

# ICT to support the transformation of Science in the Roaring Twenties

**Cees de Laat**

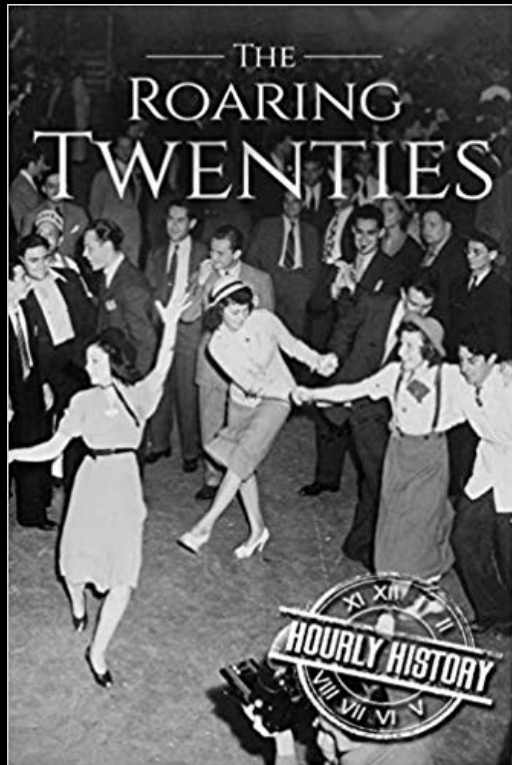
Systems and Networking Laboratory

University of Amsterdam

This trip is supported by SARNET, DL4LD and EPI projects.



# ICT to support the transformation of Science in the Roaring Twenties



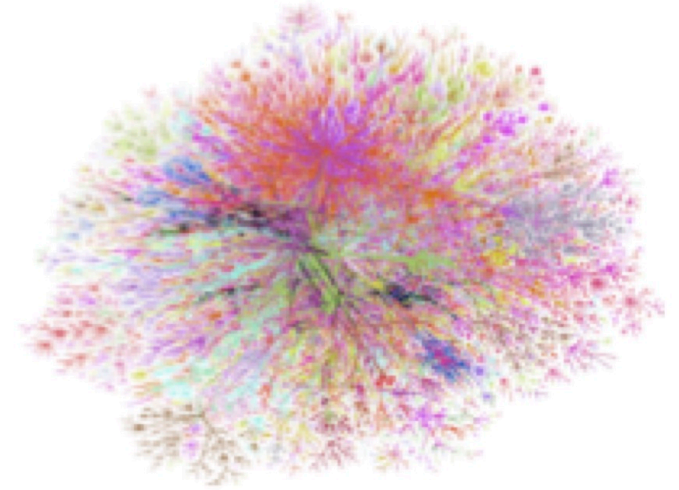
From Wikipedia: The Roaring Twenties refers to the decade of the 1920s in Western society and Western culture. It was a period of **economic prosperity** with a distinctive cultural edge in the United States and Western Europe, particularly in major cities such as Berlin, Chicago, London, Los Angeles, New York City, Paris, and Sydney. In France, the decade was known as the "**années folles**" ('crazy years'), emphasizing the era's **social, artistic and cultural dynamism**. Jazz blossomed, the flapper redefined the modern look for British and American women, and **Art Deco** peaked....

This period saw the large-scale development and use of automobiles, telephones, movies, radio, and electrical appliances being installed in the lives of thousands of Westerners. Aviation soon became a business. Nations saw **rapid industrial and economic growth, accelerated consumer demand**, and introduced significantly new changes in **lifestyle and culture**. The media focused on celebrities, especially sports heroes and movie stars, as cities rooted for their home teams and filled the new palatial cinemas and gigantic sports stadiums. In most major democratic states, women won the right to vote. The **right to vote** made a huge impact on society.

# AIM

- Observe how the art of Science is transforming with AI & ML.
- Understand how the ICT world looks like in 2030.
- Understand what hinders Science, Industry, Society to progress.
- What is needed to obtain EU leadership
  - Why?
  - Where?
  - What?

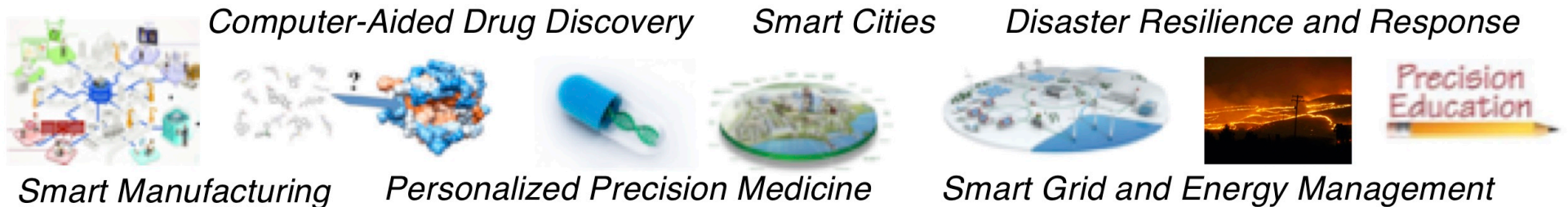
**In most applications, utilization of Big Data often needs to be combined with Scalable Computing.**



*COMPUTING AT DIVERSE SCALES*

*"BIG" DATA*

*Enables dynamic data-driven applications*



# Workflows for Data Science Center of Excellence at SDSC

[WorDS.sdsc.edu](http://WorDS.sdsc.edu)



Real-Time Hazards Management  
[wifire.ucsd.edu](http://wifire.ucsd.edu)

Data-Parallel Bioinformatics  
[bioKepler.org](http://bioKepler.org)

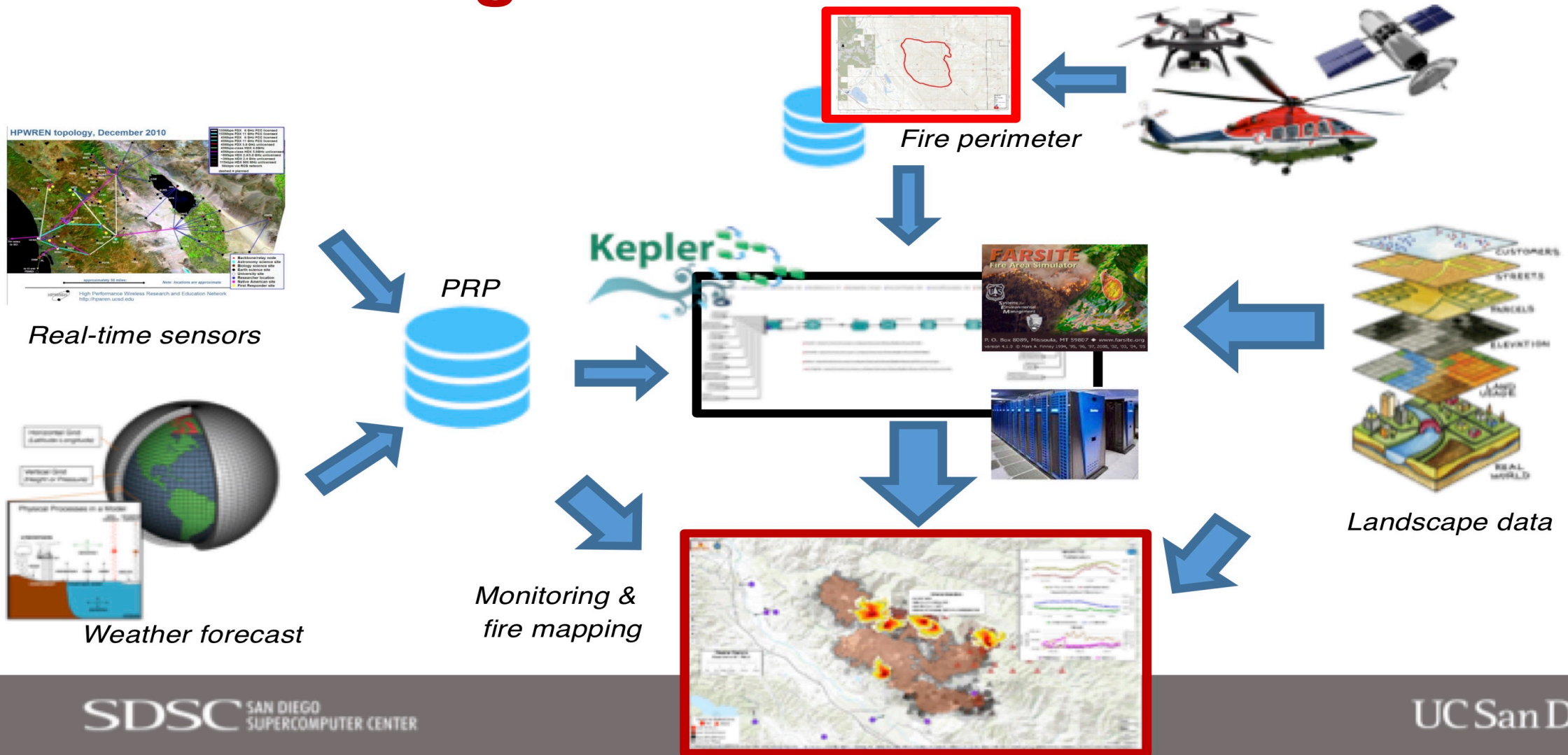
- Find, access and analyze data
- Support exploratory design
- Scale computational analysis
- Fuel reuse and reproducibility
- Save time, energy and money
- Formalize and standardize
- Train the next generation

**Goal:** Methodology and tool development to build automated and operational workflow-driven solution architectures on big data and HPC platforms.

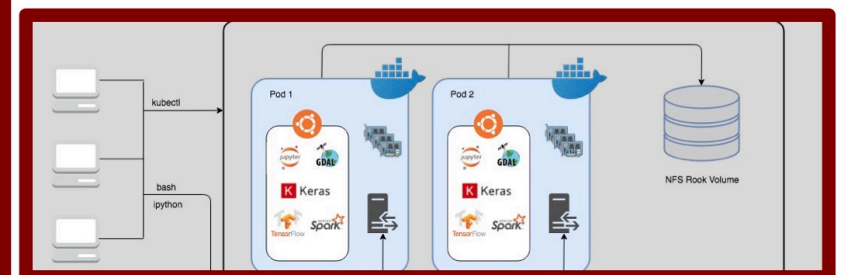
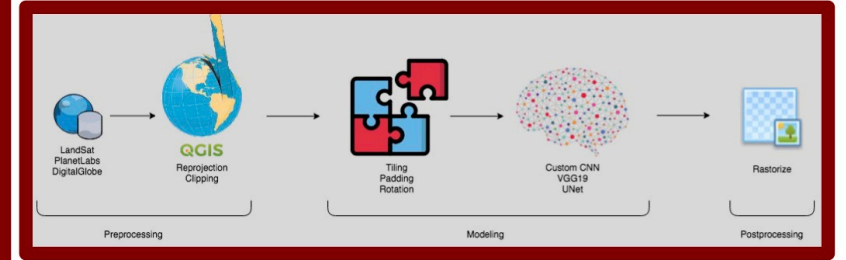
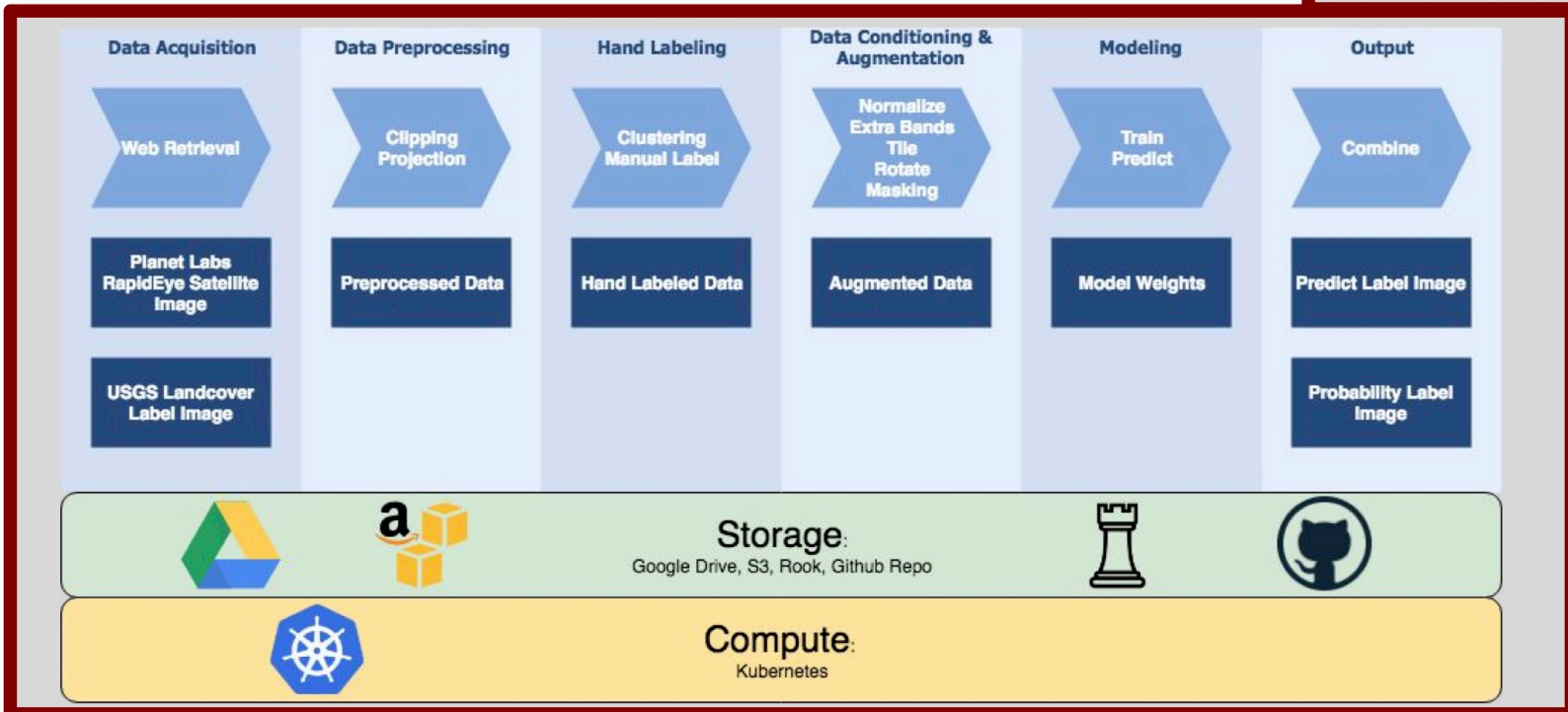
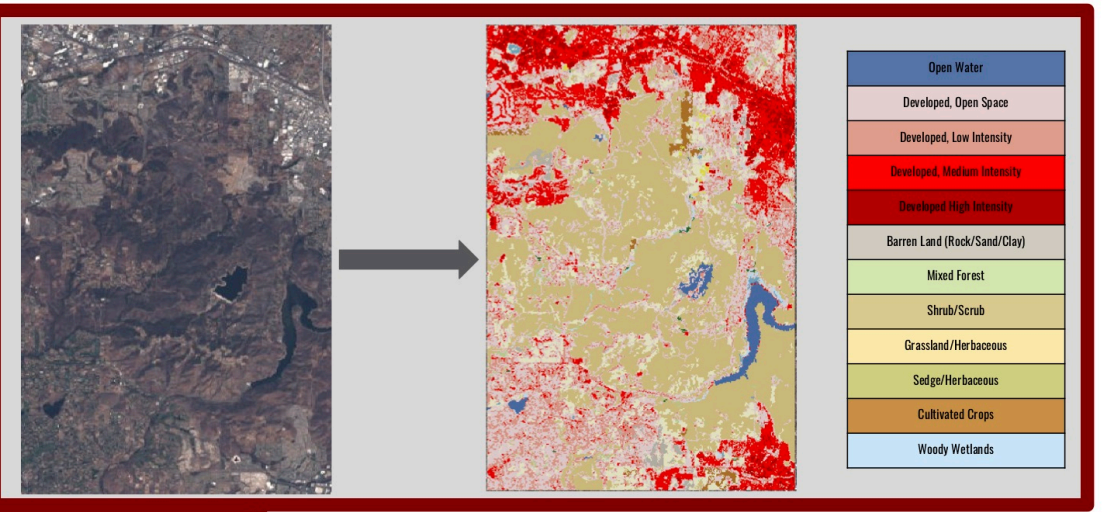
[kepler-project.org](http://kepler-project.org)

Scalable Automated Molecular Dynamics and Drug Discovery  
[nbc.ucsd.edu](http://nbc.ucsd.edu)

# Fire Modeling Workflows in WIFIRE



# One Piece of the Puzzle: Vegetation Classification using Satellite Imagery



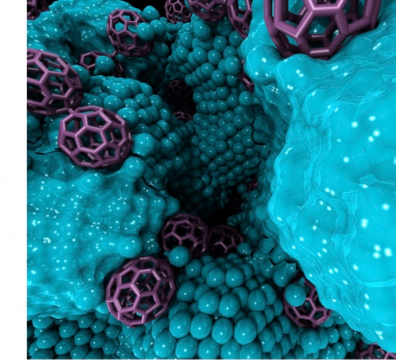
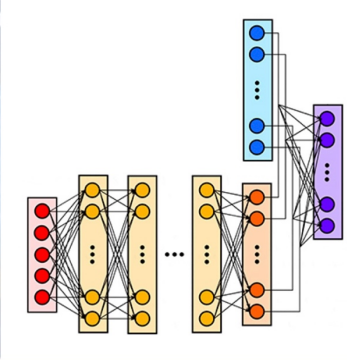
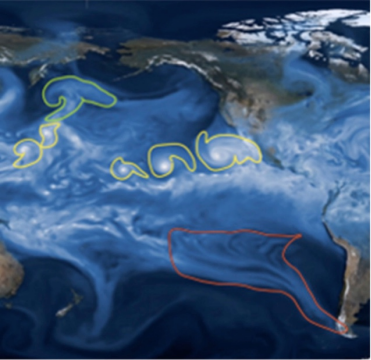
*"All for one, one for all!"*



<http://christianclipartreview.blogspot.com>

CI for AI & AI for CI





# Scientific Machine Learning & Artificial Intelligence

Scientific progress will be driven by

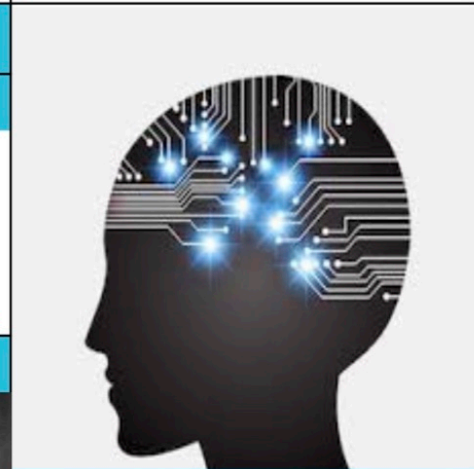
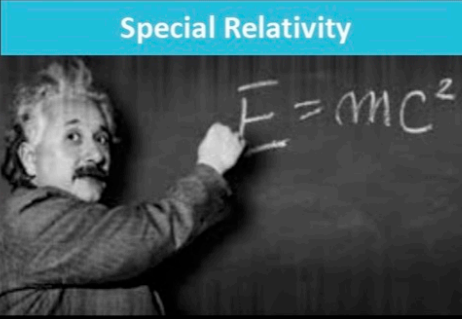
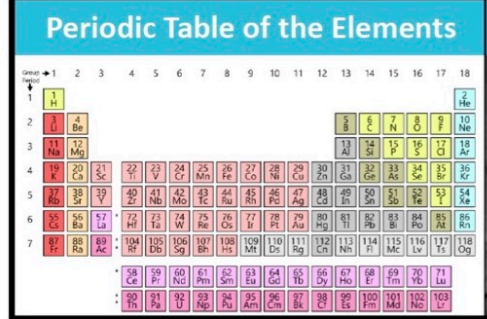
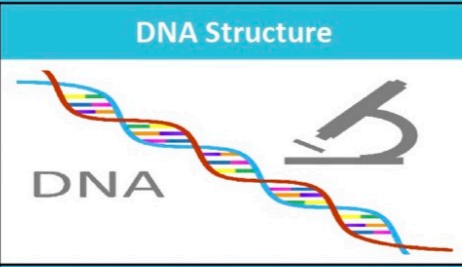
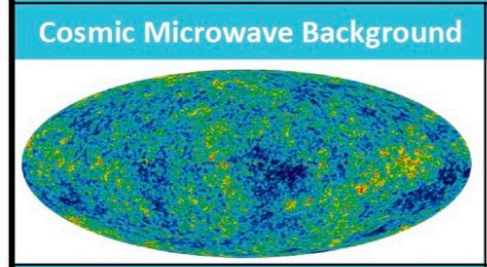
- Massive data: sensors, simulations, networks
- Predictive models and adaptive algorithms
- Heterogeneous high-performance computing

Trend: Human-AI collaborations will transform the way science is done.

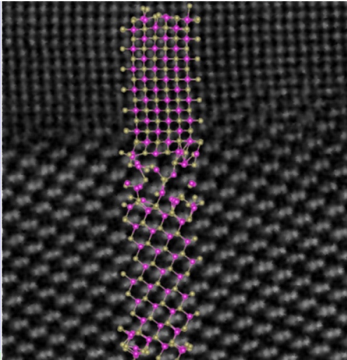
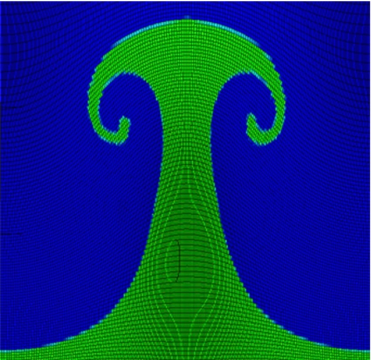
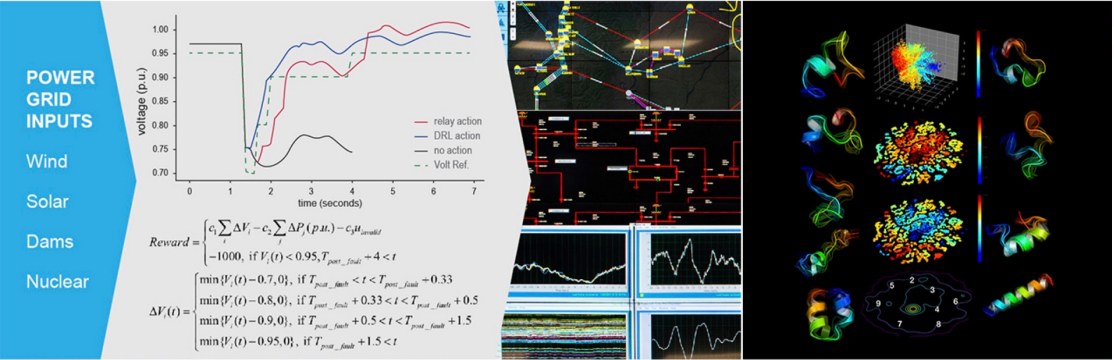
## BASIC RESEARCH NEEDS FOR Scientific Machine Learning

Core Technologies for Artificial Intelligence

### EXEMPLARS OF SCIENTIFIC ACHIEVEMENT



Human-AI insights enabled via scientific method, experimentation, & AI reinforcement learning.



Office of Science  
DOE Applied Mathematics Research Program  
Scientific Machine Learning Workshop (January 2018)

Prepared for U.S. Department of Energy Advanced Scientific Computing Research

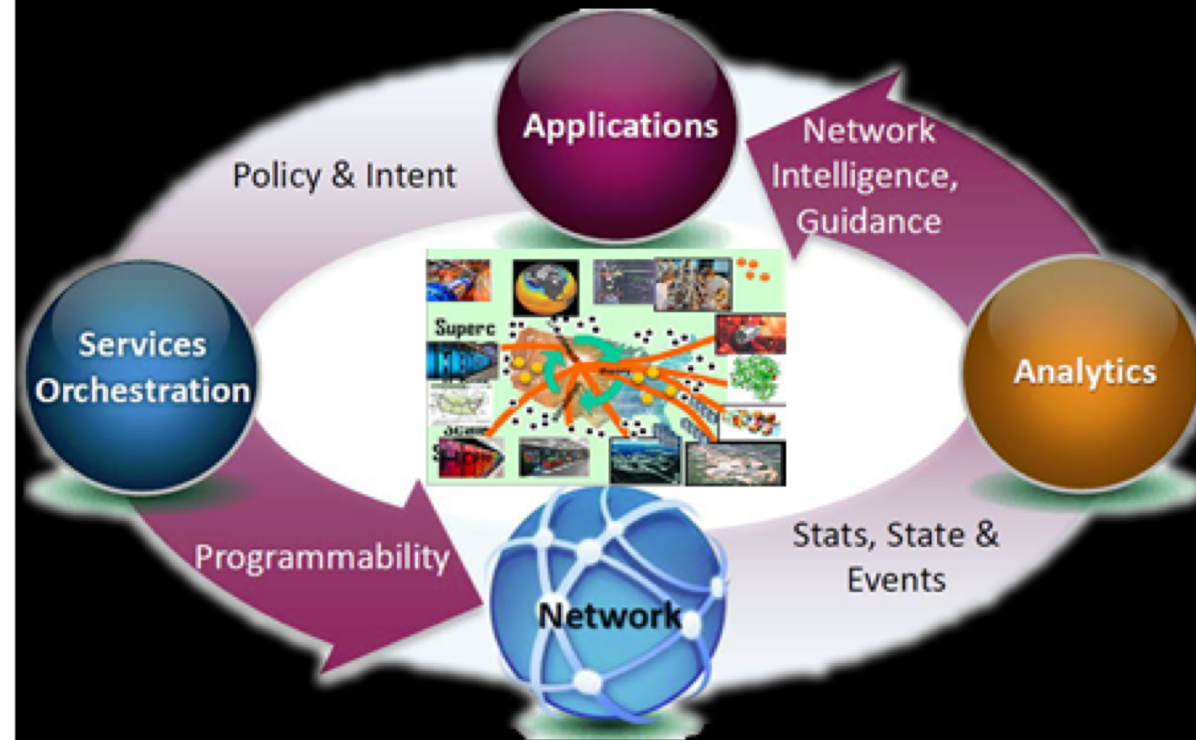
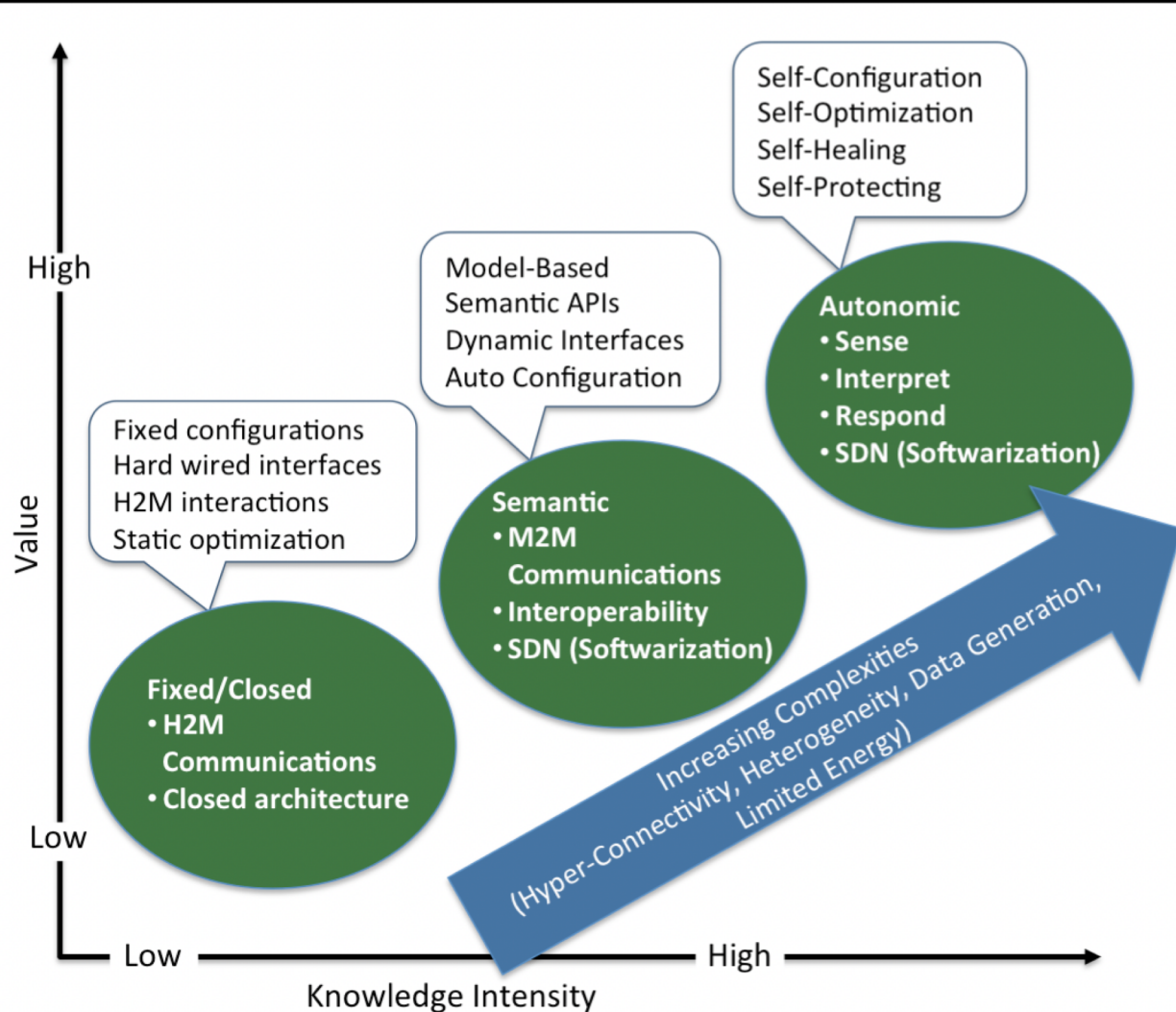


Workshop report:  
<https://www.osti.gov/biblio/1478744>

# DoE workshop on smart networks

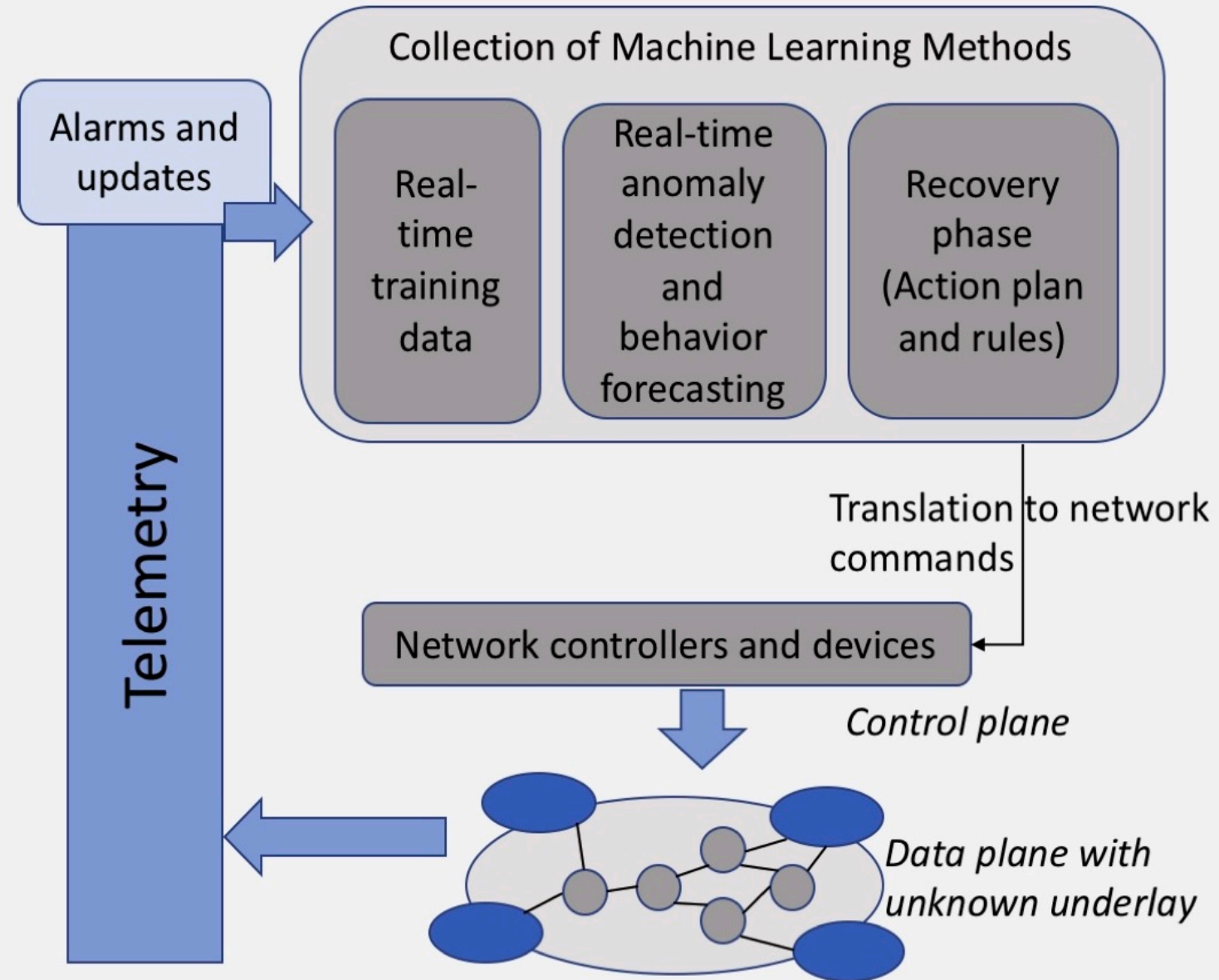
Bring AI in control plane to harness complexity

<https://www.ornl.gov/smarthp2016/>



# Example 1: Optimizing Network Traffic with Machine Learning

Exascale and increasingly complex science applications are exponentially raising demands from underlying DOE networks, such as traffic management, operation scale, and reliability constraints. Networks are the backbone to complex science workflows, ensuring data are delivered securely and on time for important computations to happen. To optimize these distributed workflows, networks are required to understand end-to-end performance needs in advance and be faster, efficient, and more proactive, anticipating bottlenecks before they happen. However, to manage multiple network paths intelligently, various tasks, such as pre-computation and prediction, must be done in near real time. ML provides a collection of algorithms that can add autonomy and assist in decision making to sup-



# Rethinking NSF's Computational Ecosystem for 21st Century Science and Engineering

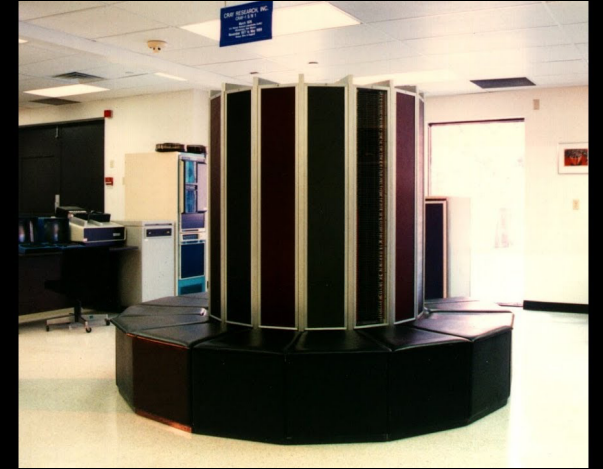
Workshop Website: <https://uiowa.edu/nsfcyberinfrastructure>

Workshop Report: <https://www.uiowa.edu/nsfcyberinfrastructure/report.pdf>

Initial debates about resource management and delivery options focused on **expert personnel as a critical component** of successful cyberinfrastructure delivery. Several examples such as Campus Champions (CC) or XSEDE's ECSS were described as critical to scientific advance but insufficient in numbers to meet demand. Regionally tasked staff might help to alleviate this shortfall. Benefits could include greater use of cloud or national resources if there was a local expert to help researchers with initial utilization. Along these lines, it was mentioned that the **NSF CC\* programs changed campus culture**, spurring local networking expertise. A similar program to promote workforce development to incentivize local computational and data scientists could, for instance, result in the integration of otherwise isolated clusters on campuses with national resources. These **key personnel**, ranging from ECSS experts and developers to CCs, are often in careers that need professionalization.

# Change in computing

- Early days a few big Supercomputers
  - Mostly science domain
- Via grid to commercial cloud
  - AWS, Azure, Google Cloud, IBM, Salesforce
  - The big five: Apple, Alphabet, Microsoft, Facebook and Amazon
  - Computing has transformed into an utility
- Data => Information is the key



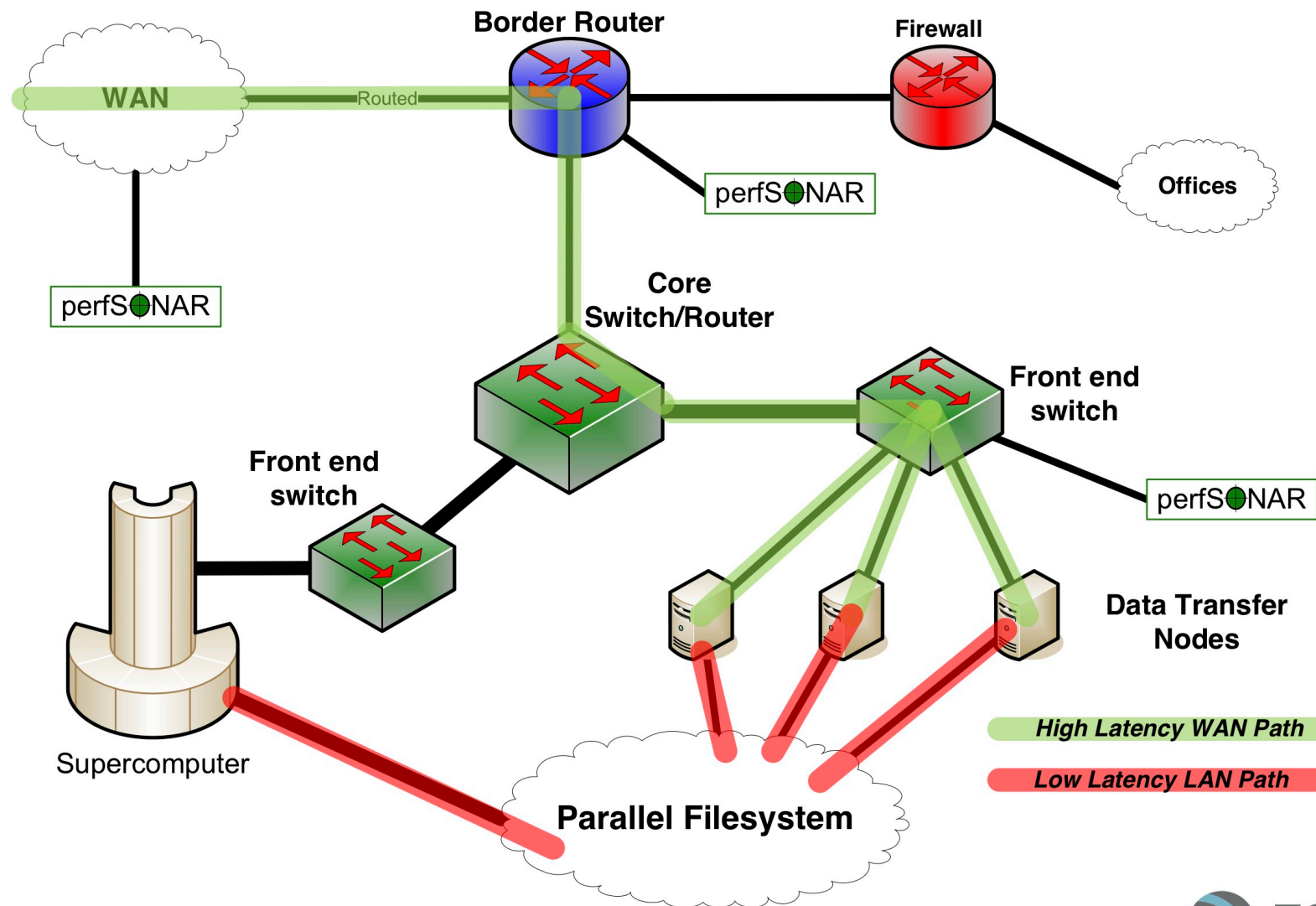
# Now, how do we get and use data?

## 2019 This Is What Happens In An Internet Minute

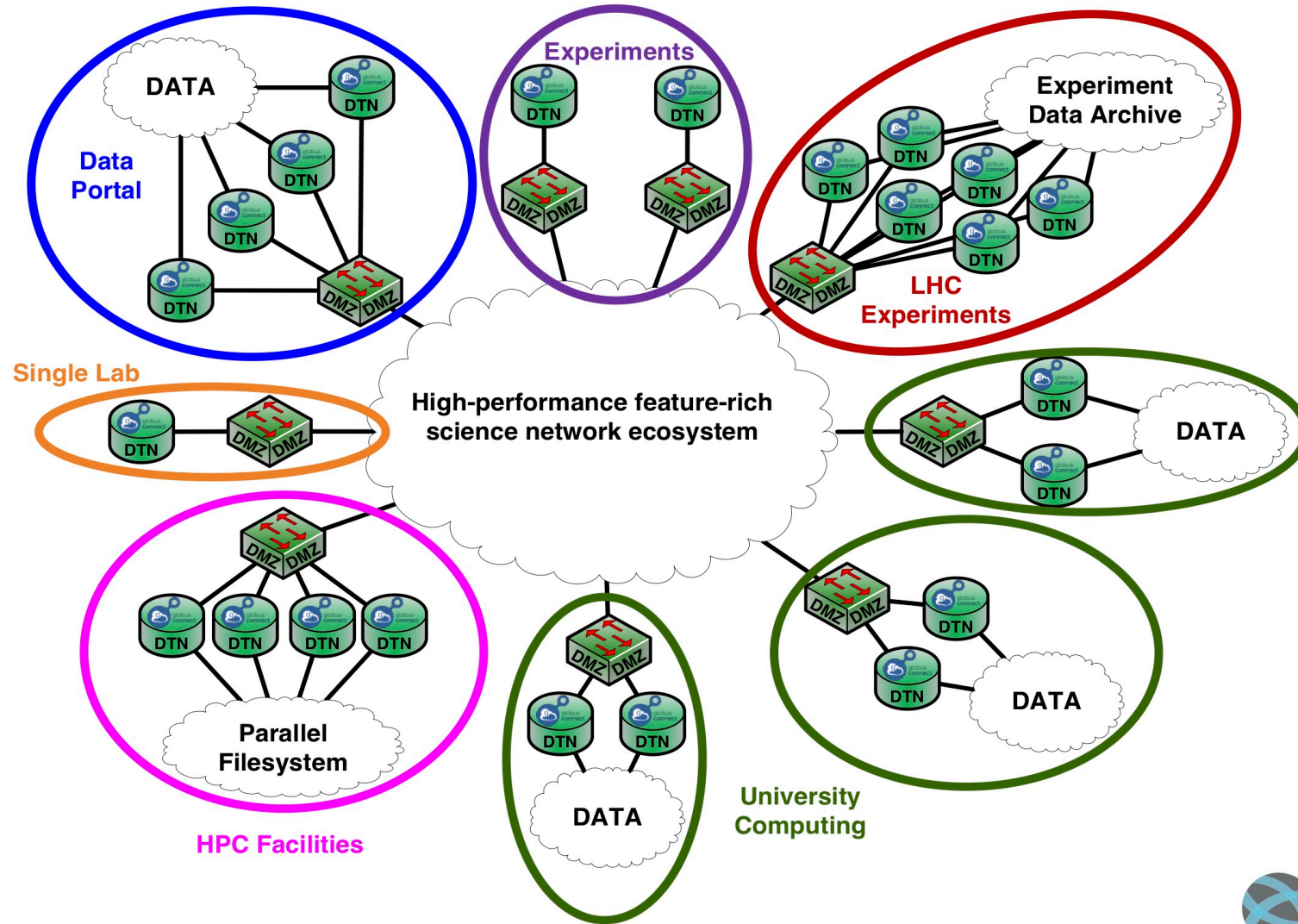


- Move towards streaming
  - Netflix
  - youtube
- Same in science world
  - SKA / LOFAR
  - Light Source
  - Environmental (Marine, Meteorology, ...)
- Data is not always huge
  - Sometimes it is very complex
  - Some example:
    - biodiversity

# Science DMZ – HPC Center DTN Cluster

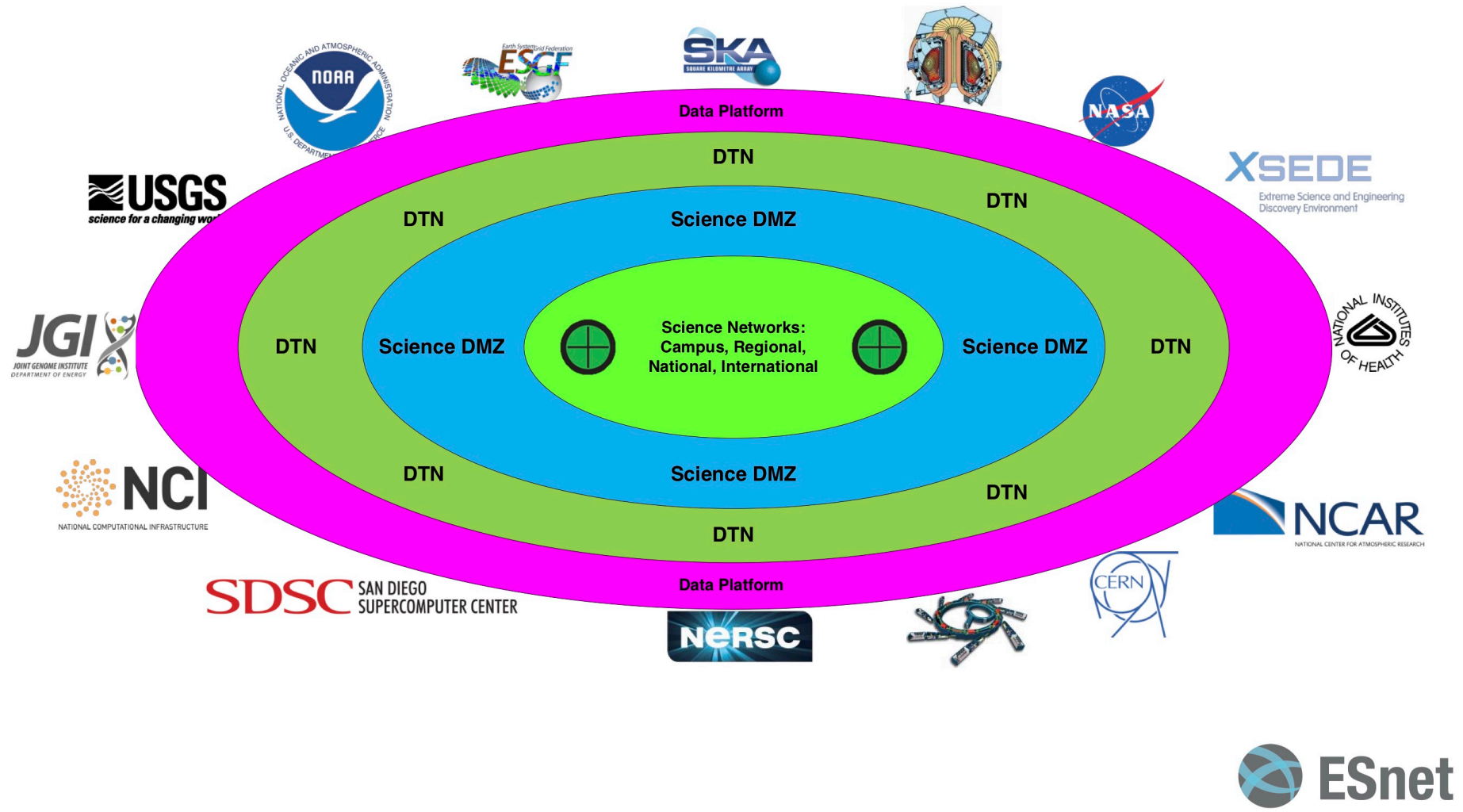


# Science DMZs for Science Applications





# Data Ecosystem – Concentric View



# DTN Cluster Performance – HPC Facilities (2017)

Petascale DTN Project

November 2017

L380 Data Set

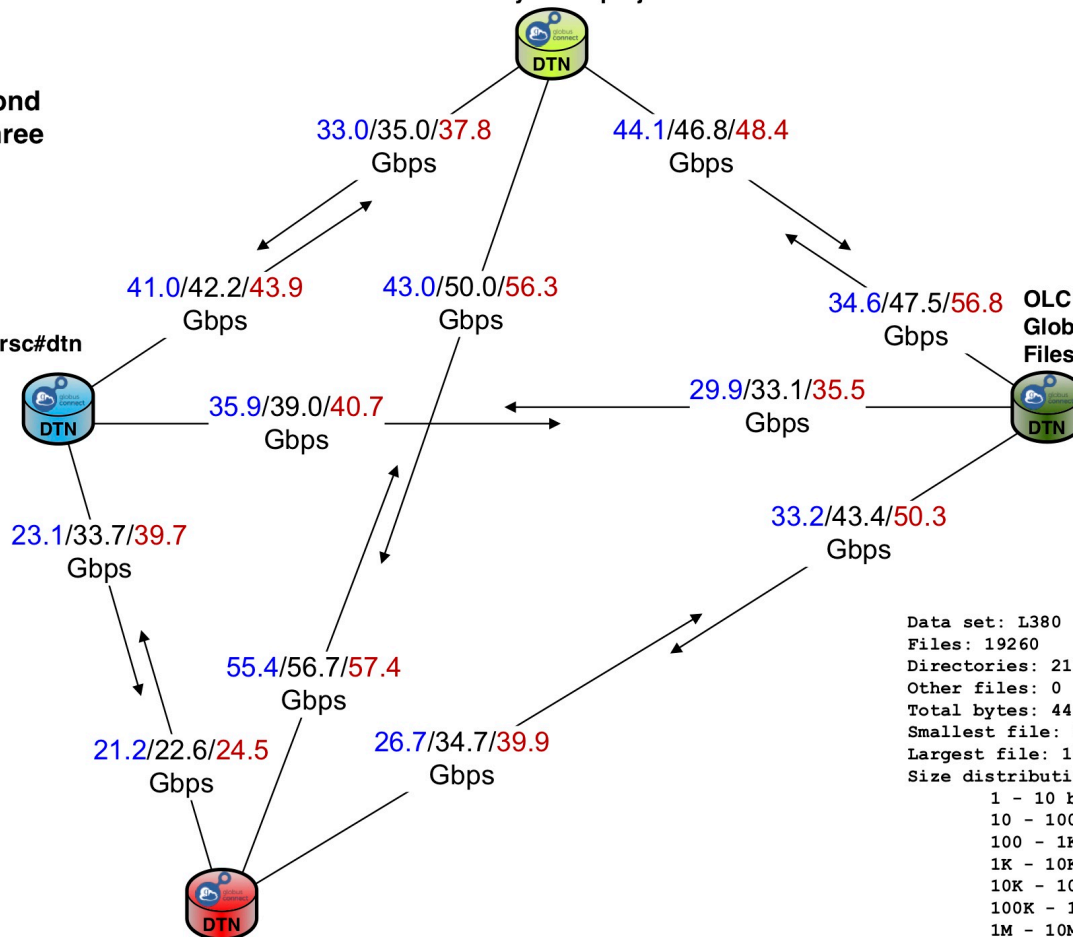
Gigabits per second  
(min/avg/max), three transfers

NERSC DTN cluster  
Globus endpoint: nersc#dtm  
Filesystem: /project

ALCF DTN cluster  
Globus endpoint: alcf#dtm\_mira  
Filesystem: /projects

OLCF DTN cluster  
Globus endpoint: olcf#dtm\_atlas  
Filesystem: atlas2

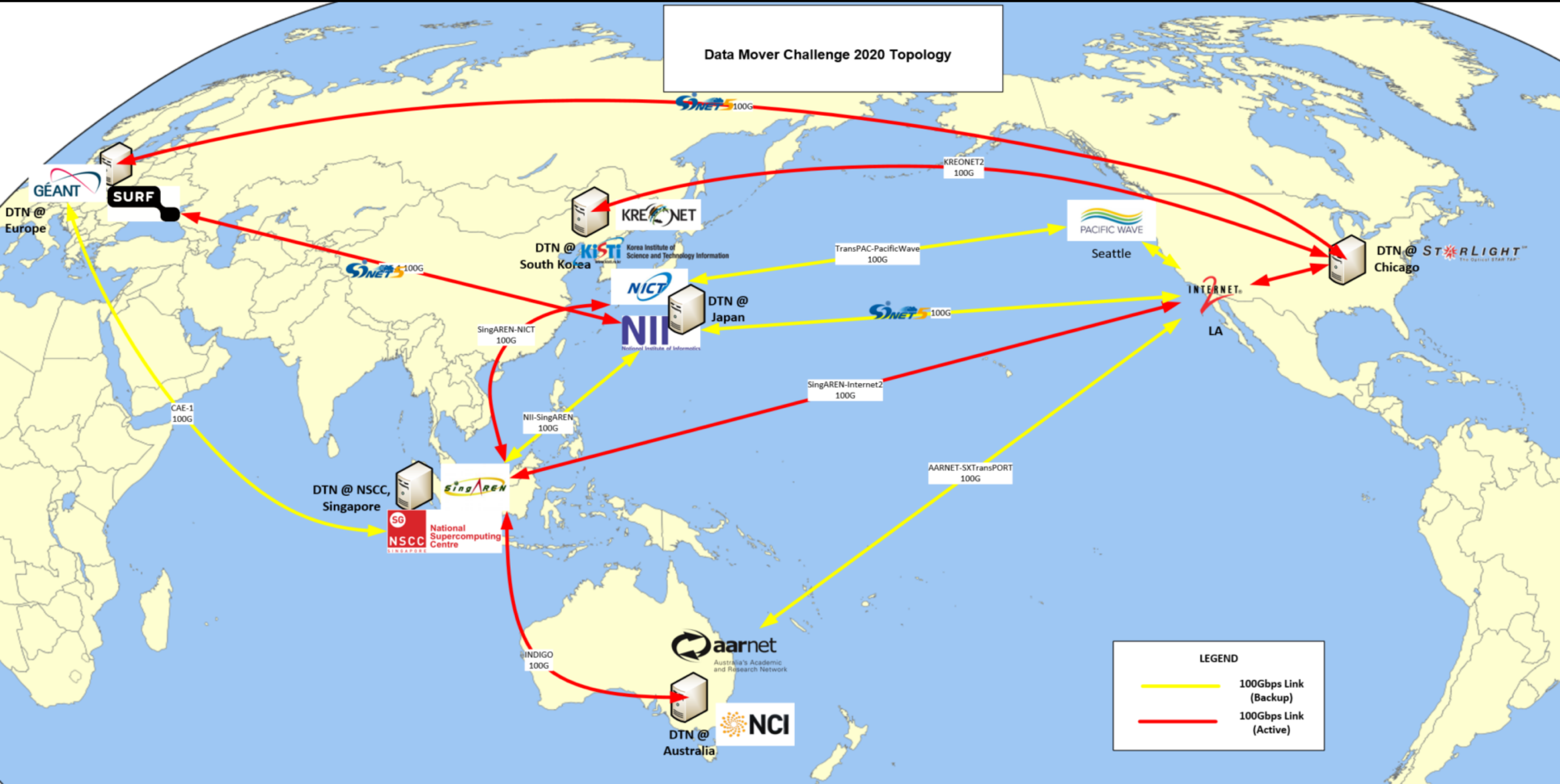
NCSA DTN cluster  
Globus endpoint: ncsa#BlueWaters  
Filesystem: /scratch



Data set: L380  
Files: 19260  
Directories: 211  
Other files: 0  
Total bytes: 4442781786482 (4.4T bytes)  
Smallest file: 0 bytes (0 bytes)  
Largest file: 11313896248 bytes (11G bytes)  
Size distribution:  
1 - 10 bytes: 7 files  
10 - 100 bytes: 1 files  
100 - 1K bytes: 59 files  
1K - 10K bytes: 3170 files  
10K - 100K bytes: 1560 files  
100K - 1M bytes: 2817 files  
1M - 10M bytes: 3901 files  
10M - 100M bytes: 3800 files  
100M - 1G bytes: 2295 files  
1G - 10G bytes: 1647 files  
10G - 100G bytes: 3 files



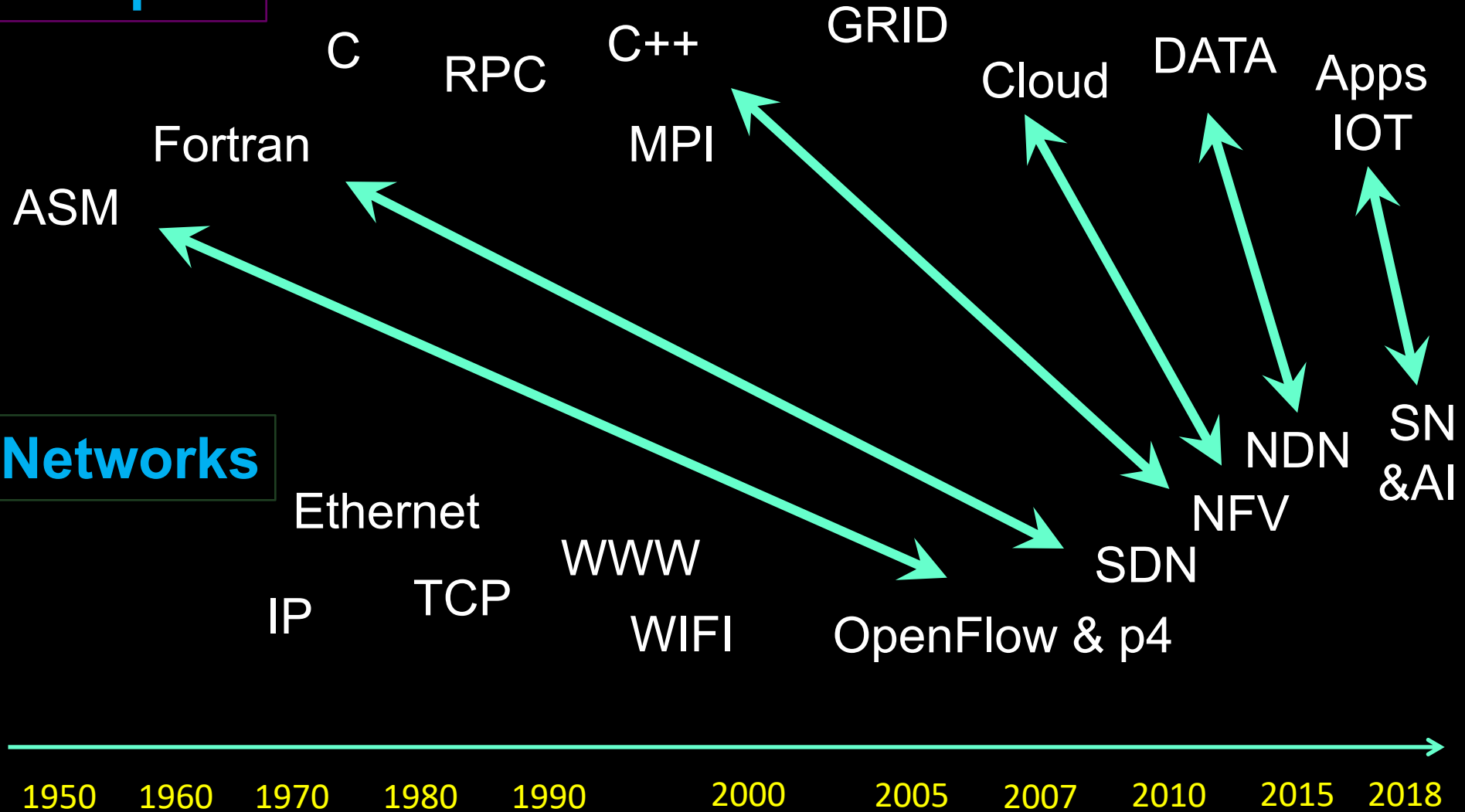
# <https://www.sc-asia.org/data-mover-challenge/>

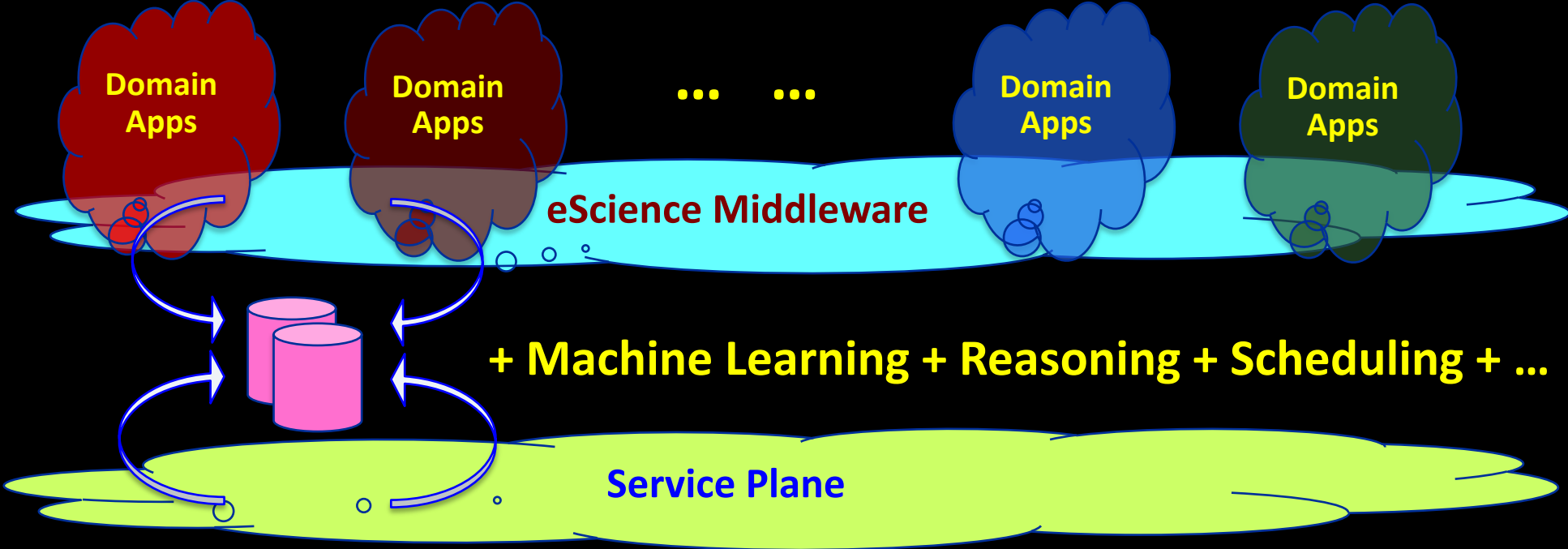


# TimeLine

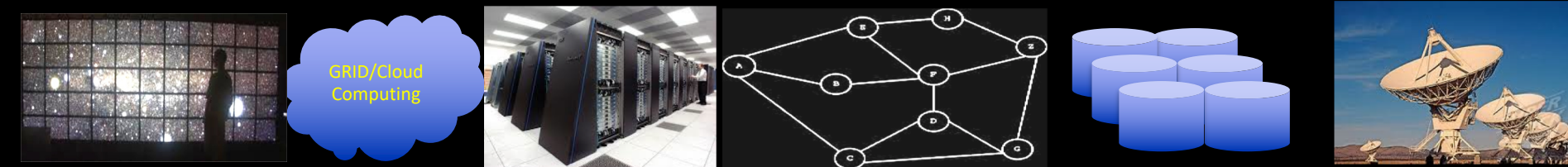
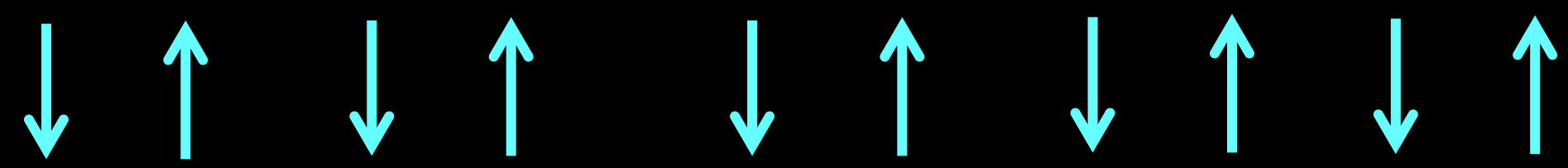
**Compute**

**Networks**





- Chromium CGLX
- SAGE MTP
- OCCI JSDL
- GIR UR
- SNMP OpenFlow SDN / NSI
- PerfSonar ICMP
- Cassandra iRODS
- Hadoop Storm
- WSRF SensorML
- WebServ INSPIRE



# The Big Data Challenge

Doing Science

ICT to enable Science

Wisdom

AI

Knowledge to act

Analytics  
Decision Support

Information

Web/OWL

Data  
a.o. from ESFRI's

Docker, VM,  
XML, RDF, rSpec, SNMP



# The Big Data Challenge

Doing Science

ICT to enable Science

Wisdom

Scientists live here!

AI

Interdisciplinary Science App Store

Knowledge to act

Analytics library / Github / etc

Analytics Decision Support

MAGIC DATA CARPET

curation - description - trust - security - policy - integrity

Information

Web/OWL

Data

a.o. from ESFRI's

Docker, VM, XML, RDF, rSpec, SNMP



# The Big Data Challenge

Doing Science

ICT to enable Science

Wisdom

AI

Scientists live here!

Interdisciplinary Science App Store

Knowledge to act

Analytics library / Github / etc

Analytics Decision Support

MAGIC DATA CARPET  
curation - description - trust - security - policy - integrity

Information

Web/OWL

Data

a.o. from ESFRI's

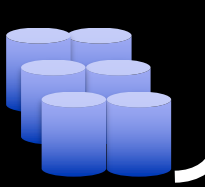
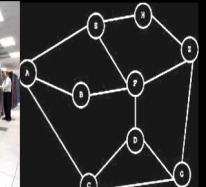
Docker, VM,

XML, RDF, rSpec, SNMP

DSC  
eScience

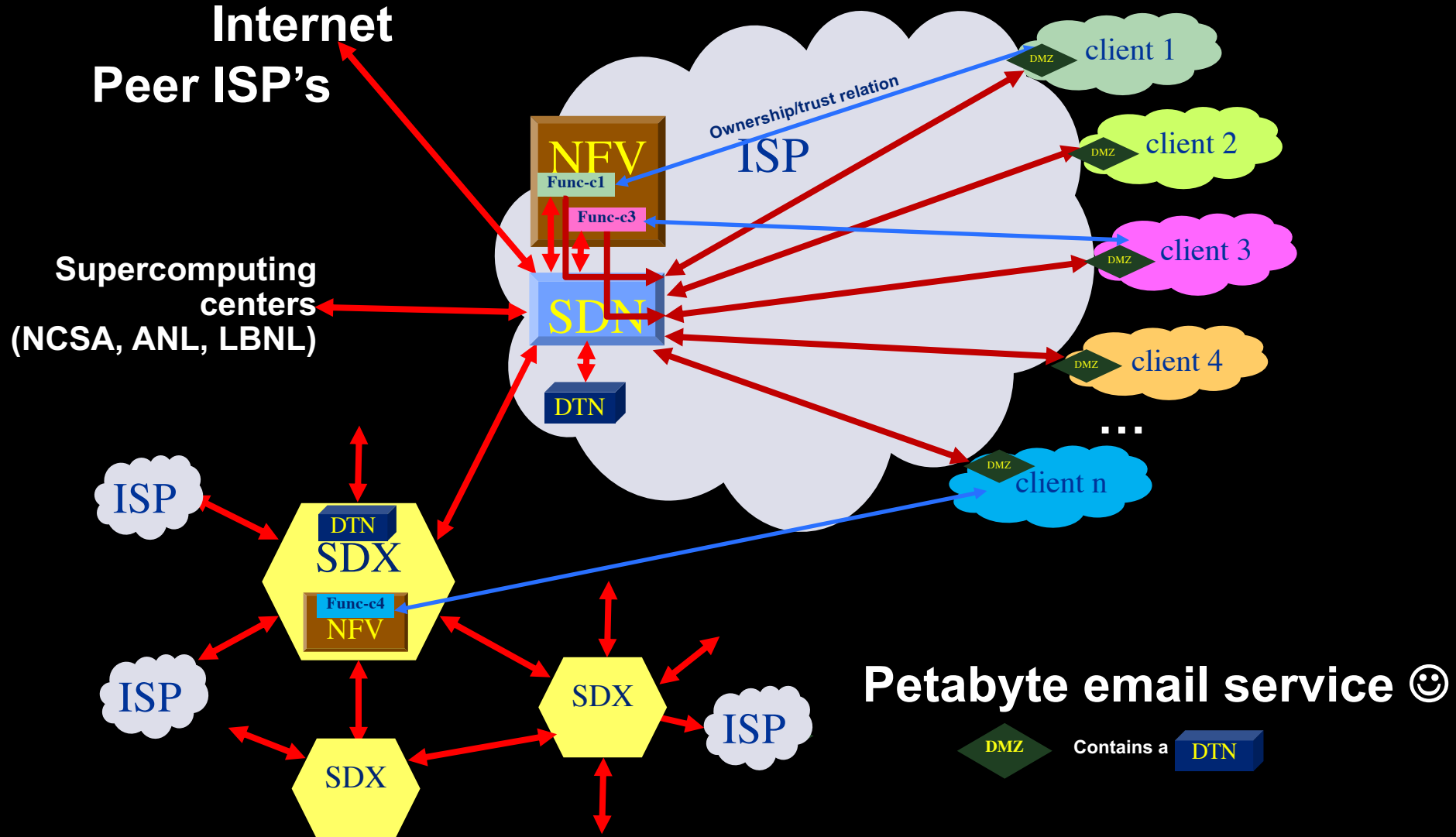
RDM/  
DANS

ICT/  
SURF

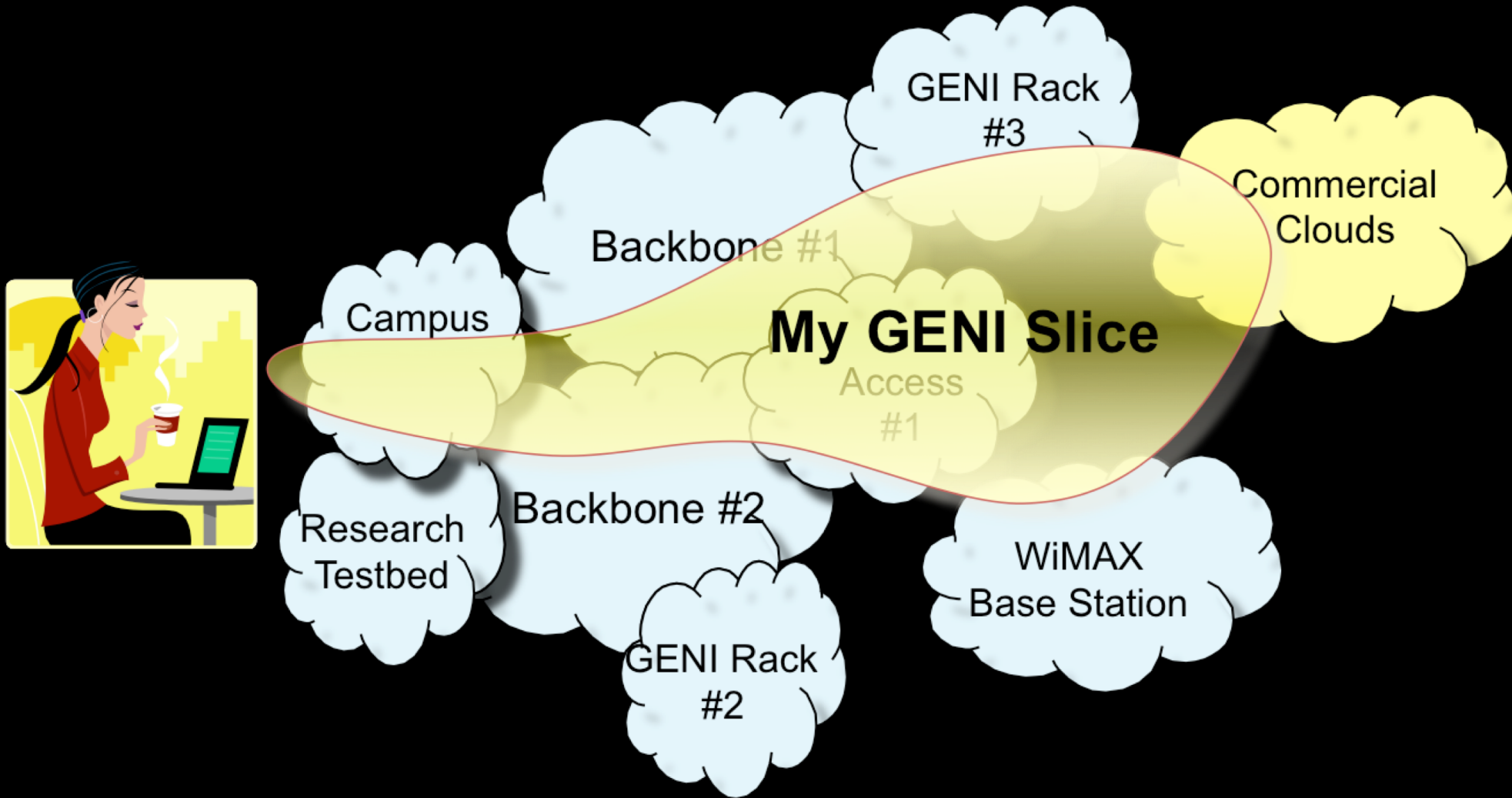




# Networks of ScienceDMZ's & SDX's



# GENI: Virtualizing CI

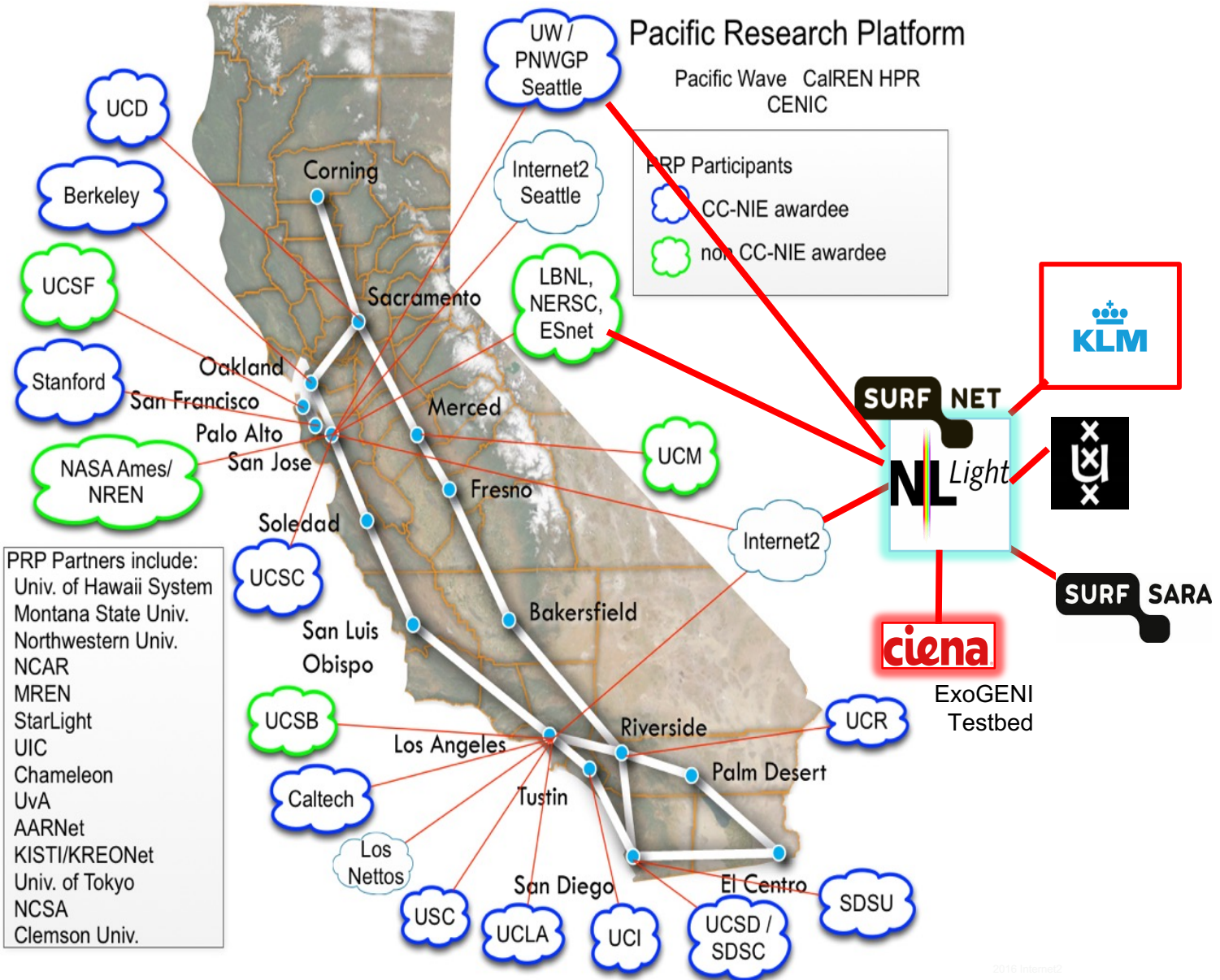


# Pacific Research Platform testbed involvement

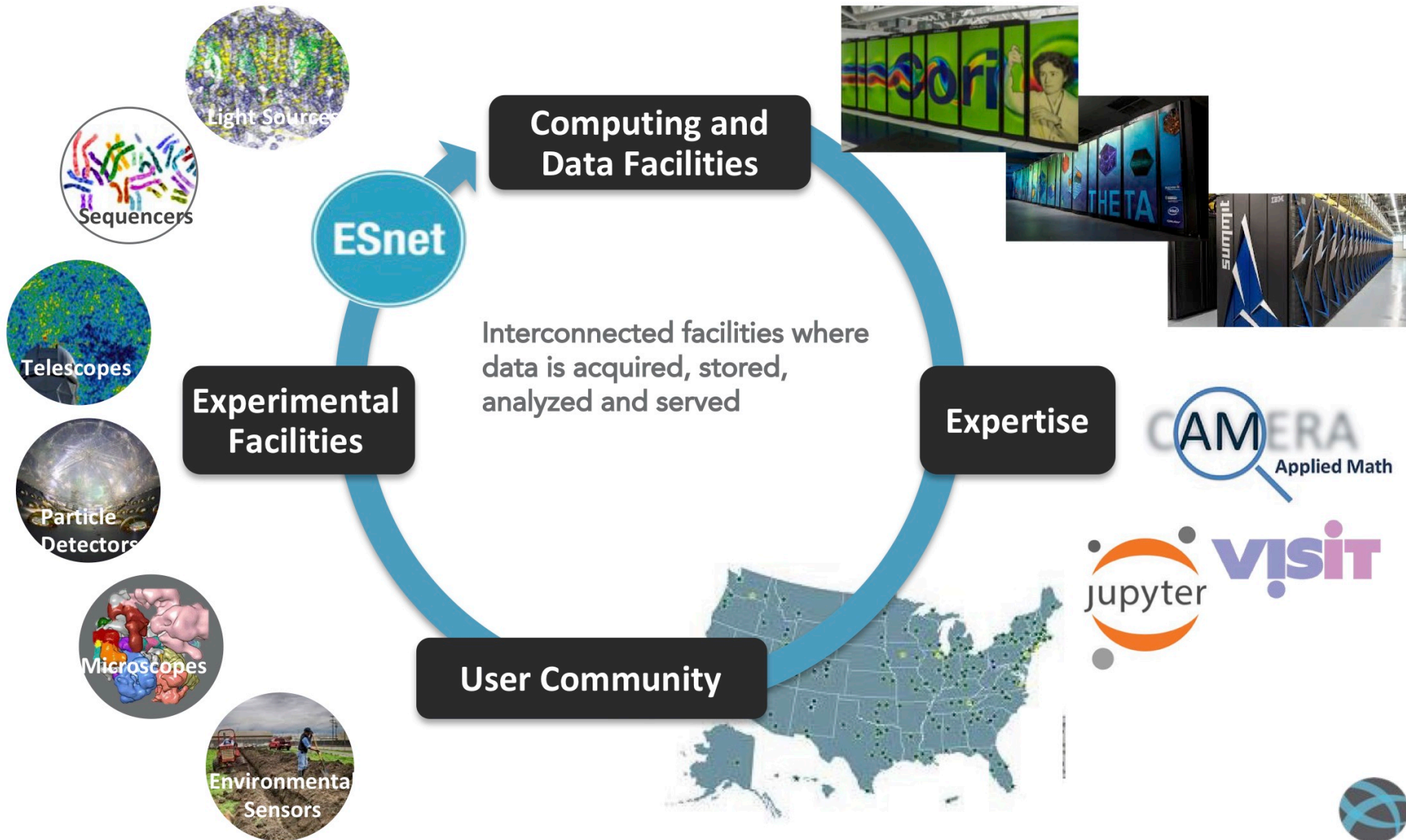
**Research goal:**  
Explore value of academic network research capabilities that enable innovative ways & models to share big data assets



PRP Partners include:  
Univ. of Hawaii System  
Montana State Univ.  
Northwestern Univ.  
NCAR  
MREN  
StarLight  
UIC  
Chameleon  
UvA  
AARNet  
KISTI/KREONet  
Univ. of Tokyo  
NCSA  
Clemson Univ.



# Superfacility Model for Productive, Reproducible Science



Extreme App Store for Science; Amazonizing infrastructure

# SARNET: Security Autonomous Response with programmable NETWORKS

Marc Lyonnais, Leon Gommans, Rodney Wilson, Lydia Meijer, Frank Fransen Tom van Engers, Paola Grosso, Gauravdeep Shami, Cees de Laat, Ameneh Deljoo, Ralph Koning, Ben de Graaff, Gleb Polevoy, Stojan Travanovski.



# Big Data: real time ICT for logistics Data Logistics 4 Logistics Data (dl4ld)

Lydia Meijer (PI), Cees de Laat (Co-PI), Leon Gommans, Tom van Engers, Paola Grosso, Kees Nieuwenhuis.



# EPI: Enabling Personalized Interventions

Cees de Laat(PI), Sander Klous (PL), Leon Gommans, Tom van Engers, Paola Grosso, Henri Bal, Anwar Osseyran, Aki Harma, Douwe Biesma, Peter Grünwald, Floortje Scheepers, Gertjan Kaspers.



# SARNET: Security Autonomous Response with programmable NETWORKS

Marc Lyonnais, Leon Gommans, Rodney Wilson, Lydia Meijer, Frank Fransen Tom van Engers, Paola Grosso, Gauravdeep Shami, Cees de Laat, Ameneh Deljoo, Ralph Koning, Ben de Graaff, Gleb Polevoy, Stojan Travanovski.



## Big Data: real time ICT for logistics Data Logistics 4 Logistics Data (dl4ld)

Lydia Meijer (PI), Cees de Laat (Co-PI), Leon Gommans, Tom van Engers, Paola Grosso, Kees Nieuwenhuis.



## EPI: Enabling Personalized Interventions

Cees de Laat(PI), Sander Klous (PL), Leon Gommans, Tom van Engers, Paola Grosso, Henri Bal, Anwar Osseyran, Aki Harma, Douwe Biesma, Peter Grünwald, Floortje Scheepers, Gertjan Kaspers.



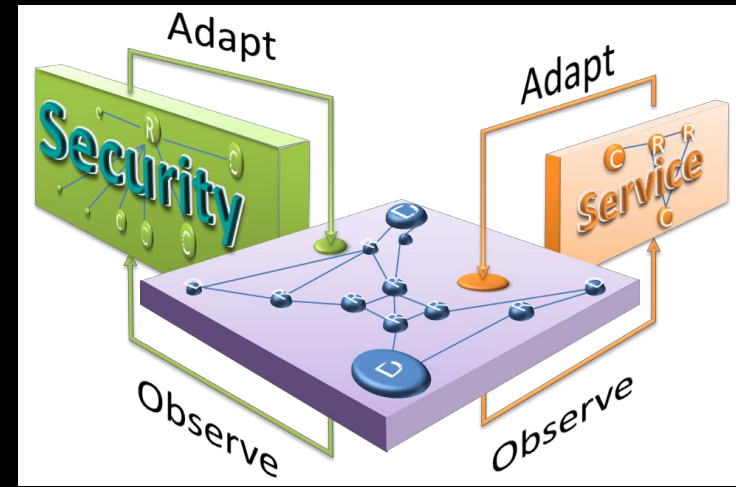
# Cyber security program

## SARNET

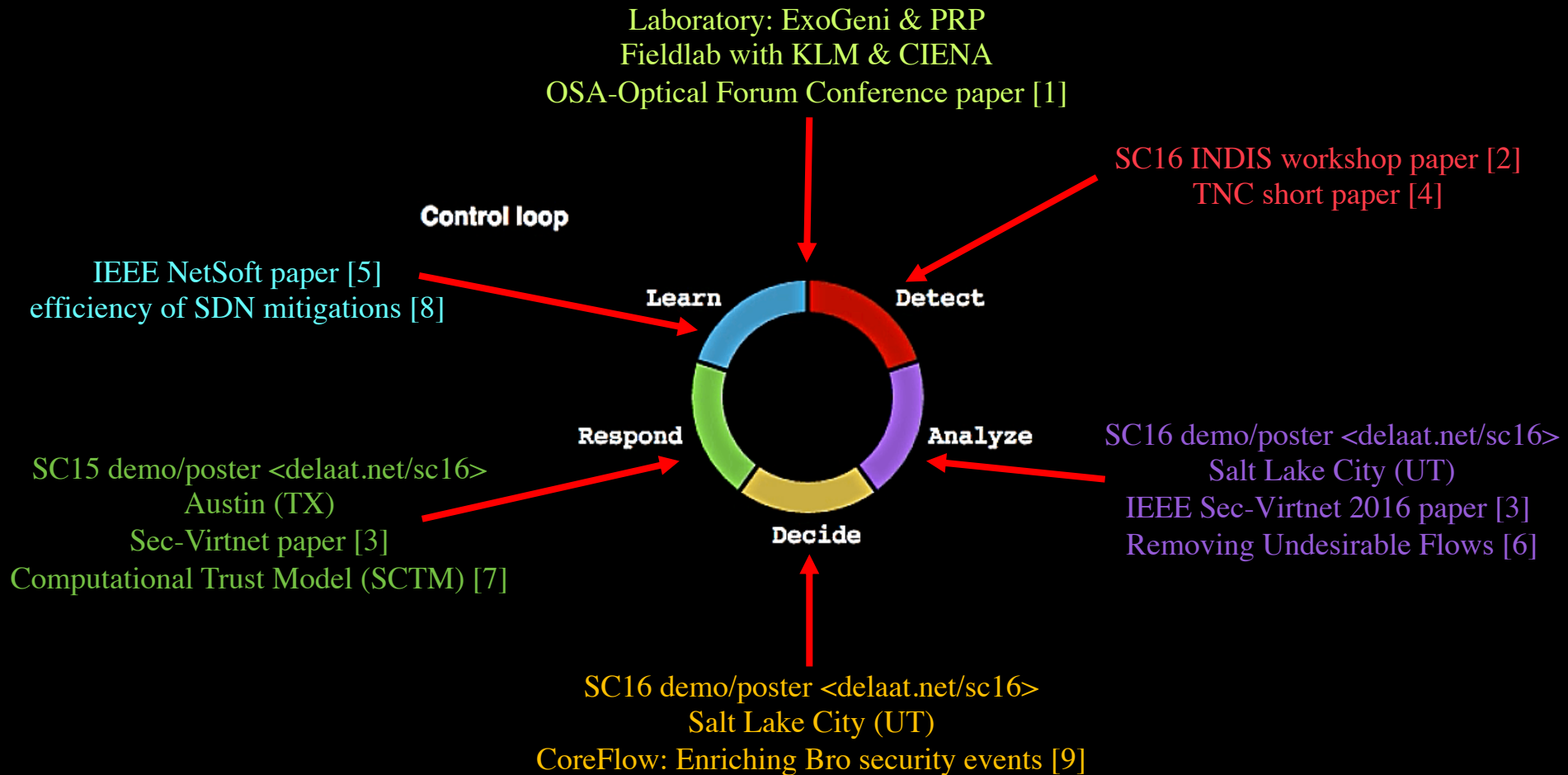
Research goal is to obtain the knowledge to create ICT systems that:

- model their state (situation)
- discover by observations and reasoning if and how an attack is developing and calculate the associated risks
- have the knowledge to calculate the effect of counter measures on states and their risks
- choose and execute one.

In short, we research the concept of networked computer infrastructures exhibiting SAR: Security Autonomous Response.



# SARNET Publications (subset)



1. Paper: R. Koning, A. Deljoo, S. Trajanovski, B. de Graaff, P. Grosso, L. Gommans, T. van Engers, F. Franssen, R. Meijer, R. Wilson, and C. de Laat, "Enabling E-Science Applications with Dynamic Optical Networks: Secure Autonomous Response Networks", OSA Optical Fiber Communication Conference and Exposition, 19-23 March 2017, Los Angeles, California.
2. Paper: Ralph Koning, Nick Buraglio, Cees de Laat, Paola Grosso, "CoreFlow: Enriching Bro security events using network traffic monitoring data.", Special section on high-performance networking for distributed data-intensive science, SC16", Future Generation Computer Systems, <accepted for publication>
3. Paper: Ralph Koning, Ben de Graaff, Cees de Laat, Robert Meijer, Paola Grosso, "Analysis of Software Defined Networking defenses against Distributed Denial of Service attacks", The IEEE International Workshop on Security in Virtualized Networks (Sec-VirtNet 2016) at the 2nd IEEE International Conference on Network Softwarization (NetSoft 2016), Seoul Korea, June 10, 2016.
4. Short paper: Nick Buraglio, Ralph Koning, Cees de Laat, Paola Grosso, "Enriching network and security events for event detection", Conference proceedings TNC2017, <https://tnc17.geant.org/core/presentation/30>.
5. Paper: Ralph Koning, Ben de Graaff, Robert Meijer, Cees de Laat, Paola Grosso, "Measuring the effectiveness of SDN mitigations against cyber attacks", IEEE Conference on Network Softwarization (Netsoft 2017 - SNS 2017), Bologna, Italy, July 3-7, 2017.
6. Paper: Gleb Polevoy, Stojan Trajanovski, Paola Grosso and Cees de Laat, "Removing Undesirable Flows by Edge Deletion.", COCOA'2018 conference, December 15 - 17, 2018, Atlanta, Georgia, USA, Springer-Verlag.
7. Paper: Ameneh Deljoo, Tom van Engers, Leon Gommans, Cees de Laat, "Social Computational Trust Model (SCTM): A Framework to Facilitate Selection of Partners". In: Proceedings of 2018 IEEE/ACM Innovating the Network for Data-Intensive Science (INDIS), Dallas, TX, USA, 2018
8. Paper: R. Koning, B. de Graaff, G. Polevoy, R. Meijer, C. de Laat, P. Grosso, "Measuring the efficiency of SDN mitigations against attacks on computer infrastructures", Future Generation Computer Systems 91, 144-156.
9. Ralph Koning, Nick Buraglio, Cees de Laat, Paola Grosso, "CoreFlow: Enriching Bro security events using network traffic monitoring data.", Special section on high-performance networking for distributed data-intensive science, SC16", Future Generation Computer Systems



# SARNET: Security Autonomous Response with programmable NETWORKS

Marc Lyonnais, Leon Gommans, Rodney Wilson, Lydia Meijer, Frank Fransen Tom van Engers, Paola Grosso, Gauravdeep Shami, Cees de Laat, Ameneh Deljoo, Ralph Koning, Ben de Graaff, Gleb Polevoy, Stojan Travanovski.



## Big Data: real time ICT for logistics Data Logistics 4 Logistics Data (dl4ld)

Lydia Meijer (PI), Cees de Laat (Co-PI), Leon Gommans, Tom van Engers, Paola Grosso, Kees Nieuwenhuis.



## EPI: Enabling Personalized Interventions

Cees de Laat(PI), Sander Klous (PL), Leon Gommans, Tom van Engers, Paola Grosso, Henri Bal, Anwar Osseyran, Aki Harma, Douwe Biesma, Peter Grünwald, Floortje Scheepers, Gertjan Kaspers.



# Harvard Business Review



Harvard Business Review

ECONOMY

## Managing Our Hub Economy


by Marco Iansiti and Karim R. Lakhani

FROM THE SEPTEMBER–OCTOBER 2017 ISSUE

WHAT TO READ NEXT

The IT Transformation Health Care Needs

SUMMARY SAVE SHARE COMMENT 3 TEXT SIZE PRINT \$8.95 BUY COPIES



THOMAS M. SCHEER/EYEEM/GETTY IMAGES

### I. The Problem

The global economy is coalescing around a few digital superpowers. We see unmistakable evidence that a winner-take-all world is emerging in which a small number of “hub firms”—including Alibaba, Alphabet/Google, Amazon, Apple, Baidu, Facebook, Microsoft, and Tencent—occupy central positions. While creating real value for users, these companies are also capturing a disproportionate and expanding share of the value, and that’s shaping our collective economic future. The very same technologies that promised to democratize business are now threatening to make it more monopolistic.

Data value creation  
monopolies



Create an equal  
playing field



Sound Market  
principles

<https://hbr.org/2017/09/managing-our-hub-economy>

# Data Sharing: Main problem statement

- Organizations that normally compete have to bring data together to achieve a common goal!
- The shared data may be used for that goal but not for any other!
- Data or Algorithms may have to be processed in foreign data centers.
  - How to organize alliances?
  - How to translate from strategic via tactical to operational level?
  - How to enforce policy using modern Cyber Infrastructure?
  - What are the different fundamental data infrastructure models to consider?

# Big Data Sharing use cases placed in airline context



**Global Scale**



Aircraft Component Health Monitoring (Big) Data  
NWO **CIMPLO** project  
4.5 FTE

**National Scale**



Cargo Logistics Data  
(C1) DaL4LoD  
(C2) **Secure scalable policy-enforced distributed data Processing**  
(using blockchain)



Cybersecurity Big Data  
NWO COMMIT/  
**SARNET** project  
3.5 FTE

**City / regional Scale**

**Campus / Enterprise Scale**

**NLIP iShare project**



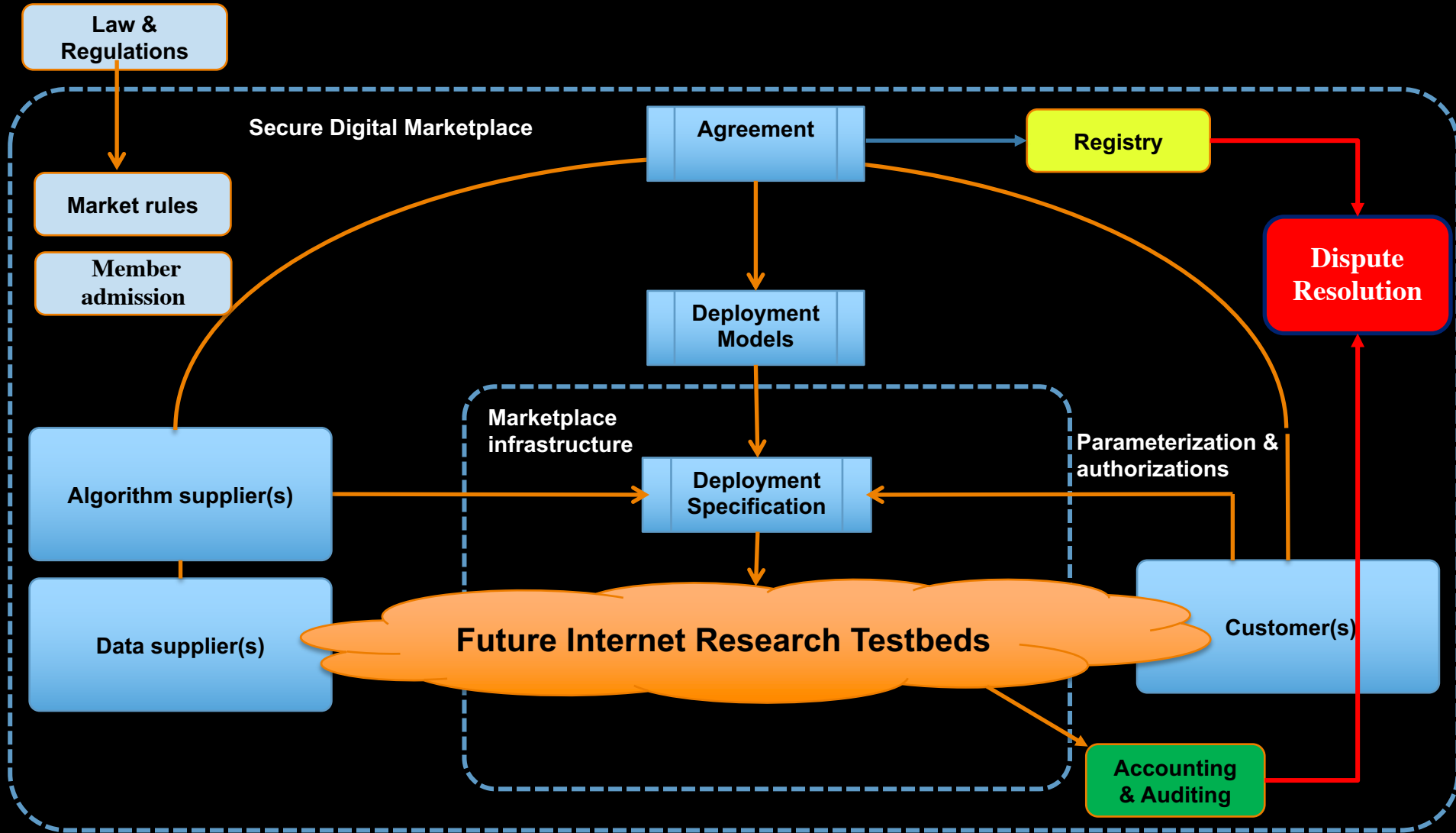
**iSHARE**  
powered by NLIP

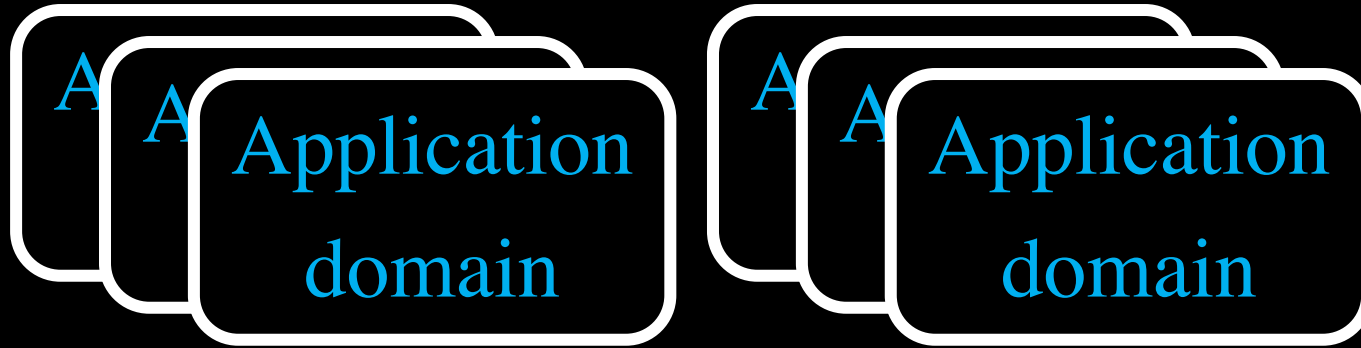
# Approach

- Strategic:
  - Translate legislation into machine readable policy
  - Define data use policy
  - Trust evaluation models & metrics
- Tactical:
  - Map app given rules & policy & data and resources
  - Bring computing and data to (un)trusted third party
  - Resilience
- Operational:
  - TPM & Encryption schemes to protect & sign
  - Policy evaluation & docker implementations
  - Use VM and SDI/SDN technology to enforce
  - Block chain to record what happened (after the fact!)



# Secure Digital Market Place Research





**AmDex**

Data objects & methods  
Data & Algorithms service

**FAIR / USE**

**AmsIX**

Routers - Internet – ISP's - Cloud  
IP packet service

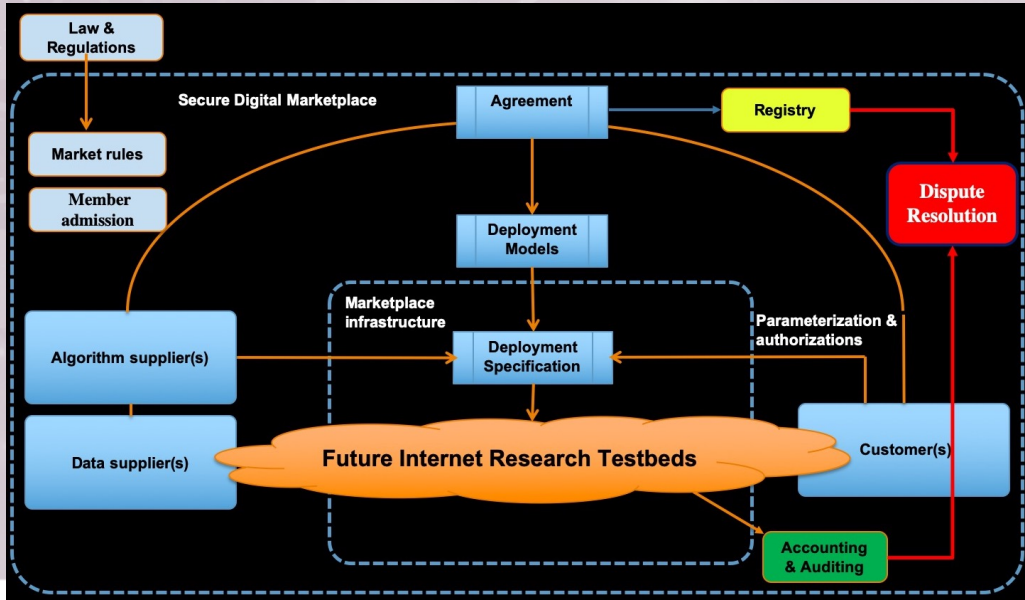
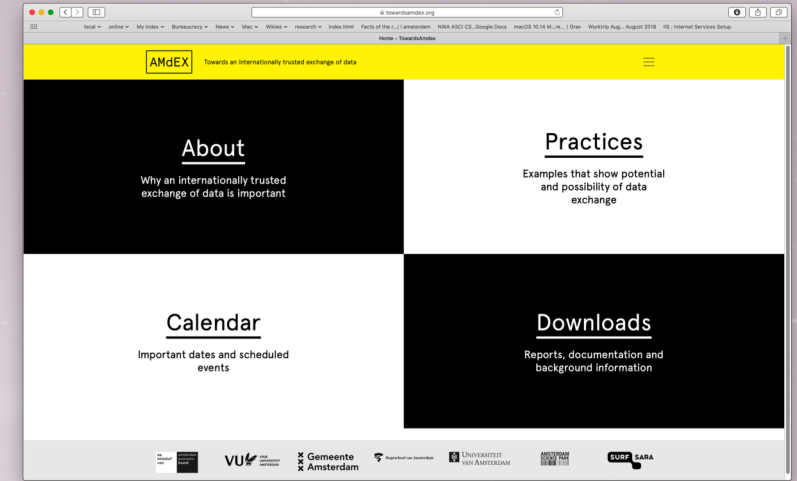
**IP / BGP**

Layer 2 exchange service  
Ethernet frames

**ETH / ST**

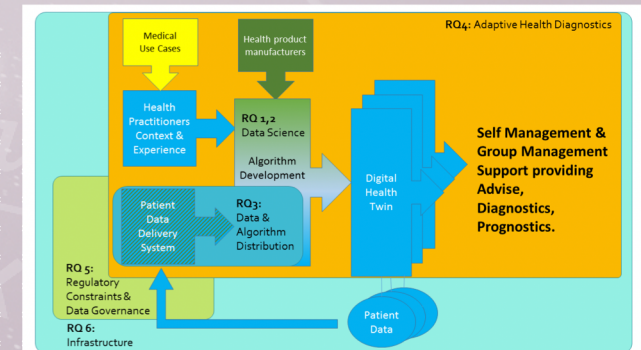
# AMdEX.eu

- Competing organisations, share data for common benefit
- Trust, Risk, data ownership & control
  - Industry: AF-KLM, Health, etc
  - Science: European Open Science Cloud
  - Society: Amsterdam Economic Board



Aircraft Maintenance AF-KLM

## Health: Enabling Personal Interventions





# SARNET: Security Autonomous Response with programmable NETWORKS

Marc Lyonnais, Leon Gommans, Rodney Wilson, Lydia Meijer, Frank Fransen Tom van Engers, Paola Grosso, Gauravdeep Shami, Cees de Laat, Ameneh Deljoo, Ralph Koning, Ben de Graaff, Gleb Polevoy, Stojan Travanovski.



## Big Data: real time ICT for logistics Data Logistics 4 Logistics Data (dl4ld)

Lydia Meijer (PI), Cees de Laat (Co-PI), Leon Gommans, Tom van Engers, Paola Grosso, Kees Nieuwenhuis.



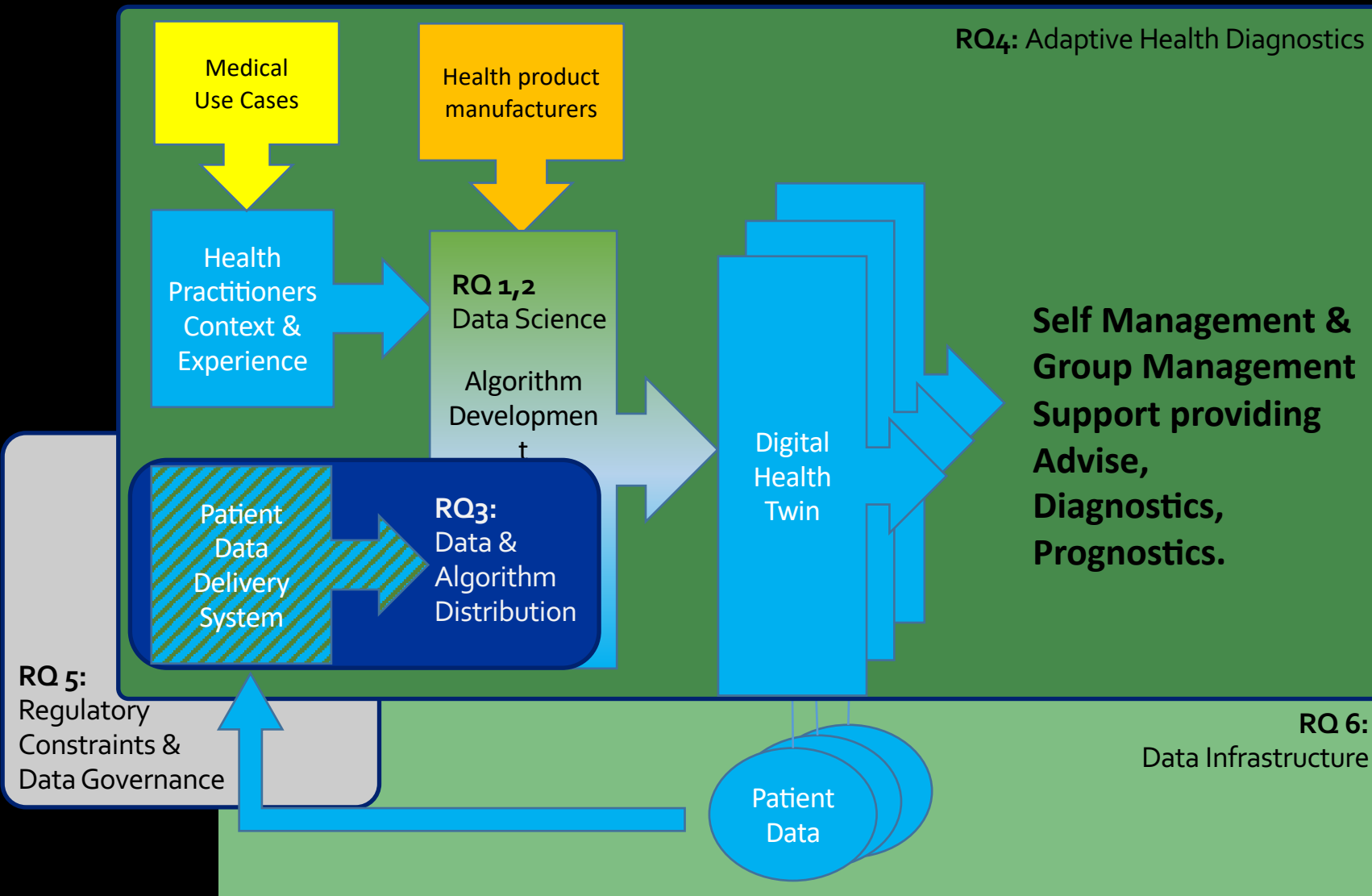
## EPI: Enabling Personalized Interventions

Cees de Laat(PI), Sander Klous (PL), Leon Gommans, Tom van Engers, Paola Grosso, Henri Bal, Anwar Osseyran, Aki Harma, Douwe Biesma, Peter Grünwald, Floortje Scheepers, Gertjan Kaspers.



# Health Use Case

## Enabling Personal Interventions



The overall aim of this project is to explore the use and effectiveness of data driven development of scientific algorithms, supporting personalized self- and joint management during medical interventions / treatments.

The key objective is to use data science promoting health practically with data from various sources to formulate lifestyle advice, prevention, diagnostics, and treatment tailored to the individual and to provide personalized effective real-time feedback via a concept referred in this proposal as a digital health twin.

# Research questions

- RQ1: **Dynamically Analyzing Interventions** based on Small Groups: how can we determine, based on as little data as possible, whether an intervention does or does not work for a small group or even an individual patient?
- RQ2: **Dynamically Personalizing the Group**: how can we identify effective intervention strategies and optimize personalization strategies applicable for different patient and lifestyle profiles via dynamic (on-line) clustering of patients? Can those clusters be adapted as new data about patients and results of interventions come in and as other data may be removed or modified?
- RQ3: **Data and Algorithm Distribution**: what are the consequences of a distributed, multi-platform, multi-domain, multi-data-source big data infrastructure on the machine learning algorithms and what are potential consequences on performance?
- RQ4: Adaptive health diagnosis leading to optimized intervention: how can we **enhance self- / joint management** by dynamically integrating updated models generated from machine learning from various data sources in state of the art health support systems that based on personal health records, knowledge of health modes and effective interventions?
- RQ5: **Regulatory constraints** and data governance: how can we create scalable solutions that meet legal requirements and consent or medical necessity-based access to data for allowed data processing and preventing breaches of these rules by embedded compliance, providing evidence trails and transparency, thus building trust in a sensitive big data sharing infrastructure?
- RQ6: **Infrastructure**: how can the various requirements from the use-cases be implemented using a single functional ICT-infrastructure architecture?

# SC16 Demo

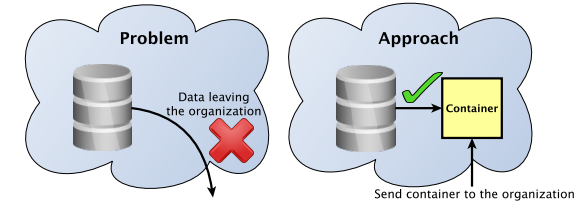
# DockerMon Sending docker containers with search algorithms to databases all over the world.

<http://sc.delaat.net/sc16/index.html#5>

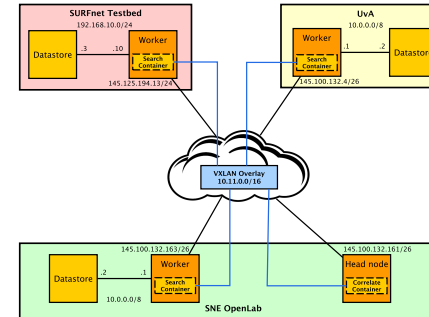
## Container-based remote data processing

### Problem Description

- Scientific datasets are usually made publicly available  
...but data cannot always leave the organization premises
- On-site data processing can be challenging because of incompatibility of systems or lack of manpower
- Can a container-based system perform remote on-site data processing efficiently?
- What are the networking issues to solve?



### Underlay and Overlay



#### Main features:

- Networked containers
- VXLAN overlay
- Containers that perform data retrieval and computation
- Containers built on-demand
- On-site data processing
- Distributed data source
- Multiple sites with datasets

### The Game

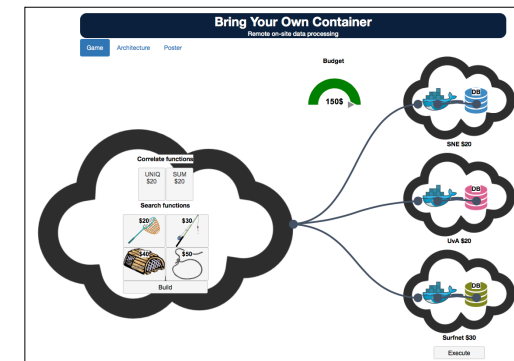
Our SC16 demo is a gamification of the remote dataset processing architecture.

How many different animal species can you find? You have a fixed budget and each function and processing will cost you money!

In our game you will:

- Select a correlate function to combine the results of the different sites.
- Pick different search functions, represented as tools, to find animals in the remote datasets.
- Build containers with the search and correlate functions.
- Execute the containers on the sites of your choice.

Will you have the best score?



# SC17 Posters and proof of concepts & demo's

<http://sc.delaat.net/sc17>

## Unlocking the Data Economy via Digital Marketplaces

Researching governance and infrastructure patterns in airline context

### Use Case: Sharing Aircraft Data to develop a Maintenance Credit System



- A **Digital Twin** estimates time before maintenance is needed after data is received from a corresponding aircraft system.
- Algorithm quality increases when data, owned by different airline operators, can be shared during its development.
- **Sharing data assets carries risk** (e.g. non-compliance).
- **Research Question:** "Can Digital Marketplace concepts organize trust amongst its stakeholders to enable common benefits no single organization can achieve, whilst observing economic principles?"

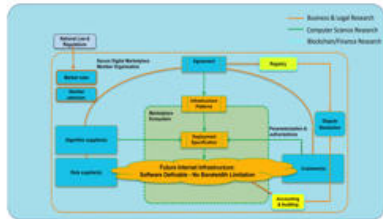
### Digital Marketplace as a means to organize trusted data asset sharing

A Digital Marketplace is a membership organization identified by a common goal: *Share data to enable development of a Maintenance Credit System.*

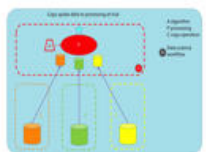
Membership organization is institutionalized to create, implement and enforce membership rules.

Market members create **digital agreements**.

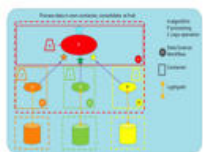
Agreements are translated into different software defined infrastructures using **infrastructure patterns** offered by a **Digital Marketplace Ecosystem**.



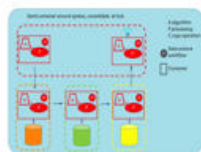
### Examples of infrastructure patterns offered by a Digital Marketplace



Public cloud model



Container model



Turntable model

## DEMONSTRATION: LIGHT PATHS AND DATA TRANSFER NODES FOR AIRCRAFT MAINTENANCE

Air France-KLM uses a 100 Gbit/s link, connected to Netherlight, to research an aircraft maintenance industry use case. Via this open exchange, Data Transfer Nodes (DTNs) of Air France-KLM in the Netherlands and iCAIR - present in Chicago at StarLight - connect to each other using light paths over their links. In this demonstration, users at SC17 in Denver will experience the difference in file transfer rates with and without using DTNs.

### USE CASE: AIRCRAFT MAINTENANCE

Besides people and luggage, aircrafts transport data they generate, like flight information, technical statistics and sensor readings. These data tell pilots and engineers if the aircraft's critical systems are doing their job safely. When data are transferred and analyzed rapidly, defects can be solved more quickly, possibly even while the aircraft is waiting at the gate. When receiving the data within minutes, expert engineers in a remote airport can readily verify with the home base engineers if an engine vibration warning was caused by the engine or by a falling sensor.

### INTERNET VS LIGHT PATH

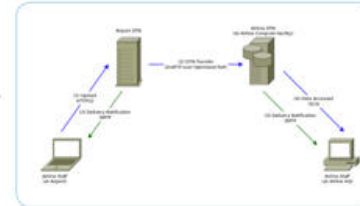
Air France-KLM uses a 100 Gbit/s light path and researches its benefits. Using light paths, you can transport huge amounts of data at high speed and with a guaranteed bandwidth between 2 points. When using high volumes of data, the current Air France-KLM's internet connections are not private or fast enough to transfer the data within the requested time frame. Transferring a terabyte of engine data via the current internet connections would take around 30 hours, with a 100 Gbit/s light path this could take less than 2 minutes.

### AIR FRANCE-KLM CONNECTED TO NETHERLIGHT

Ciena and SURF facilitate the connection from Air France-KLM at Schiphol to NetherLight, SURF's European hub for international light paths in Amsterdam. SURF provides the 100 Gbit/s light path from Air France-KLM via NetherLight to the aircraft's destination. For this demonstration, the location is StarLight in Chicago, a hub similar to NetherLight.

### DATA TRANSFER NODES

Data Transfer Nodes are high-performance systems that are optimized to transfer huge amounts of data. The interconnects between these systems exist of high-capacity dedicated bandwidth, removing network bottlenecks within the mesh of global DTNs. To date, DTNs are present on a small scale, e.g. a couple per continent. By copying a file from an end user system directly into the nearest DTN, the global DTN system sends the file to the DTN nearest to the final file destination, optimizing the process of high-latency international transfers.



### LIVE DEMONSTRATION

In this demonstration, there are two end user systems, one in Amsterdam and one in Chicago. Neither system will be optimized for long range transfers, however each will have access to a nearby Data Transfer Node. Visitors are allowed to transfer pre-prepared datasets between the systems via the DTNs with graphs showing various performance metrics. As a comparison, the performance of a direct connection between the two systems - without using DTNs - will also be shown. The intention is to show that systems not optimized for long distance transfers can benefit from using nearby DTNs to facilitate the transfer and decreasing file transfer time.

### RESEARCH IN OTHER INDUSTRIES

In addition to the aircraft research, high bandwidth, low latency light paths offer possibilities for research in other industries as well. For example, fundamental research on data transfer protocols suitable for these bandwidths can also help excel diagnosis by doctors when they can have access to terabytes of patient and other related research data within minutes, instead of days or weeks. Imagine what this would enable other research disciplines to do too. Possibilities are almost infinite!

More information: [www.surf.nl/en/100-G-Air-France-KLM](http://www.surf.nl/en/100-G-Air-France-KLM)



## Data Transfer Node (DTN) Workflows

Joseph Hill, Gerben van Malenstein, Cees de Laat, Paola Grosso, Leon Gommans

### Why Data Transfer Nodes (DTNs)

- DTNs can act as an interface to a high performance link
- Configured to maximize performance for a given workflow
- Simplifies configuration of client systems
- Multiple clients may share a DTN
- DTNs strategically placed to best benefit clients
- DTNs can be compared to specialized high speed transport systems of the past

Pneumatic Tube Messaging System, 1943



United States Library of Congress's Press and Photograph Division (Digital ID 55a-82308-1)

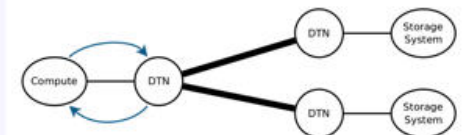
### Example: Entry Point for High Speed Transport

A typical use case for DTNs is as a high speed file transfer service. A computer system's configuration may allow for the utilization of all available bandwidth in a LAN environment. However, it is often the case that in a WAN environment with high latency or packet loss the same system performs poorly. A DTN could be tuned to maximize performance on a high latency path. It could also use specialized transfer protocols to mitigate high packet loss. The DTN may also have access to an optimized path such as a light path. Files destined for a distant receiver would be first sent to a DTN located on the same LAN as the sender. That DTN would then forward it to a DTN near the receiver. That DTN would then forward it to the final destination.



### Example: Storage Access Point

Another possible use case for DTNs is to be used to access distributed data from remote locations. In this scenario a system located at a compute facility requests the data from the local DTN as it is required. That DTN would then transparently retrieve the data from multiple remote sites as needed. In contrast to the first example here block level access is provided by the DTNs. To the system performing the computations the nearby DTN appears to be the actual and only storage system. This hides both the remote and distributed nature of the data. While the compute side DTN may perform some caching, there need not be permanent storage of data at the compute facility.

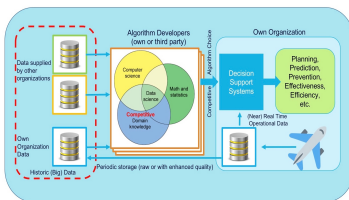


# SC18 – Dallas TX

## Training AI/ML models using Digital Data Marketplaces

Creating value and competition by enabling access to additional big data owned by multiple organizations in a trusted, fair and economic way

### The more data - the better: an aircraft maintenance use-case



- AI/ML algorithm based Decision Support Systems create business value by supporting real-time complex decision taking such as **predicting the need for aircraft maintenance**.
- Algorithm quality increases with the availability of aircraft data.
- Multiple airlines operate the same type of aircraft.
- **Research Question:** "How can AI/ML algorithm developers be enabled to access additional data from multiple airlines?"
- **Approach:** Applying Digital Data Marketplace concepts to facilitate trusted big data sharing for a particular purpose.

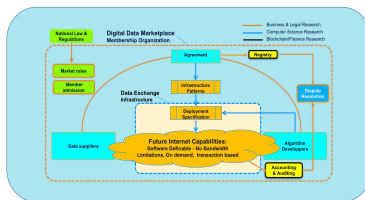
### Digital Data Marketplace enabling data sharing and competition

A **Digital Data Marketplace** is a membership organization supporting a common goal: e.g. enable data sharing to increase value and competitiveness of AI/ML algorithms.

Membership organization is institutionalized to create, implement and enforce membership rules organizing **trust**.

Market members arrange **digital agreements** to exchange data for a **particular purpose** under specific conditions.

Agreements subsequently drive data science transactions creating processing infrastructures using infrastructure patterns offered by a Data Exchange as **Exchange Patterns**.

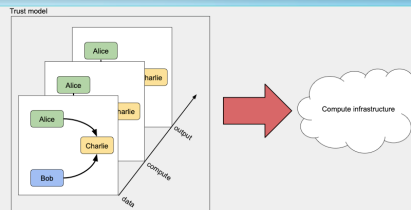


## Dataharbours: computing archetypes for digital marketplaces

Reginald Cushing, Lu Zhang, Paola Grosso, Tim van Zalingen, Joseph Hill, Leon Gommans, Cees de Laat, Vijaay Doraiswamy, Purvish Purohit, Kaladhar Voruganti, Craig Waldrop, Rodney Wilson, Marc Lyonnais

### The problem

How can competing parties share compute and data? The architecture of a digital marketplace is an active research field and has many components to it. Here we investigate a federated computing platform which is molded into different **archetypes** based on **trust** relationships between organizations.



### The components

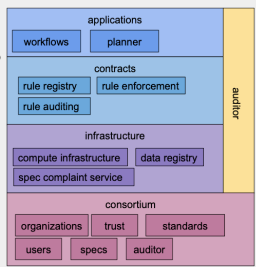
**Consortium:** is an initial document which brings together organizations that wish to collaborate. It defines static information such as keys to identify parties.

**Infrastructure:** A single domain organization infrastructure that securely hosts data, compute containers and, optionally, compute infrastructure. We dub this infrastructure a **data harbour**. A harbour implements a set of protocols that allows it to interact with other harbours.

**Contracts:** Are a set of rules that are shared amongst participating harbours which describe how objects (data, compute) can be traded between harbours and who can process data. In its simplest form is a 7-tuple which binds a user, data object, compute container, contract, consortium, harbour, and expiry date.

**An application:** Is a distributed pipeline which can make use of several contracts. The combination of application and contract defines the archetype of the computation i.e. how data and compute are moved to effect computation.

**Auditor:** A trusted entity that collects audit trails for use in litigation of policy violations.

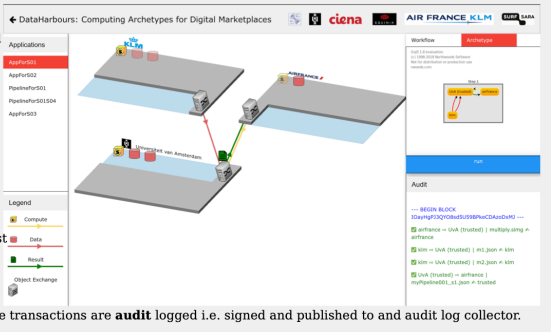


### In action

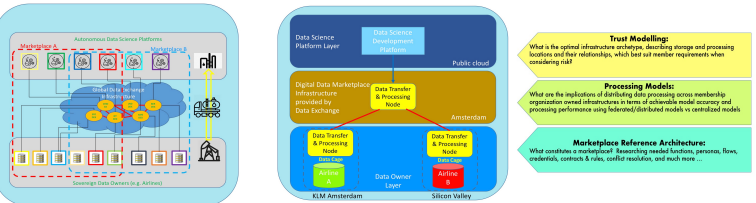
Federated computing on 3 distributed data harbours. Here we illustrate one archetype where KLM and Airfrance do not trust each other and employ a trusted 3rd party to send the data and compute for processing.

For the scenario to succeed the different harbours need to effect several transactions which are governed by contractual rules.

The transaction **protocol** involves first identifying both parties are who they say they are through pub/priv key challenges and secondly, that at least a **contract rule** is matched to allow the transaction. Important steps of the transactions are **audit** logged i.e. signed and published to and audit log collector.



### Researching Exchange Patterns to support Digital Data Marketplaces



The screenshot shows the SC18 website with the URL 'http://sc.delaat.net'. It features a navigation bar, a list of presentations (e.g., 'Data Harbours: A compute infrastructure for data marketplaces', 'Building User-friendly Data Transfer Nodes'), and a list of booth members (e.g., AIR FRANCE KLM, TRANSFIDES, evofenedex, ciena, THALES, SCinet). A group photo of the SC18 booth is also visible.



# CONCLUSIONS

- Overall advice
  - It is about people & knowledge
  - Base on society relevant applications
  - Get faculty drivers from each campus
  - Governance model is essential
  - align with education (soft&hard money)
- Applications
  - Health
  - Instrumenting IOT
  - Energy transition/critical infrastructures IT
  - CyberSecurity

# CONCLUSIONS

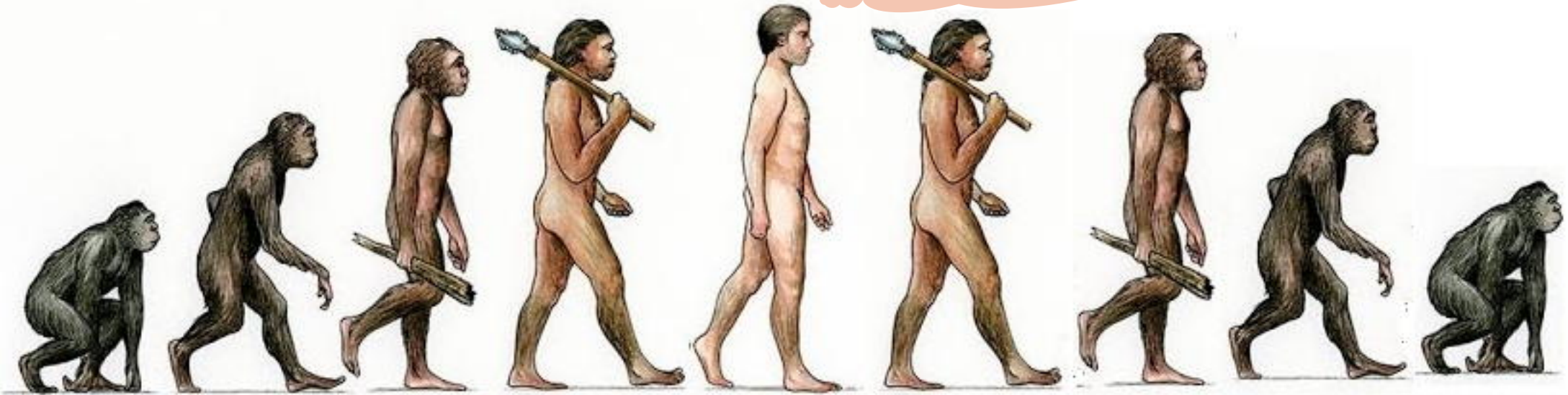
- Themes
  - global data & methods ecosystem supporting applications
  - Explainable AI to aid managing CI
  - Security
  - Super-facility
  - revisiting Internet standards with current technology in mind
  - Quantum compute and networking
- Quotes:
  - Wouter Los: Any future data infrastructure should accommodate the preferred governance model related to the cultural dimension. What are checks and balances, and who decides (has power) on what? How is this framed in the context of (self regulating) micro markets, when billions of agents interact.
  - Tom Defanti: ML is like training your dog without knowing how the dog works.
  - Larry Smarr: Manage the exponential.
  - Mike Norman: It is not about hardware, it is about the people!!!
  - Inder & me: The kids of today are the decision makers tomorrow and have no feeling for classic CI.



# Past & future ICT research infrastructures

- TEN34 / TEN155
  - Geant testbed & JRA's
  - FIRE
  - Grid5000 (FR)
  - DAS1-5 (NL)
-  Was connected by LightPath around 2010!
- Need for breakable CS oriented testbed
  - Must include: Programmable networks, Cloud, Exascale SC, DTN's, streaming, access to public services, IOT, Wireless
  - Must include work on AI & ML, fundamental data security
  - At Scale → SILECS - <https://www.silecs.net>

# AI forking off



Artificial Intelligence

NOW

# Conclusions, Info, Acknowledgements, Q&A

- Need for Network to Data level experimental Infrastructure, Europe's own DTN infra, CC program, CI Ambitions, Data at scale.
- Data hindered by risk of unexpected use, lack of trust; Using market principles, enforcement and determining incentives and value in the data life cycle to make data flow
- More information:
  - <http://delaat.net/dl4ld>      <http://delaat.net/sarnet>      <http://delaat.net/epi>      <https://towardsamdex.org>
  - <https://www.esciencecenter.nl/project/seconnet>

P.S. I did not mention Quantum Compute & Networking; See:

- <https://www.orau.gov/quantumnetworks2018/default.htm>
- [https://science.energy.gov/%7E/media/ascr/pdf/programdocuments/docs/2019/QNOS\\_Workshop\\_Final\\_Report.pdf](https://science.energy.gov/%7E/media/ascr/pdf/programdocuments/docs/2019/QNOS_Workshop_Final_Report.pdf)
- <https://delaat.net/qn>

