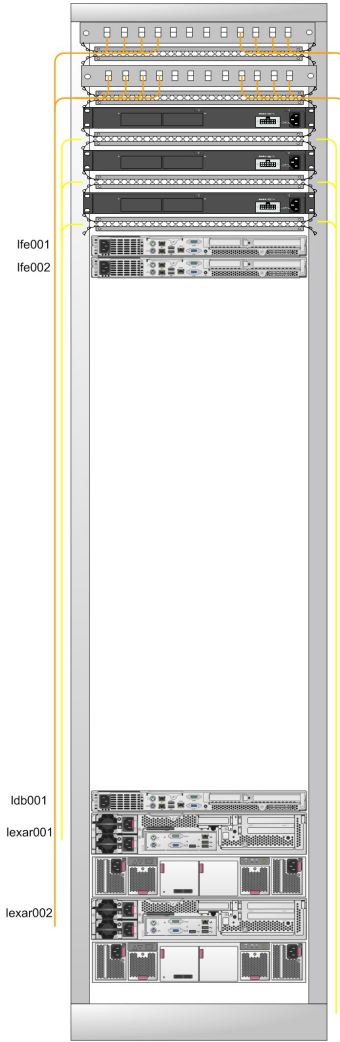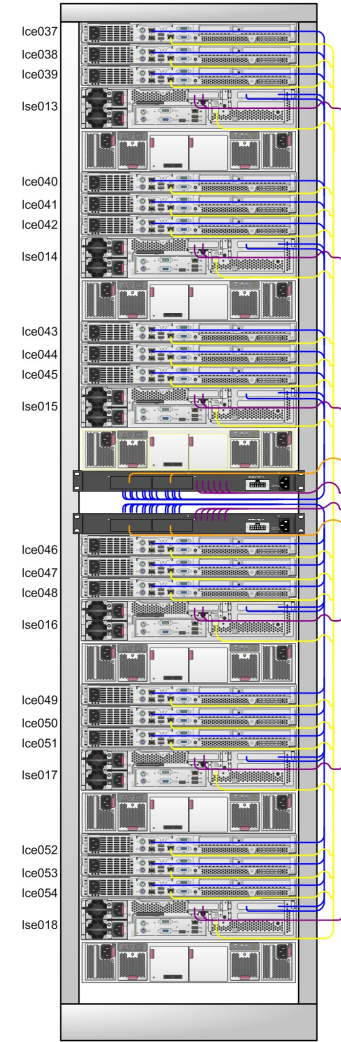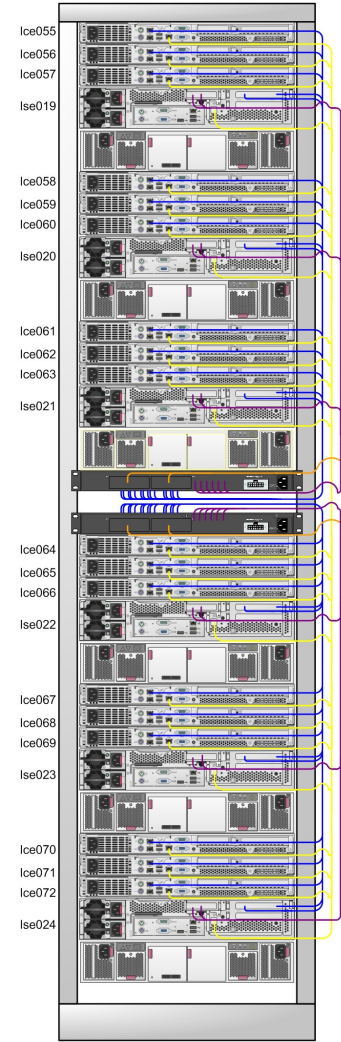
# Lessons learned from phase I

- Storage was/is too slow
    - RAID controllers quite slow
    - Simultaneous read/write virtually impossible
    - CPUs on the storage nodes mostly idle
- Sub-clusters mostly worked out
    - But moving data from storage nodes is not ideal
    - Good for commissioning, not for production
- Gigabit Ethernet is not ideal
    - Every storage node uses 4x GbE
    - 10GbE to 1GbE boundary performance sucked
- We need to Archive data

# Lessons learned from phase I

- Storage was/is too slow

  - Reduce the capacity of each node

  - Add lots more nodes (more controllers)

  - Require minimum write capacity 3.1 Gbps/node
    - Minimum required for uncompressed pulsar observations

  - Bring compute power to the data

|  | Read (MB/s) | Write (MB/s) | Read/Write (MB/s) |
| --- | --- | --- | --- |
| HP P800 | 228 | 259 | 118 / 278 |
| Areca ARC1880 | 1052 | 919 | 418 / 670 |

ASTRON          LOFAR          NWO

# Lessons learned from phase I

- Subclusters not suited for production
  - We need a place for commissioners and developers
  - But production needs most/all of the new cluster
  - Reuse phase I cluster for commissioning and development
- Gigabit Ethernet is not ideal
  - Infiniband for cluster nodes
  - Requires a bridge between the LOFAR Ethernet network and the Infiniband interconnect
  - BridgeX offered and rejected, PC based solution

ASTRON          LOFAR          NWO

# Lessons learned from phase I

- We need to archive data
  - Was not fully considered when building phase I
  - Reuse storage nodes phase I as staging cluster
  - Target project in Groningen
    - Provide high bandwidth archive facility
    - For end products and raw data
    - Possible to stream directly to Target from phase II cluster
    - Hopefully will provide significant additional compute power too
  - LTA's in Amsterdam and Juelich
    - Currently less bandwidth available (but we're ready for more)
    - Archive only end products here

ASTRON    LOFAR    NWO

# Headlines (1)

- 100 hybrid storage / compute nodes
  - 2x AMD Opteron 6172 CPUs per node
    - 12 Cores, 2.1 GHz
    - 24 Cores total, 201,6 GFlops
  - 64 GB main memory per node
    - 2.67 GB/core
  - Areca ARC-1880ixl RAID controller
  - 12x 2TB 7200 rpm disks
    - ~20 TB storage capacity available per node
  - QDR Infiniband interconnect
    - 32 Gbp theoretical maximum throughput (~20 Gbps realistic)

**ASTRON**   **LOFAR**   **NWO**

# Headlines (2)

- Storage capacity
  - 12x 2 TB disks in RAID5
  - 2 TB scratch space per node
  - 20 TB available
  - 2 PB total capacity
- Computational capacity
  - 100x hybrid compute/storage nodes
  - 20.6 TFlops Peak performance

ASTRON  LOFAR  NWO

# What will I notice

- A lot more cores per node
  - Parallelism even more important
- A bit more memory per core; a lot more per node
- Data is now (hopefully / mostly) stored locally
  - No / fewer NFS issues to contend with
- RAID performance should be much better
- High performance interconnect
  - No fully non-blocking, but much better than phase I cluster
- Bandwidth to and from cluster via bridge
  - Performance as yet unknown

ASTRON    LOFAR    NWO

# Phase II cluster – Bandwidths (Gbps)

| From \ To | Stations | BG/P | Cluster | Staging | Target | LTA's |
|-----------|----------|------|---------|---------|--------|-------|
| Stations | X | 200 | 80-160 | 80 | - | - |
| BG/P | - | X | 80-160 | 80 | - | - |
| Cluster | - | 80-160 | X | 80 | 80 | - |
| Staging | - | - | 80 | X | 80 | 10-80 |
| Target | - | - | - | - | X | - |
| LTA's | - | - | - | - | - | X |