



International
Centre for
Radio
Astronomy
Research



SDP Data Management Challenges

Andreas Wicenec

on behalf of

the Data Intensive Astronomy team at ICRAR

&

others

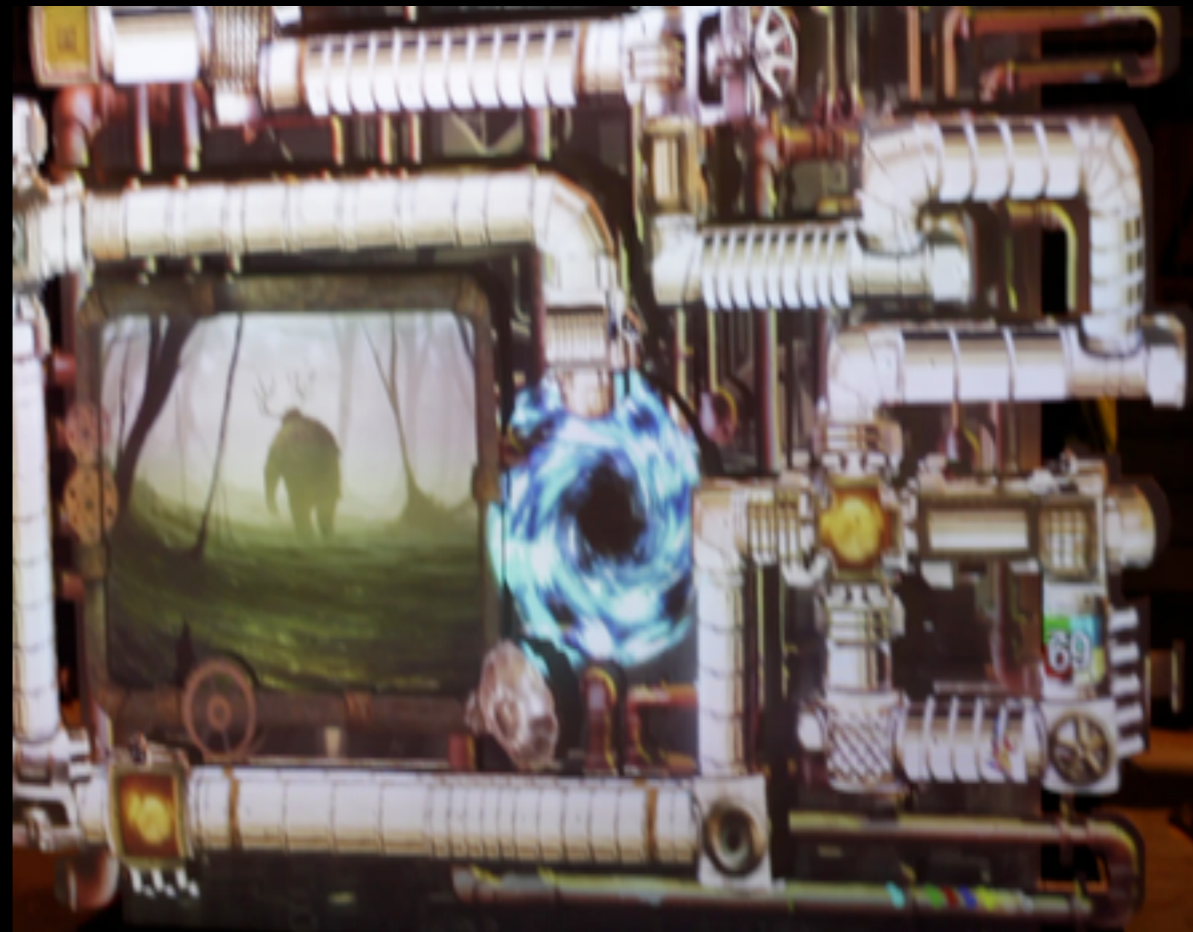
CONTEXT

Boundary Conditions

- Ingest rates ~ 400 GB/s per site
- Multiple concurrent observing projects
- Data sharing between projects
- Capital and operational budget limited
 - Power, Cooling
 - Acquisition, maintenance & software development costs
- Data parallelism: Millions of related tasks on thousands of nodes.
- Throughput: produce ~0.2-10 Tera Voxels/second
- Automatic 23/7 type of operation

Pipeline processing

- New ones are ‘invented’ for almost every single project.
- Very often based on scripting together modules from CASA, Miriad and AIPS as well as homegrown modules into an unmaintainable monster, that only a few people understand.
- Very dynamic and project specific. One of those monsters does not work for the next project.



We are at the limit!

★ SKA, ASKAP, MWA, LOFAR, MeerKAT

- are producing very high data volumes at very high rates
- are a challenge for currently available compute infrastructures (at least at affordable costs)

★ that means

- just throwing more hardware at the problem won't do the trick anymore.
- we need to use existing hardware more efficiently.
- at least the SKA requires significant innovation in order to approach the science potential of the arrays.

PROPOSAL & PROTOTYPE:

... A TINY BIT OF
INNOVATION

DALiNGE!

Data Activated

Graph

...think about
deluge!

flow
Engine

SEPARATION OF CONCERNS

Separation of concerns

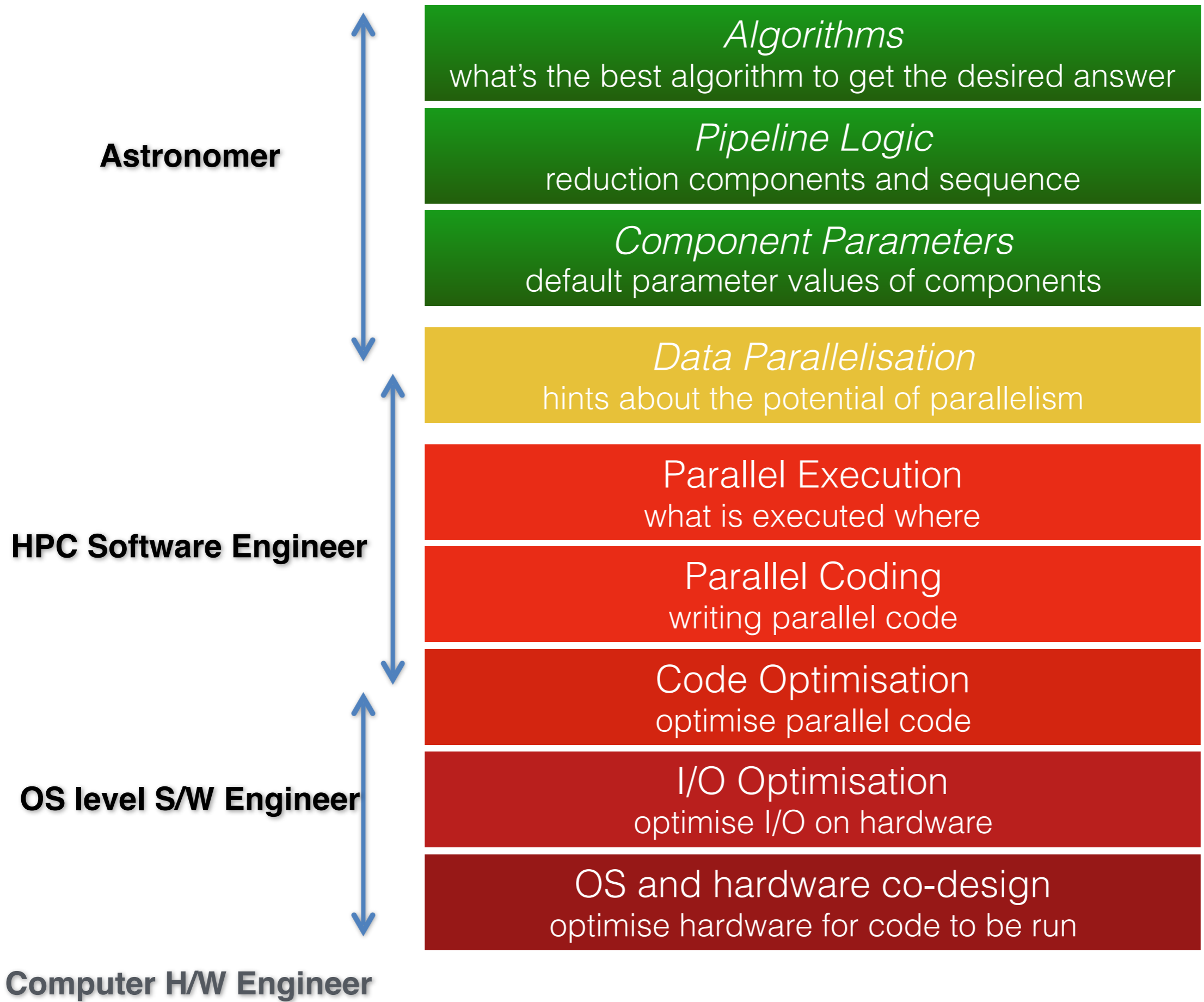
★ Let software engineers think about and write software:

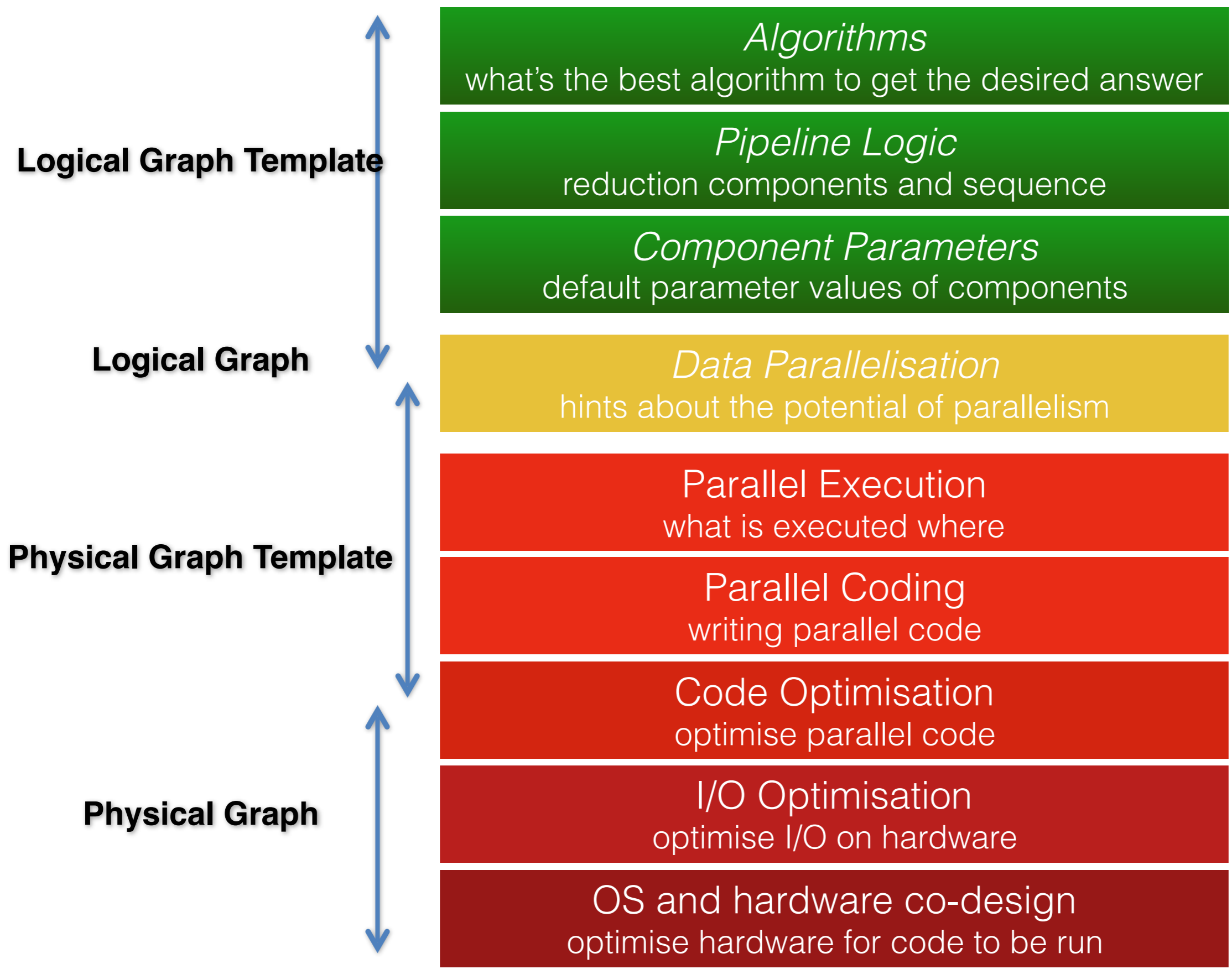
- Optimised code using the most appropriate language
- Novel ways of using latest hardware
- Using modern I/O techniques.
- Using advanced DB technologies
- Parallel code (even only a few software engineers can do this well!)
- HPC coding and optimisation (even less people can do this well!)
-

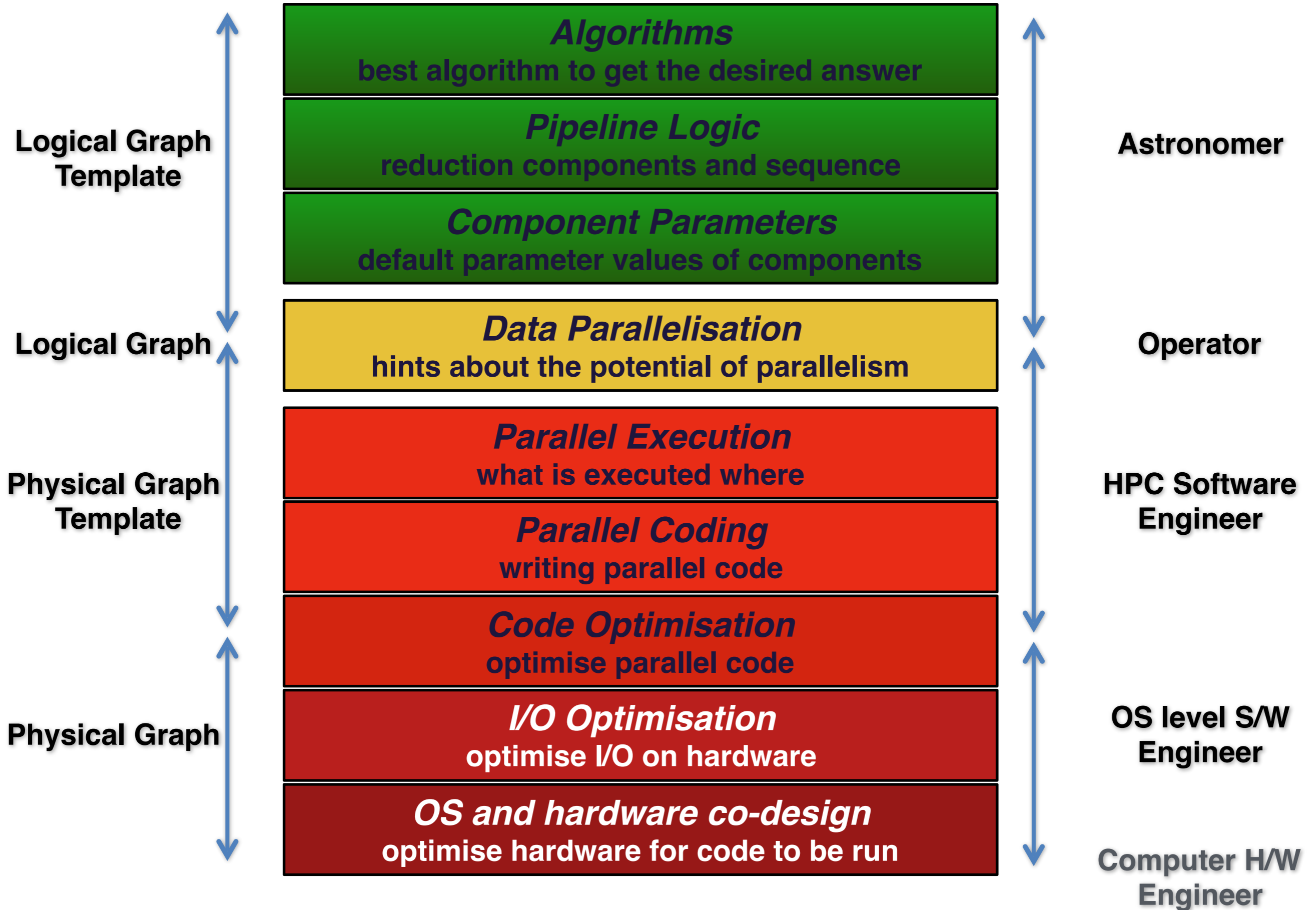
Separation of concerns

★ Let astronomers think about and do astronomy:

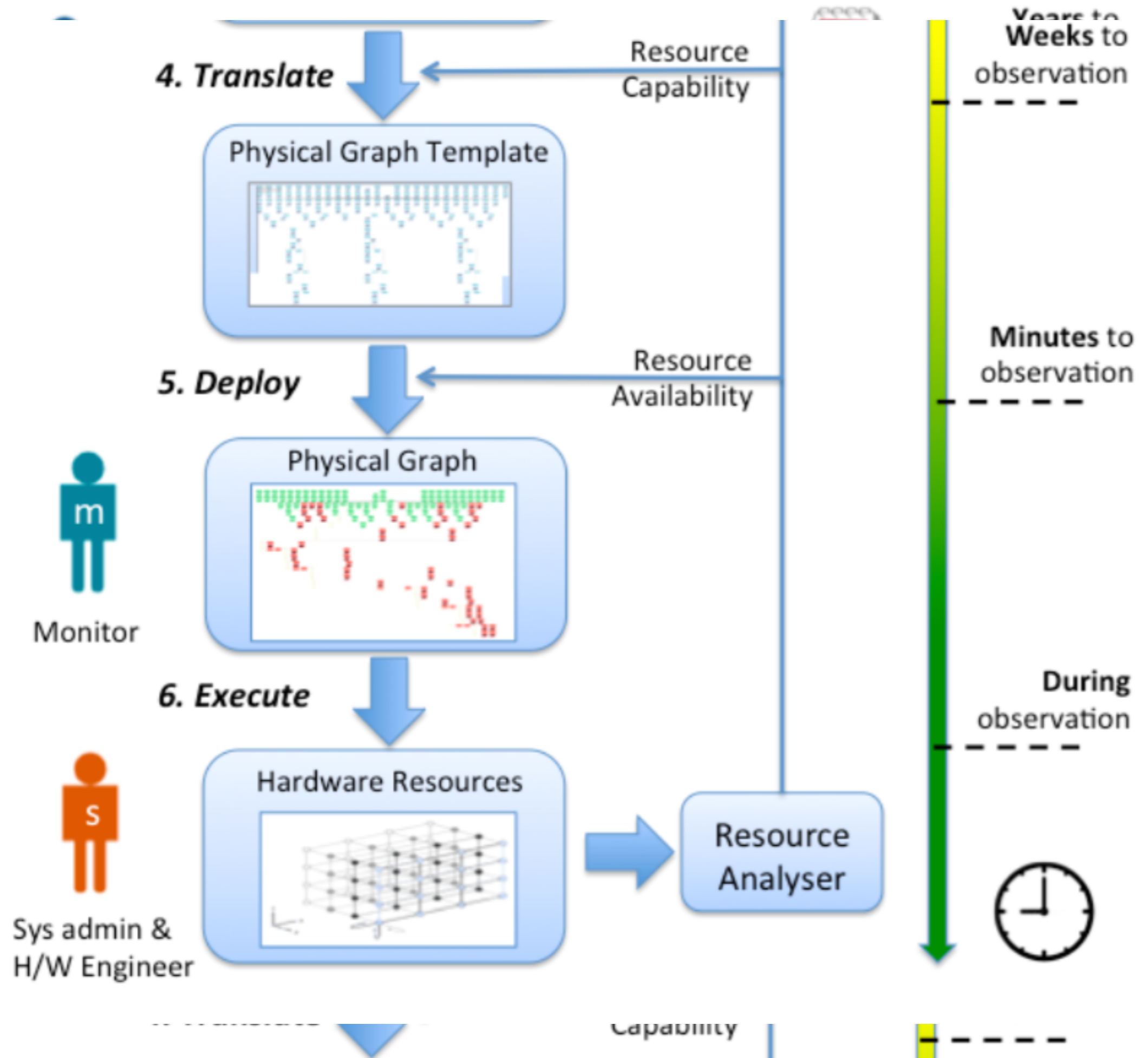
- Astronomical algorithms
- Pipeline logic
- Novel ways of extracting science
- New science
- Interpretation of extracted information
- Training of AI methods
-

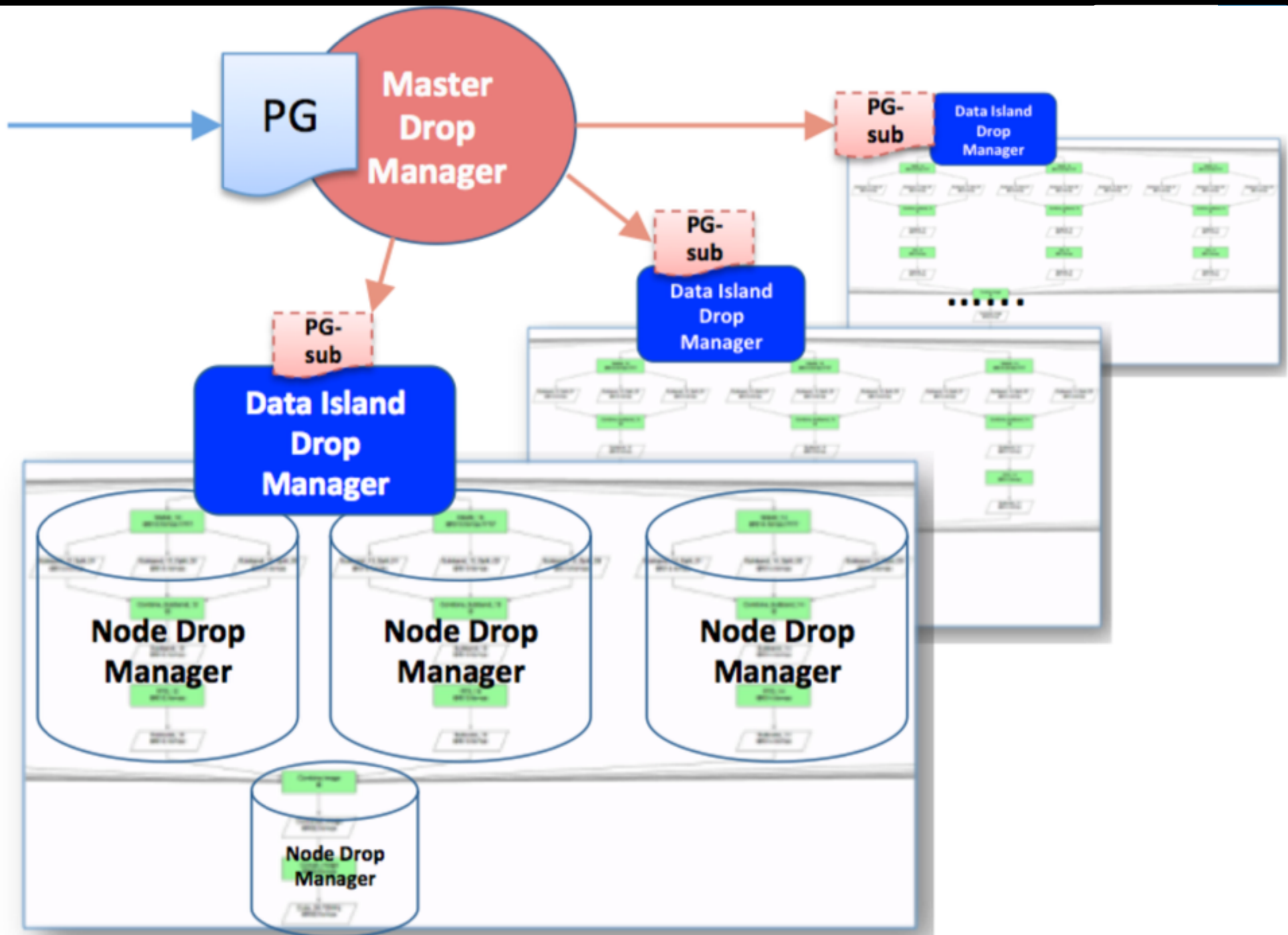






DEVELOPMENT, DEPLOYMENT & ASYNCHRONOUS EXECUTION



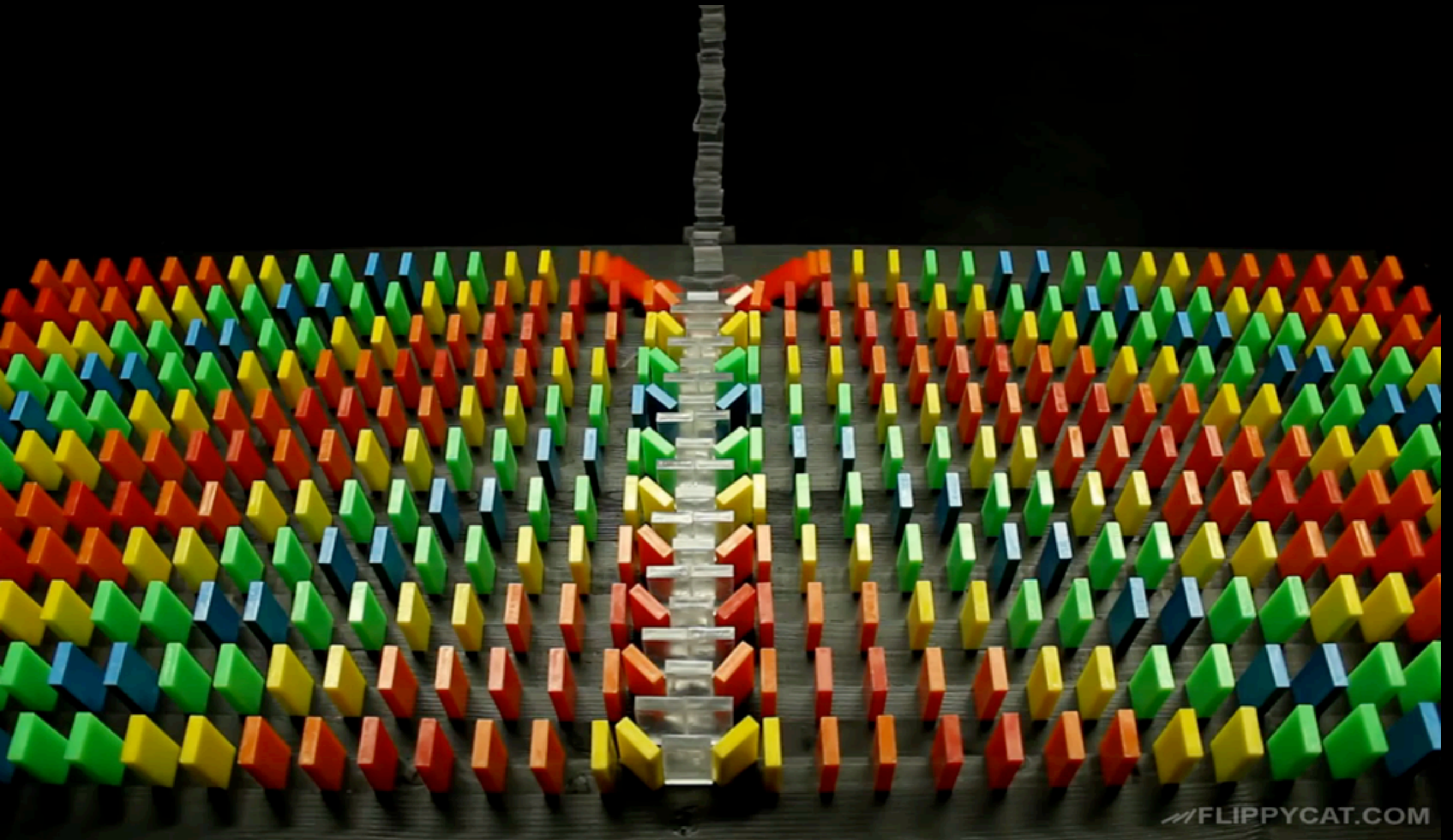


Execution (1)

- **DALiuGE's** graph implementation translates Data constructs into *DataDrop* objects.
- A **Drop** is an active software object wrapper referring to a payload.
- Typically the data is represented as a URL reference in a Drop object.
- A Drop payload can be anything, actual data but also executable binaries (*ApplicationDrop* objects)
- Enables the implementation of a Drop life-cycle management system for data and applications.
- Enables the deployment of Data

Execution (2)

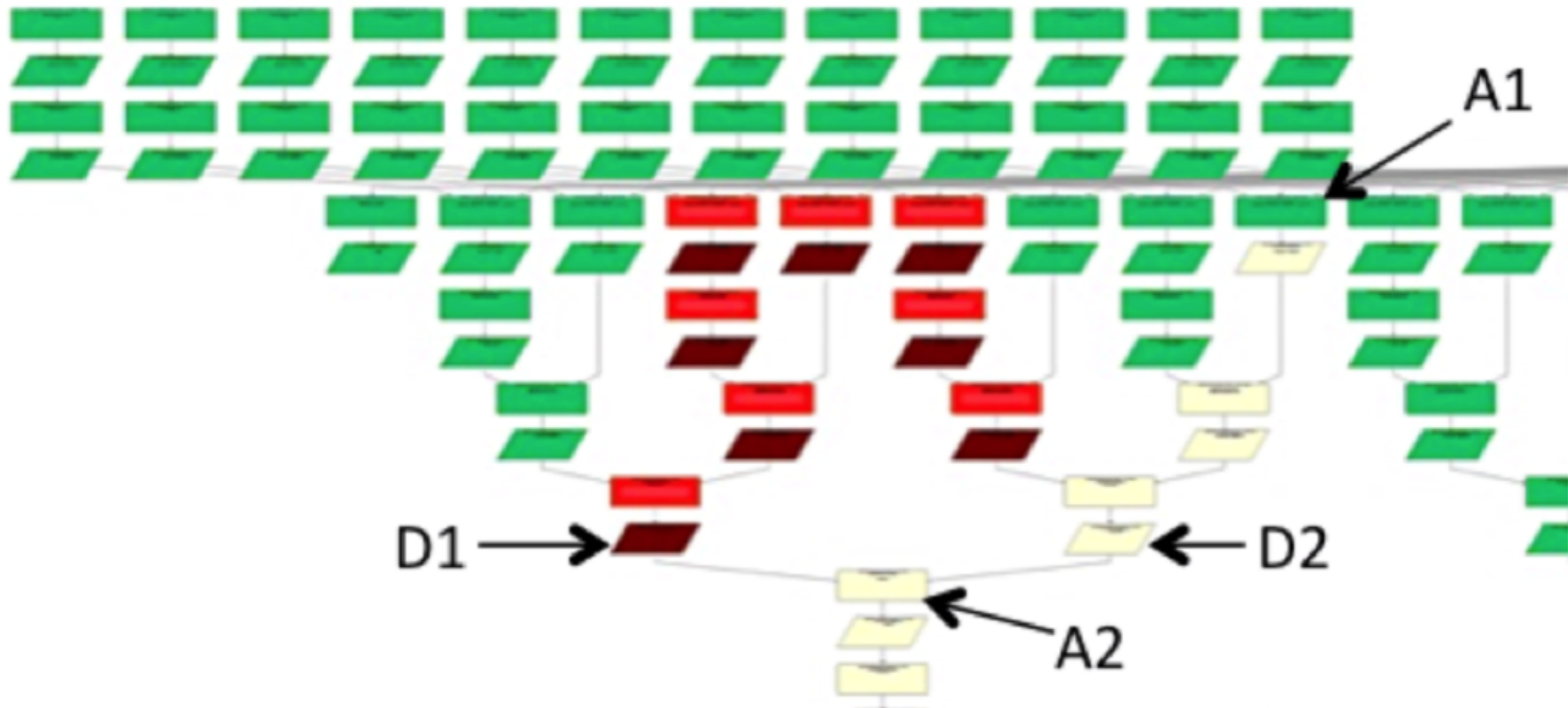
- Drop objects also define event channels, which are implemented as ZeroMQ PUB/SUB sockets.
- Drop objects implement a state machine.
- **Combine these:**
 - When a DataDrop object transitions to the ‘completed’ state it fires an event to its subscribers.
 - This triggers one or more consumer ApplicationDrops to act on the data payload.
- That means a single root DataDrop on a PG can trigger the execution of the whole rest of the PG.
- Very similar to a complex domino challenge.



ERROR HANDLING

Errors

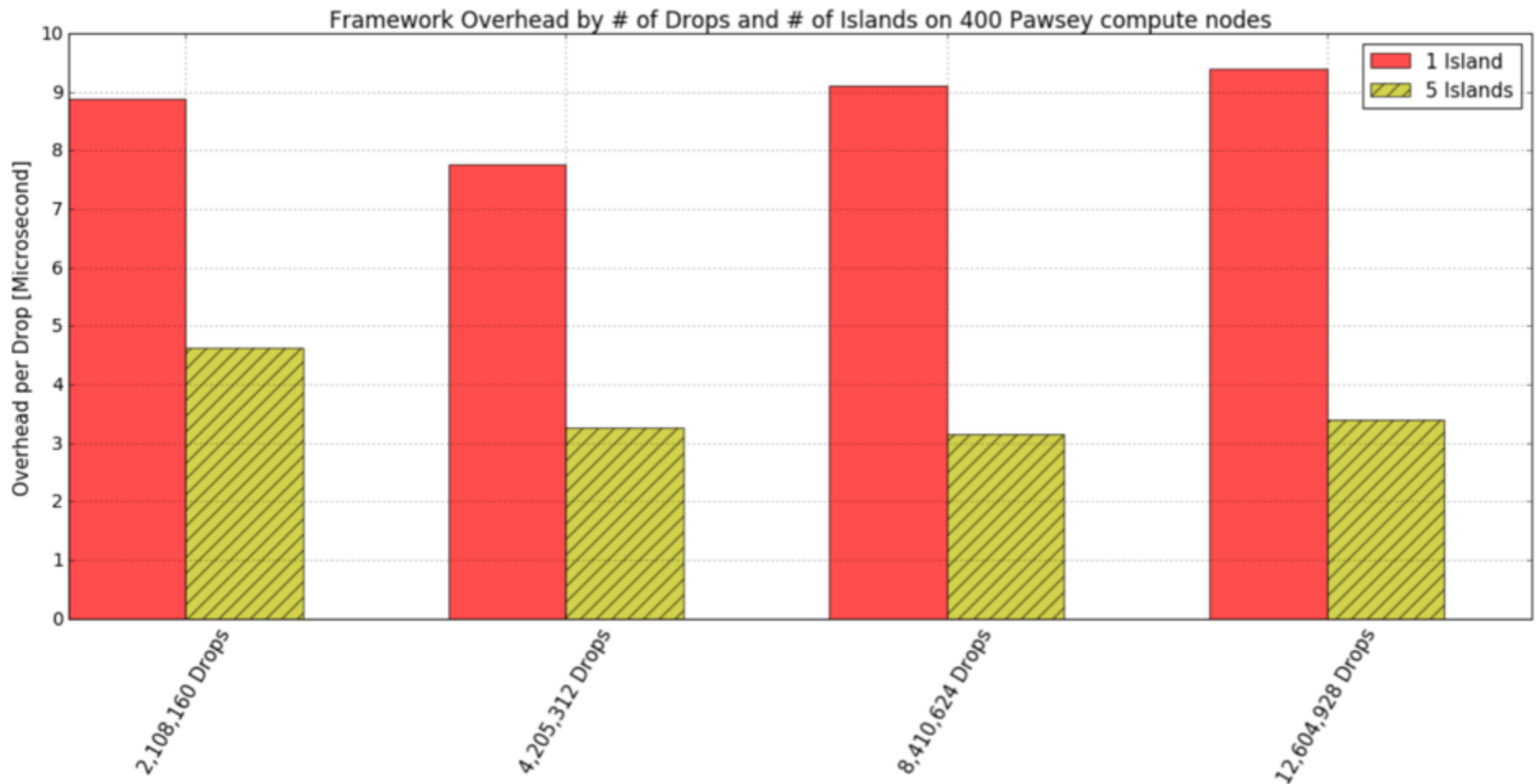
- Like with dominos, some of the chains WILL fail.
- We have constructs with thresholds to deal with that....



SCALABILITY

How many Drops?

- We have executed a wide series of scalability tests on various platforms, including Tianhe-2 in China.
- Measured plain overhead during execution time.



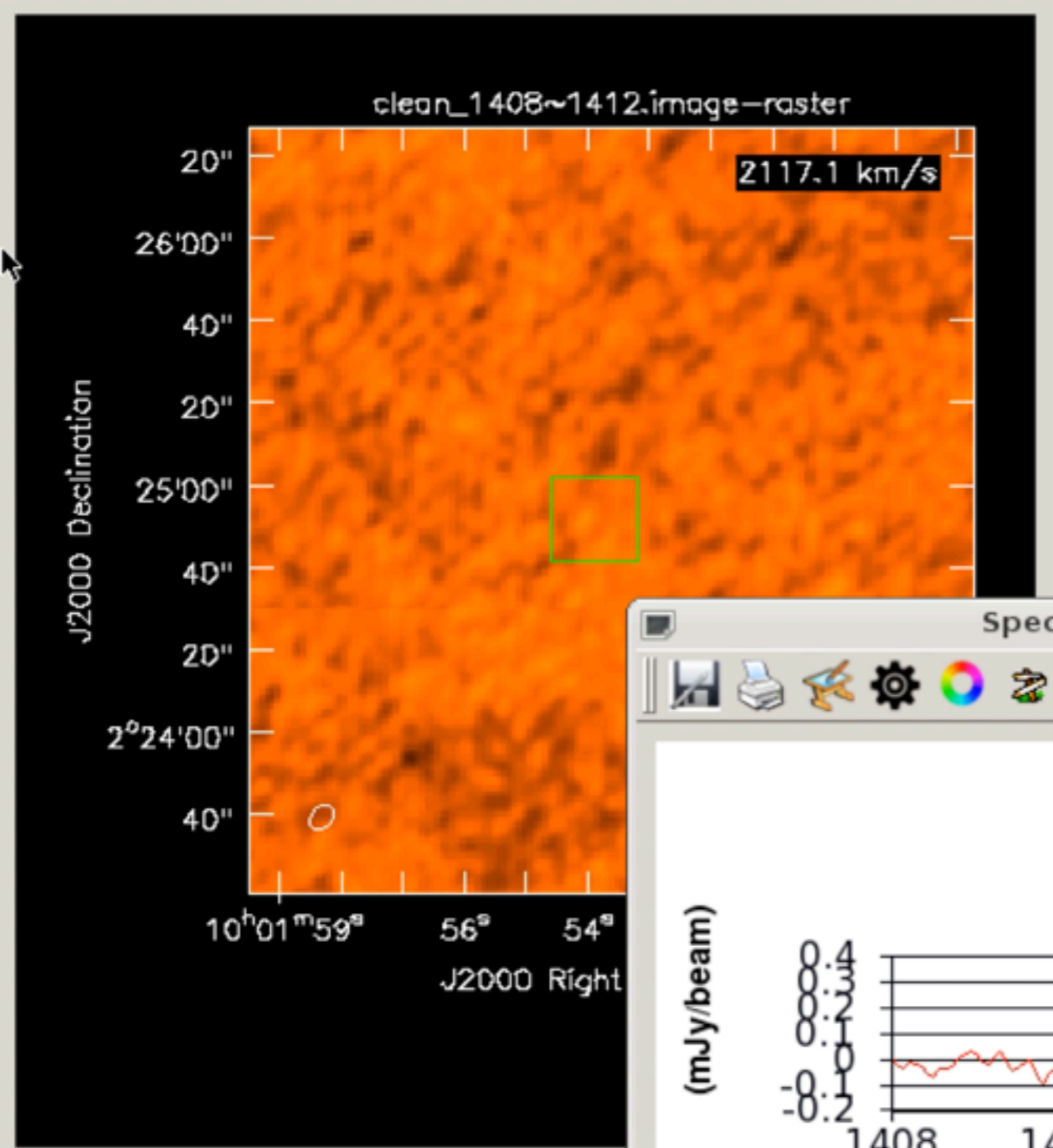
REAL WORLD EXAMPLES

CHILES

- **DALiuGE** has been verified using CHILES data on AWS, in-house cluster, Magnus and Galaxy.
- The current version of the code creates 40+ Node managers all running on separate heterogeneous AWS instances; with a single Data Island Manager controlling them.
- These graphs contain 7,000 ~ 88,000 Drops
- The graph generator knows the AWS instance types and can deploy more CPU/IO intensive tasks to more powerful nodes.
- The CasaPy tasks are all run from within Docker containers controlled by the **DALiuGE**



Display



Animators

Channels

Navigation icons: left, right, up, down, zoom in, zoom out, zoom reset

Rate: 5

Jump 152 256

Slider: 150 to 200

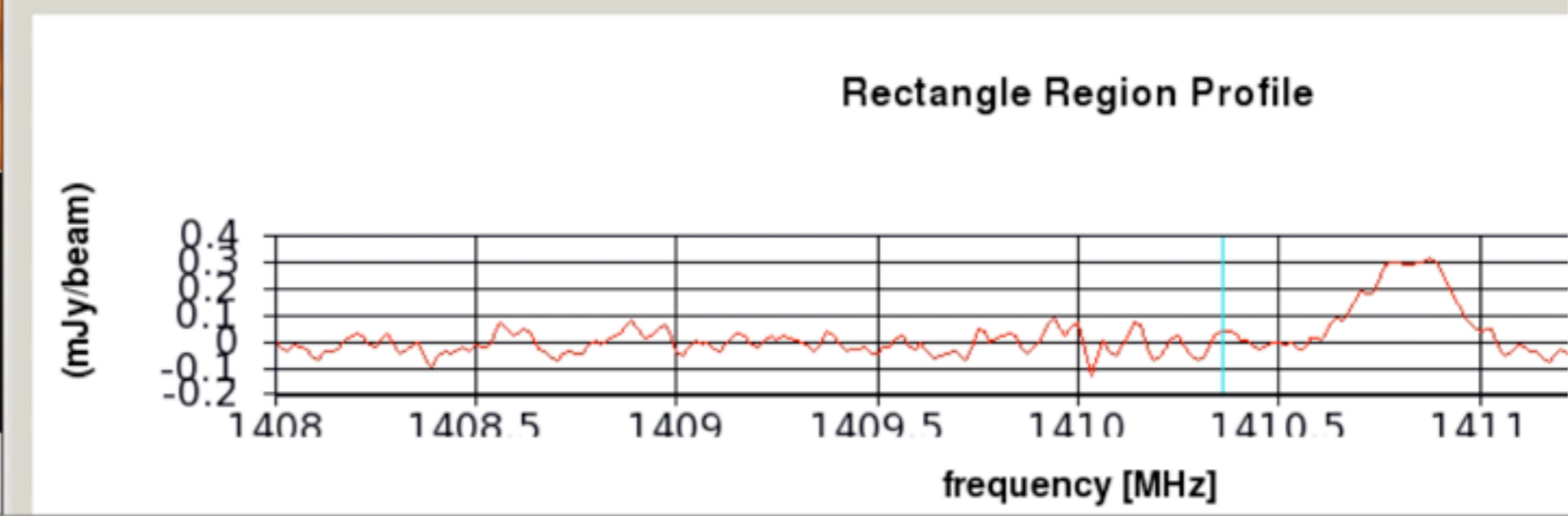
Images

Cursors

clean_1408~1412.image-raster

```
+3.19447e-05 Jy/beam Pixel: 642 1253 0 153
10:01:55.850 +02.25.46.455 I 2113.8 km/s (lsrk/radio velocity)
```

Spectral Profile - clean_1408~1412.image-raster (on epeius.icrar.org)



...and more

- we have ported the basic version of a MWA GLEAM pipeline to Daliuge.
- ASTRON is working to port and run the LOFAR pipeline.
- Kunming University ported their existing MUSER pipeline.
- we are also integrating OSKAR2 to simulate and reduce SKA-scale data sets on Tianhe-2.
- code is available on SKA SDP github.
- documented and fully tested code (continuous integration with loads of test code)
- Graph translation and optimised scheduling is a really hard problem...

Daliuge run on Tianhe2
66473 Drops

