



OPENFABRICS
ALLIANCE

12th ANNUAL WORKSHOP 2016

EXPERIENCES WITH LARGE-SCALE MULTI-SUBNET INFINIBAND FABRICS

David Southwell, CVO

Obsidian Strategics Inc.

[April 7th, 2016]





OPENFABRICS
ALLIANCE

TAKING INFINIBAND FURTHER

VISION

Enabling broad InfiniBand adoption through expanded capabilities



OBSIDIAN is about:

STRATEGICS

- **Range extension – global reach using standard WANs (2005)**
- **In-line AES cryptography – encryption and authentication (2008)**
- **Native IB routing – multiple subnets, compound topologies (2013)**
- **Industrial-strength fabric management (2015)**

- **Robust telco-grade hardware platforms – FPGA**
- **Vertical technology integration**
- **Commitment to open source, open standards, interoperability**



OPENFABRICS
ALLIANCE

INFINICORTEX & INFINICLOUD

INFINICORTEX

A global network for beneficially aggregating data sources, HPC, storage and analysis

**“A Galaxy of Supercomputers” was initial motivation at A*STAR:
(Agency for Science Technology and Research - Singapore)**



Marek Michalewicz, Ph.D

CEO: A*STAR Computational Resource Centre



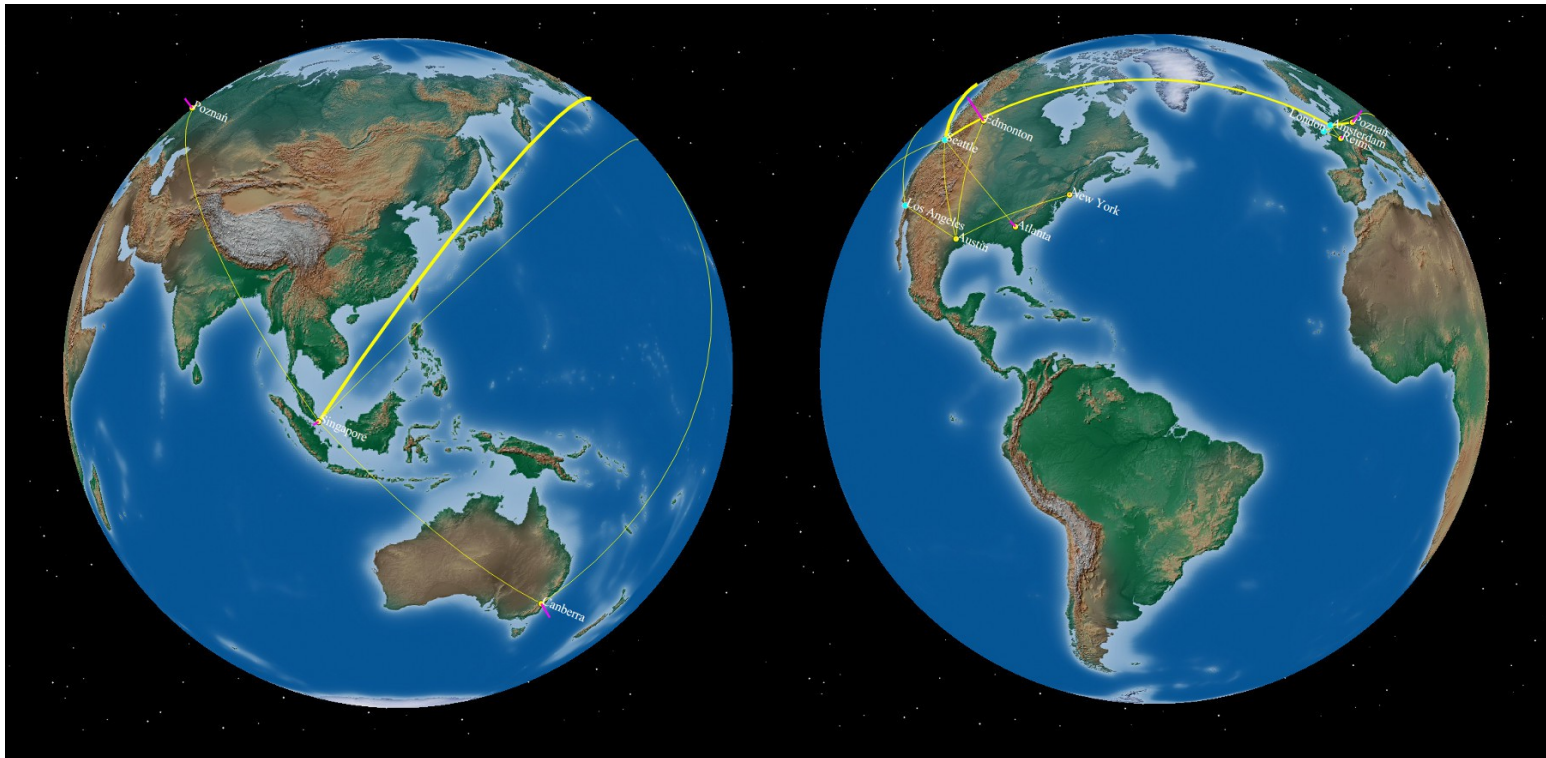
Tin Wee TAN, Ph.D

Chairman: A*STAR Computational Resource Centre
Director: National SuperComputing Centre (NSCC)

INFINICORTEX

A global network for beneficially aggregating data sources, HPC, storage and analysis

- **Concept developed by A*STAR CRC, Singapore**
- **An infrastructure for novel HPC workflows**
- **Entirely InfiniBand based**



INFINICORTEX

A global network for beneficially aggregating data sources, HPC, storage and analysis

Team Singapore:



Team Europe:



Team Canada:



Team Australia:

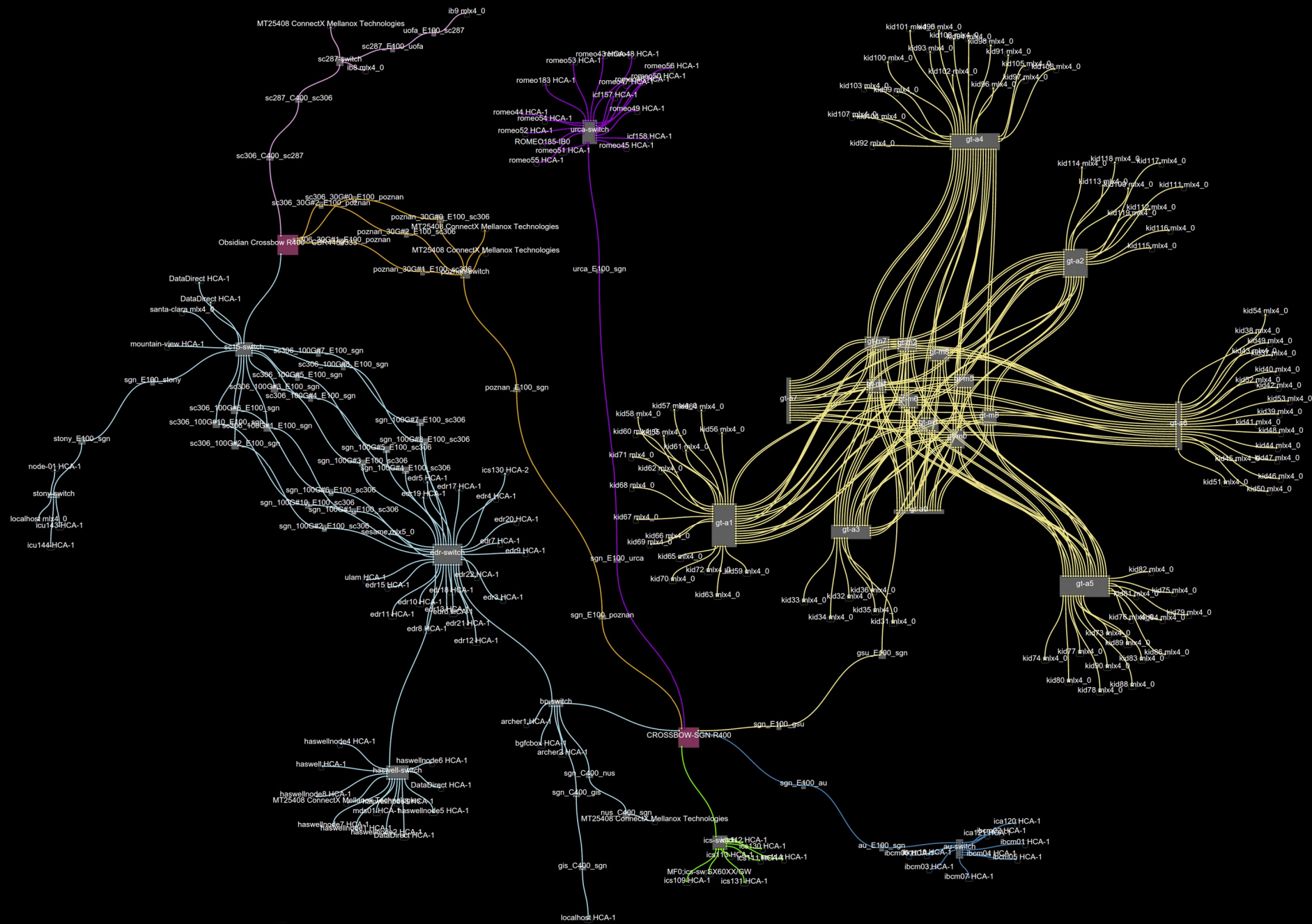


- **Bandwidth efficient global storage migration (asynchronous)**
- **Low latency metro area storage mirroring (synchronous)**
- **Cluster aggregation (MPI)**
- **Direct connect to streaming data sources (sequencers, physics...)**
- **Project remote high fidelity interactive visualisation**
- **Globally dispersed HPC stream processing (geo-pipelining)**

INFINICORTEX

A global network for beneficially aggregating data sources, HPC, storage and analysis





Edmonton



Reims



Austin



Poznań



Atlanta



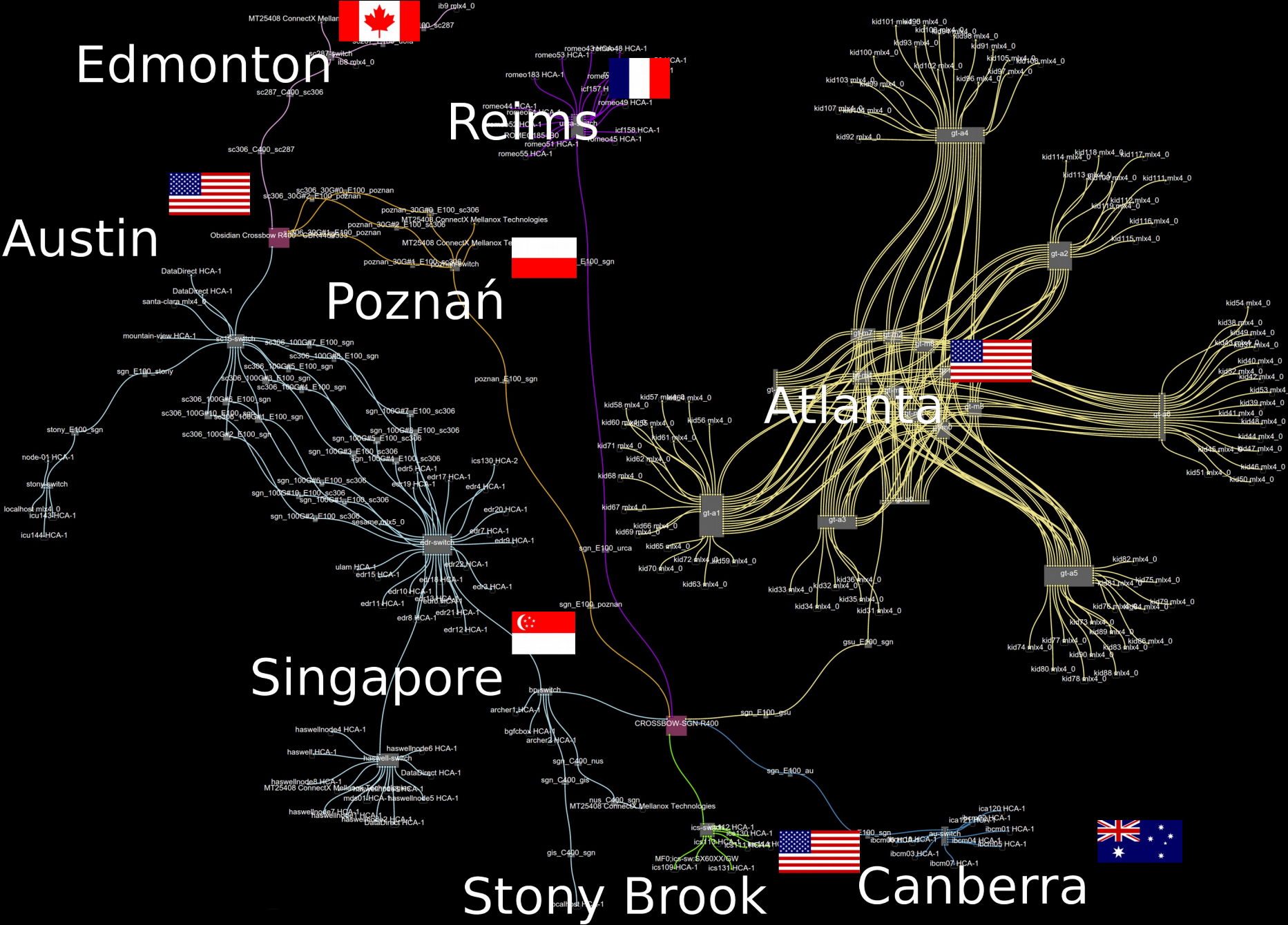
Singapore



Stony Brook

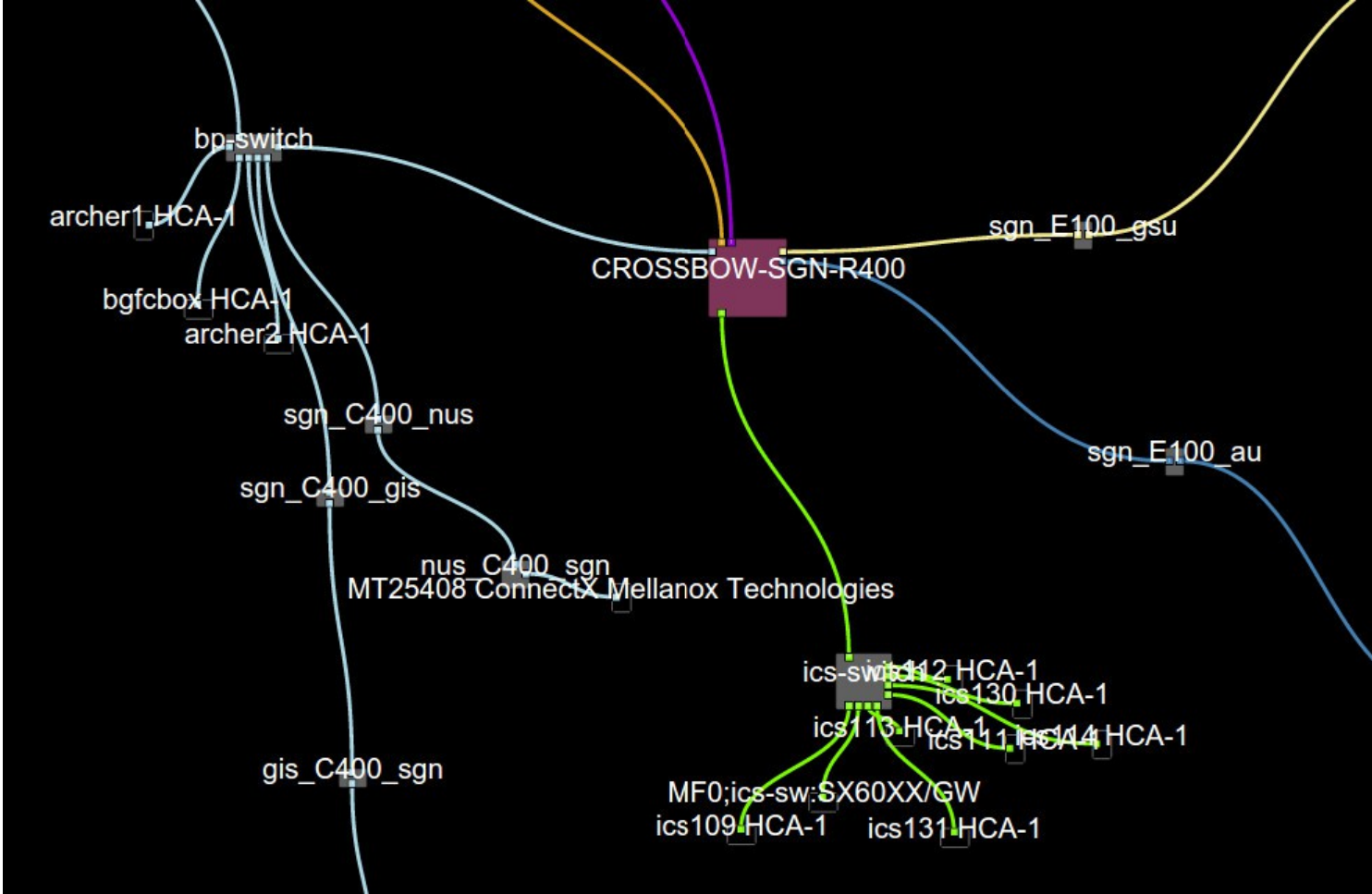


Canberra



INFINICORTEX

A global network for beneficially aggregating data sources, HPC, storage and analysis



INFINICORTEX

A global network for beneficially aggregating data sources, HPC, storage and analysis

Management Functions:

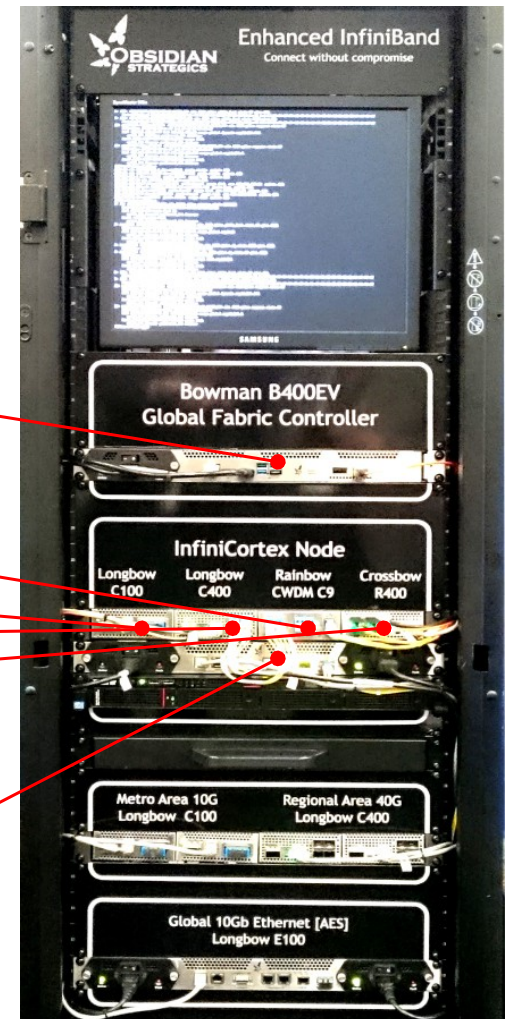
- Bowman B400

Interconnect Functions:

- Optical Processing (Rainbow C9)
- Metro Range (Longbow C100)
- Regional Range (Longbow C400)
- Routing (Crossbow R400-6)

Security Functions:

- Secure Global Range (Longbow E100)





OPENFABRICS
ALLIANCE

INFINICORTEX APPLICATIONS

INFINICORTEX APPLICATIONS

Enabling broad InfiniBand adoption through expanded capabilities



InfiniCortex Demonstrations at Supercomputing 2016

Applications

- High-speed file transfers among PSNC, ROMEO HPC Center, and A*CRC Fusionopolis data centre
- University of Lille: Asynchronous linear solvers running between Singapore and ROMEO HPC Center
- ICM & A*CRC: distributed weather models
- Remote visualization in Singapore of CFD application in Poznań
- Distributed bio-informatics
- University of Alberta: distributed genomic computations and visualization

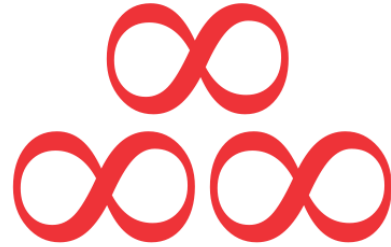
Connectivity

- Singapore — SingAREN, TIEN, Yata, Internet2, ESNET
- Poland — Górnik, Géant, AN200, Internet2
- Canada — Canarie, Cybersix
- Australia — AARNet
- USA — ESNET, Internet2/CAIR

Participants

- A*STAR CRC, Singapore
- ROMEO HPC Center
- University of Reims
- University of Lille
- University of Alberta
- NCI, Canberra
- GA Tech/ORNL
- Stonybrook University
- Poznań Supercomputing and Networking Center

INFINICORTEX APPLICATIONS



InfiniCloud

High performance cloud computing platform
for data-intensive workflows

Kenneth Ban

Dept of Biochemistry, NUS
& IMCB, A*STAR

Jakub Chrzesczyk

National Computational
Infrastructure, ANU

INFINICORTEX APPLICATIONS

Breaking out of the traditional HPC mold



Specialized applications
High performance



Flexible
Virtualization overhead



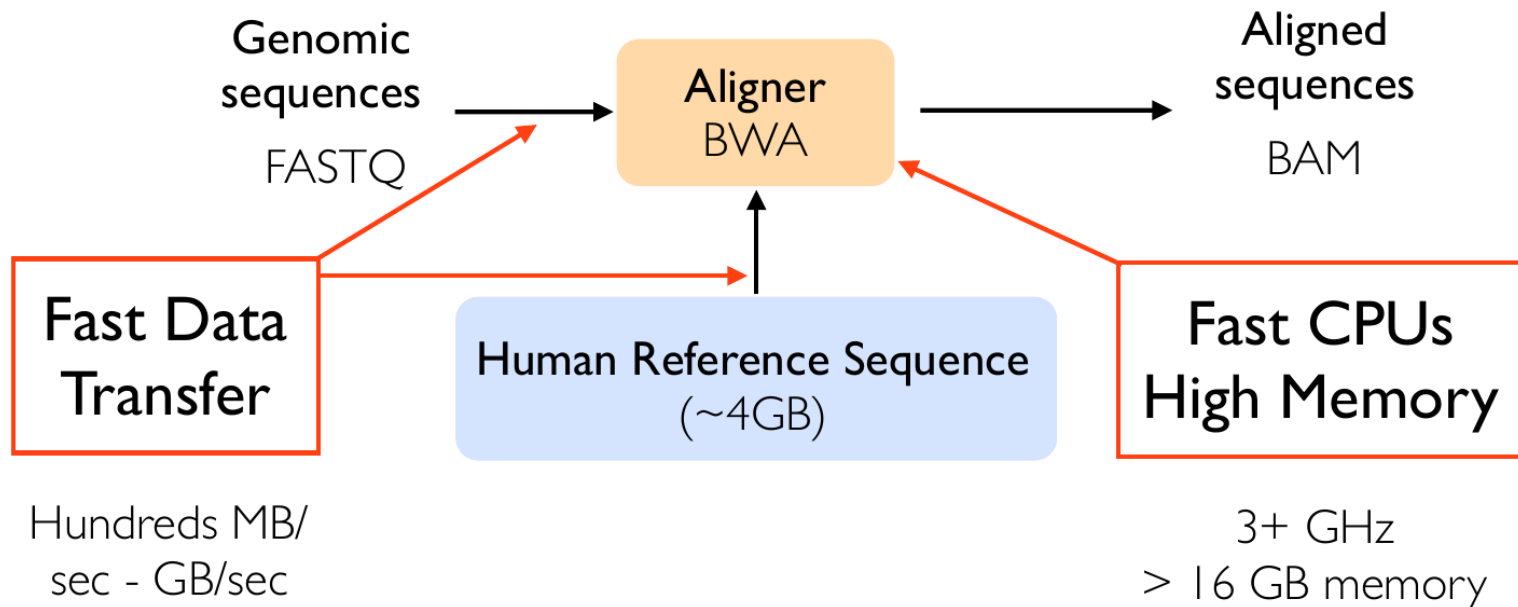
InfiniCloud



INFINICORTEX APPLICATIONS

What components do we need?

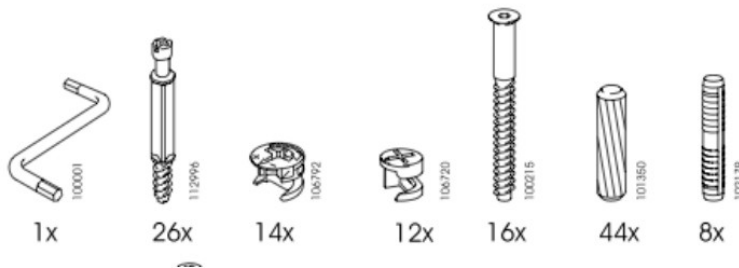
- **High CPU/memory and network performance** for rapid analysis of large datasets



INFINICORTEX APPLICATIONS

What components do we need?

- **Reproducible** and well documented **workflows** that can be run on different hardware platforms



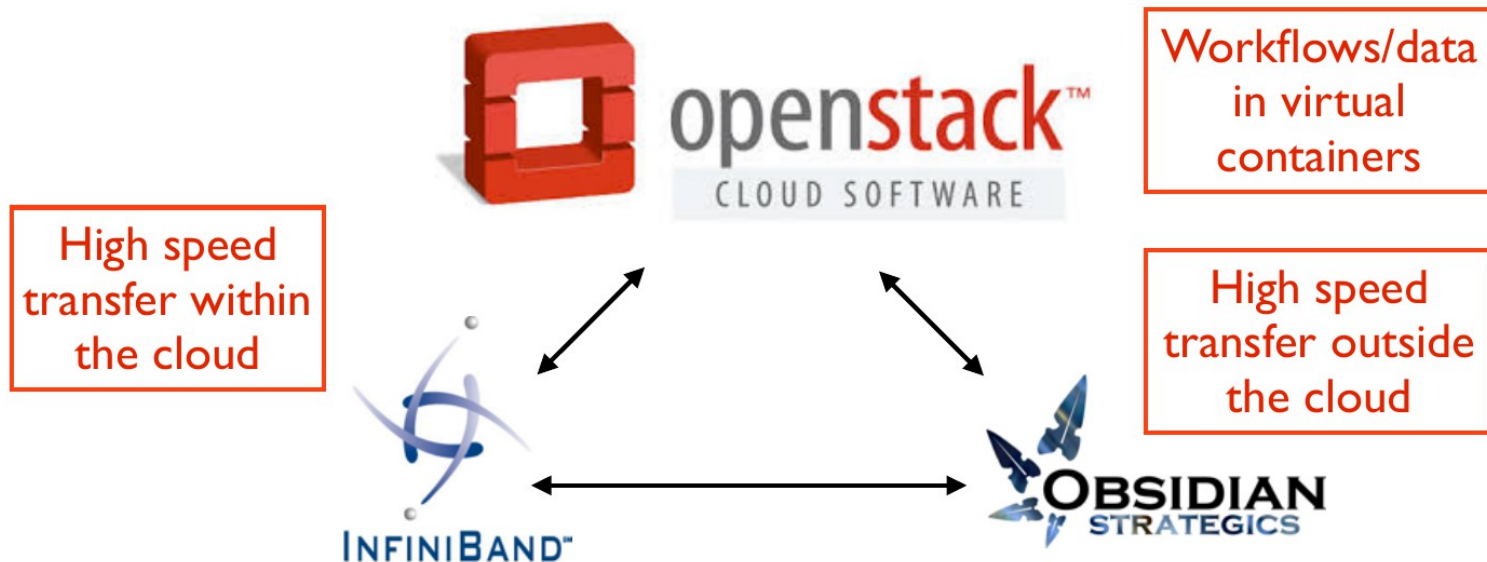
There are many parts
(and different versions)
in an analytical pipeline



Fitting them
together properly
can be challenging

INFINICORTEX APPLICATIONS

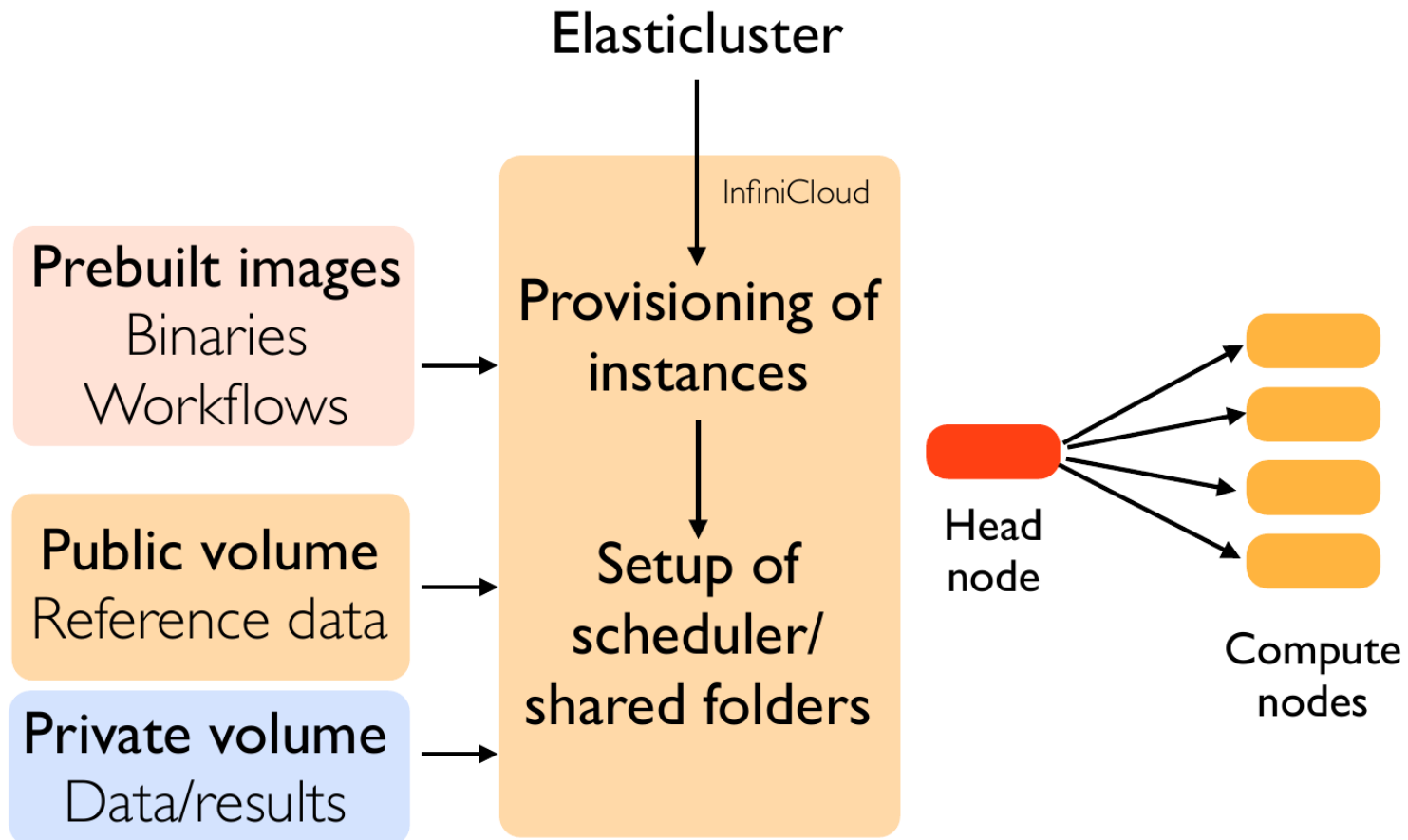
InfiniCloud: a flexible high performance cloud computing platform



- **Cloud** infrastructure for **flexible** computing
- High speed/low latency **Infiniband** interconnect
- **Long-range** Infiniband (**global** reach)

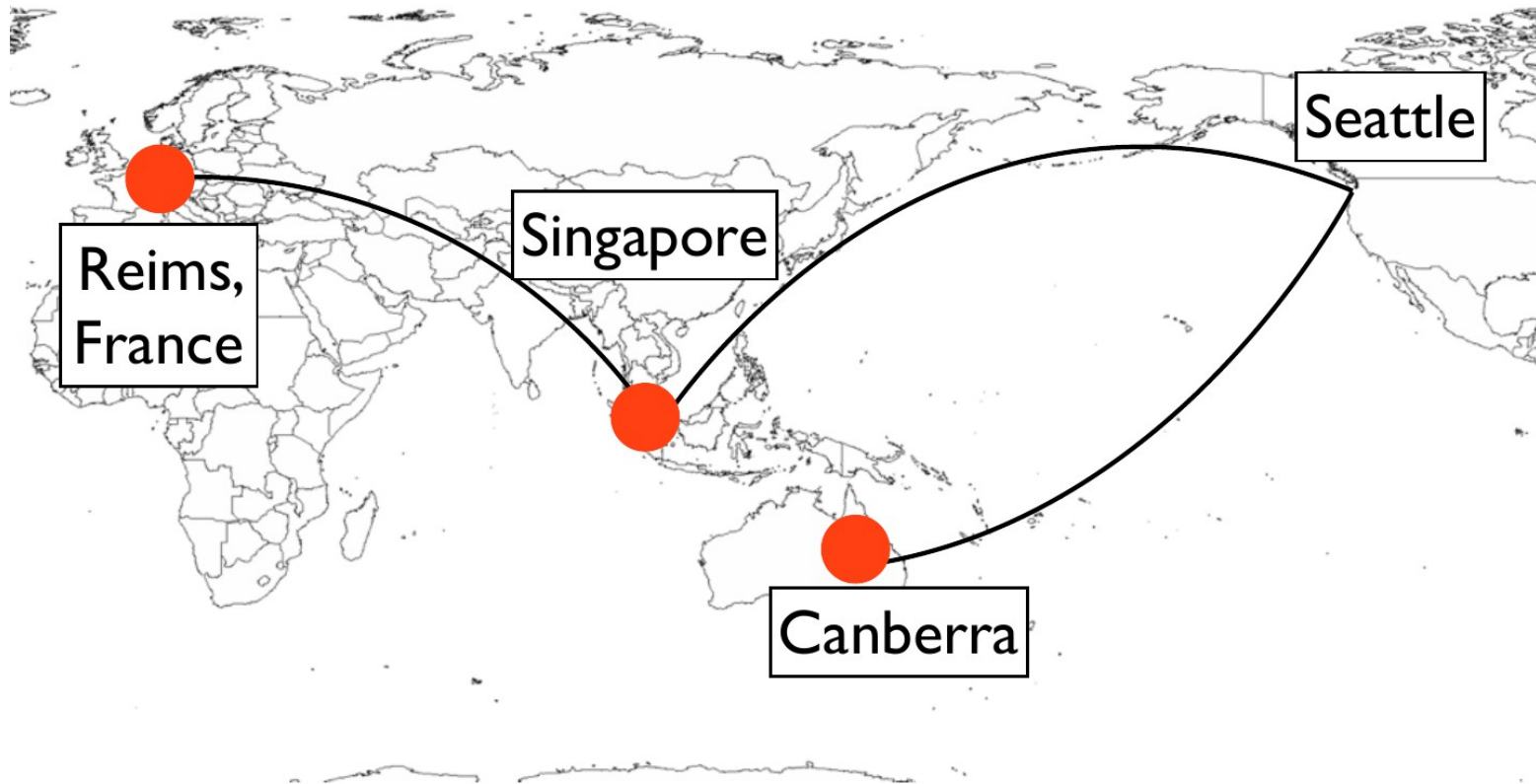
INFINICORTEX APPLICATIONS

On-the-fly provisioning of virtual clusters for distributed computing



INFINICORTEX APPLICATIONS

A geo-distributed virtual cluster
connected by long range InfiniBand



INFINICORTEX APPLICATIONS

Setup of geo-distributed virtual cluster

InfiniCloud™

[Home](#)

[Resources](#)

[Status](#)

[Images](#)

[Instances](#)

[Log In](#) ▾

Status of Virtual Machines

Show entries

Search:

Instance Name	IP address	Zone	Status
geopipeline-compute001	10.2.1.95	singapore	ACTIVE
geopipeline-compute002	10.2.1.96	australia	ACTIVE
geopipeline-compute003	10.2.1.98	europa	ACTIVE
geopipeline-compute004	10.2.1.97	australia	ACTIVE
geopipeline-frontend001	10.2.1.94	nova	ACTIVE

Showing 1 to 5 of 5 entries (filtered from 14 total entries)

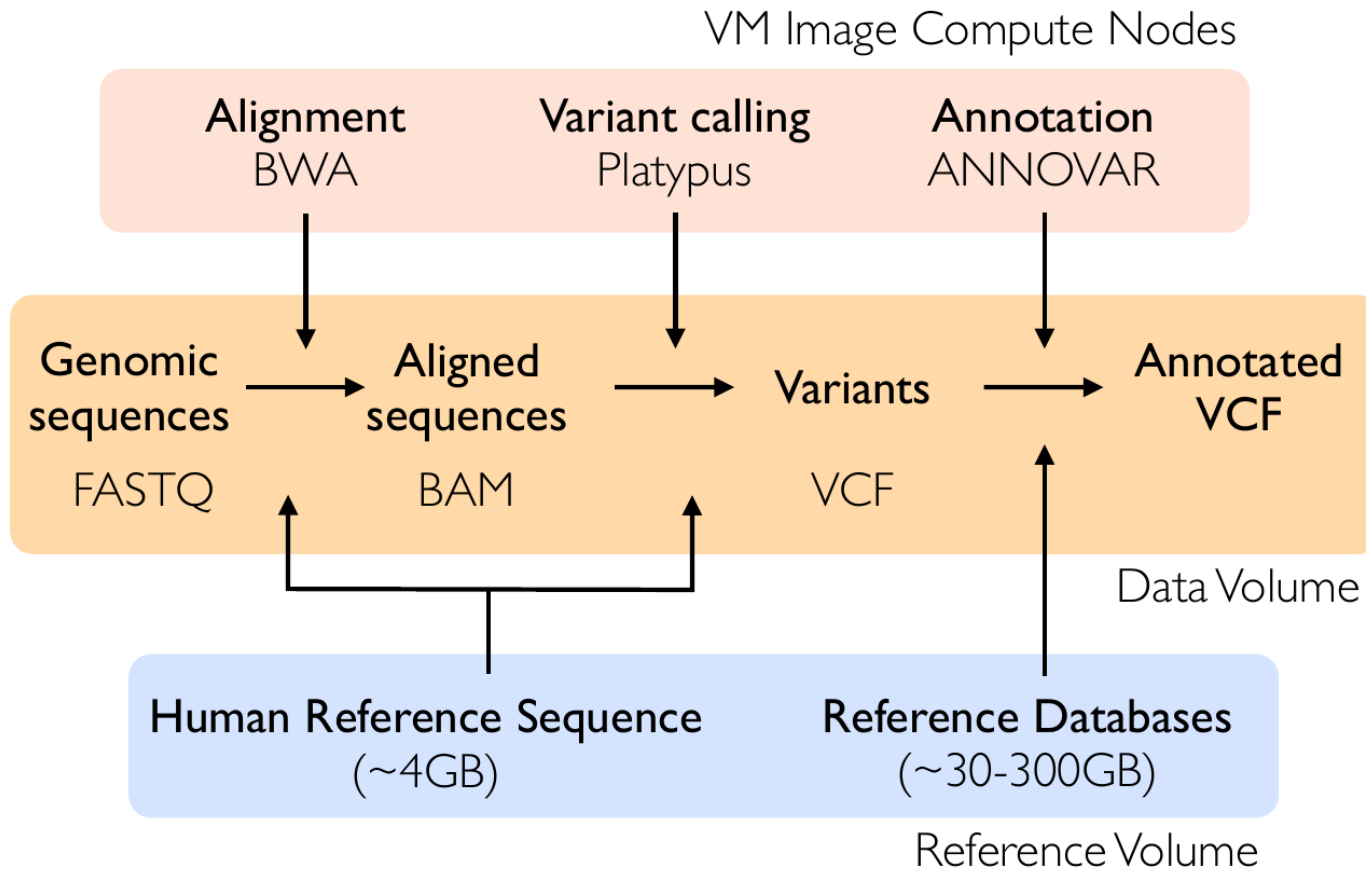
[Previous](#)

[1](#)

[Next](#)

INFINICORTEX APPLICATIONS

Demo pipeline for variant calling

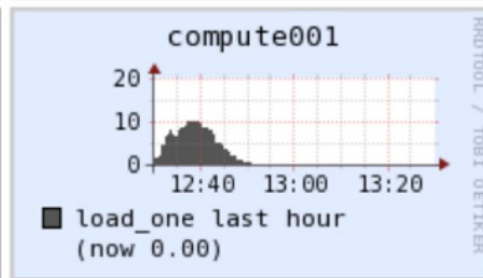
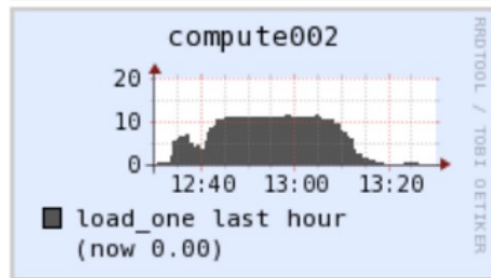
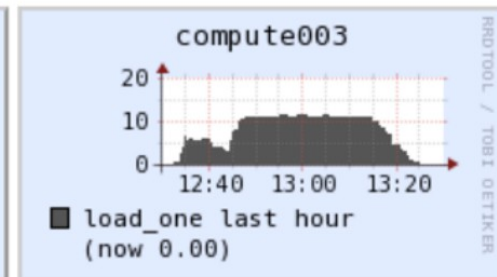
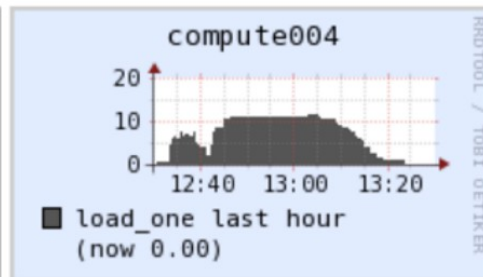
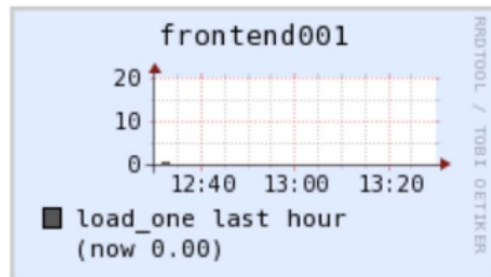


INFINICORTEX APPLICATIONS

Geo-distributed pipeline for identification of mutations in cancer samples

Australia

France



Australia

Singapore

INFINICORTEX APPLICATIONS

Acknowledgements



Computational
Resource Centre



Institute of
Molecular and
Cell Biology



Yong Loo Lin School of Medicine

INFINICORTEX APPLICATIONS

Enabling broad InfiniBand adoption through expanded capabilities

InfiniCloud is interesting for several reasons:

- **An InfiniBand based HPC cloud overlay on InfiniCortex**
- **Virtualised HPC nodes through OpenStack containers (SRIOV)**
- **Spans multiple InfiniBand subnets**
- **All WAN links are hardware AES encrypted (bioinformatics...)**
- **All WAN link end points are hardware authenticated**
- **Implements high performance data flow pipeline across the globe!**
- **Fully automated across heterogeneous metal**
- **Sustained wire speed operation (simultaneous streaming and computation)**



OPENFABRICS
ALLIANCE

THE HEART OF THE FABRIC: BGFC

BGFC

Bowman Global Fabric Controller

Middleware like InfiniCloud requires a fast, transparent, secure, scalable, segmentable, highly stable and manageable fabric.

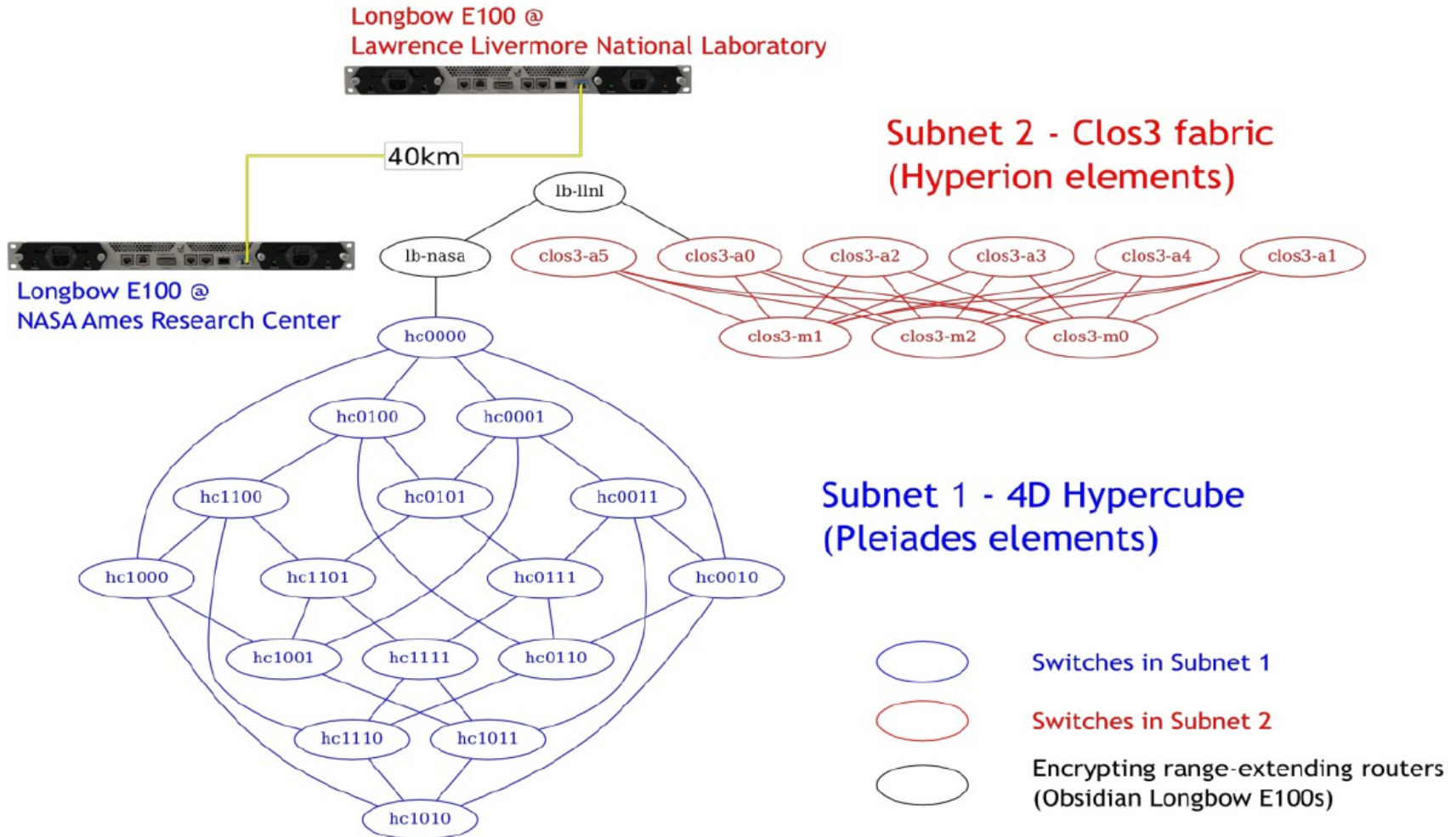
Such a fabric needs a controller that did not exist before BGFC.

Matt Leininger (LLNL) and Bob Ciotti (NASA Ames) approached Obsidian in 2011 describing similar challenges with their LAN, CAN and WAN InfiniBand deployments.

Not seeing OpenSM as a viable platform upon which to build, Obsidian responded with a green field fabric controller architecture that would address immediate requirements and many more besides...

BGFC

Bowman Global Fabric Controller



BGFC

Bowman Global Fabric Controller

Now 5 years into the adventure, Obsidian has completed phase I (LLNL and NASA) and phase II (A*STAR) of the original development program.

Phase III is imminent (TRL-9 testing at scale), and BGFC will be open-sourced thereafter.

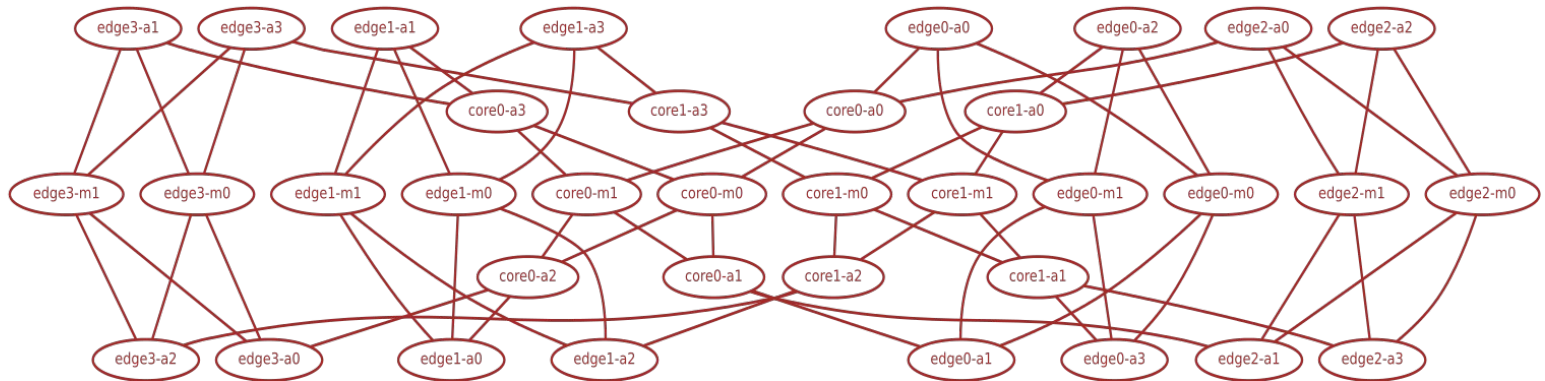
Built for a much grander mission than simple subnet management, BGFC is aimed at complex multi-domain fabrics supporting international traffic while preserving sovereign administrative domains.

BGFC

Bowman Global Fabric Controller

Exact and Mathematically Perfect LFTs

Graph theory-based subgraph isomorphism
Guaranteed deadlock-free routing
Reliably supports very large subnets
Deterministic QoS
Python scripts for topology descriptions



BGFC

Bowman Global Fabric Controller

Python topology prescription examples:

Simple, direct from templates ...

```
from bgfc.template import *  
topology = {"my-topo": Hypercube(9)}
```

Customised from existing templates ...

```
from bgfc.template import *  
class MyNetwork(ClosTree5):  
    def __init__(self):  
        ClosTree5.__init__(self, core=(36,18), cores=2, edge_conns=18//2)  
    def disperse_edges(self):  
        [.. insert site-specific wiring function ..]  
topology = {"my-topo": MyNetwork()}
```

, or arbitrarily complex by deriving from the `Topology` class.

BGFC

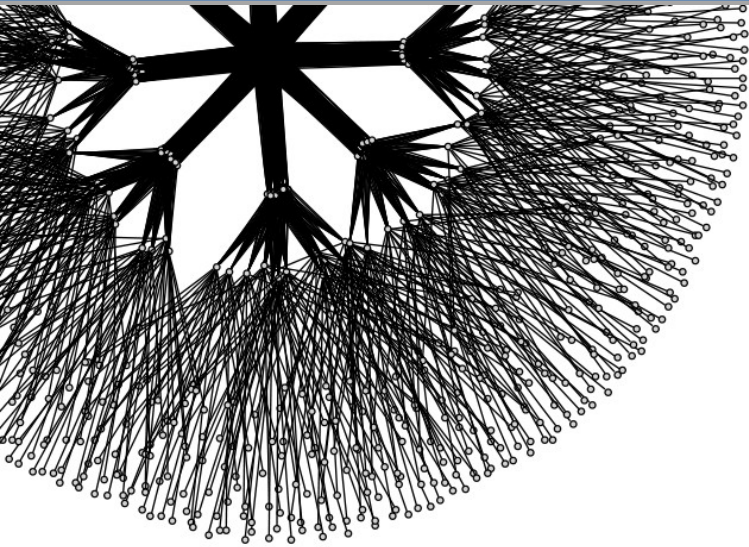
Bowman Global Fabric Controller

Precise topology definitions allow simple mathematical routing functions, but also provide a powerful means of detecting unintended deviations.

BGFC uses persistent topology and LFT solution databases to make it easy to ensure a fabric initialises the same way if required, increasing reproducibility of fabric behaviour and drastically reducing initialisation time.

BGFC

Bowman Global Fabric Controller



Core routing concept:

Optimal dead lock free IB routing is a NP problem with complexity related to the number of buffers, and no apparent easy shortcuts.

Subgraph isomorphism is a NP problem with complexity related to the number of switches, with well known shortcuts.

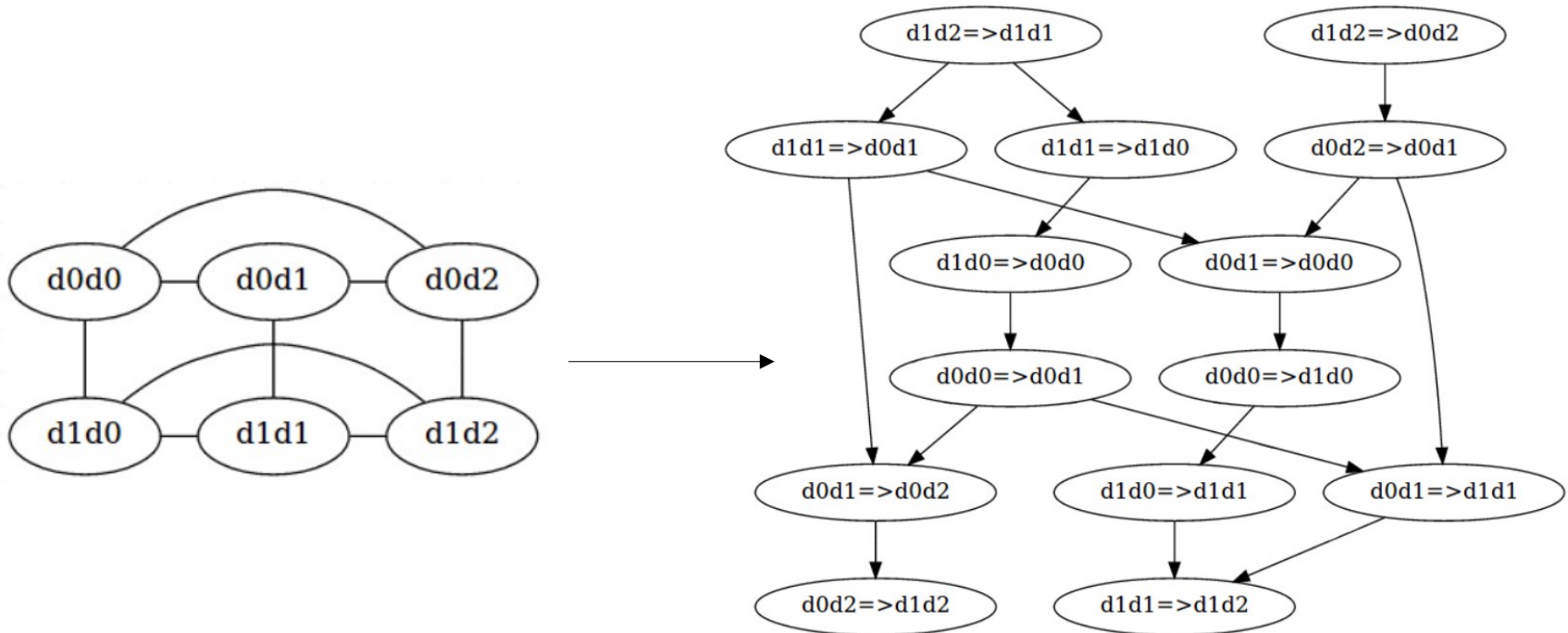
So, solve subgraph isomorphism and then get routing, rather than try to solve routing directly!

BGFC: 624 lines of Python
OpenSM: 17,200 lines of C

BGFC

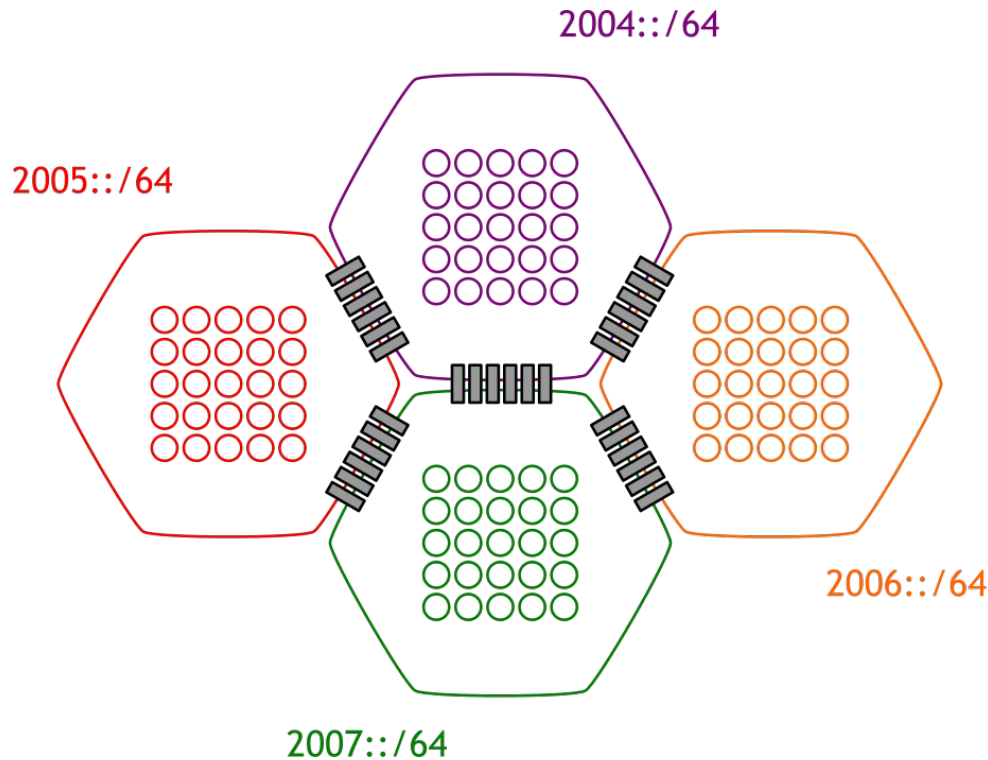
Bowman Global Fabric Controller

Example; a torus(2,3) and its acyclic flow group graph:



BGFC

Bowman Global Fabric Controller



Multi-subnet Native InfiniBand Routing

- Administrative demarcation
- Fault isolation
- Performance at scale
- Complex topologies
- Inter-site separation

BGFC

Bowman Global Fabric Controller

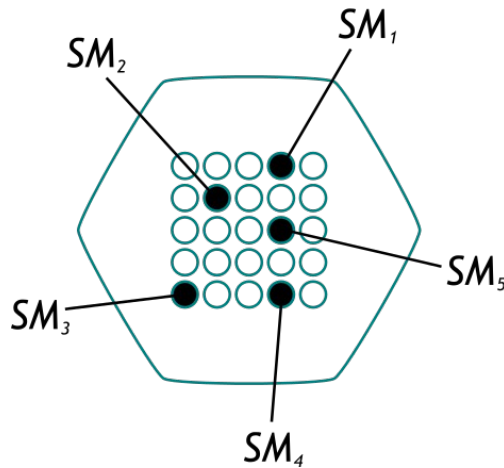
Why is it called an InfiniBand 'subnet'...

...if there is only one in your cluster?

BGFC

Bowman Global Fabric Controller

N-way Clustered Active SMs



- Parallel host-based subnet managers
- High performance at scale
- Lockless WODB architecture
- Decentralised control
- Extreme fault tolerance
- C++11 implementation

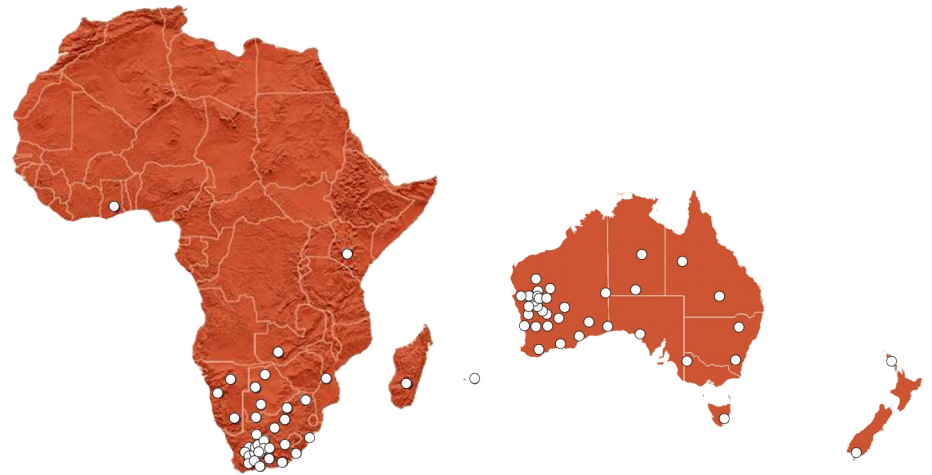
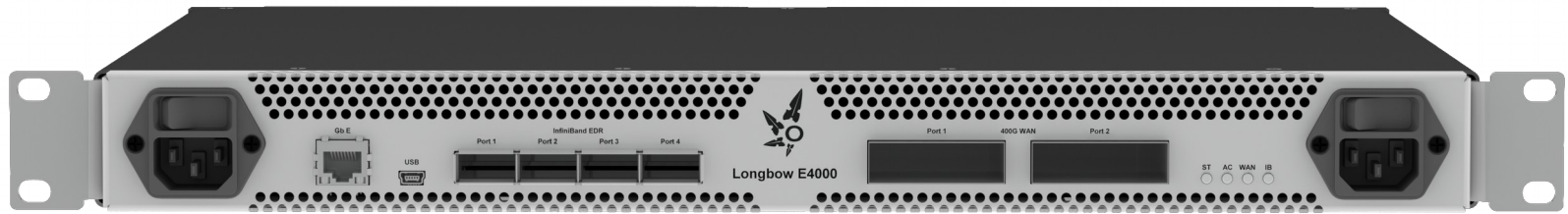


OPENFABRICS
ALLIANCE

FUTURES...

FUTURES

100 and 400Gbits/s Thresholds





OPENFABRICS
ALLIANCE

12th ANNUAL WORKSHOP 2016

THANK YOU

David Southwell, CVO

Obsidian Strategics Inc.

