





# Lecture 1011

## Federation & scaling approaches for exascale data

Lecturer: Dr Jeremy Coles (jc574)



# Lecture 1011

(1011)<sub>2</sub> = (1 × 2<sup>3</sup>) + (0 × 2<sup>2</sup>) + (1 × 2<sup>1</sup>) + (1 × 2<sup>0</sup>)

## Federation & scaling approaches for exascale data

Lecturer: Jeremy Coles (jc574)



Credit: SKAO



# Before we start ...

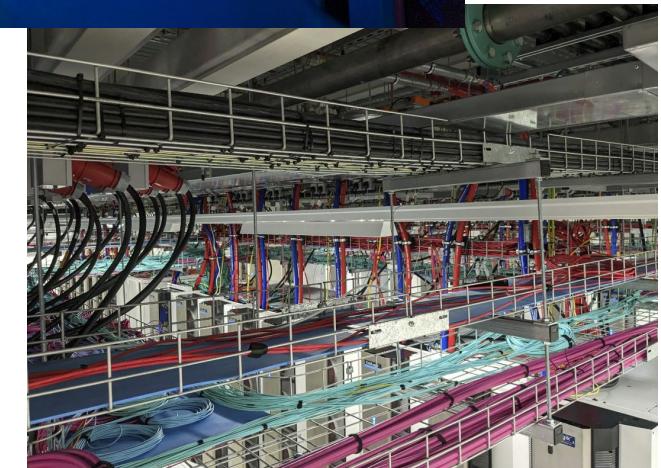
Data Centre tours. 5 signed up for Monday 26th.

We're also running some visits on Wednesday 28th. But **you must be signed-up in advance AND bring your photo ID on the day.**

CSD3 - still in top 500:  
<https://www.top500.org/site/49143/>

**Dawn** - The Cambridge Open Zettascale Lab is hosting Dawn, the UK's fastest artificial intelligence (AI) supercomputer. <https://www.hpc.cam.ac.uk/d-w-n>

6PB Lustre system ...

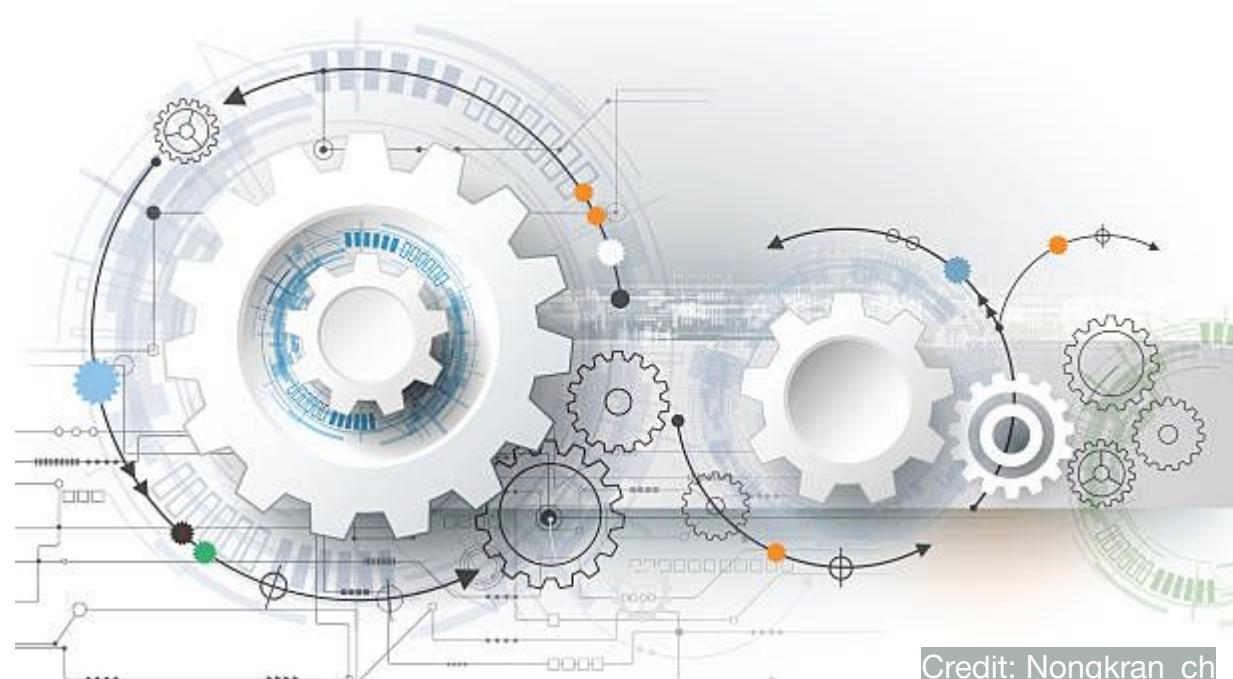


Credit: Joe Bishop



# Conclusions

- There are lots of moving parts in large data projects
- This lecture can only cover some aspects
- Existing computing solutions cover specific use-cases. For federating we often evolve them.
- Project constraints and capabilities require constant trade-offs
- High-Performance Computing (HPC), and federated computing, is a research topic in itself!



*Spoiler alert*



# Why bother?

<http://www.>



Credit: Paul Clark



# Be curious!



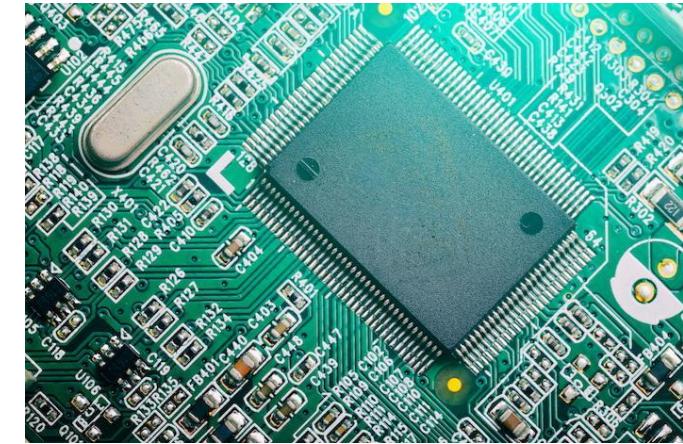
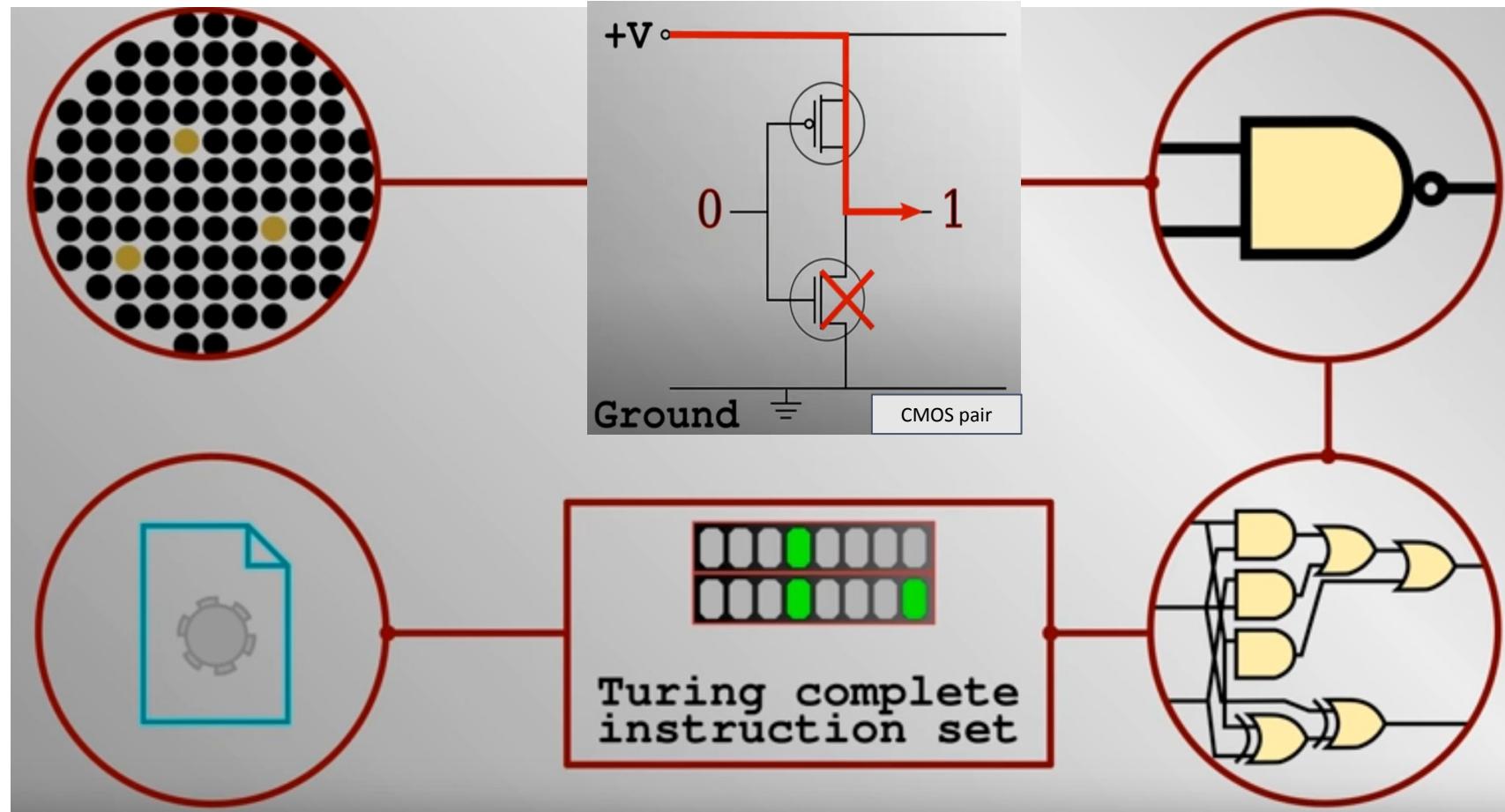
## Infrastructure

*Infra-* means "below;"

*Structure* From [Middle French \*structure\*](#),  
from [Latin \*structūra\*](#) ("a fitting together...")

The infrastructure is the "underlying structure" ... the fixed installations that it needs in order to function.

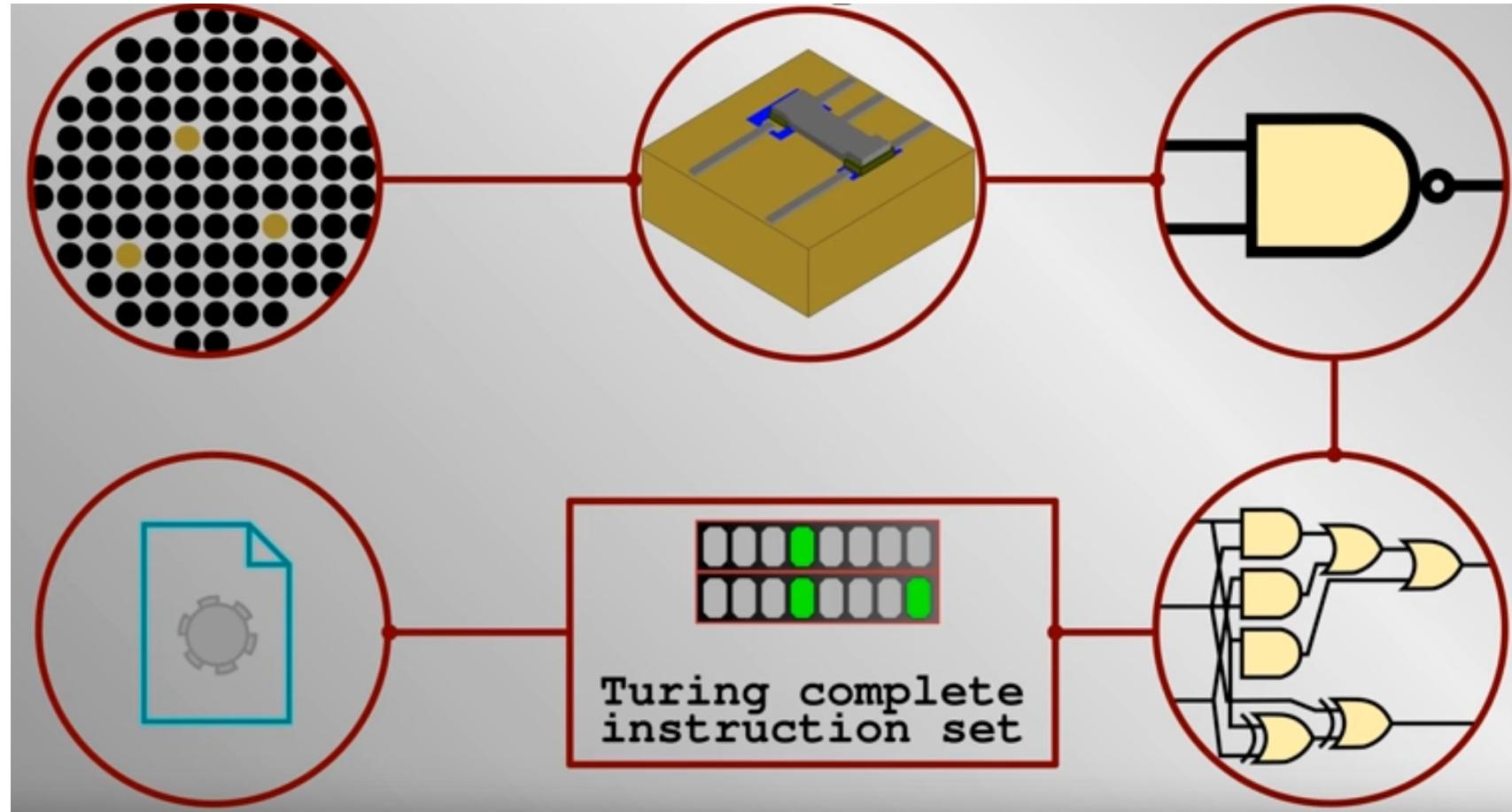
# How compute, storage, switch works



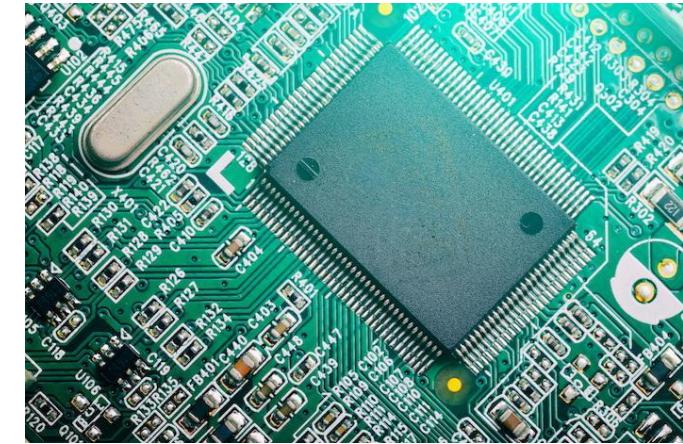
# Applying voltages drives switching

Credit: IM

# How compute, storage, switch works

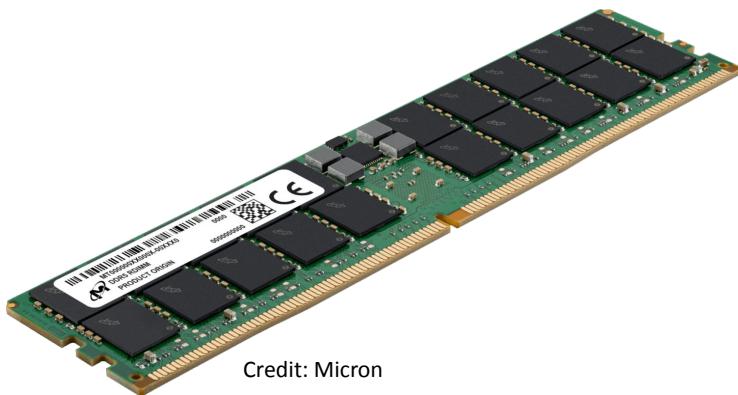


Credit: IM





# Core components



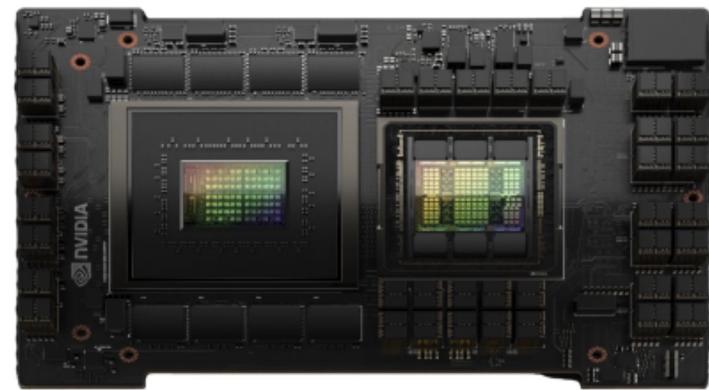
Dynamic Random Access Memory (DRAM) - billion capacitor memory cells (17 nanoseconds)



Solid State Drive (SSD) 3D trillion memory cells



Hard disk drive



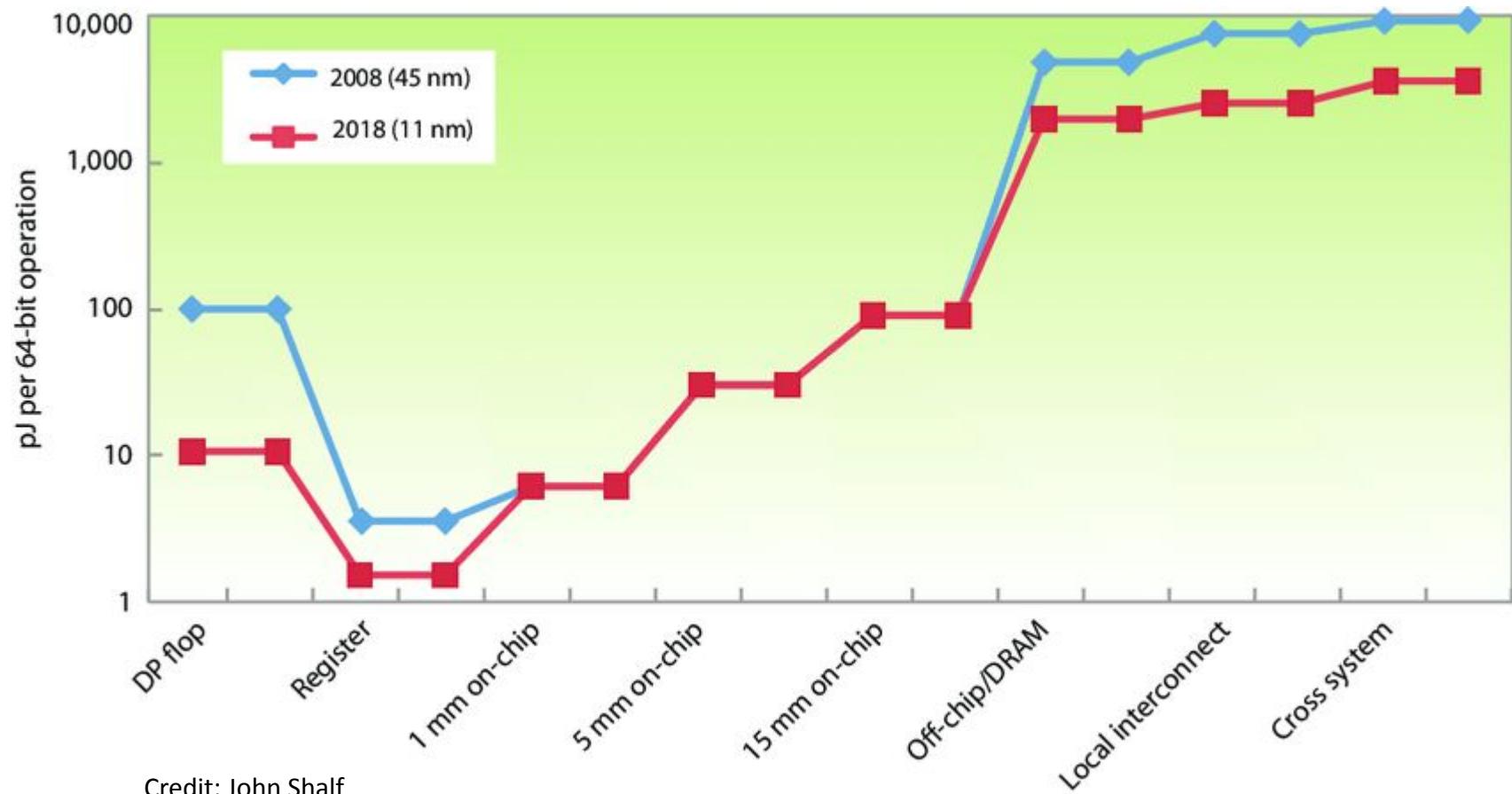
Credit: nVidia

Memory: Cache | registers | DRAM | SSD | HDD. Speed vs capacity - cost

Processing: CPU vs GPU and combinations!

*Each consumes energy to work*

# Bits, energy and cost



*Every bit has a cost*

# How does this impact Data Science activities?

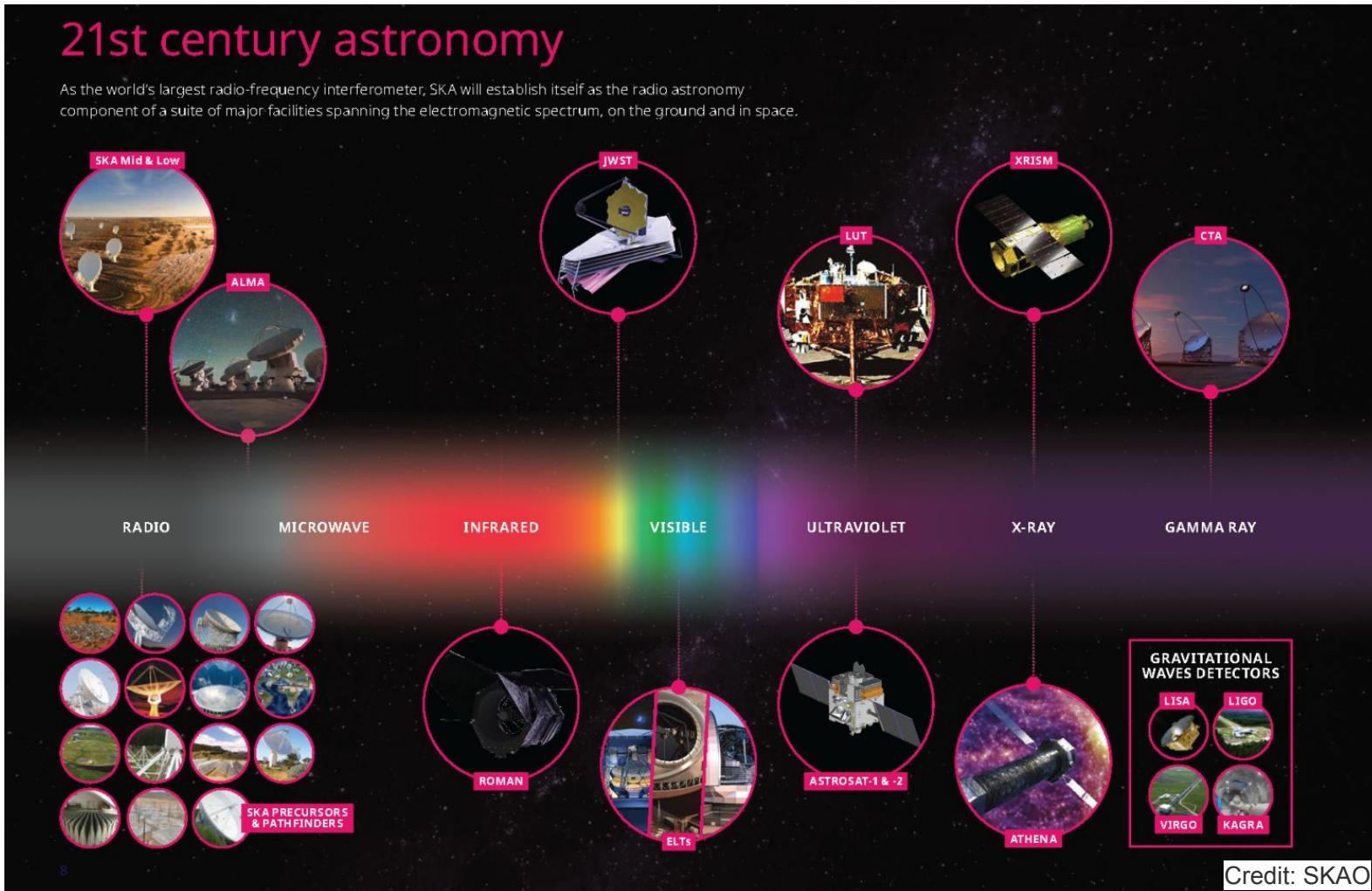


- Maintenance
- Problem solving
- Capabilities
- Performance
- Connectivity
- Optimisation
- Total and hidden **costs**
- ...

*Informed tradeoffs are needed*



# A reminder



Radio and visible wavelengths are less absorbed by the Earth's atmosphere



# A reminder

The screenshot shows a news article from **nature astronomy** about a black hole. The article title is **The accretion of a solar mass per day by a 17-billion solar mass black hole**. The authors listed are Christian Wolf<sup>1,2</sup>, Samuel Lai<sup>1</sup>, Christopher A. Onken<sup>1</sup>, Neelesh Amrutha<sup>1</sup>, Fuyan Bian<sup>3</sup>, Wei Jeat Hon<sup>1,4</sup>, Patrick Tisserand<sup>1,5</sup> & Rachel L. Webster<sup>1,4</sup>. The article was received on 19 April 2023, accepted on 4 January 2024, and published online on 19 February 2024. A 'Check for updates' button is present. The article discusses the cataloging of quasars and the properties of the most luminous ones. At the bottom, there is a collage of various space and astrophysics projects, including the SKA, Roman, ELTs, AstroSat-1 & -2, ATHENA, LISA, LIGO, VIRGO, and KAGRA. A red annotation on the right side of the slide reads: *Lots of exciting discoveries to come!*



# The SKA context

## SKA1-Mid

the SKA's mid-frequency telescope



Location: South Africa

Credit: SKAO



Frequency range:  
**350 MHz**  
to  
**15.4 GHz**  
with a goal of 24 GHz



**197 dishes**  
(including 64 MeerKAT dishes)



Maximum baseline:  
**150km**

## SKA1-Low

the SKA's low-frequency telescope



Location: Australia



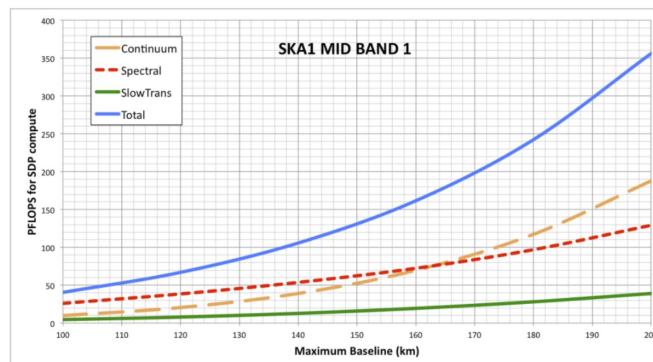
Frequency range:  
**50 MHz**  
to  
**350 MHz**



**131,072**  
antennas spread between  
512 stations



Maximum baseline:  
**~74km**

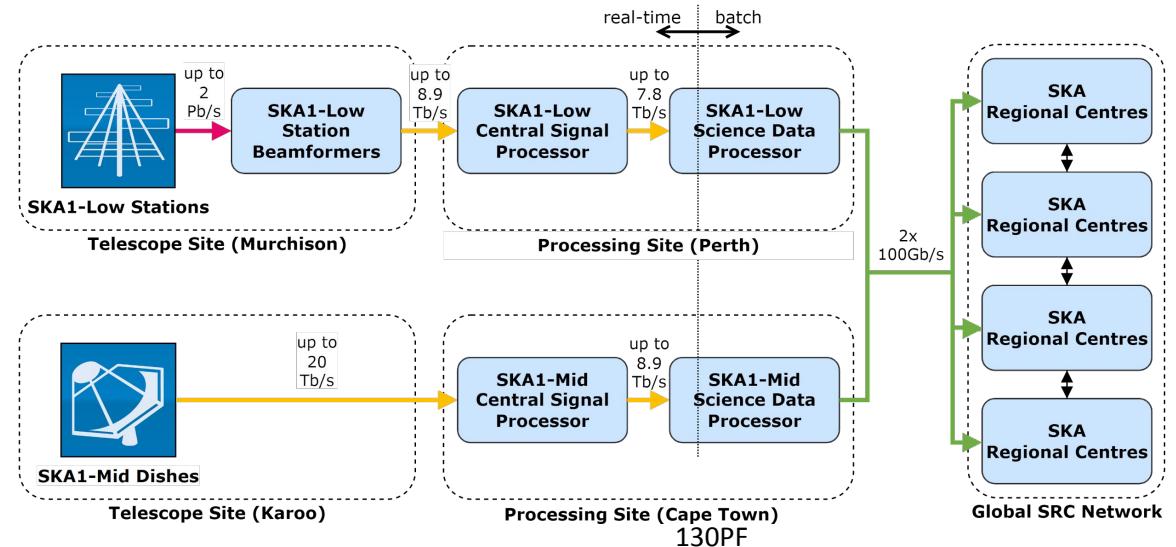
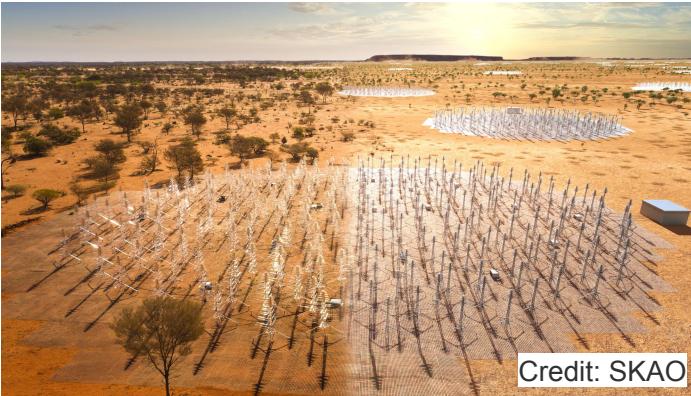


One petaFLOPS is equal to 1,000,000,000,000 (one quadrillion) floating-point operations per second!

A floating-point operation is an addition or multiplication of two real (or floating-point) numbers represented in some machine-readable and manipulatable form.



# SKA - End-to-End data movement



SKA Project budget: \$2.2B - up to Science Data Processor (SDP)

Lower limit for each telescope just for data in year 1:

Commercial cloud storage: \$0.021/GB → \$2.73M (\$4M to buy)

Data put/retrieve: \$0.01/GB → \$1.30M

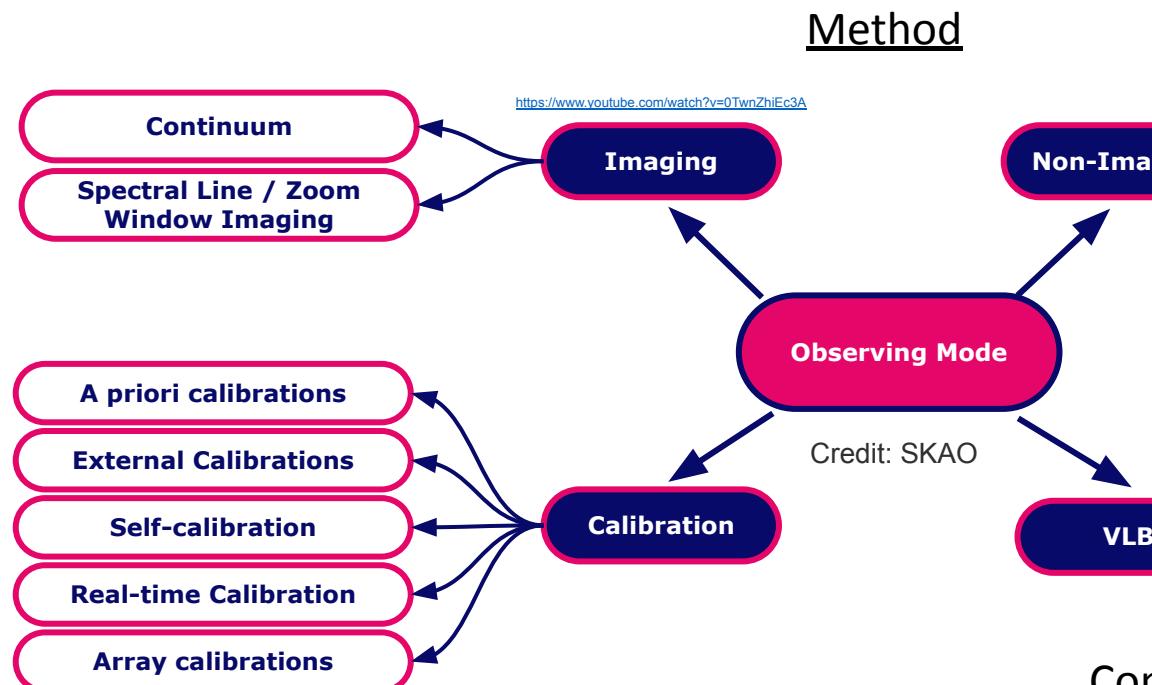
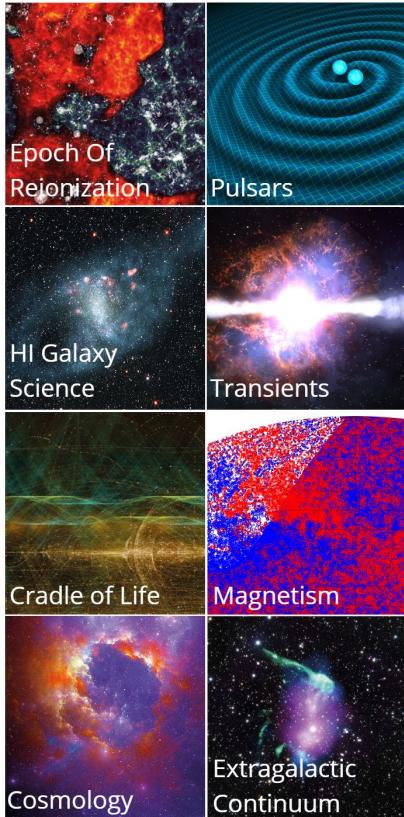
Cost vs energy vs performance vs science vs ....

*Project tradeoffs*

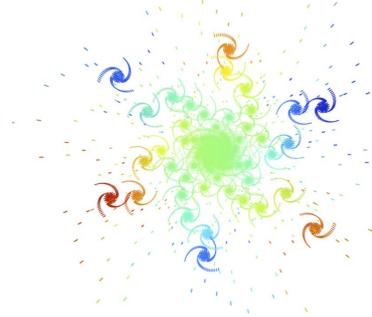


# Some of our inputs

## Science



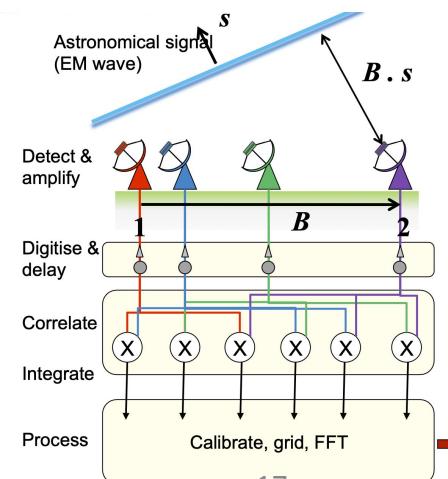
## Complexity



*Considerations are complicated!*

**VLBI** = Very-long-baseline interferometry is a type of astronomical interferometry used in radio astronomy. In VLBI a signal from an astronomical radio source, such as a quasar, is collected at multiple radio telescopes on Earth or in space.

## Scale



# Functional vs Non-Functional Requirements



- Functional Requirements direct the system on a functional level.
- Non-Functional Requirements (NFRs) define the **Quality Attributes**.
- Well-defined non-functional requirements:
  - are measurable,
  - can be verified (tested),
  - provide teams with ability to test the outcome of small task, feature, project or process against a defined set of success factors.
- NFRs define what does the success look like before detailed work begins.



*NFRs define the Quality Attributes*

# Science Data Processor (SDP)

## Quality Attributes



UNIVERSITY OF  
CAMBRIDGE

### Performance & Scalability

- ❖ Compute, I/O & Storage
  - >10 Pflop/s effective
  - ~0.77 TB/s ingest rate
  - ~4 TB/s into processing
  - >40 PB tiered buffer
- ❖ Need to scale
  - Trivial (e.g. Ingest) and expensive (e.g. ICAL) workflows co-exist
  - Early Array Assembly workloads will be unrepresentative
  - SKA “>1” will be even harder on SDP

### Modifiability & Maintainability

- ❖ Long lifespan (>50 yrs)
- ❖ Software changes
  - Execution Engines
  - Science Workflows
  - Processing Components
  - Data Models
- ❖ Hardware changes
  - Processing
  - Storage
  - Network

### Buildability, Affordability

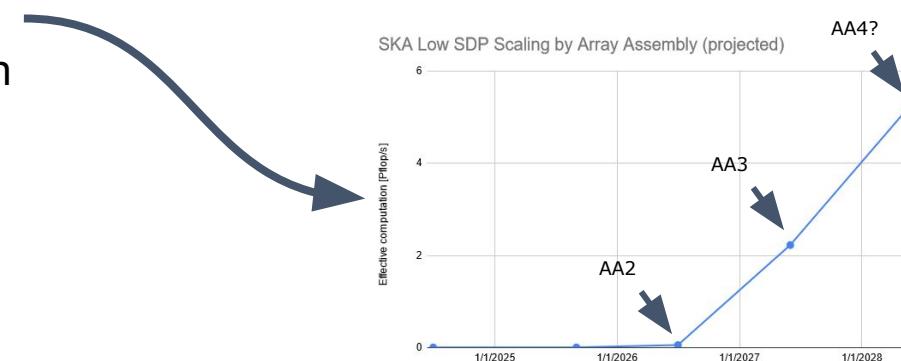
- ❖ COTS components
  - Get going quickly
  - Externalise/share maintenance
- ❖ Support agile development
  - Small functional increments
  - Parallel work of different teams

### Testability

- ❖ Continuous Integration
  - Test outside SDP
  - Support different development speeds
- ❖ Ensure scientific validity
  - Must trust pipelines with autonomous analysis

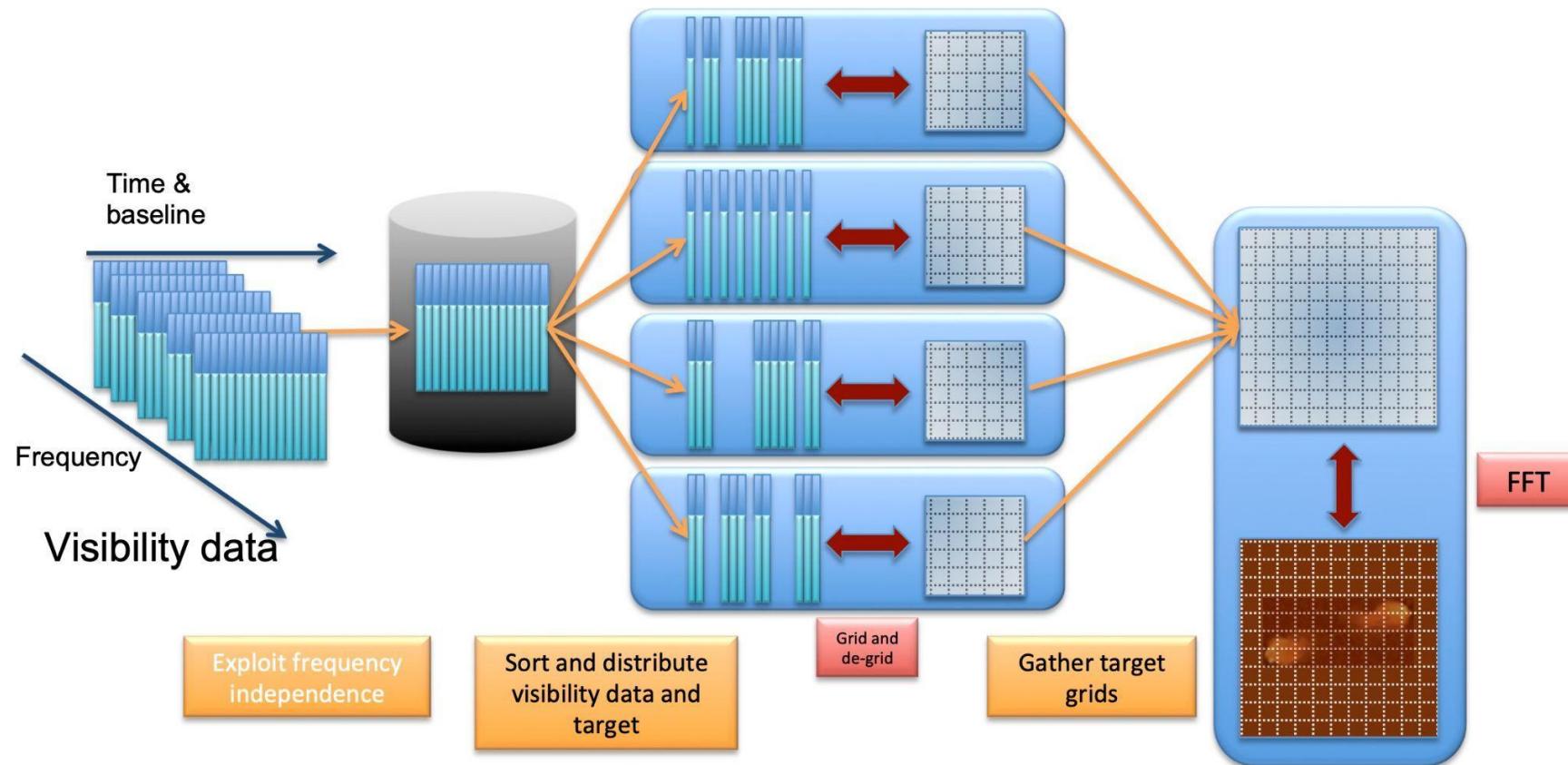
### Portability

- ❖ SKA Low, SKA Mid, SRCs



*Key considerations for SDP*

# Mapping



- Parallelize by frequency and time steps
- Can be used for load balancing

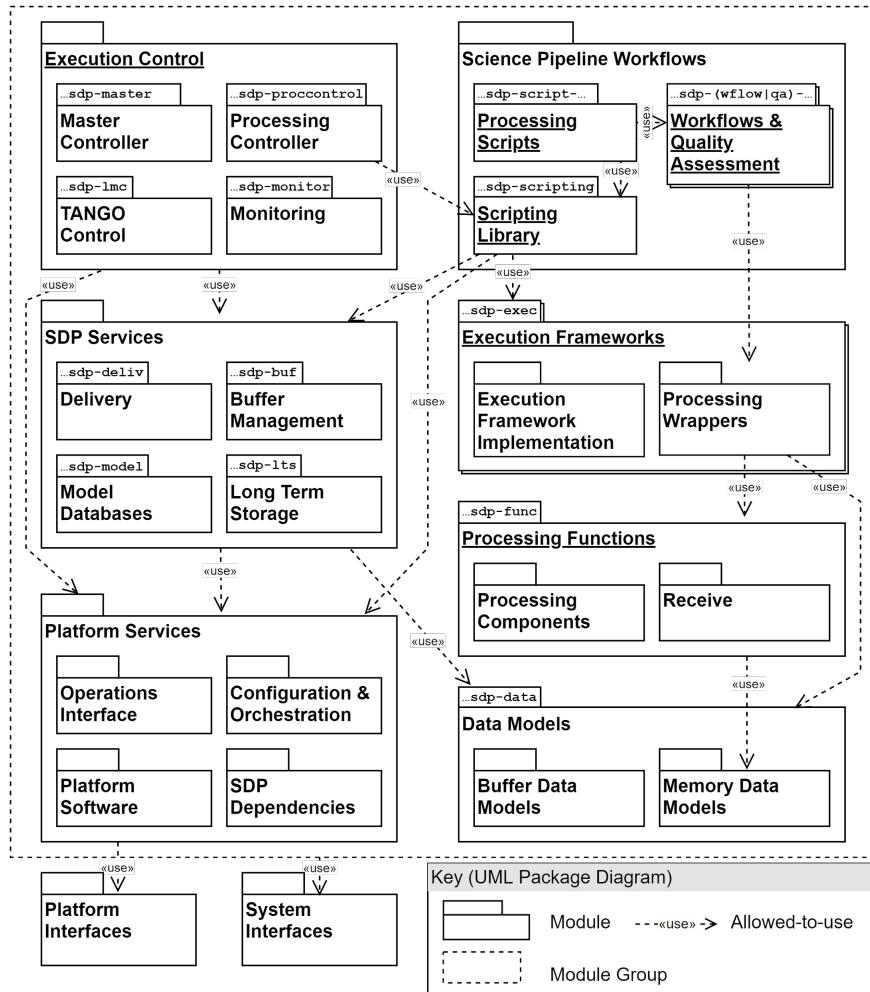
*Exploit natural parallelism in the data*

FFT = Fast  
Fourier  
Transform

# SDP Module View



UNIVERSITY OF  
CAMBRIDGE

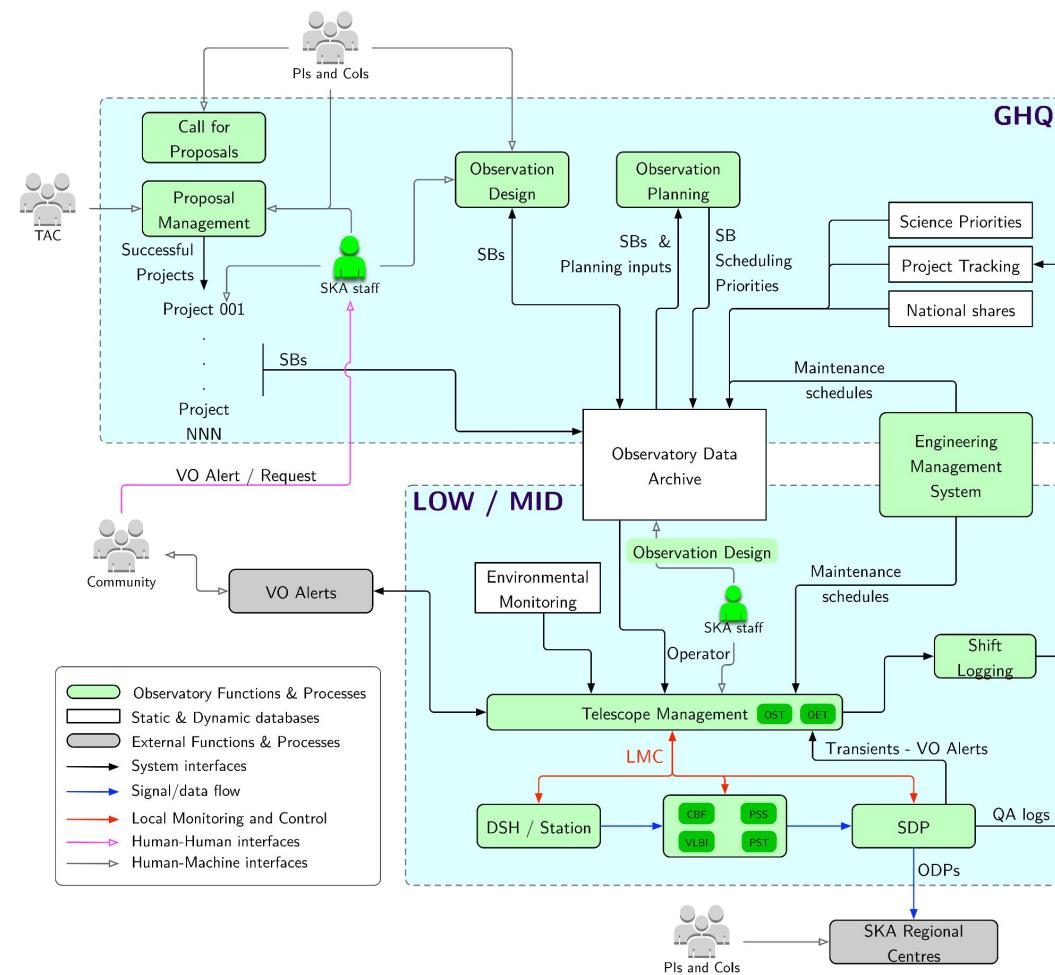


Layered processing architecture:

- Workflows  
(Interface of processing for rest of system)
- Execution Frameworks  
(Coordinate data movement and processing)
- Processing Functions  
(Algorithmic kernels)
- Data Models  
(Standards and libraries for data exchange)

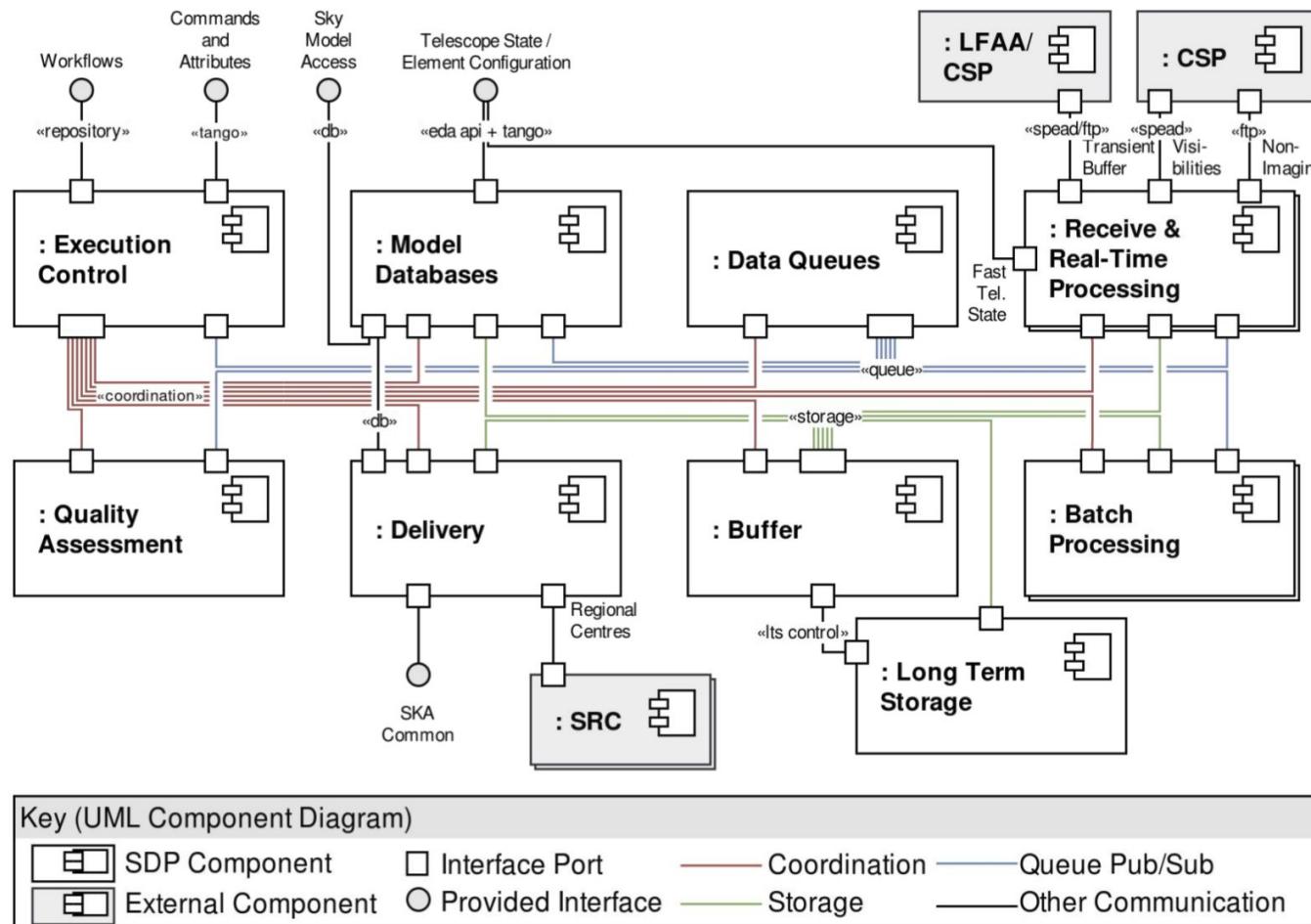
*Architecture views help us reason about the software*

# Aside: The broader integration challenge



The SKAO software has various interfaces beyond SDP

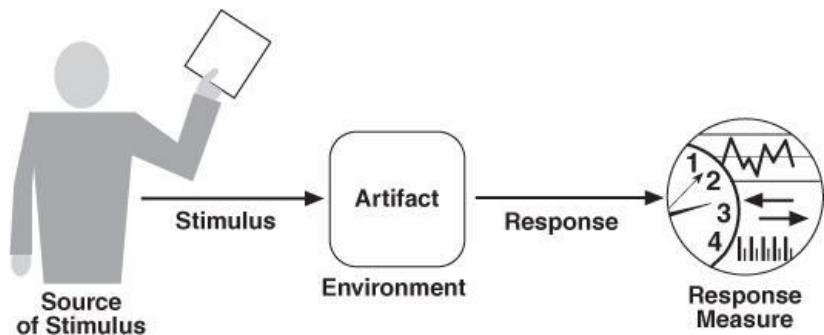
# High-level Component & Connector view



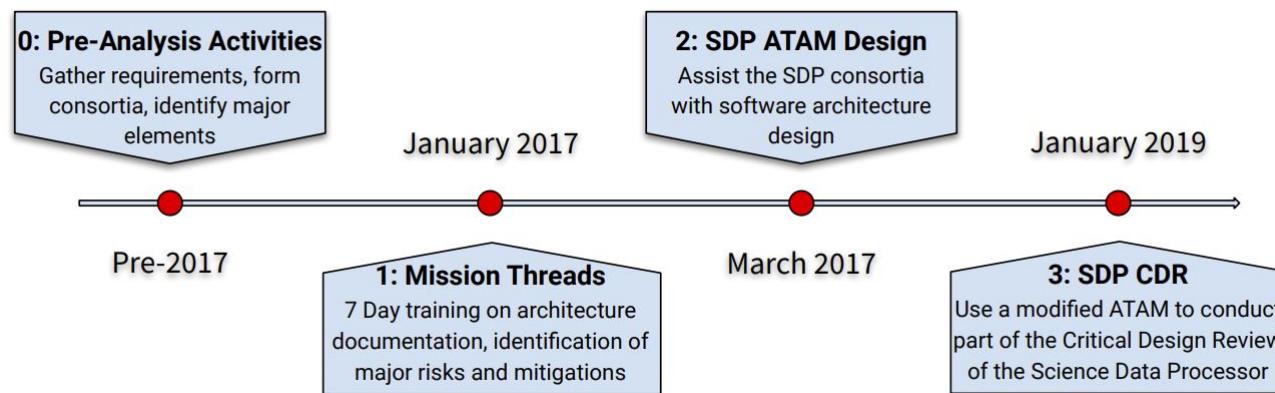
We need to map the  
data interfaces



# A way to approach trade-offs



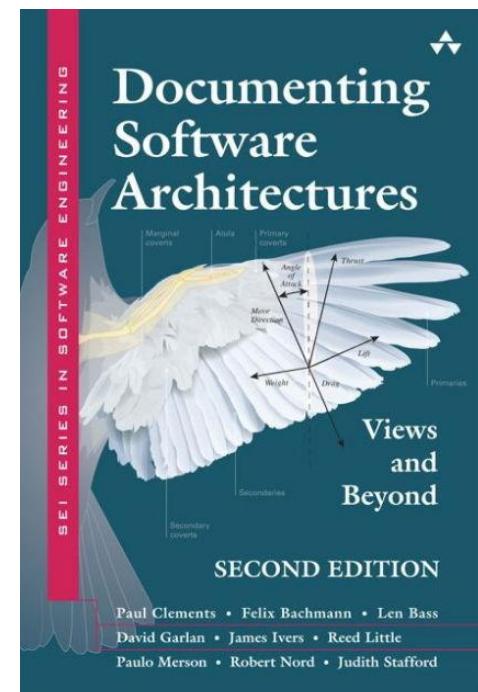
A scenario stresses a **quality** of the system ... the Non-functional Requirements



*There are methods to help us validate the architecture*

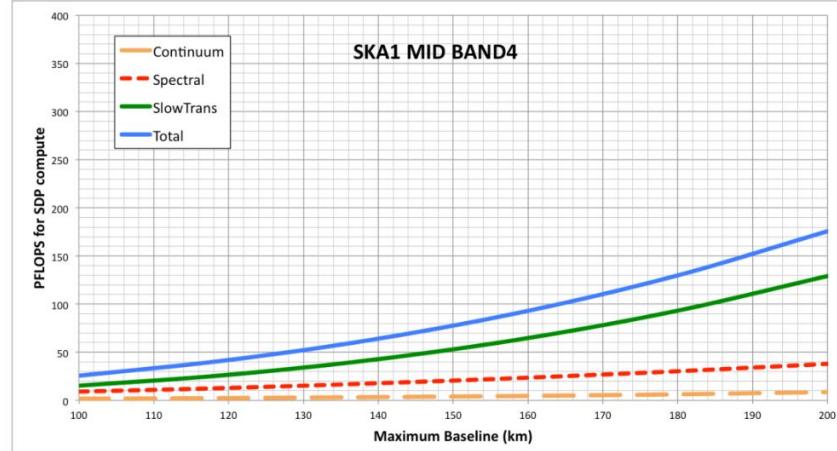
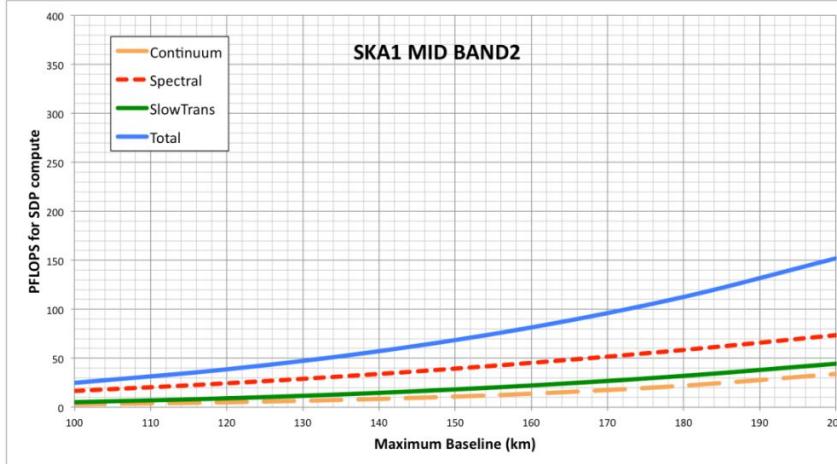
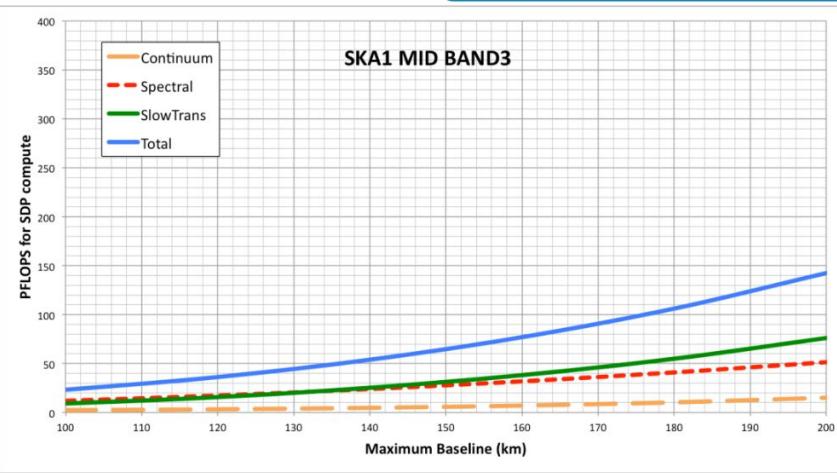
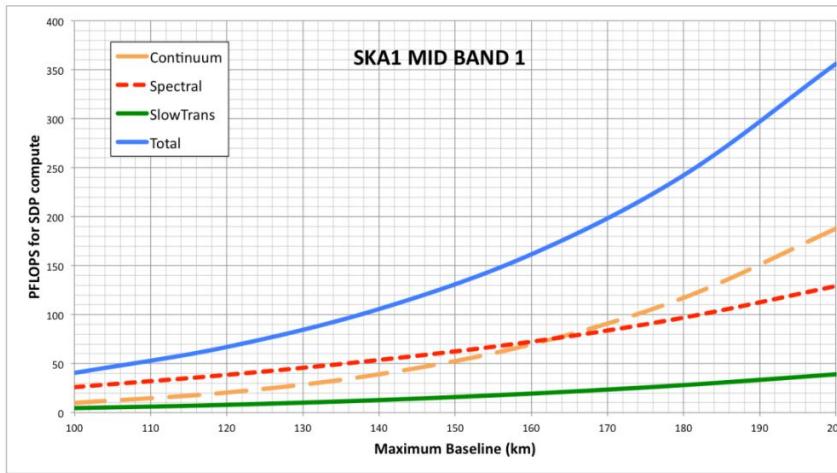
*"Software Architecture refers to the high-level structures of a software system and the discipline of creating such structures and systems. Each structure comprises software elements, relations among them, and properties of both elements and relations."*

## Modules, Components and Allocations





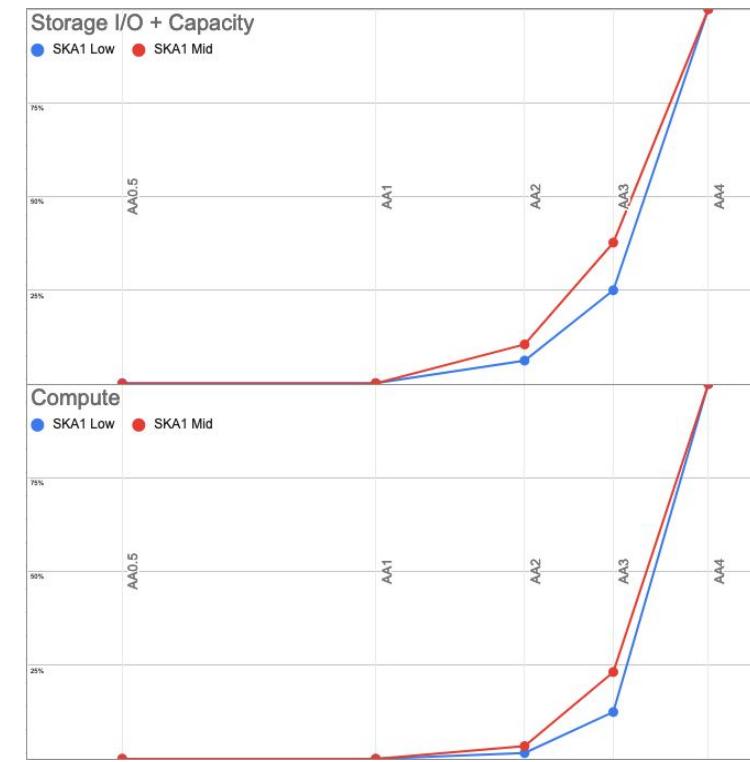
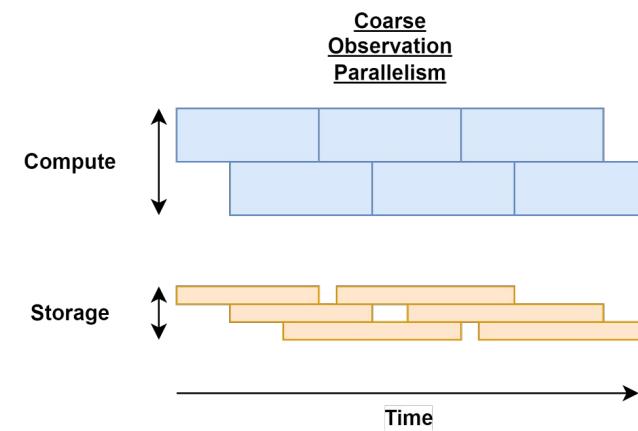
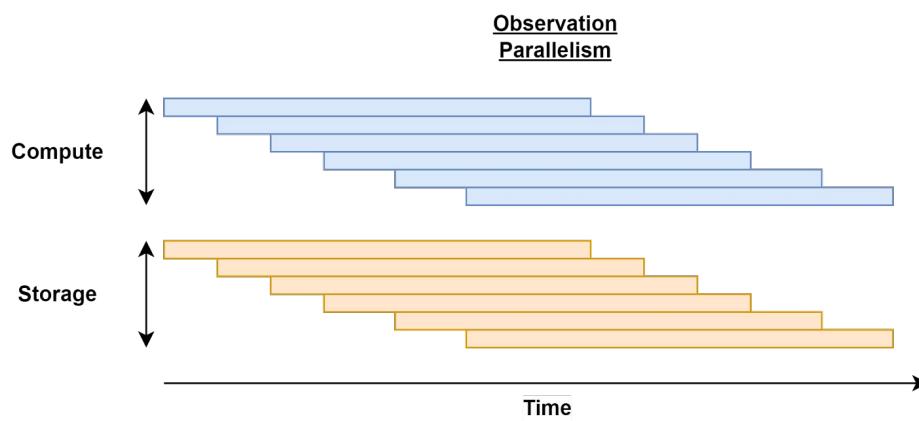
# Energy budgets (4-10MW)



We are limited by  
funding not science  
opportunities



# The challenge of data volumes

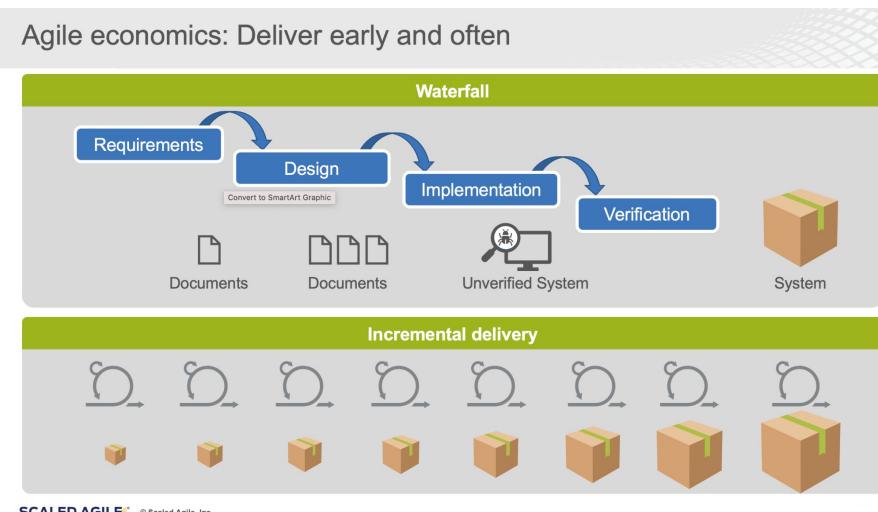


Storage needs/costs grow rapidly

# Improvements for key pipeline(s)

Scaling	I/O efficiency	Optimisations
<p><b>Scaling</b></p> <p>Scale to additional nodes</p> <ul style="list-style-type: none"> <li>Frequency parallelism (main focus for AA2)</li> <li>Snapshot parallelism</li> <li>Facet parallelism (main focus for AA*)</li> </ul> <p><b>Potential:</b> ~x10 improvement? (3 nodes → 30 nodes for AA2? More later...)</p>	<p><b>I/O efficiency</b></p> <p>Minimise I/O overheads</p> <ul style="list-style-type: none"> <li>Minimise inter-node I/O (main focus for AA2 - current bottleneck preventing scaling)</li> <li>Work entirely in-memory (i.e. not even write model visibilities)</li> </ul> <p><b>Potential:</b> ~x5 improvement? (assuming 10 facets, half of time spent on I/O?)</p>	<p><b>Optimisations</b></p> <p>Optimise processing functions (intra-node)</p> <ul style="list-style-type: none"> <li>Continue following algorithmical advances</li> <li>Use of vectorisation and accelerators</li> </ul> <p><b>Potential:</b> ~x2 improvement? (typical GPU lead per-cost, algorithmic advances harder to predict)</p>

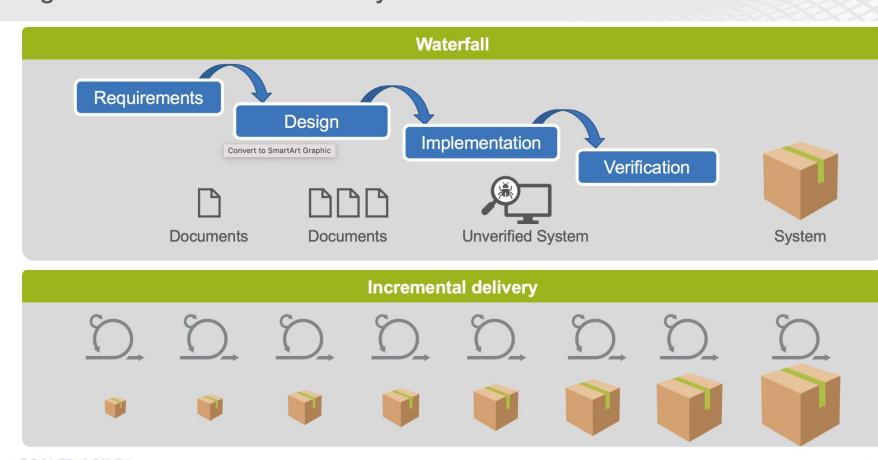
# Aside: How do we develop the software and systems?



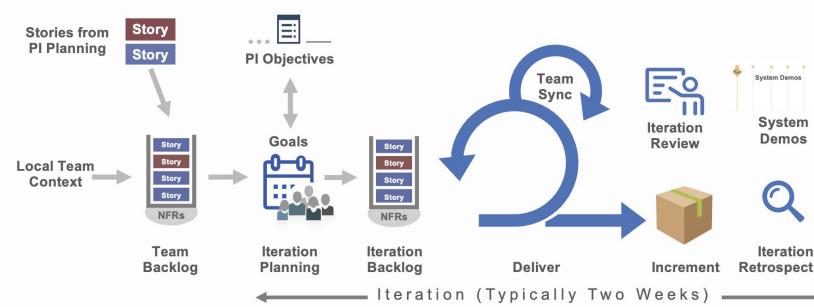


# Aside: How do we develop the software and systems?

Agile economics: Deliver early and often



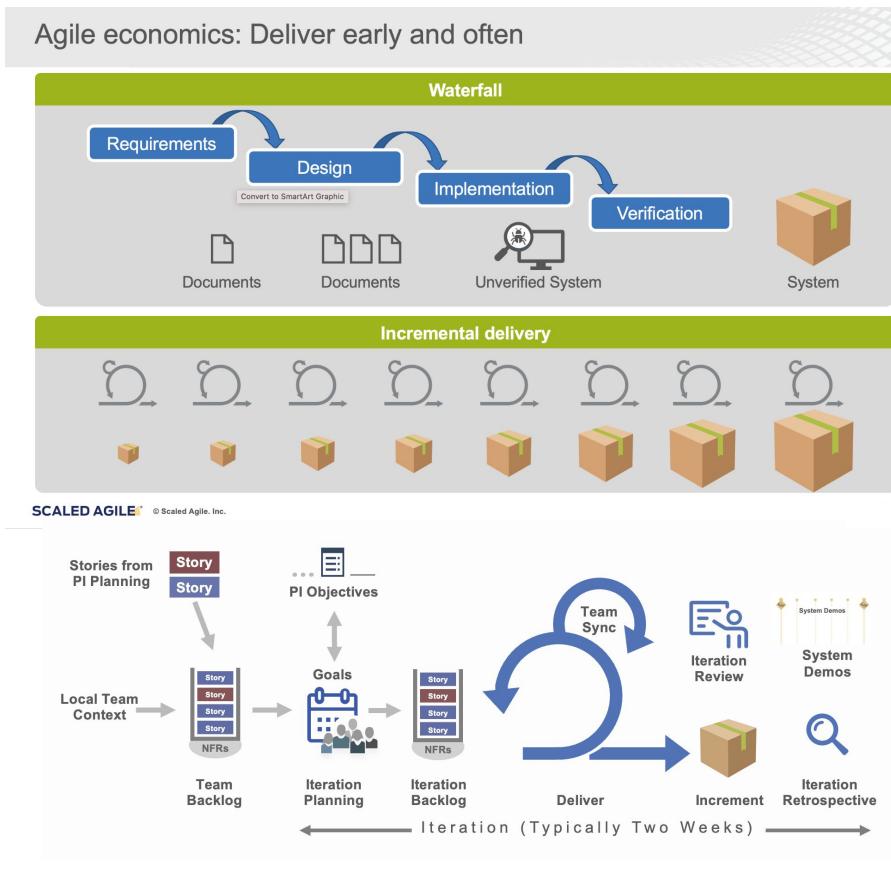
SCALED AGILE® © Scaled Agile, Inc.



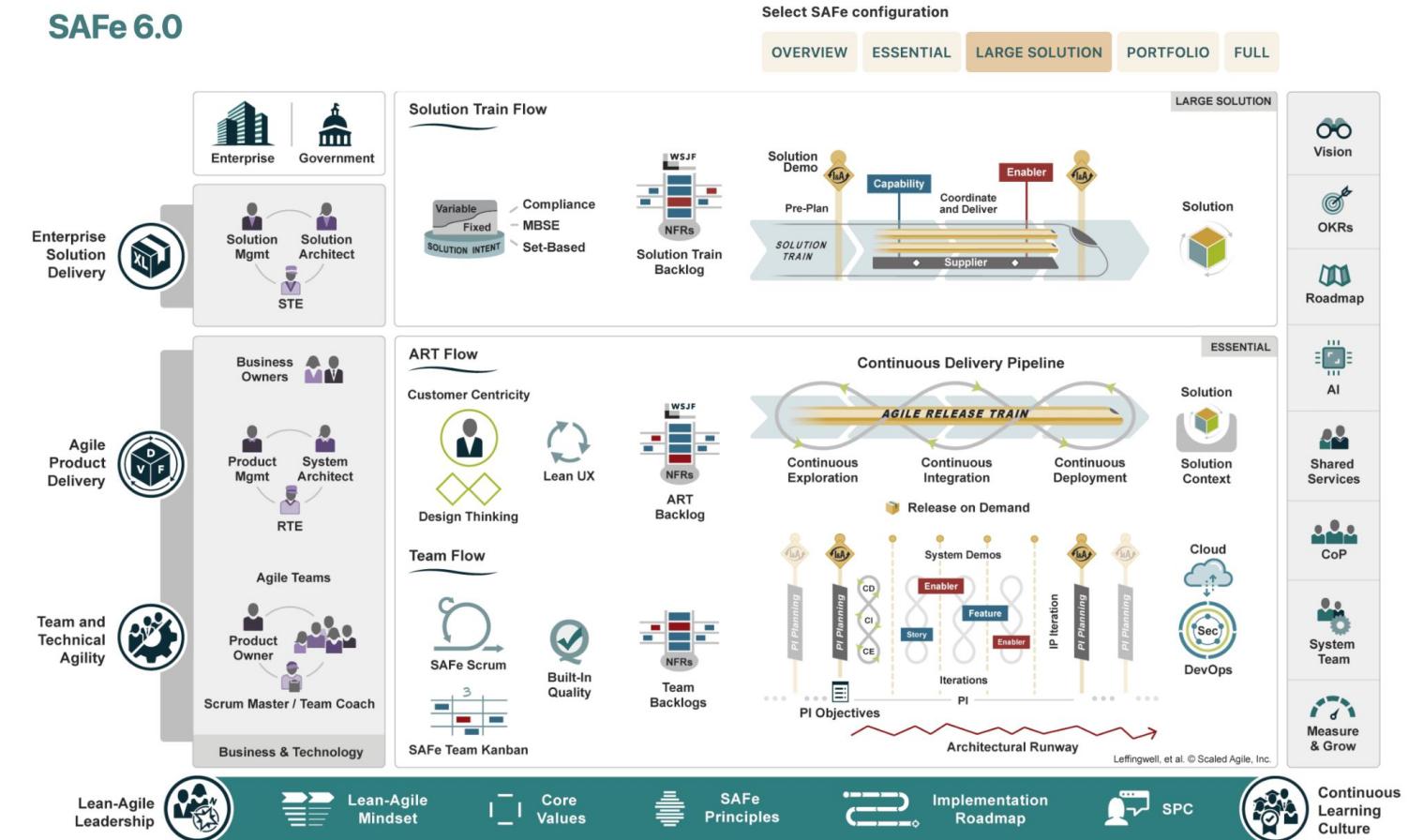


# Aside: How do we develop the software and systems?

<https://scaledagileframework.com/#largesolution>



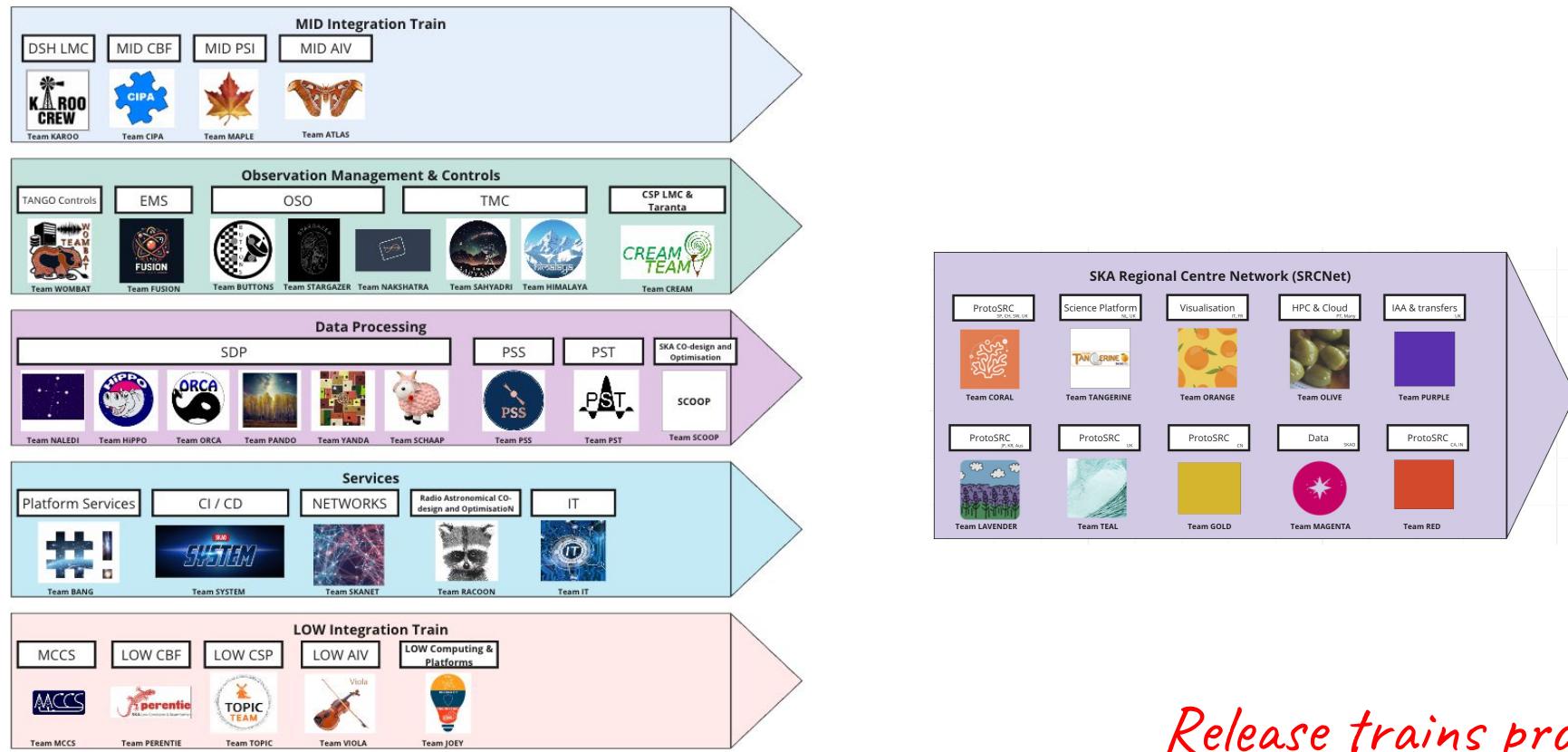
## SAFe 6.0



Agile techniques are being used

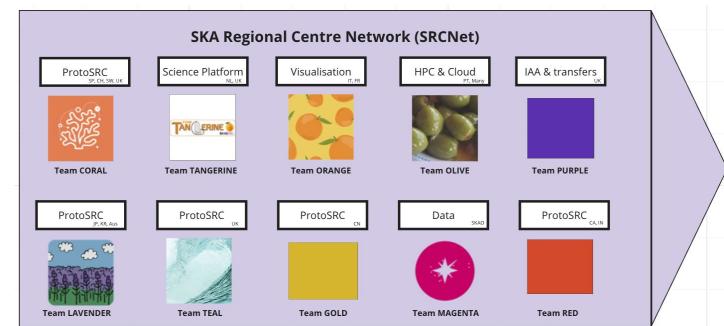
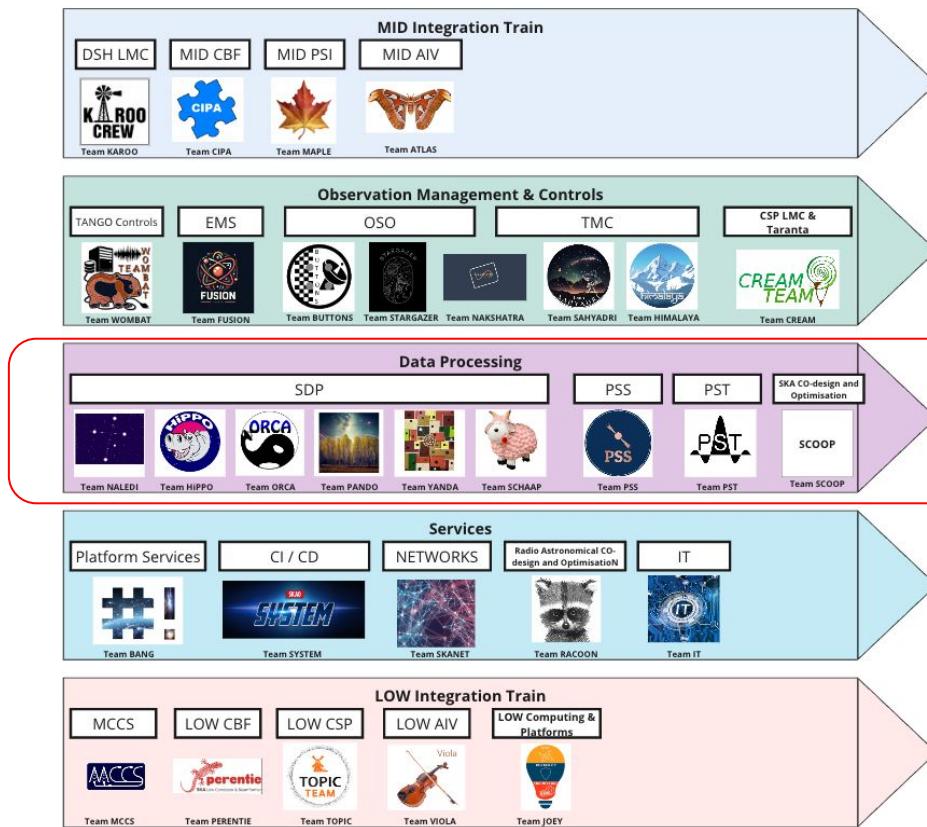


# The team of teams





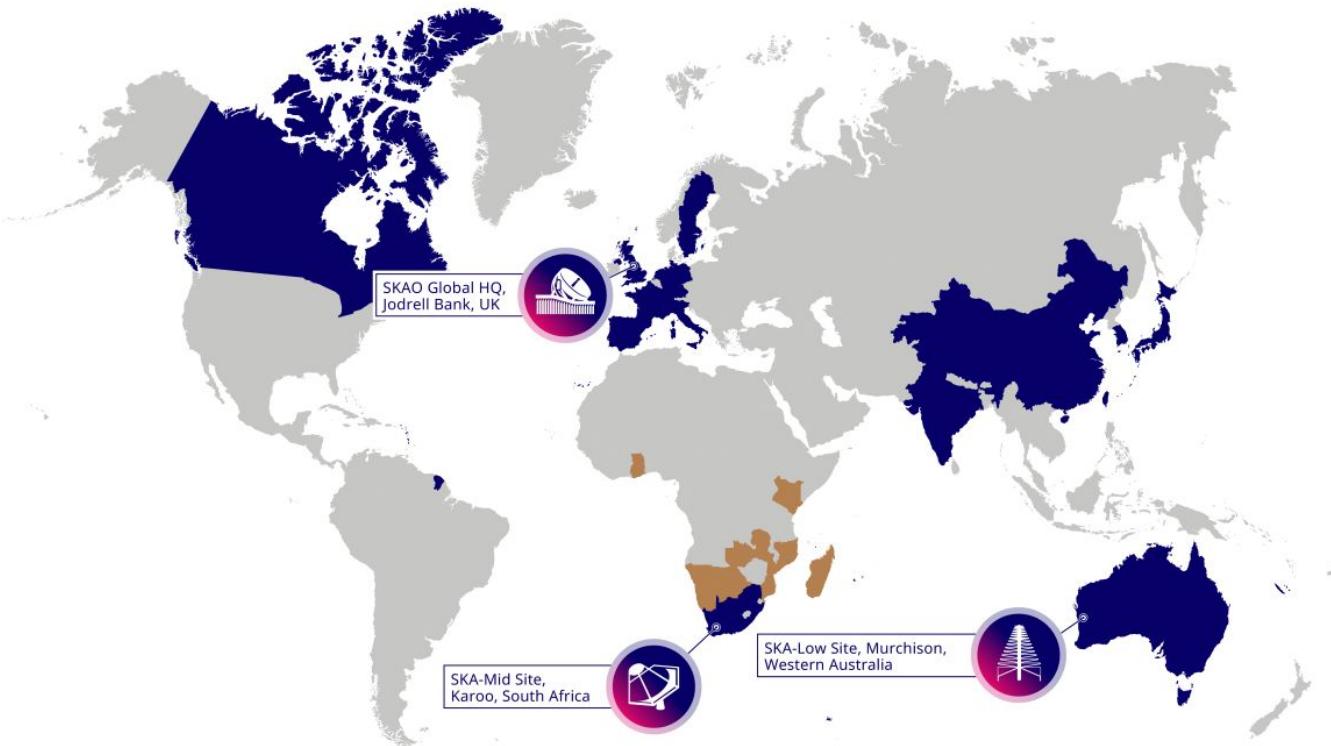
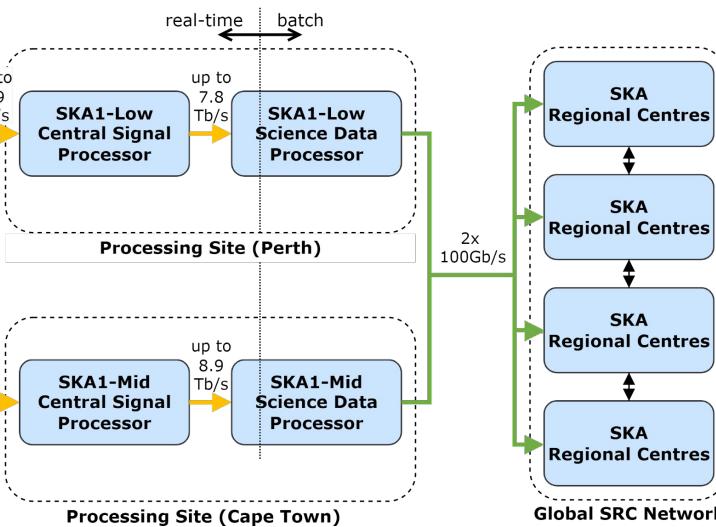
# The team of teams



Next come the federating considerations



# Where does the data analysis happen?



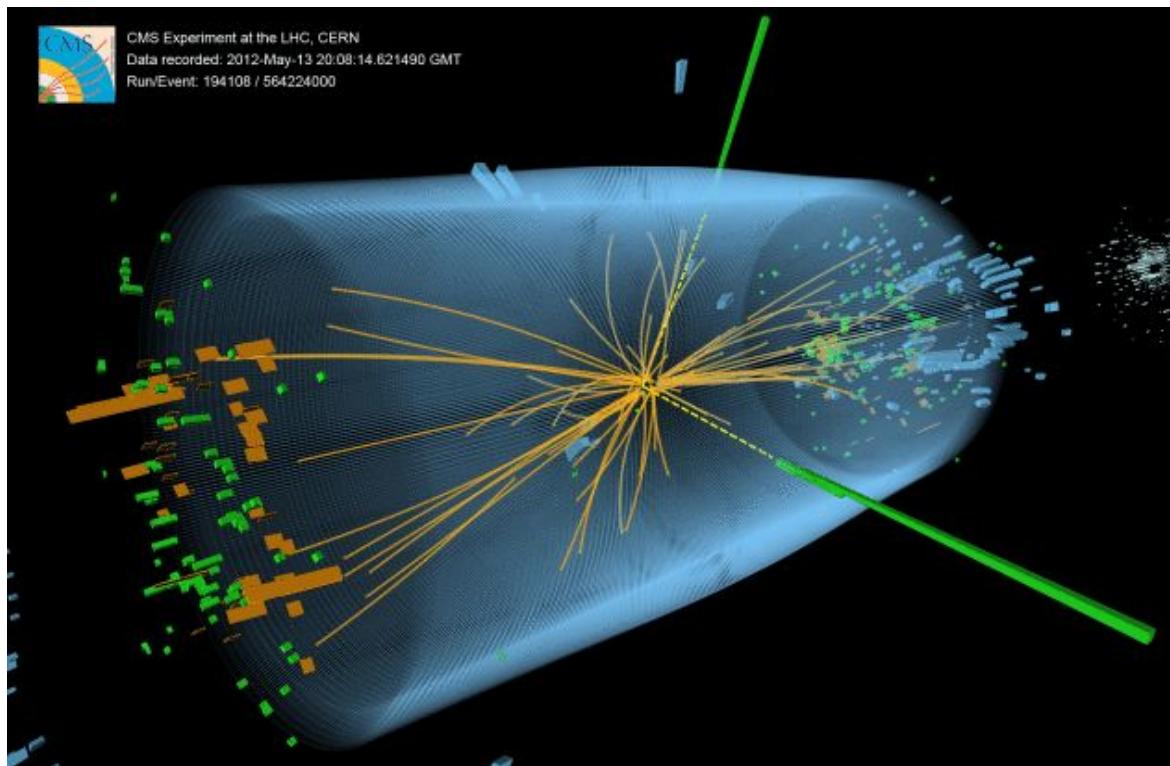
Credit: SKAO

Data is to be distributed and further processed globally

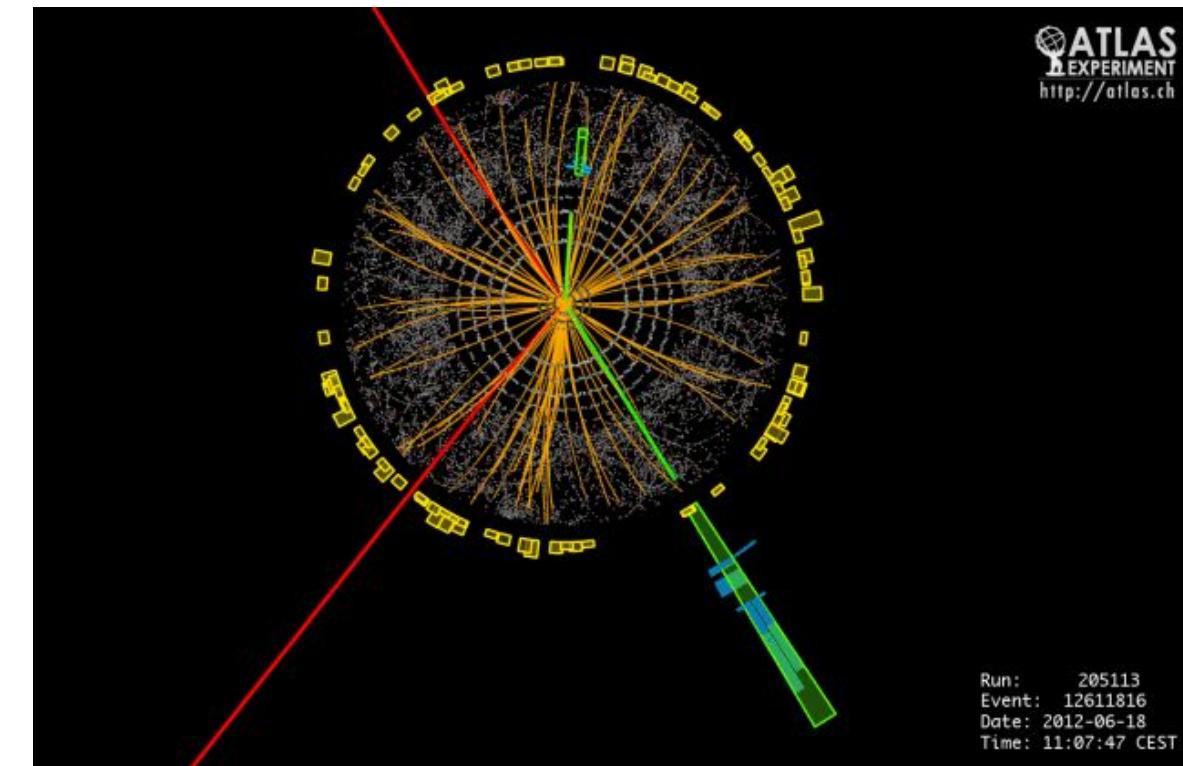


# We've been in a similar situation

Credit: CMS, CERN



Credit: ATLAS, CERN



Event recorded with the CMS detector in 2012 at a proton-proton centre of mass energy of 8 TeV. The event shows characteristics expected from the decay of the SM Higgs boson to a pair of photons (dashed yellow lines and green towers).  
(Image: CERN)

Event display of a  $H \rightarrow 2e2\mu$  candidate event with  $m(4l) = 122.6$  (123.9) GeV without (with) Z mass constraint. The masses of the lepton pairs are 87.9 GeV and 19.6 GeV. The event was recorded by ATLAS on 18-Jun-2012, 11:07:47 CEST. (Image: CERN)

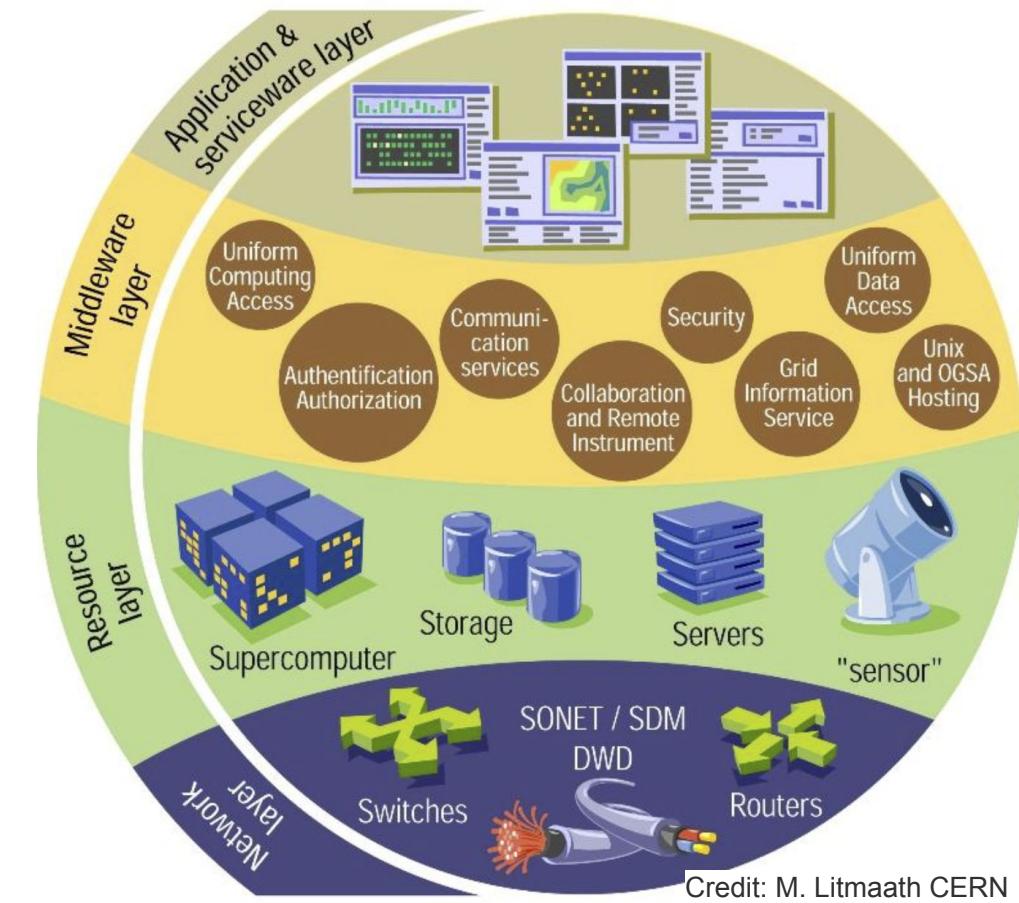
*The Higgs discovery depended on distributed computing*

# Worldwide LHC Computing Grid

- What is a computing grid: facilitates collaboration between members of a distributed community.
- Middleware: multiple data centres appear as a single system:
- Monitoring & accounting
- Job management
- Data management
- Information system
- Security
- Operations

Challenges include:

- Local policies & priorities
- Systems
- User communities ...



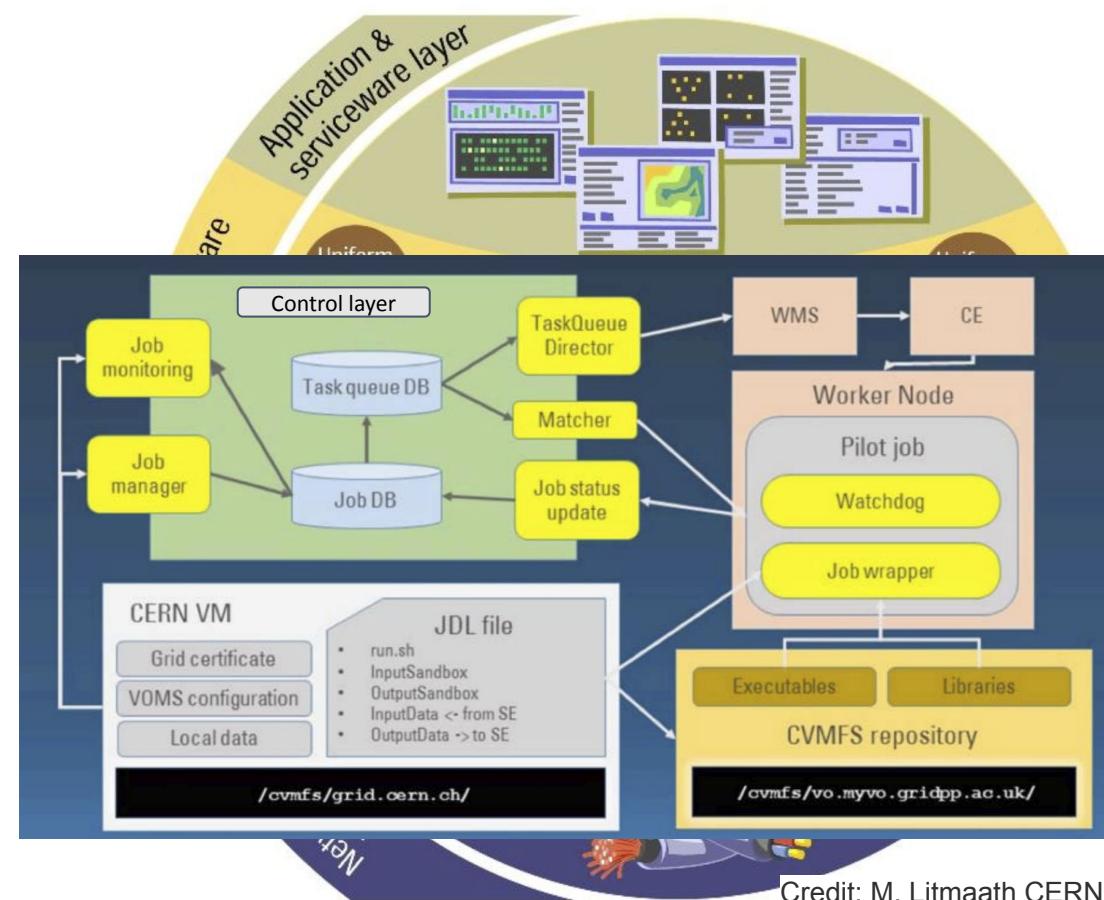
*Grids came before clouds*

# Worldwide LHC Computing Grid

- What is a computing grid: facilitates collaboration between members of a distributed community.
- Middleware: multiple data centres appear as a single system:
- Monitoring & accounting
- Job management
- Data management
- Information system
- Security
- Operations

Challenges include:

- Local policies & priorities
- Systems
- User communities ...



Credit: M. Litmaath CERN

*Centralised job scheduling depended on various middleware developments*



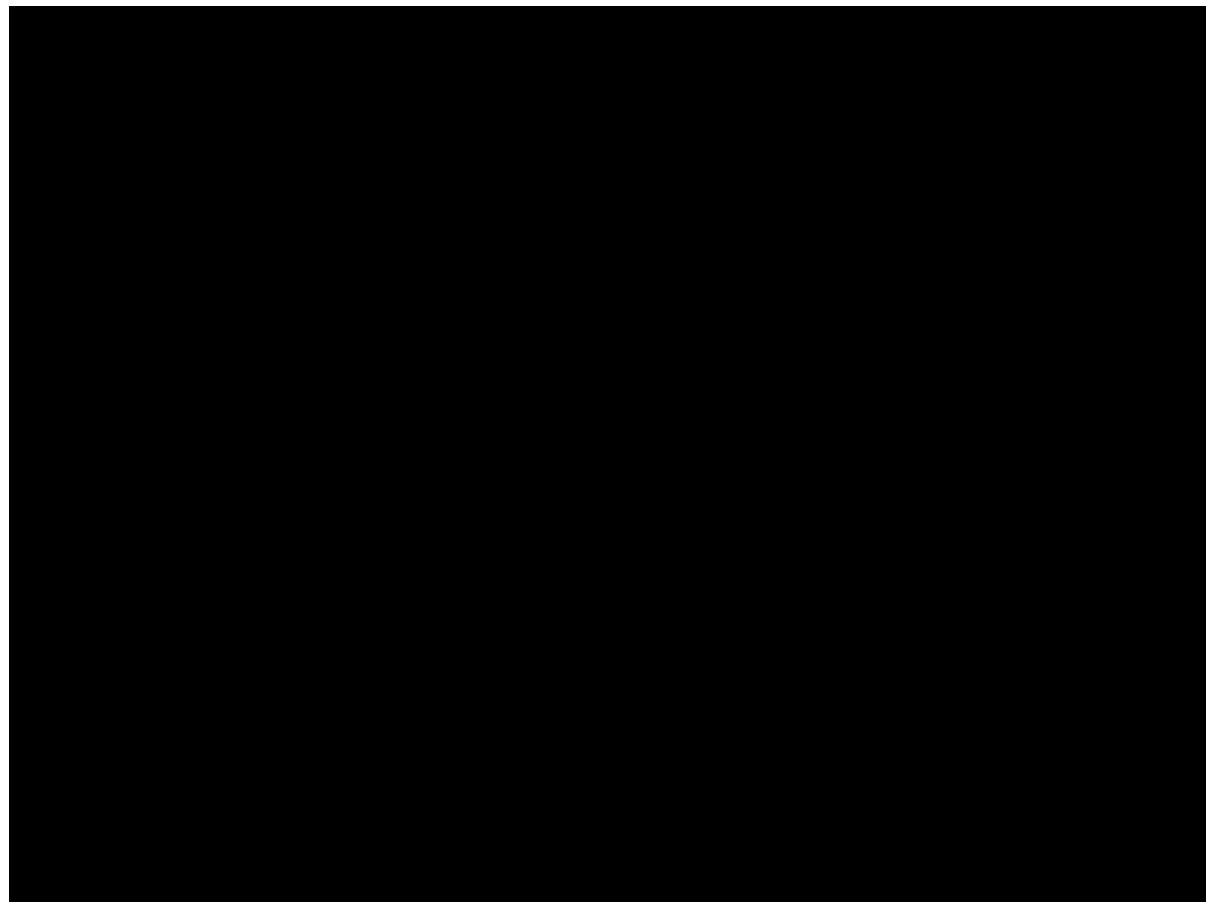
# Job flows in WLCG

WLCG provides global computing resources for the storage, distribution and analysis of the data generated by the LHC.

WLCG combines about **1.4 million computer cores and 1.5 exabytes of storage** from over 170 sites in 42 countries.

This massive distributed computing infrastructure provides more than 12 000 physicists around the world with near real-time access to LHC data, and the power to process it.

It runs over 2 million tasks per day and, at the end of the LHC's LS2, global transfer rates exceeded 260 GB/s.



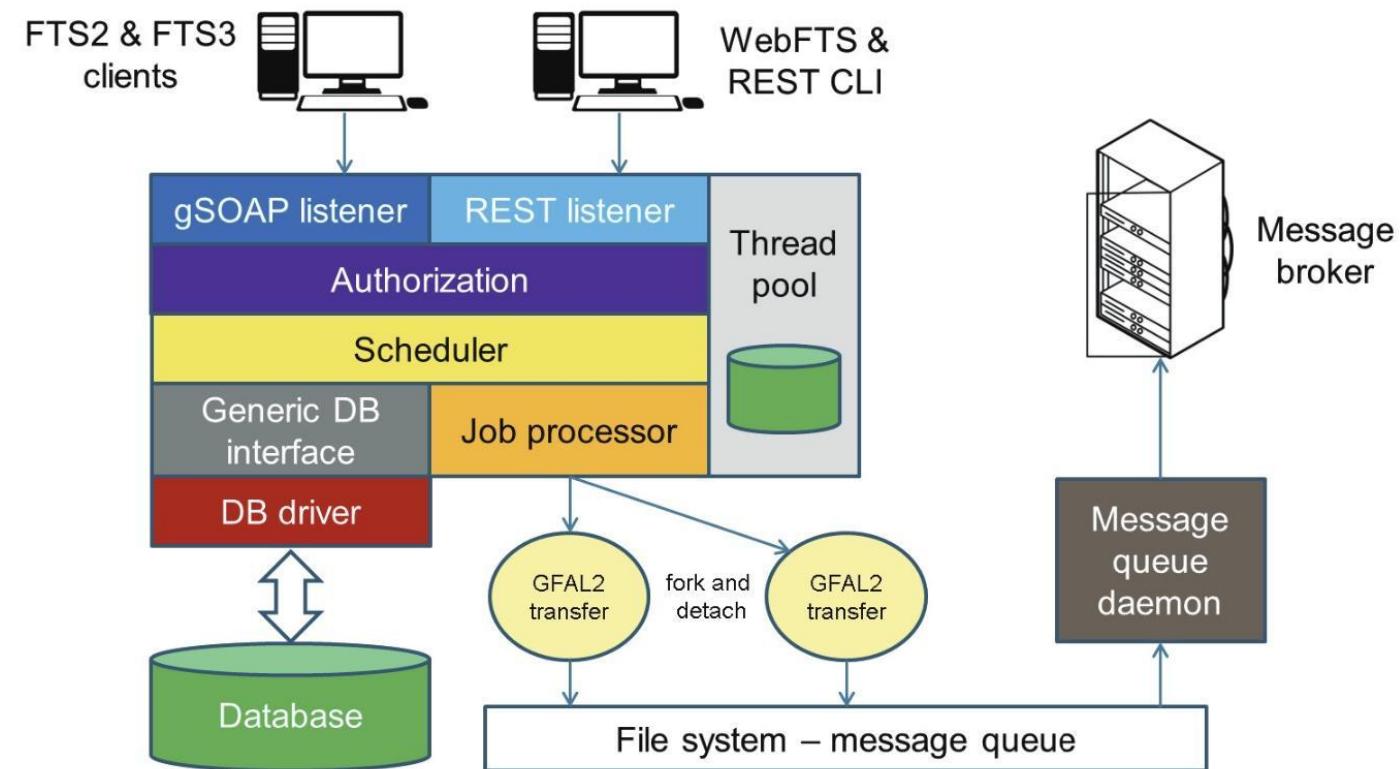
Credit: CERN

# File Transfer Service

One of the most important tasks is reliable file replication. It is a complex problem, suffering from transfer failures, disconnections, transfer duplication, server and network overload, differences in storage systems, etc.

Example capabilities:

- transfer auto-tuning / adaptive optimization;
- endpoint-centric VO configuration;
- transfer multi-hop;
- VO activity shares;
- multiple replica support;
- REST-style interface for transfer submission and status retrieval;
- retry failed transfers mechanism;
- staging files from archive;
- bulk deletions;
- support for Oracle and MySQL database back-ends



*Reliable 'at scale' data transfer is essential*

# The story so far ...

*It all depends on paired transistors*

*Each consumes energy to work*

*Every bit has a cost*

*Informed tradeoffs are needed*

*Radio and visible wavelengths are less absorbed by the Earth's atmosphere*

*Considerations are complicated!*

*NFRs define the Quality Attributes*

*Exploit natural parallelism in the data*

*Architecture views help us reason about the software*

*Software has various interfaces*

*There are methods to help us validate the architecture*

*We are limited by funding not science opportunities*

*We are exploring ways to improve the compute performance*

*Agile techniques are being used*

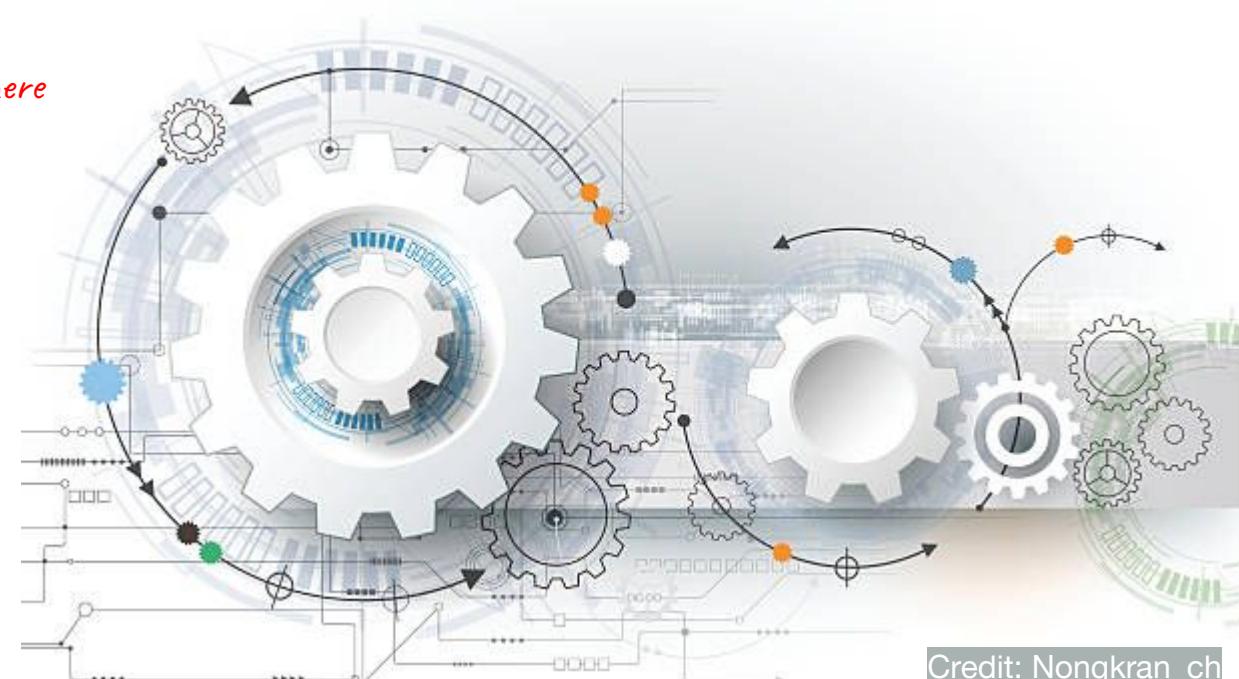
*Release trains provide alignment*

*SKA data is to be distributed and further processed globally*

*The Higgs discovery depended on distributed computing*

*Centralised job scheduling depended on various middleware developments*

*Reliable 'at scale' data transfer is essential*



Credit: Nongkran\_ch