



Future Infrastructure for Data-Intensive Science

David Schade
Canadian Astronomy Data Centre
National Research Council Canada & University of
Victoria

November 22, 2017 Vienna



National Research
Council Canada

Conseil national de
recherches Canada

Canada 

Missing recommendation

The UN recognizes that governments have invested hundreds of millions of dollars to create the present-day network of astronomy data services. The powerful science capabilities provided by this network are the foundation upon which Open Universe will operate.

These investments must continue and increase in order to support the types of services proposed by the Open Universe.

Fresh new ideas and approaches

- Few new ideas in the Open Universe initiative
- The new factor is the involvement of the UN

A fresh new approach would be:

A substantial transfer of the benefits of astronomy data and supporting infrastructure to the public through education, outreach, and citizen science with a focus on developing nations as the highest priority.

Scale of CADC 2016

CADC

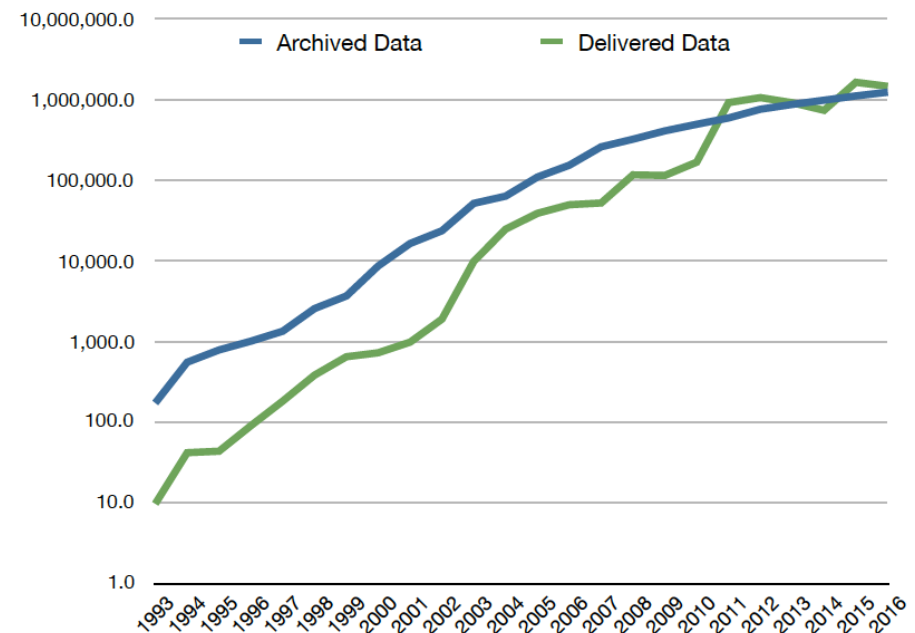
- was created in 1996 and parallels Hubble Space Telescope
- has 21 staff: scientists, programmers, operations
- 1 billion files
- 2.6 Petabytes

Data flows

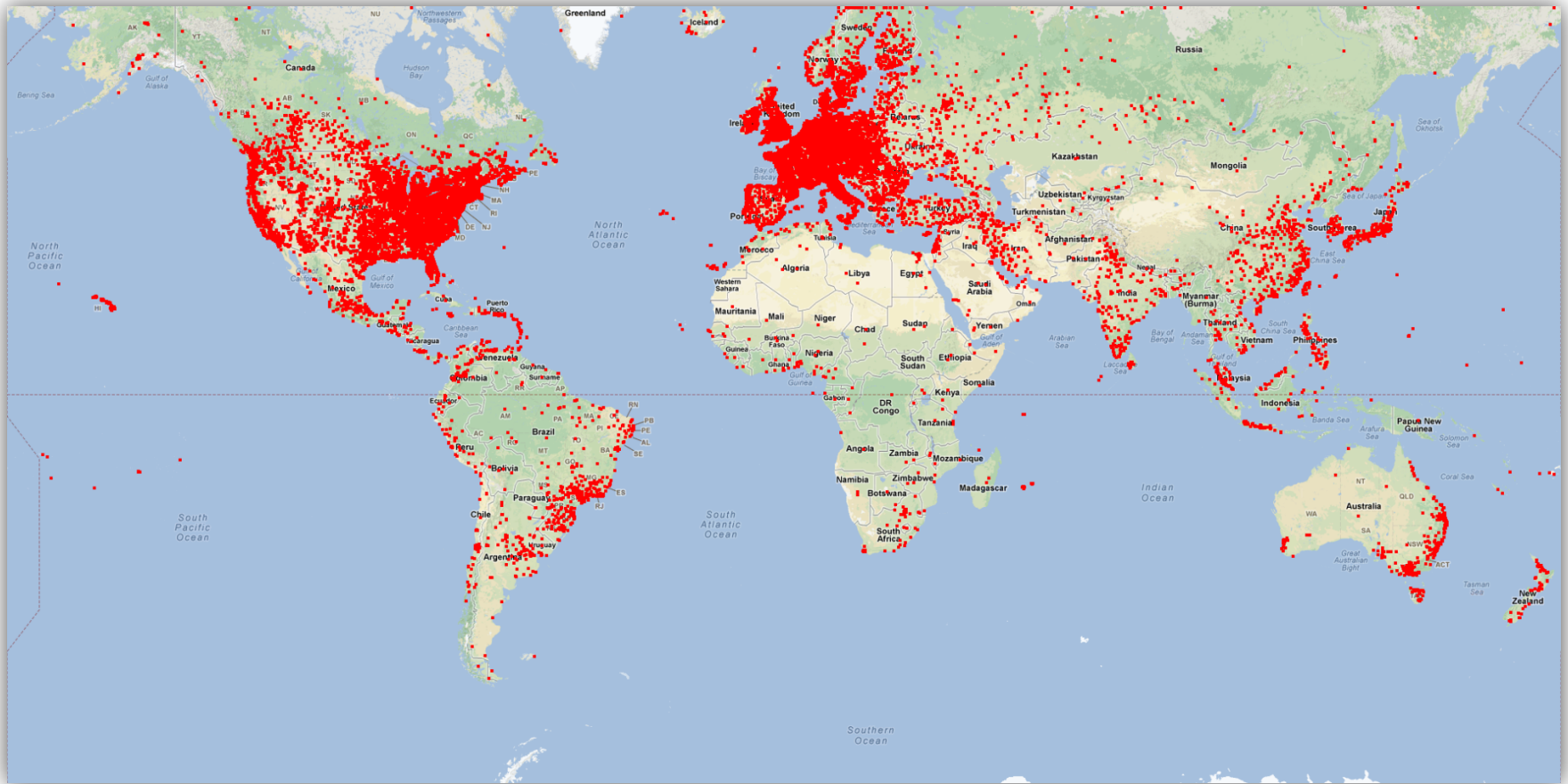
- 1.4 Petabytes of data out
 - 75 million individual calls
- 300 Terabytes put back into CADC system
 - 15 million calls

Processing

- 3,671,737 jobs in batch mode
- 387 interactive Virtual Machines
- 460 core years of processing used



CADC data delivery



The Future: Two Themes

Integration of data resources

- Integration within data centres
- Integration across data centres

Integration of data with computing infrastructure

- Integration within Canada
- Integration internationally



Advanced Search

Search Results Error ADQL Help

Search

Reset

Integration of data from 115 instruments

Click on ? for explanations

Observation Constraints

Observation ID ?
P.I. Name ?
Proposal ID ?
Proposal Title ?
Proposal Keywords ?
Data Release Date ?

Science and Calibration data

Spatial Constraints

Target ?
Pixel Scale ?
☐ Do Spatial Cutout

Temporal Constraints

Observation Date ?
Integration Time ?
Time Span ?

Spectral Constraints

Spectral Coverage ?
Spectral Sampling ?
Resolving Power ?
Bandpass Width ?
Rest-frame Energy ?
☐ Do Spectral Cutout

Additional Constraints

Band	Collection	Instrument	Filter	Calibration Level	Data Type	Observation Type
All (7) Gamma-ray Infrared Millimeter Optical Radio UV Unknown	All (21) CFHT CFHTMEGAPIPE CFHTTERAPIX CFHTWIRWOLF HST HSTHLA GEMINI JCMT JCMTLS DAO DAOPLATES	All (115) ACS Apogee USB/Net COS CPAPIR Cassegrain Spectrograph Cassegrain Spectropolarime Direct image ESPaDOnS F2 FTS2-SCUBA-2 Fabry image	All (2152) 0.35MB 0.35um 0.45MB 0.45um 0.75um 0.85um 1.083 um 1.210 um 1.282 um 1.3um 1.4um	All (5) (3) Product (2) Calibrated (1) Raw Standard (0) Raw Instrumental Unknown	All (6) catalog cube image Other spectrum timeseries	All (57) ACQUIRE ALIGN ARC ASTAR BIAS CAL CALIB COMPARISON DARK DIM DOME_FLAT

Search

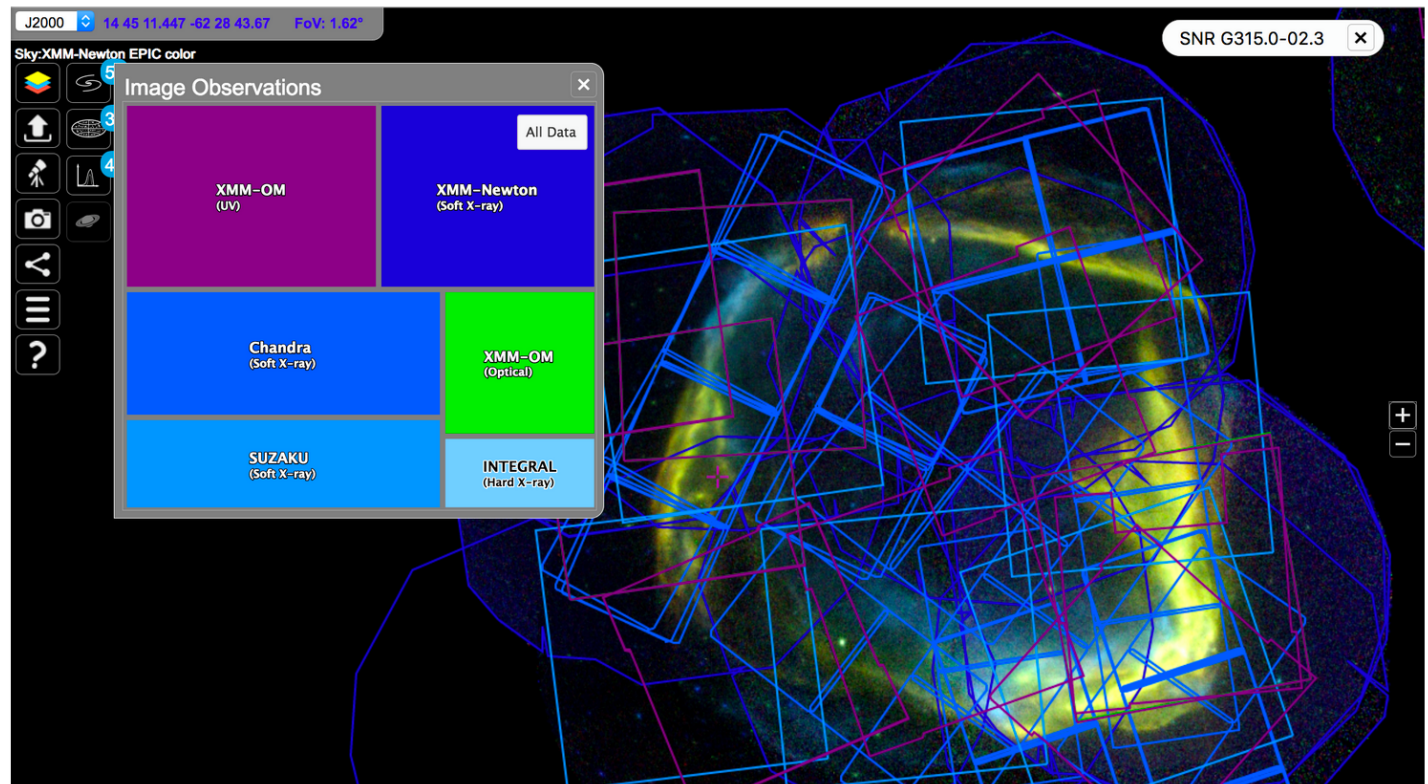
Reset

ESASKY

ESDC » ESASky » ESASky Help

ESASKY HELP

Home
About ESDC
Science Archives
Archive Image Browser
ESASky
Videos
Scientific Tutorials
Publications
VOSpec
Euro-VO Registry
Contact Us



University
of Victoria



University of
British Columbia

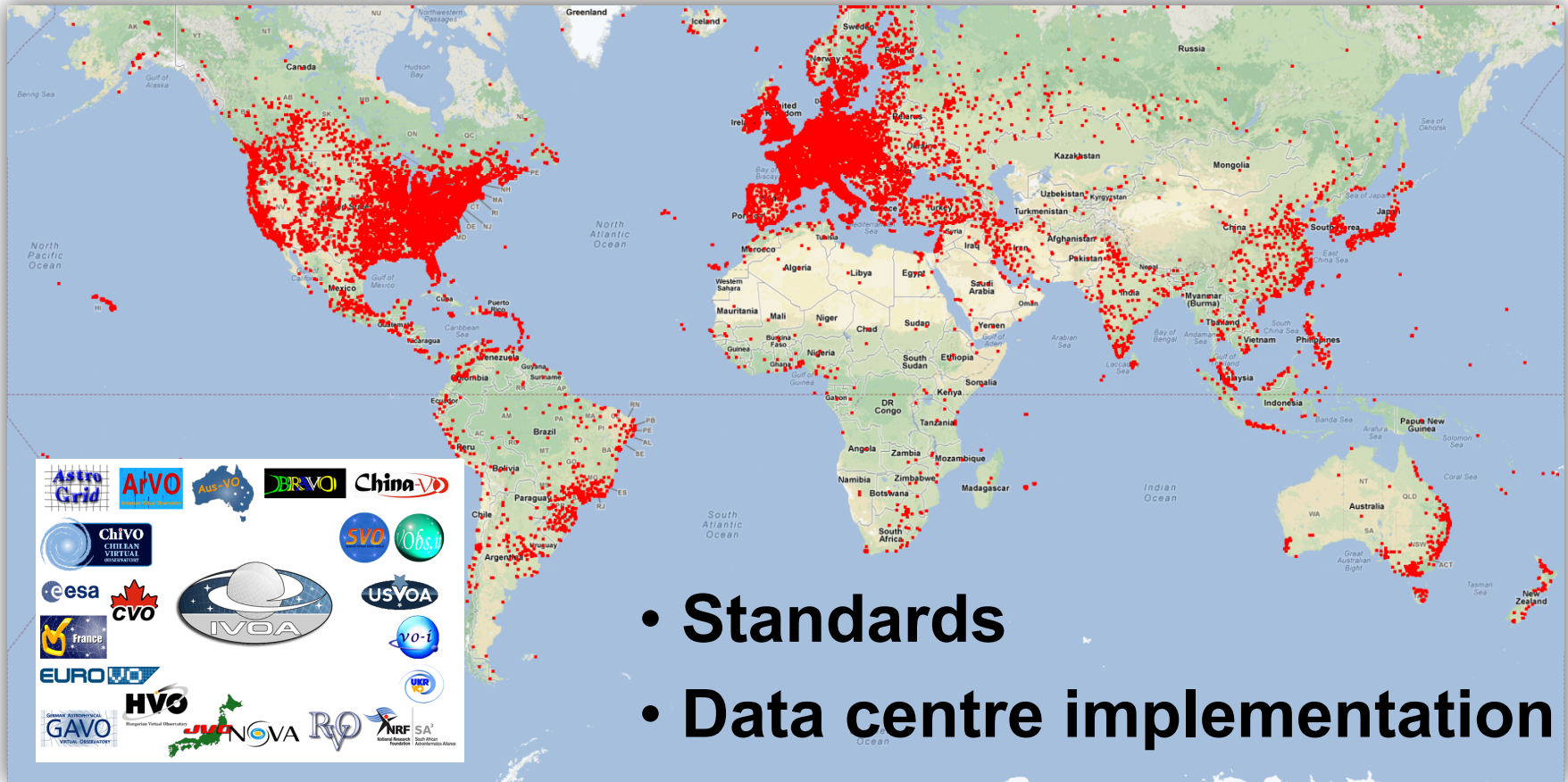
canarie



compute + calcul
CANADA



International Data Integration



Two Themes

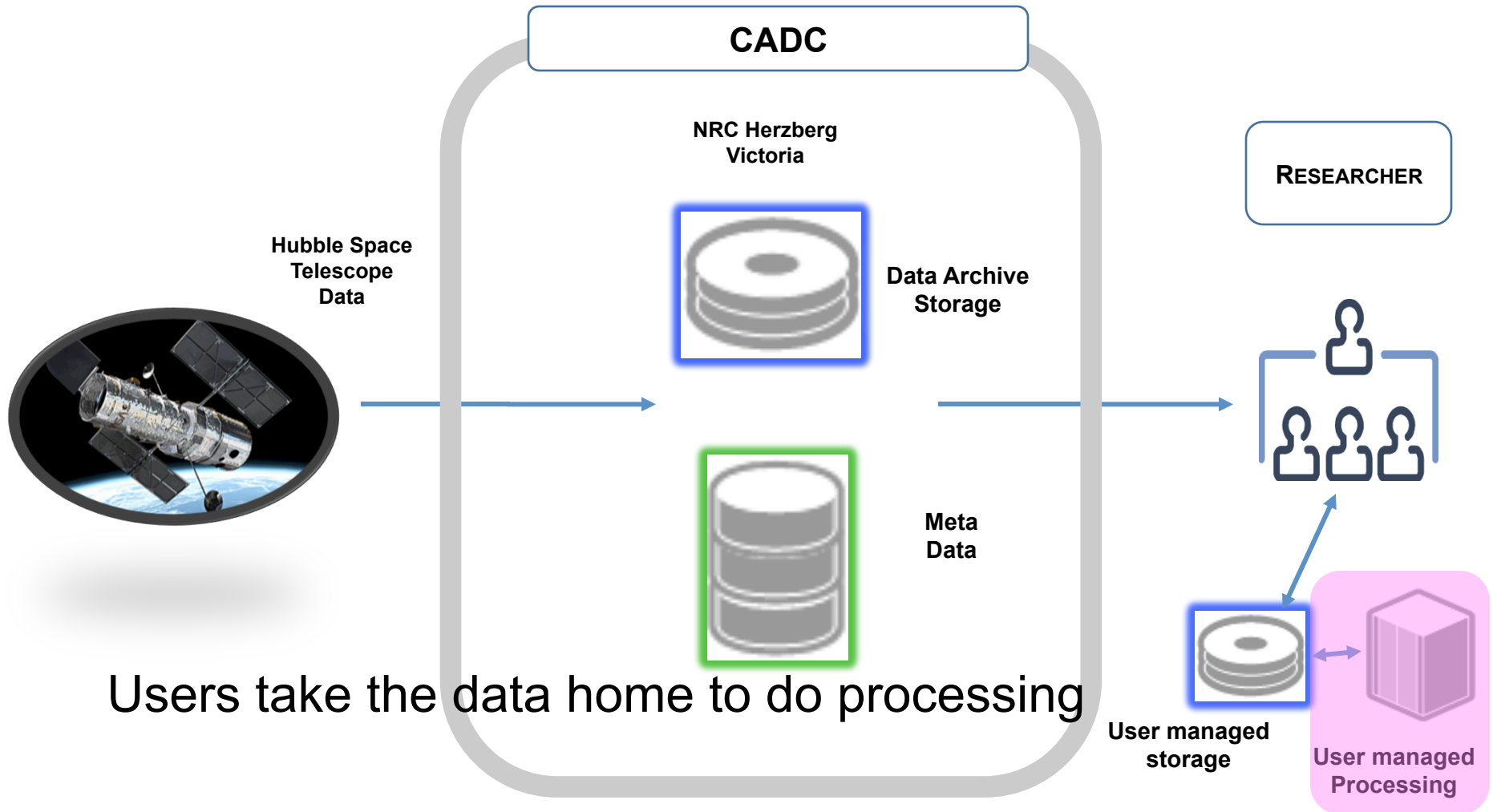
Integration of data with computing infrastructure

- Integration within Canada
- Integration internationally

Driven by:

- Large data volumes
- Government funding policy
- Science practice

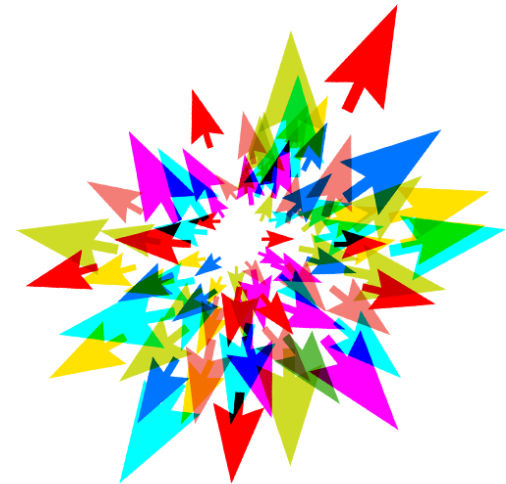
Past practice



CADC operates an integrated system of resources

- A cloud ecosystem for data intensive astronomy
- User services
 - Store and share data
 - Create and configure VMs
 - Run interactive VMs
 - Run persistent VMs
 - Batch processing with VMs
- Using research cloud resources
 - Compute Canada
 - CADC
- Integrated authentication and authorization

compute | **calcul**
canada | canada



University
of Victoria



University of
British Columbia

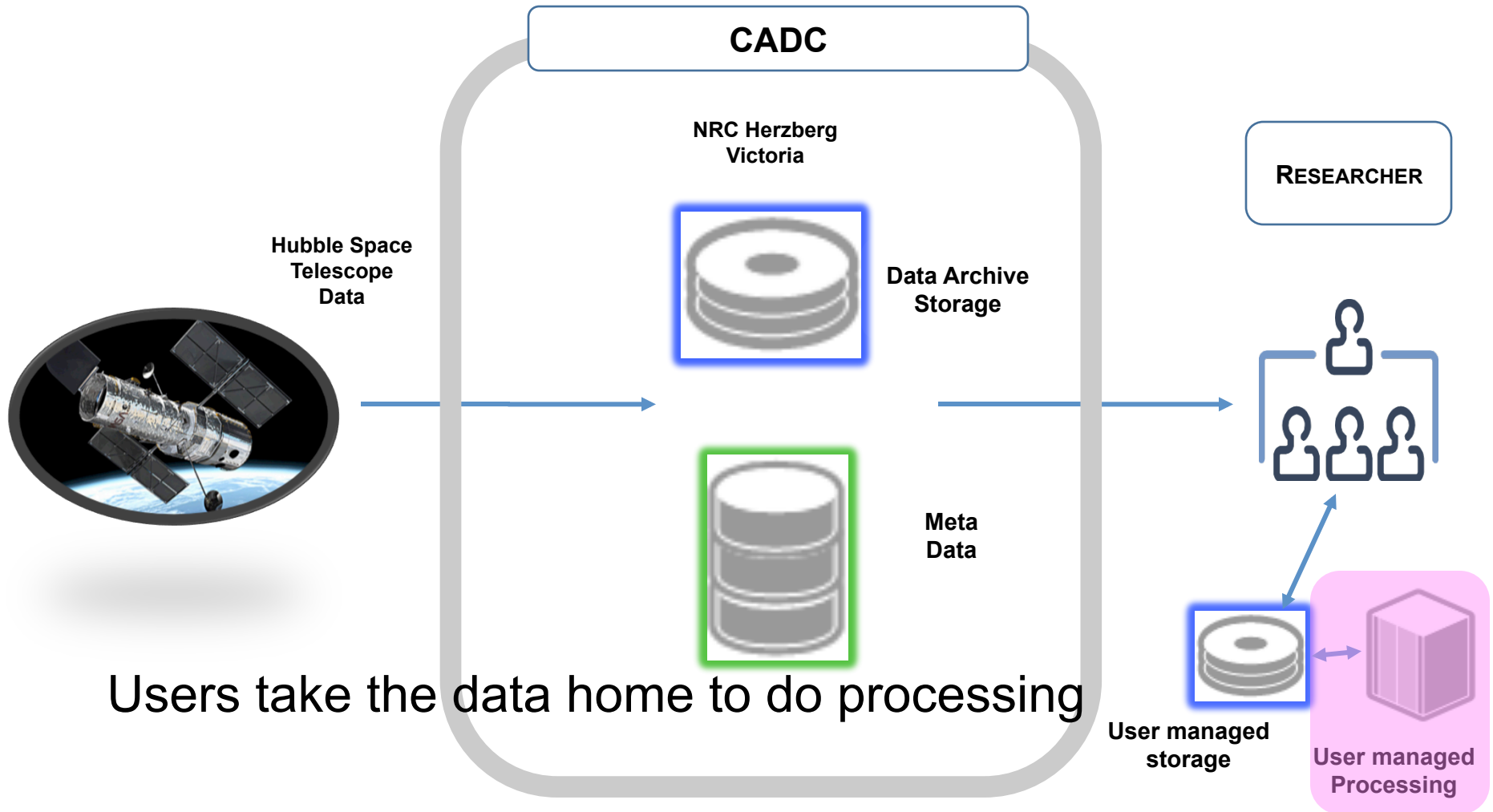
canarie



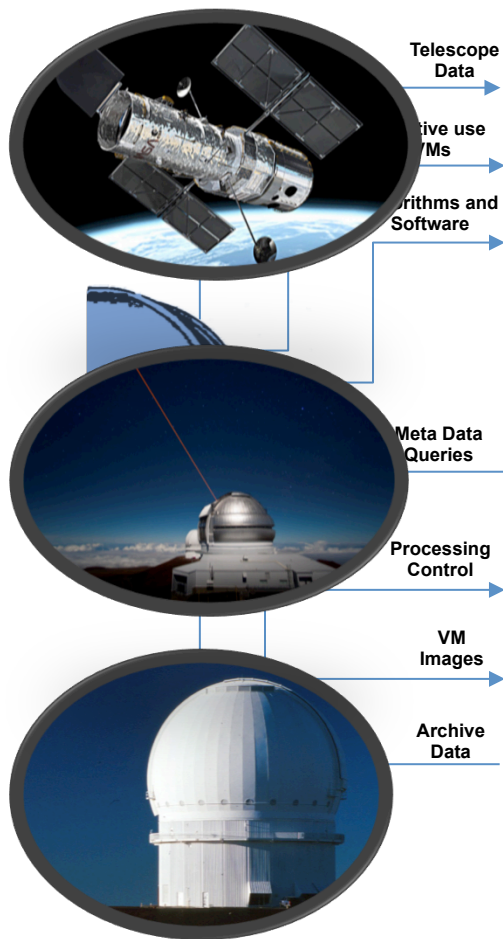
compute • calcul
CANADA



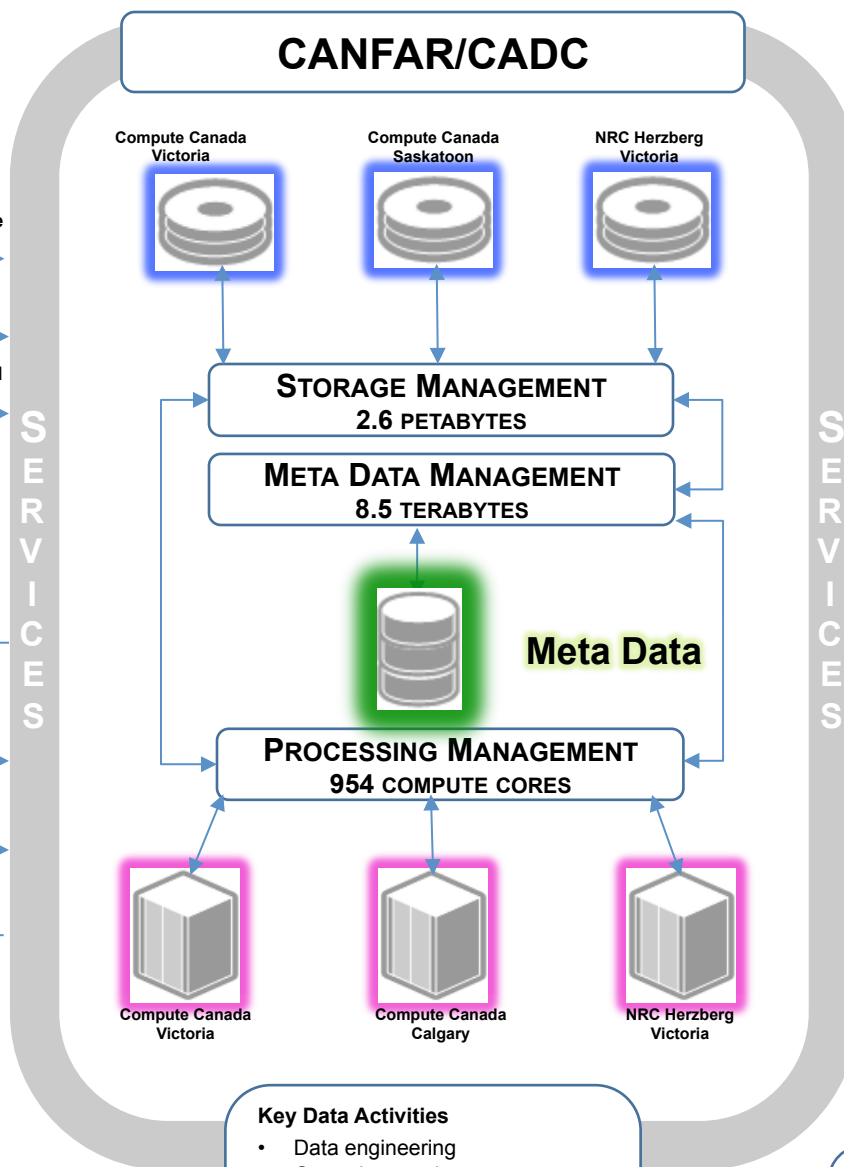
Past practice



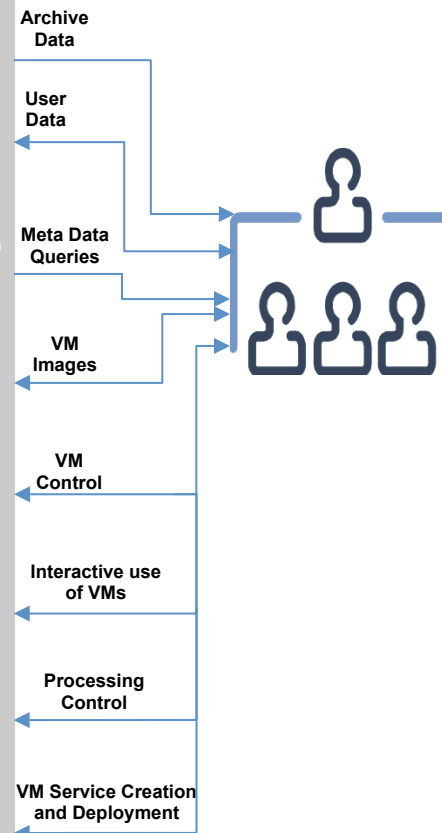
TELESCOPE CLIENT



CANFAR/CADC



UNIVERSITY RESEARCHER CLIENT



Key Data Activities

- Data engineering
- Operations and user support
- Software development
- Software integration
- Data processing
- Data management
- User web services
- User web interfaces

University researchers and telescope staff have privileges to upload data, create VMs and install science applications, run interactive VM sessions, submit batch processing jobs to VMs, share their VMs, control the life-cycle for their VMs, offer software-as-a-service applications in their VMs.

Definition: VM – Virtual Machine

	Data In		Data Out	
	# of files	Terabytes	# of files	Terabytes
Peak per day	2,169,190	8.0	648,093	16.8
Avg per day	130,952	0.4	99,253	2.6

CADC's role has changed radically

We were:

- Managers/curators/distributors of data collections

We are now:

- Managers of an an integrated system of services for data-intensive astronomy



University
of Victoria



University of
British Columbia

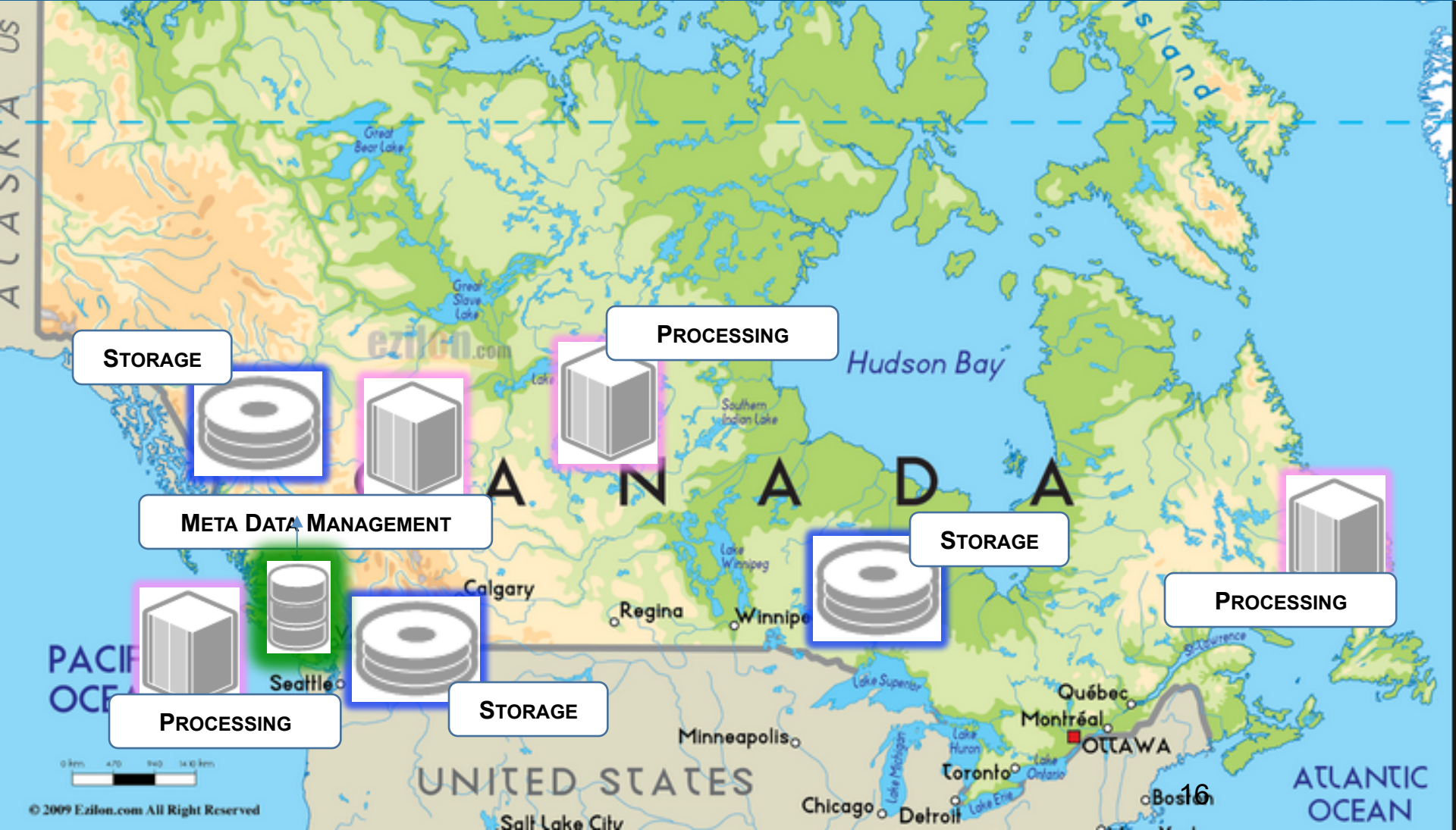
canarie



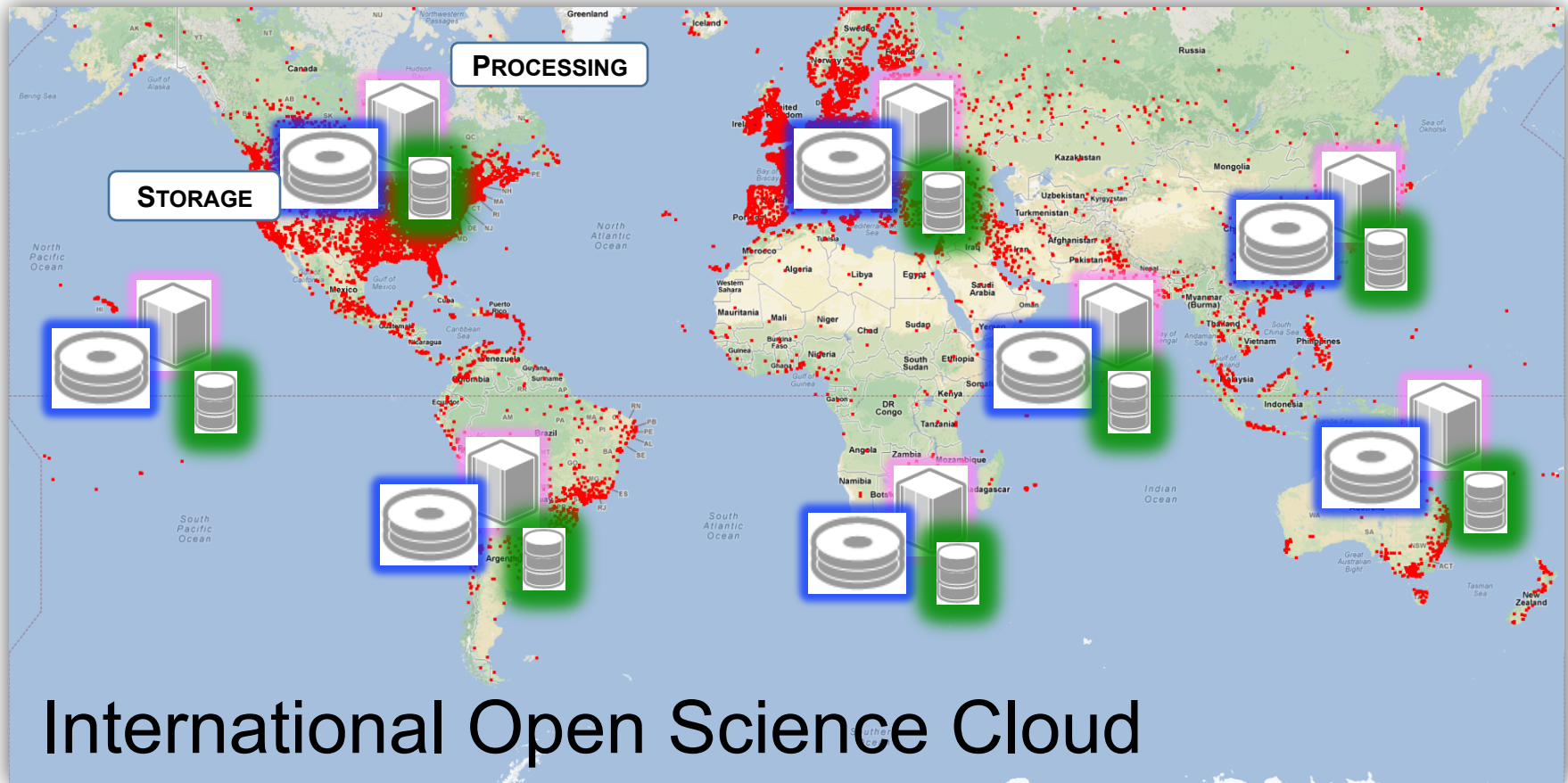
compute • calcul
CANADA



Canadian distributed astronomy platform



Shared international platform



Why INTERNATIONAL shared computing platforms?

Science practice is international

Reciprocity

- for data
- for computing infrastructure
- For services supporting data-intensive science

The Open Universe (whatever it turns out to be)

Open Universe

will be based on

IVOA standards

that support the operation of

Astronomy Data Centres

that are integrated into

Open Science Clouds

Shared infrastructure for data-intensive science

This new infrastructure creates opportunities for those who have limited access to resources

- Equalizes access for professional scientists in developing countries
- Provides new capabilities for teachers and the public

Example: Graduate student in Bangladesh

Missing recommendation

The UN recognizes that governments have invested hundreds of millions of dollars to create the present-day network of astronomy data services. The powerful science capabilities provided by this network are the foundation upon which Open Universe will operate.

These investments must continue and increase in order to support the types of services proposed by the Open Universe.

Fresh new ideas and approach

- Few new ideas in the Open Universe initiative
- The new factor is the involvement of the UN

A fresh new approach would be:

A substantial transfer of the benefits of astronomy data and supporting infrastructure to the public through education, outreach, and citizen science with a focus on developing nations as the highest priority.