# LOFAR CEP Design & Performance

Chris Broekema
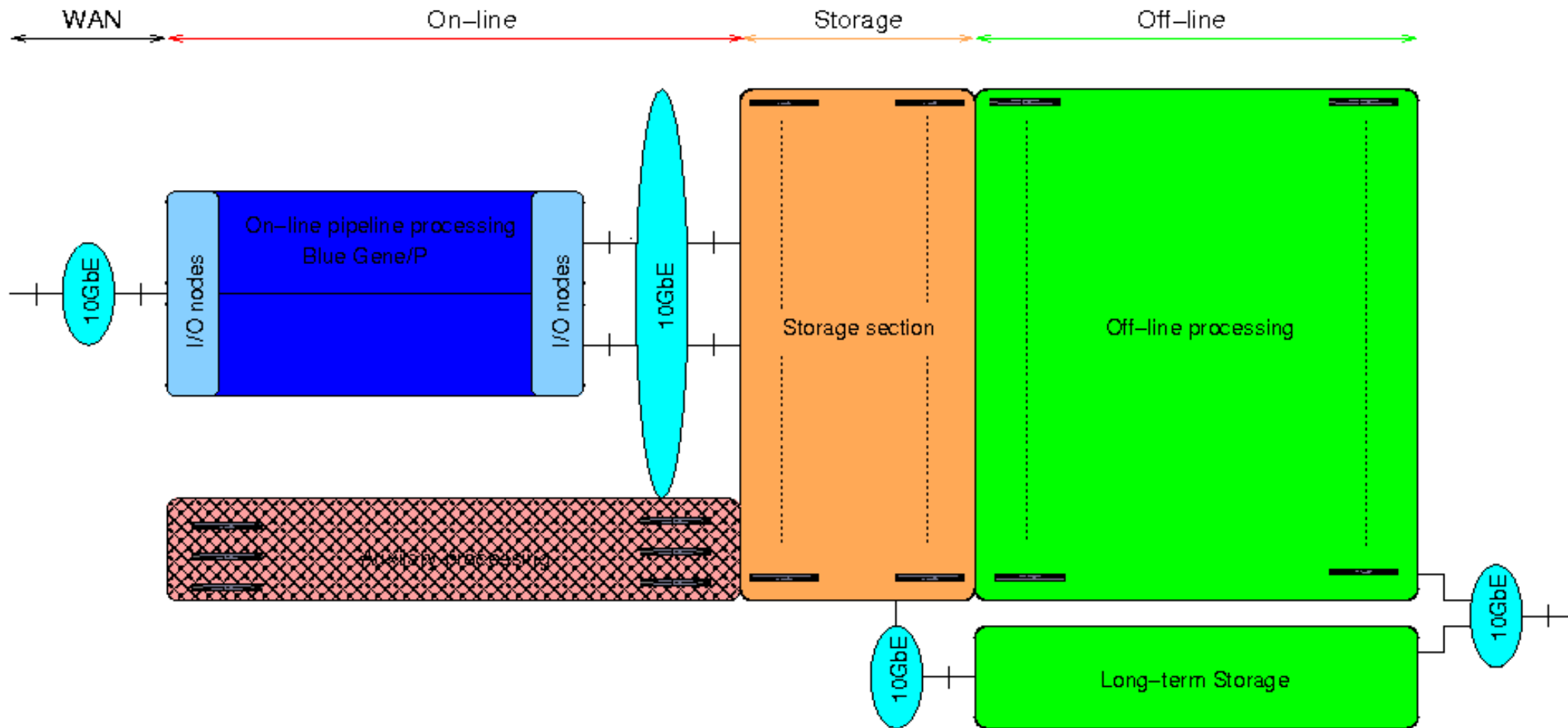
ASTRON

# Outline

- The LOFAR Central Processor

  - Top level design

  - Current hardware

- Current status and recent results

  - Standard imaging mode

  - Tied-array beamforming (pulsar mode)

- The offline processor
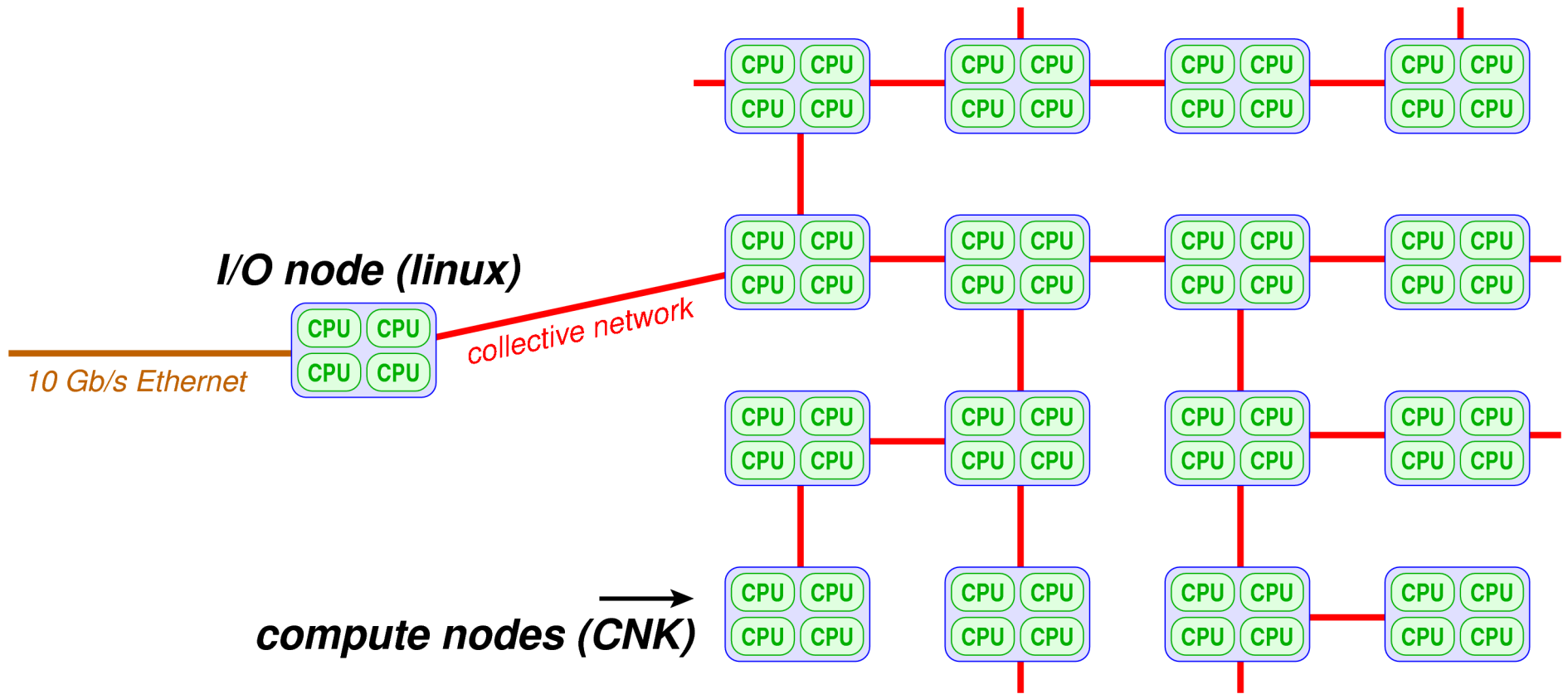
  - Performance requirements

  - Design

- Summary

# The LOFAR central processor

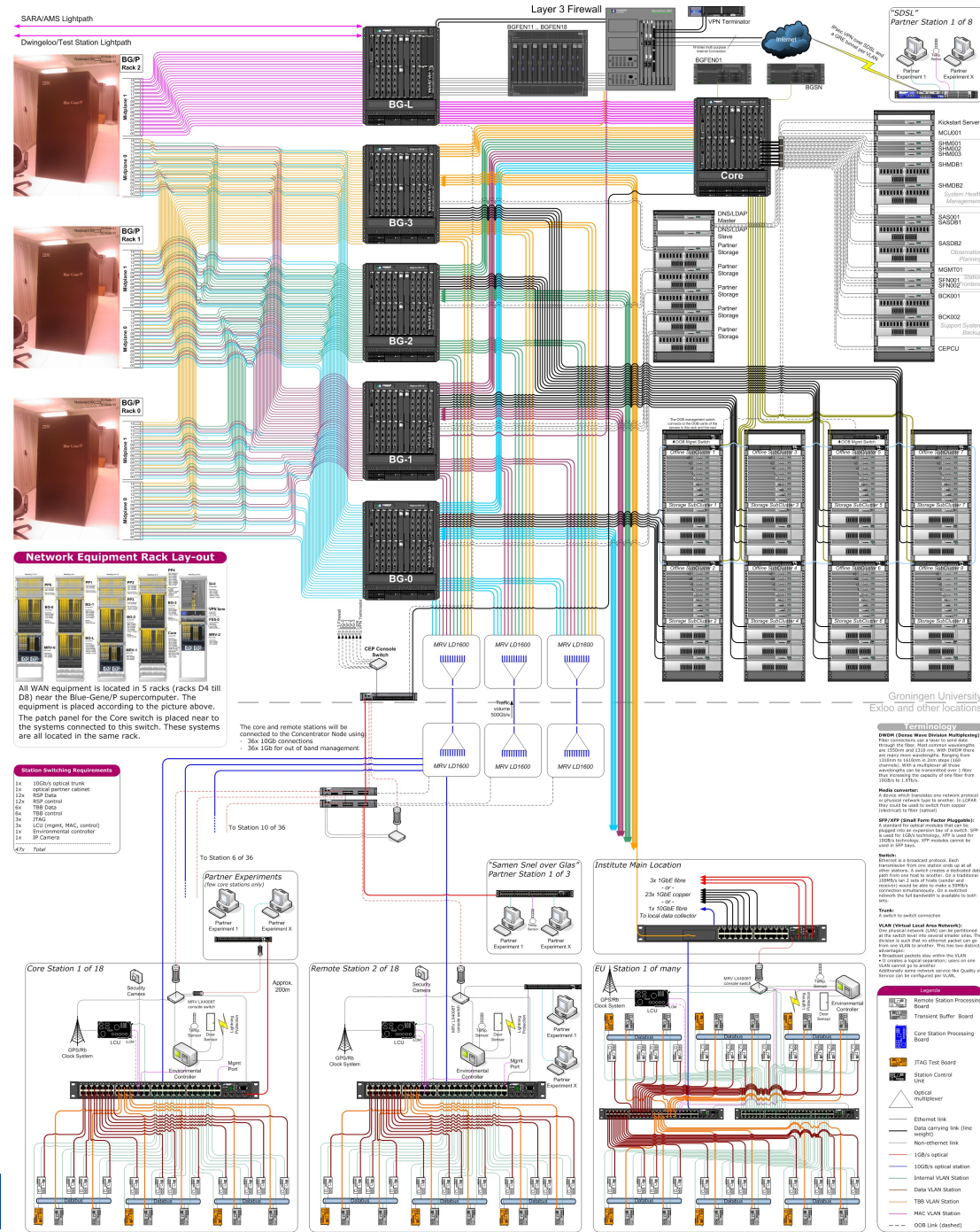# Central processor

- 3 rack IBM Blue Gene/P

  - #75 in the Top 500 (11-2008)

  - Peak performance 41.8 TFlops

    - (actually 44.4 TFlops, including I/O nodes)

  - 13056 PowerPC cores @ 850 MHz

    - Quad core system-on-chip CPUs
    - Double FPU
    - exceptional complex number support

  - ~6 TiB memory

  - 192 10GbE links

  - Several dedicated internal networks (torus, tree)
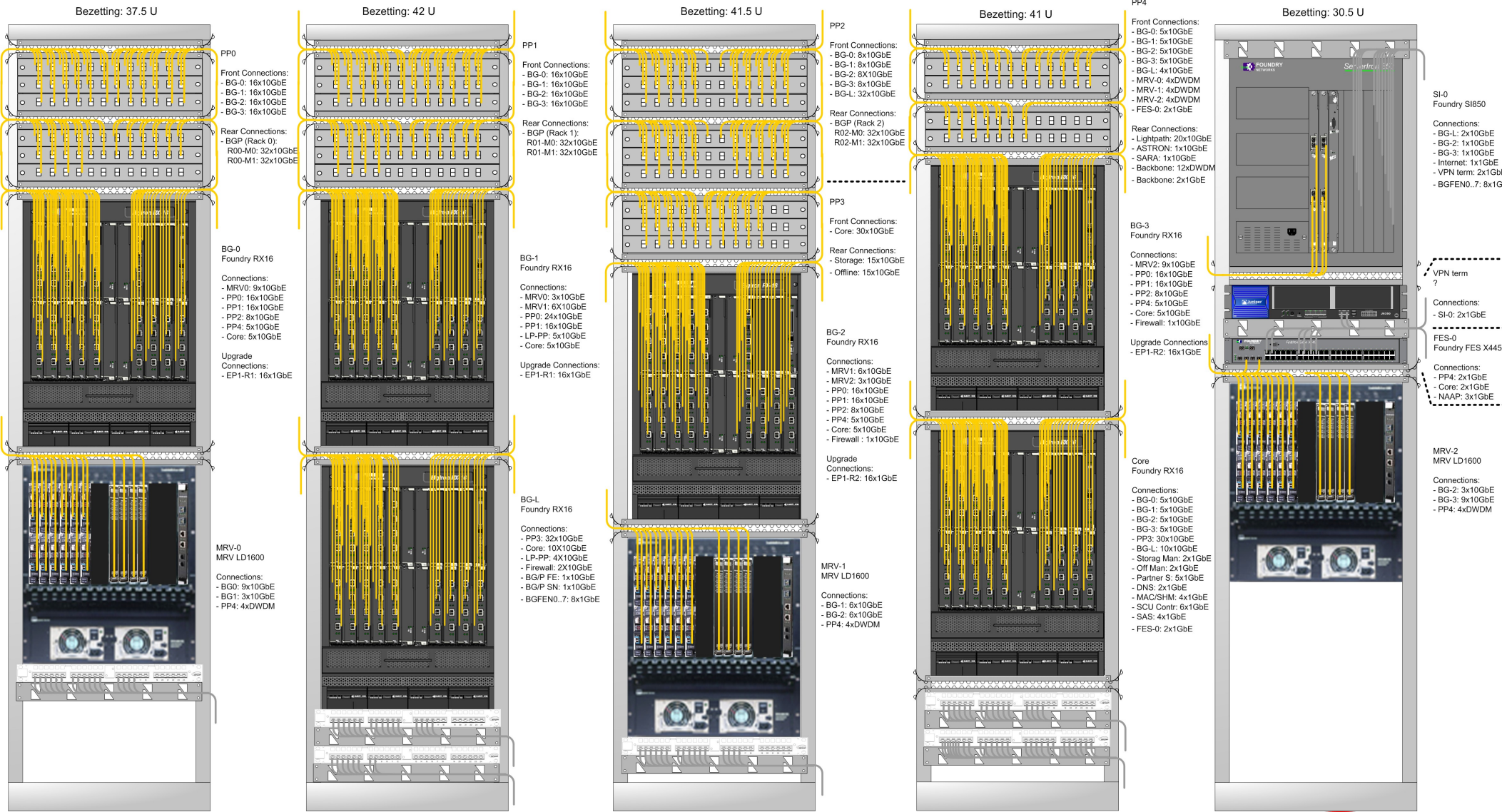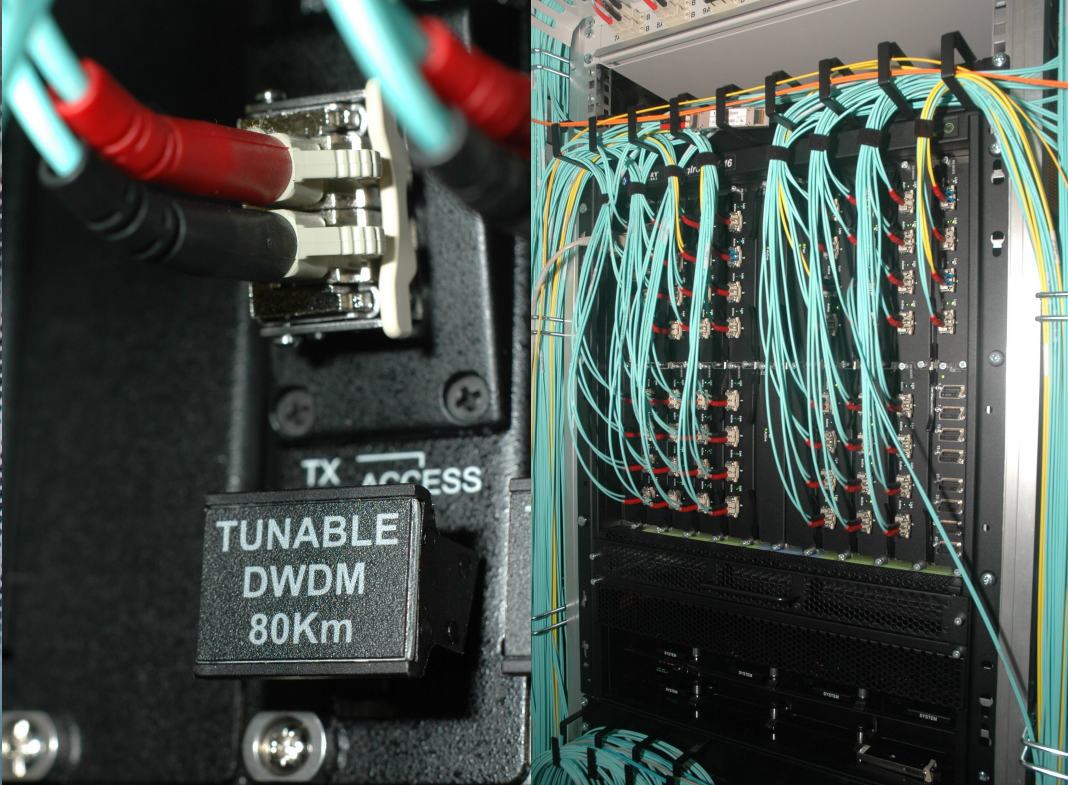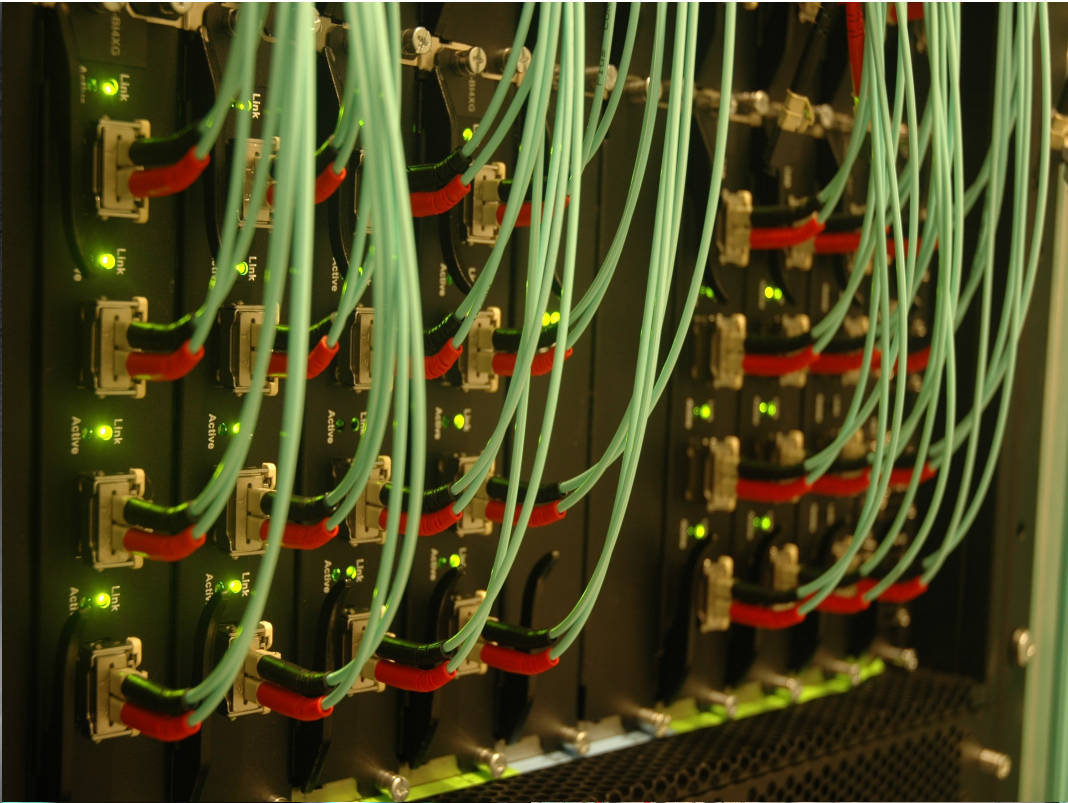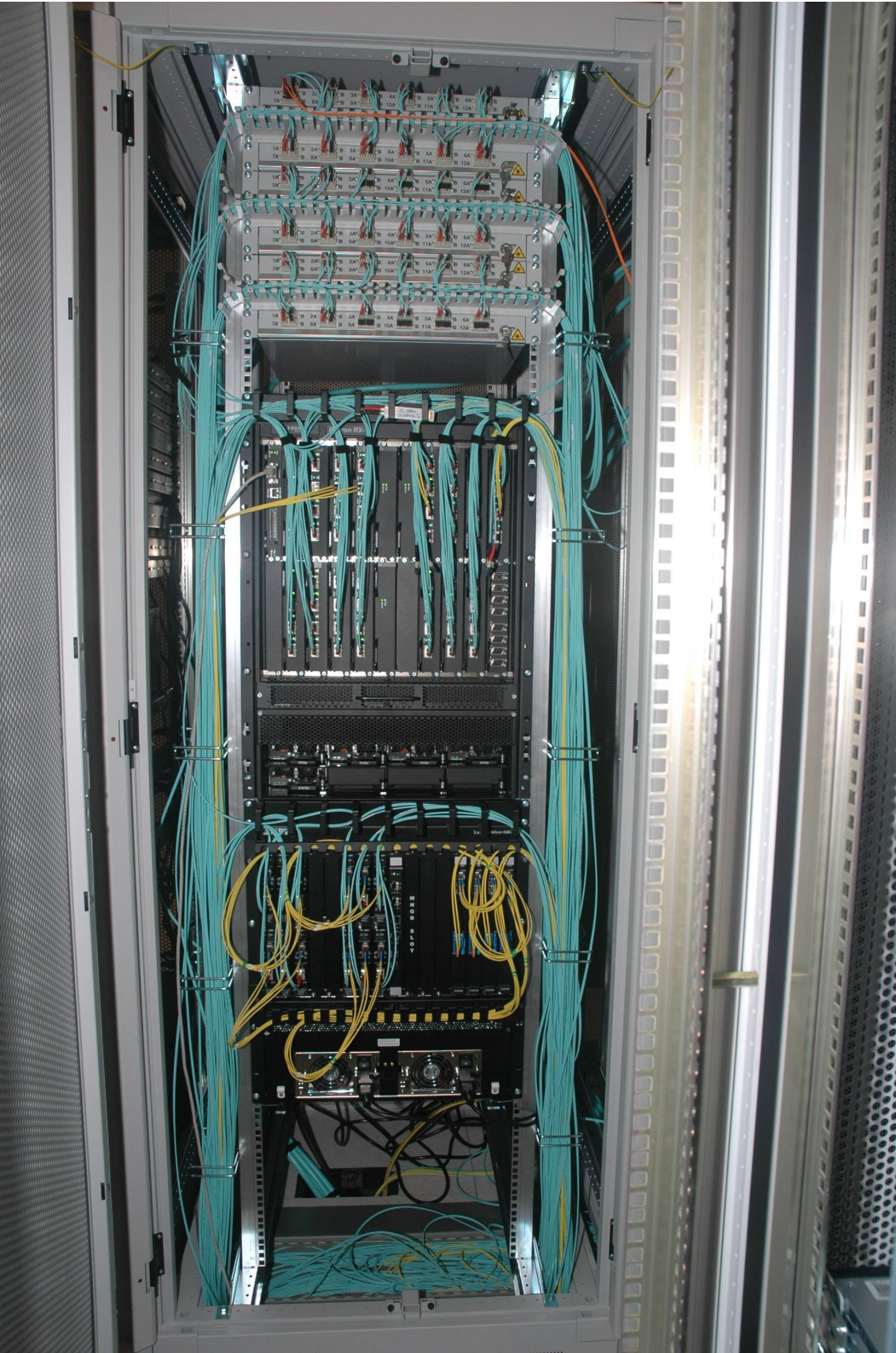
# Blue Gene/P pset

# Central processor

- 6 Foundry BigIron RX16 switching frames
    - 1 core, 4 leafs and 1 infrastructure
- ~350 10 GbE ports
    - 192 BG/P, ~70 stations, ~70 uplink, ~10 science
- ~300 GbE ports
- Dataflow optimized network design
    - keep dataflow within one switching frame
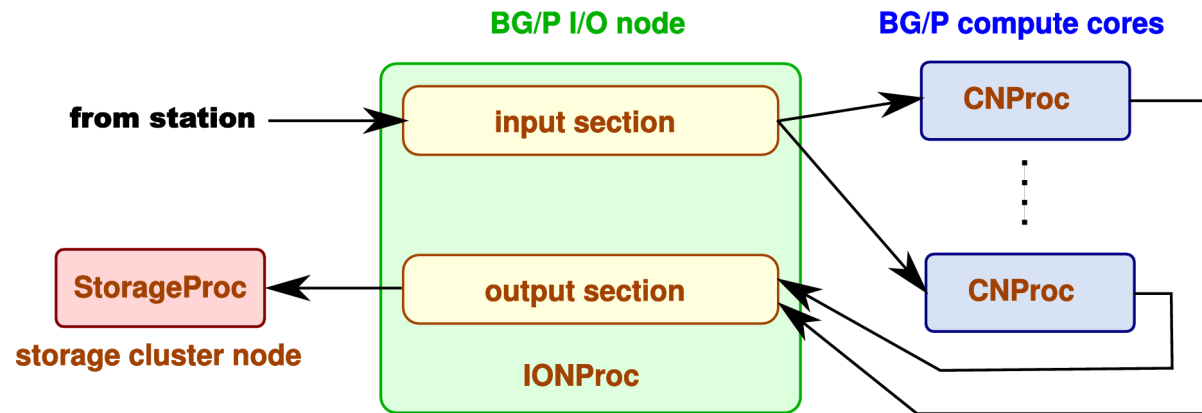    - Bandwidth between frames limited (~50 Gbps)

# LOFAR
## *Network Topology*

ASTRON · LOFAR

# Network Equipment Rack Lay-out

## Bezetting: 37.5 U

PP0

Front Connections:
- BG-0: 16x10GbE
- BG-1: 16x10GbE
- BG-2: 16x10GbE
- BG-3: 16x10GbE

Rear Connections:
- BGP (Rack 0):
  R00-M0: 32x10GbE
  R00-M1: 32x10GbE

BG-0
Foundry RX16

Connections:
- MRV0: 9x10GbE
- PP0: 16x10GbE
- PP1: 16x10GbE
- PP2: 8x10GbE
- PP4: 5x10GbE
- Core: 5x10GbE

Upgrade
Connections:
- EP1-R1: 16x1GbE

MRV-0
MRV LD1600

Connections:
- BG0: 9x10GbE
- BG1: 3x10GbE
- PP4: 4xDWDM

## Bezetting: 42 U

PP1

Front Connections:
- BG-0: 16x10GbE
- BG-1: 16x10GbE
- BG-2: 16x10GbE
- BG-3: 16x10GbE

Rear Connections:
- BGP (Rack 1):
  R01-M0: 32x10GbE
  R01-M1: 32x10GbE

BG-1
Foundry RX16

Connections:
- MRV0: 3x10GbE
- MRV1: 6X10GbE
- PP0: 24x10GbE
- PP1: 16x10GbE
- LP-PP: 8x10GbE
- Core: 5x10GbE

Upgrade Connections:
- EP1-R1: 16x1GbE

BG-L
Foundry RX16

Connections:
- PP3: 32x10GbE
- Core: 10X10GbE
- LP-PP: 4X10GbE
- Firewall: 2X10GbE
- BG/P FE: 1x10GbE
- BG/P SN: 1x10GbE
- BGFEN0..7: 8x1GbE

## Bezetting: 41.5 U

PP2

Front Connections:
- BG-0: 8x10GbE
- BG-1: 8x10GbE
- BG-2: 8X10GbE
- BG-3: 8x10GbE
- BG-L: 32x10GbE

Rear Connections:
- BGP (Rack 2)
  R02-M0: 32x10GbE
  R02-M1: 32x10GbE

PP3

Front Connections:
- Core: 30x10GbE

Rear Connections:
- Storage: 15x10GbE
- Offline: 15x10GbE

BG-2
Foundry RX16

Connections:
- MRV1: 6x10GbE
- MRV2: 3x10GbE
- PP0: 16x10GbE
- PP1: 16x10GbE
- PP2: 8x10GbE
- PP4: 5x10GbE
- Core: 5x10GbE
- Firewall : 1x10GbE

Upgrade
Connections:
- EP1-R2: 16x1GbE

MRV-1
MRV LD1600

Connections:
- BG-1: 6x10GbE
- BG-2: 6x10GbE
- PP4: 4xDWDM

## Bezetting: 41 U

BG-3
Foundry RX16

Connections:
- MRV2: 9x10GbE
- PP0: 16x10GbE
- PP1: 16x10GbE
- PP2: 8x10GbE
- PP4: 5x10GbE
- Core: 5x10GbE
- Firewall: 1x10GbE

Upgrade Connections:
- EP1-R2: 16x1GbE

Core
Foundry RX16

Connections:
- BG-0: 5x10GbE
- BG-1: 5x10GbE
- BG-2: 5x10GbE
- BG-3: 5x10GbE
- PP3: 30x10GbE
- BG-L: 10x10GbE
- Storag Man: 2x1GbE
- Off Man: 2x1GbE
- Partner S: 5x1GbE
- DNS: 2x1GbE
- MAC/SHM: 4x1GbE
- SCU Contr: 6x1GbE
- SAS: 4x1GbE
- FES-0: 2x1GbE

## Bezetting: 30.5 U

PP4

Front Connections:
- BG-0: 5x10GbE
- BG-1: 5x10GbE
- BG-2: 5x10GbE
- BG-3: 5x10GbE
- BG-L: 4x10GbE
- MRV-0: 4xDWDM
- MRV-1: 4xDWDM
- MRV-2: 4xDWDM
- FES-0: 2x1GbE

Rear Connections:
- Lightpath: 20x10GbE
- ASTRON: 1x10GbE
- SARA: 1x10GbE
- Backbone: 12xDWDM
- Backbone: 2x1GbE

SI-0
Foundry SI850

Connections:
- BG-L: 2x10GbE
- BG-2: 1x10GbE
- BG-3: 1x10GbE
- Internet: 1x1GbE
- VPN term: 2x1GbE
- BGFEN0..7: 8x1GbE

VPN term
?

Connections:
- SI-0: 2x1GbE

FES-0
Foundry FES X445

Connections:
- PP4: 2x1GbE
- Core: 2x1GbE
- NAAP: 3x1GbE

MRV-2
MRV LD1600

Connections:
- BG-2: 3x10GbE
- BG-3: 9x10GbE
- PP4: 4xDWDM

ASTRON

LOFAR

NWO

TX ACCESS

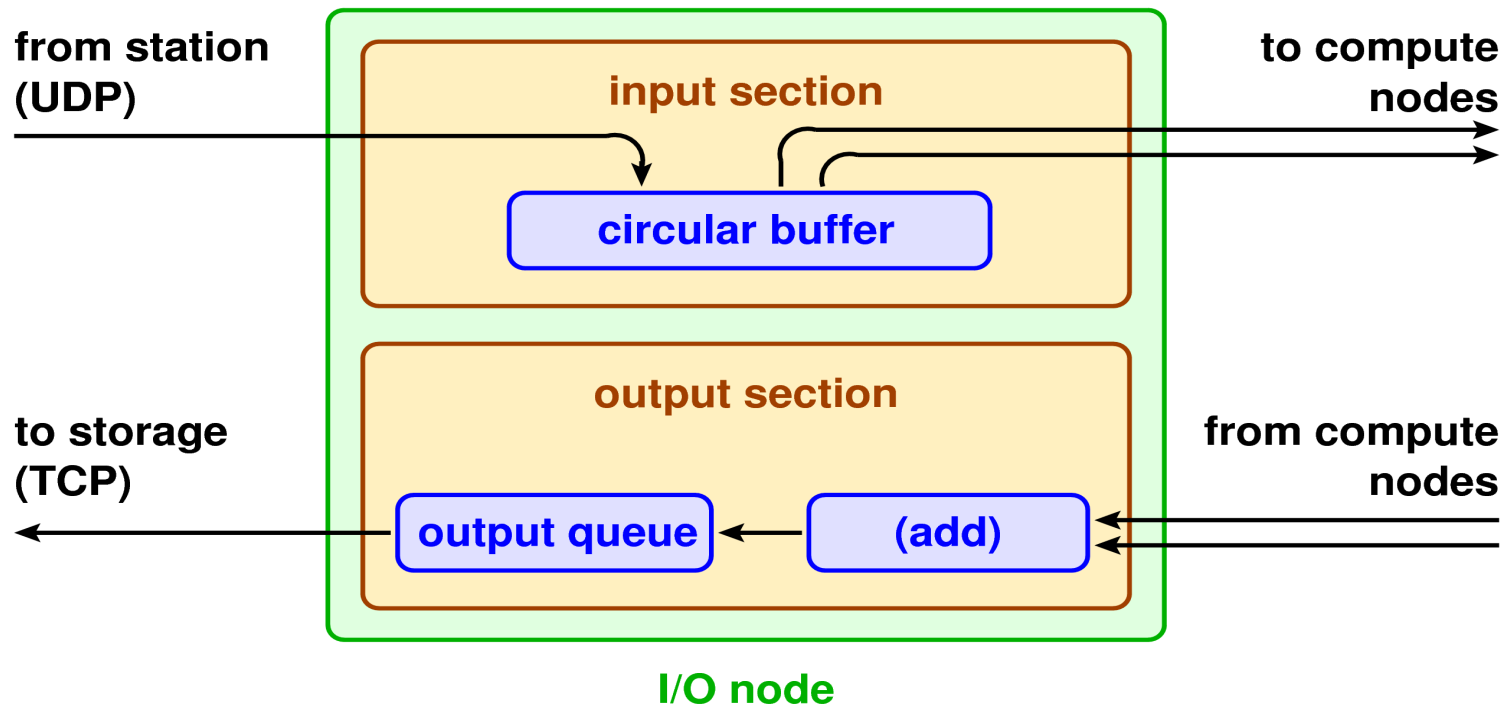TUNABLE
DWDM
80Km

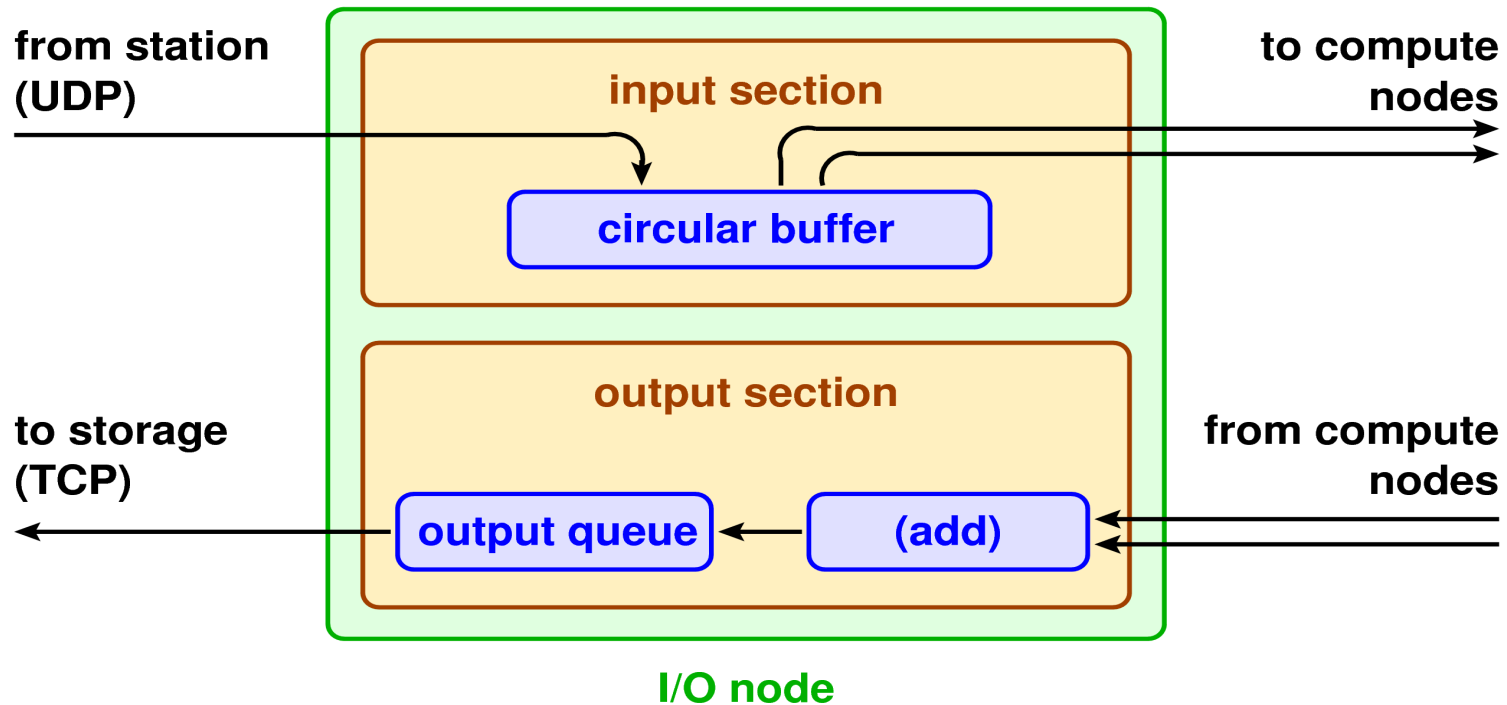# The Blue Gene/P Correlator



- three distributed applications/platforms
  - BG/P I/O nodes
  - BG/P compute nodes
  - external storage nodes
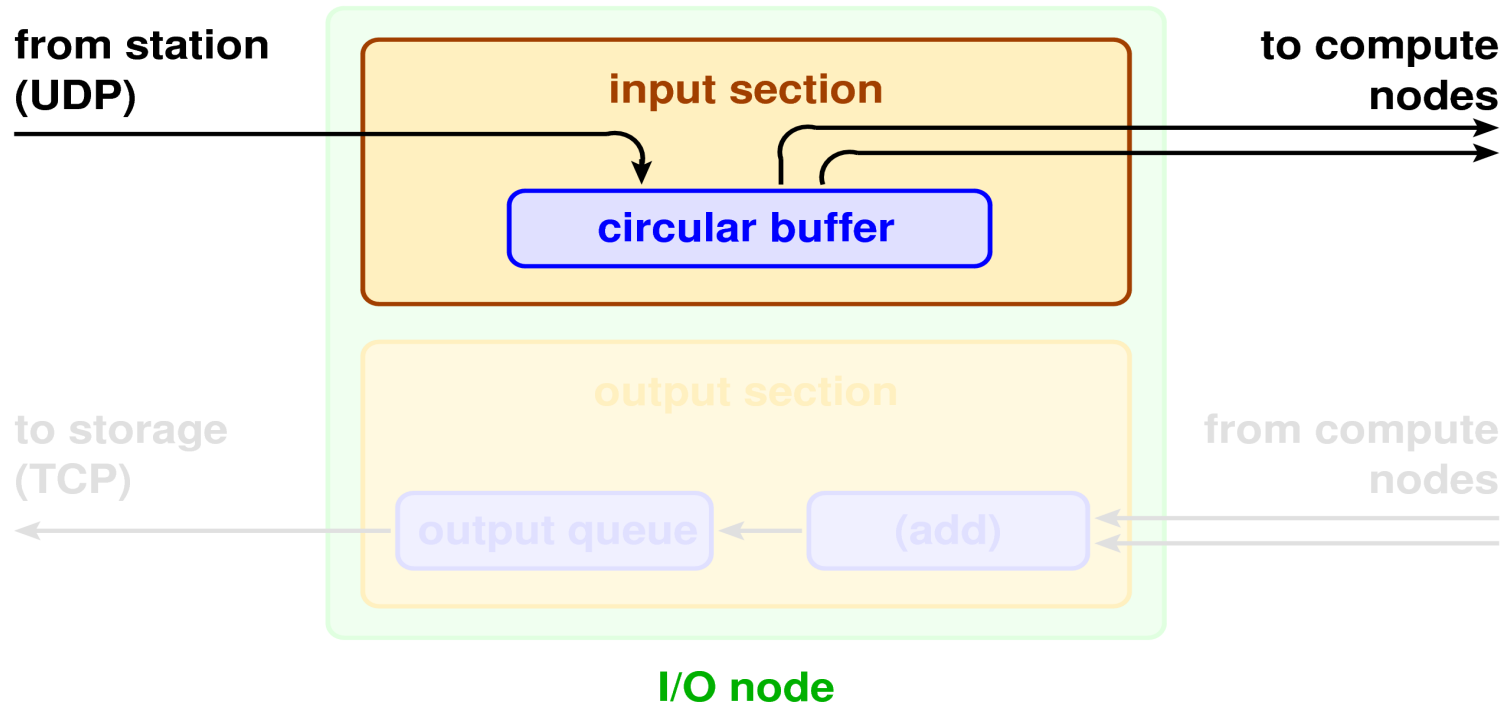
# I/O node processing



- application on I/O nodes

    – more efficient & flexible

    – BG/L: saved costs for input cluster

    – BG/L: major system software changes (ZOID) [PPoPP'08]

# I/O node processing



- Two sections
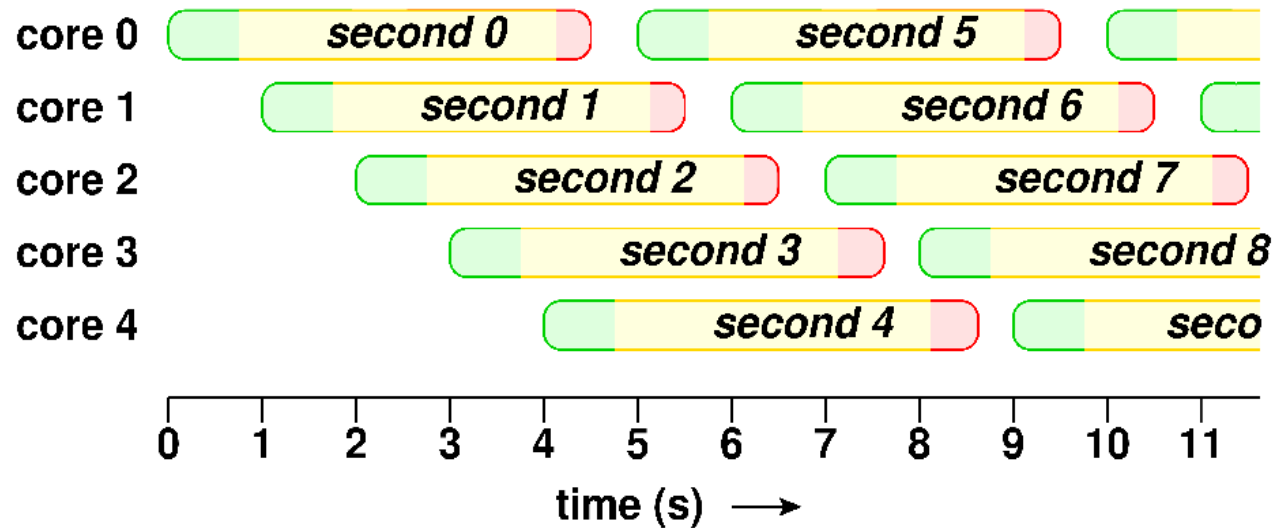    - Input section
    - Output section
- Heavily threaded & optimized

# I/O node input section

**from station (UDP)**

**input section**

**circular buffer**

**to compute nodes**

**to storage (TCP)**

**output section**

**output queue** ← **(add)** ← **from compute nodes**

**I/O node**

- one station per I/O node
- 48,828 pkt/s
- circular buffer (~2.5 s):
  - WAN delays
  - set observation direction
  - handle hiccups
- handles missing data
- wall-clock trigger

# Work distribution

- $O(100)$ independent data chunks

  - 1 second, 1 subband, all stations

  - needs > 1 second processing time

- distribute round robin over cores

  - receive, process, send, idle

# Compute node processing (1)



- Exchange (transpose)

    - All subbands; 1 station → all stations 1 subband

    - asynchronous

- Polyphase filter creates channels

- Phase correction to point accurately

# Compute node processing (2)



- Correct station-introduced bandbass
- Beam form (add) to create "Super Station" (optional)
- Correlate station samples pair-wise

# Bandbass correction

- 2 single dipoles
- 58.6 MHz
- 30 minute integration



no correction
corrected

# I/O node output section



- adds correlations (optional)
- best-effort queue
    - ensures real-time continuation of correlator

# Std imaging mode performance



- 1 rack BG/P used as correlator
- 1 rack BG/P generates simulated station data
  - Up to 64 stations @ 3.1 Gbps each
- ½ rack BG/P receives (and dumps) visibilities

# Std imaging mode performance

| observation mode | A | B | C |
|---|---|---|---|
| #stations | 64 | 64 | 48 |
| #subbands | 248 | 496 | 992 |
| #bits/sample | 16 | 8 | 4 |
| obs. bandwidth (MHz * #beams) | 48.4 | 96.9 | 194 |
| input bandwidth (Gb/s) | 64 * 3.1 | 64 * 3.1 | 48 * 3.1 |
| output bandwidth (Gb/s) | 62 * 0.58 | 62 * 1.2 | 62 * 1.3 |
| CPU load compute nodes | 35% | 70% | 85% |
| CPU load I/O nodes | 67% | 81% | 80% |
| data loss | ~ .0001% | ~ 0.01% | ~ 0.01% |

ASTRON          LOFAR          NWO

# Std imaging mode performance

- This is representative for full LOFAR

  - Up to 64 stations

- In two new observation modes (8 bit & 4 bit)

- At 150% of the specified bandwidth

- With half the designed resources

- Without significant data loss

- EoR mode can be done on 1 rack BG/P

  - (6 Racks BG/L originally)

# Std imaging mode performance
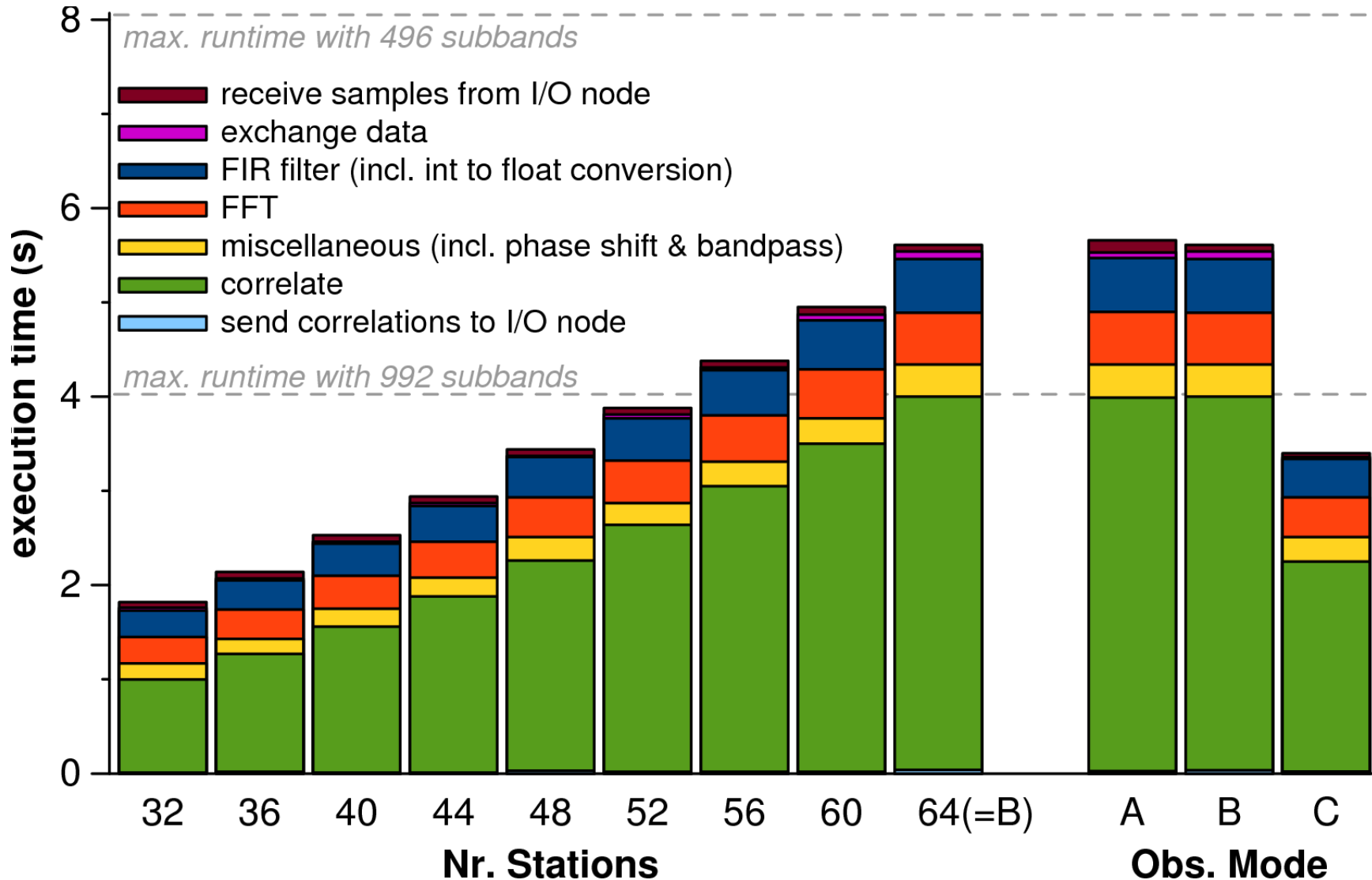## Blue Gene/P I/O node load

# Std imaging mode performance
## Blue Gene/P I/O optimizations

- Heavily modified I/O node Linux kernel

  - Avoid TLB misses

  - Optimize network stack buffer sizes

- Low overhead protocol to Compute nodes

- Optimum scheduling of threads in application

    - Use Linux real-time threads

- Use of assembler where appropriate

# Std imaging mode performance
## Blue Gene/P Compute node load

# Std imaging mode performance
## Blue Gene/P Compute node optimizations

- Heavy use of assembler in hot spots

  - Correlator (96% of peak FPU performance)

  - FIR filter (86% of peak FPU performance)

  - FFT (43% of peak FPU performance)

- Rewrite transpose to use DMA engine

  - Uses asynchronous send/recv instead of MPI_Alltoallv()
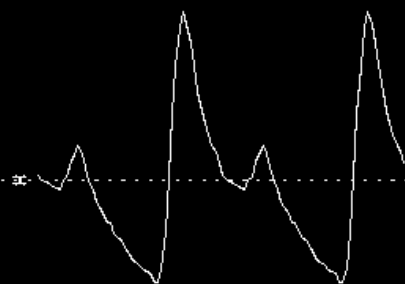
  - Hides transpose time completely

- Low overhead protocol to (and from) I/O nodes

**ASTRON**          **LOFAR**          N$\mathcal{W}$O

# Tied-Array beamforming

- Reference implementation available
  - Real-time
  - Capable of creating multiple close beams
- Complex voltages
- Stokes I
- Stokes I, Q, U, V
- Incoherent
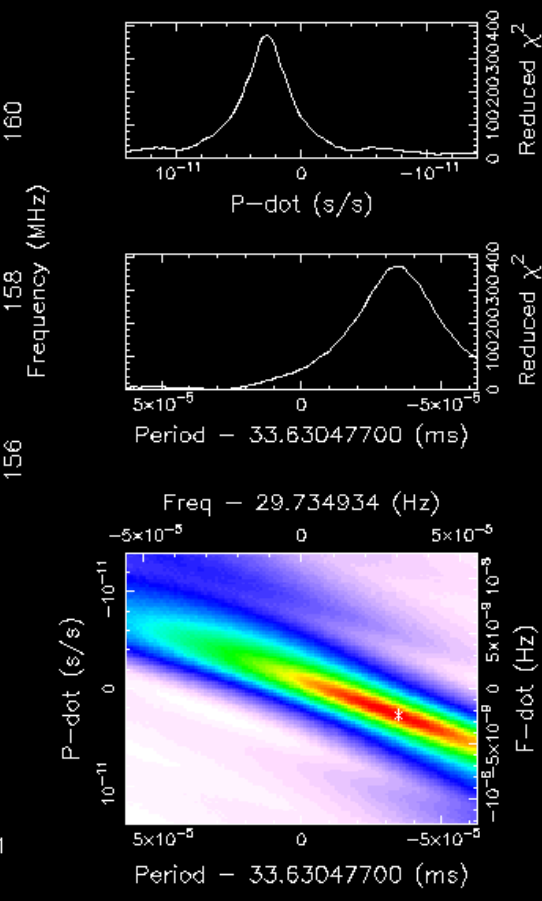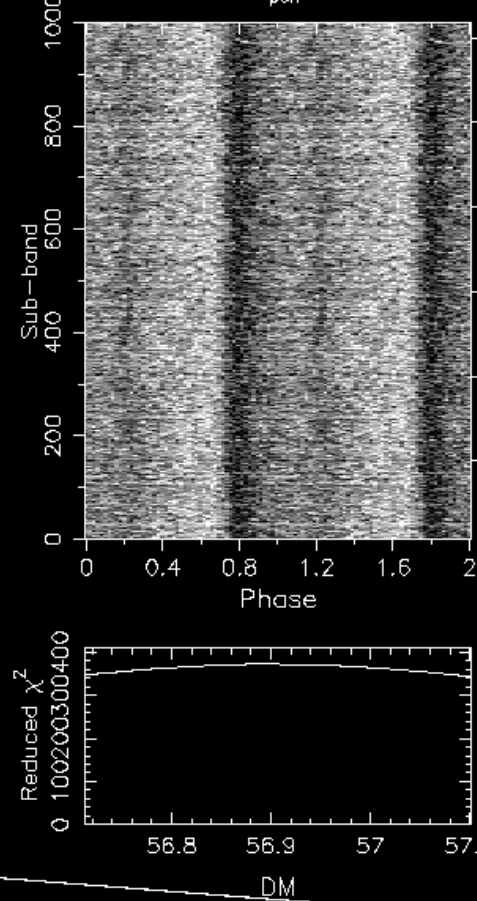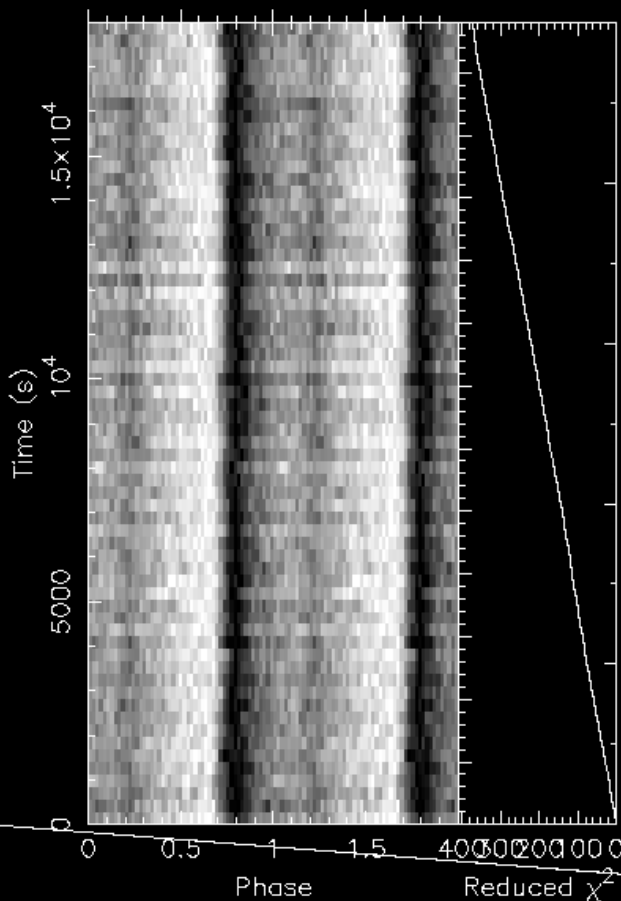- Online integration over time

# Tied-Array beamforming

- Runs correctly and stable

  - Several successful pulsar observations done

  - Multiple pencil beams, up to 64 hours

- Still way too slow

  - Can only form a few beams in real-time

  - Reference C++ implementation written for clarity

  - Optimized (assembler) version available

    - unverified, not finished

    - But observed to run ~30 times faster

# The offline processor Phase 1 specification

- Temporary storage
  - ~500 TB
  - ~15 Gbps input (continuous)
  - ~30 Gbps output (burst)
- Compute cluster
  - Flagging, calibration, imaging, source finding
  - ~5 TFlops
  - Needs to keep up with the correlator
    (although not necessarily in real-time)

# Phase 1 hardware (1)

- 24 storage nodes
  - 2 Quad-core low-power Intel Xeon CPUs
  - 16 GiB main memory
  - 24 x 1TB disks each → ~20 TB usable capacity
  - 4 GbE interfaces

- 72 compute nodes
  - 2 Quad-core low-power Intel Xeon CPUs
  - 16 GiB main memory
  - 1 TB local storage (2x 500 GB in RAID-0)
  - 2 GbE interfaces

ASTRON    LOFAR    NWO

# Phase 1 hardware (2)

- 8 GbE data switches
  - One for each sub-cluster
  - 20 Gbps uplink to Core infrastructure
- 2 frontend nodes
  - 2 Quad-core low-power Intel Xeon CPUs
  - 16 GiB main memory
  - ~2 TB storage capacity in RAID-5

# Bandwidth optimized sub-clusters

- Offline cluster does mostly batch processing

- Inter node communication is limited

- Huge data volumes

    - Communication needs to be limited

    - Necessary communication needs to be optimized

    - Cache locally to avoid unnecessary transport

- Divide cluster resources into 8 sub-clusters with optimum connectivity

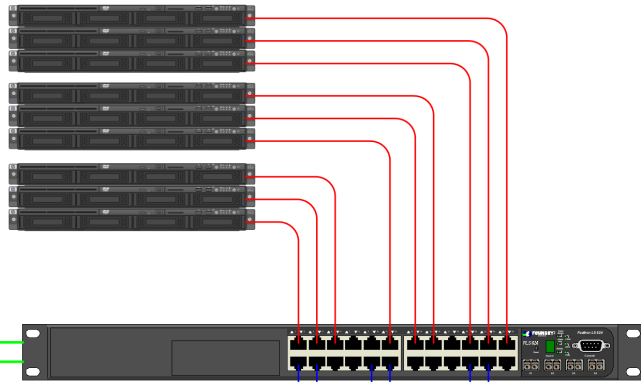# Subcluster configuration
# (INTERNAL USE ONLY)

## Legend

— 10 GbE uplink

— GbE BG/P to storage (vlan 1000)

— GbE storage out (vlan 1001)

— GbE offline cluster (vlan 1001)

Phase 1 LOFAR CEP Hardware:

8 24 port switches with 2 10 GbE uplinks each
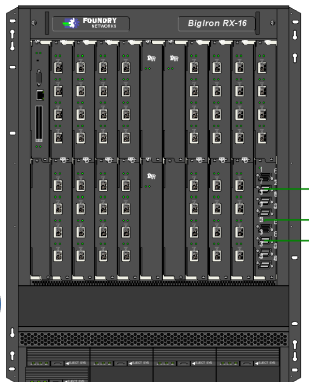24 Storage nodes, each with ~24 TB disks
72 Computational nodes

| | |
|---|---|
| Total storage capacity: | ~480 TB |
| Total input bandwidth: | ~24 Gbps |
| Total output bandwidth: | ~48 Gbps |

Each node is connected using GbE, but switch capacity is available for 2 x GbE with bonding.
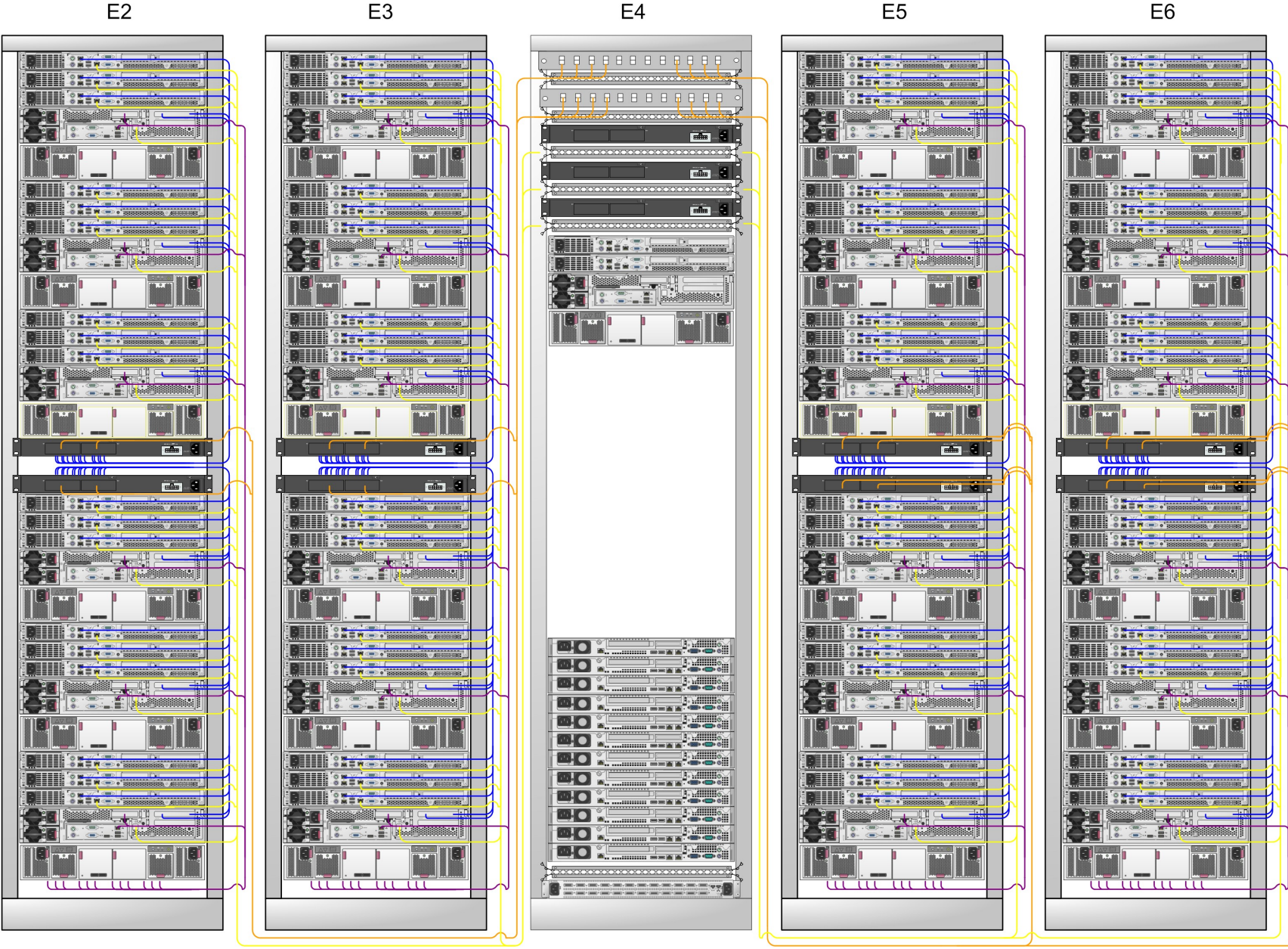
To RX16 'Core'
To RX16 'Core'

Phase 1 of the LOFAR CEP design will require eight of these subclusters. Each of the four Blue Gene/P switches will host two of these subclusters, connected to the copper blades.

Storage nodes require two GbE outputs to reach specified output bandwidths. Each storage node is connected to vlan 1000 using 1 GbE and vlan 1001 with 2 GbE devices.

ASTRO

NWO

**LOFAR Phase 1 cluster**
**Rear view**

E2  E3  E4  E5  E6

Legend
- 10GbE Uplink
- 1 GbE Data line (from Blue Gene)
- 1 GbE offline data link
- Management network

# Phase 1 hardware

- Delivery scheduled this week

- Installation until beginning of June

- Commissioning as soon as possible

    - Operating system, infrastructure & applications

    - Staggered roll-out per subcluster

    - Subclusters may be temporarily reassigned

    - Currently available cluster OS migration

# Phase 2 hardware

- Q4 2009 – Q1 2010

- Storage component grows to 2 PB
    - Input b/w ~50 Gbps (sustained)
    - Output b/w ~100 Gbps (burst)

- Offline processing cluster
    - At least ~10 TFlops
    - May not be enough

# Summary & Conclusions (1)

- The LOFAR central processor is ready
    - We can handle full LOFAR in std imaging mode
    - At 150% of designed bandwidth
    - Tied-array beamforming is coming along nicely
- The phase 1 offline processor to be built shortly
- Phase 2 specifications to be defined next
    - Using LOFAR-20 experiences
    - Probably dominated by calibration

ASTRON        LOFAR        NWO

# Summary & Conclusions (2)

- Getting to this point required specialists
  - Linux kernel hacking on BG/P I/O nodes
  - Assembler kernels for computational hotspots
  - Detailed hardware design optimized for application
- Computation cannot be separated from I/O
  - Network → node
  - Memory → cache or CPU
  - Cache → CPU
  - Many-core architectures complicate this problem

# Acknowledgements

John W. Romein

Jan David Mol

IBM:

Rob van Nieuwpoort

Bruce Elmegreen

Todd Inglet

Argonne National Lab:

Tom Leibsch

Kazutomo Yoshii

Andrew Taufener

Kamil Iskra

**ASTRON**

**LOFAR**

**NWO**

# Relevant publications

- John W. Romein, P. Chris Broekema, Jan David Mol, and Rob V. van Nieuwpoort, Processing Real-Time LOFAR Telescope Data on a Blue Gene/P SuperComputer, Under review

- Kazutomo Yoshii, Kamil Iskra, P. Chris Broekema, H. Naik, and Pete Beckman, Characterizing the Performance of Big Memory on Blue Gene Linux, International Workshop on Parallel Programming Models and System Software for High-End Computing (P2S2'09), Vienna, Austria, September, 2009

- John W. Romein, FCNP: Fast I/O on the Blue Gene/P, Parallel and Distributed Processing Techniques and Applications (PDPTA'09), Las Vegas, NV, July, 2009

- Rob V. van Nieuwpoort and John W. Romein, Using Many-Core Hardware to Correlate Radio Astronomy Signals, ACM International Conference on SuperComputing (ICS'09), New York, NY, June, 2009

- Kamil Iskra, John W. Romein, Kazutomo Yoshii, and Pete Beckman, ZOID: I/O-Forwarding Infrastructure for Peta-Scale Architectures, ACM Symposium on Principles and Paradigms of Parallel Programming (PPoPP'08), Salt Lake City, NV, February, 2008

- John W. Romein, P. Chris Broekema, Ellen van Meijeren, Kjeld van der Schaaf, and Walther H. Zwart, Astronomical Real-Time Signal Processing on a Blue Gene/L SuperComputer, ACM Symposium on Parallel Algorithms and Architectures (SPAA'06), Cambridge, MA, July, 2006