

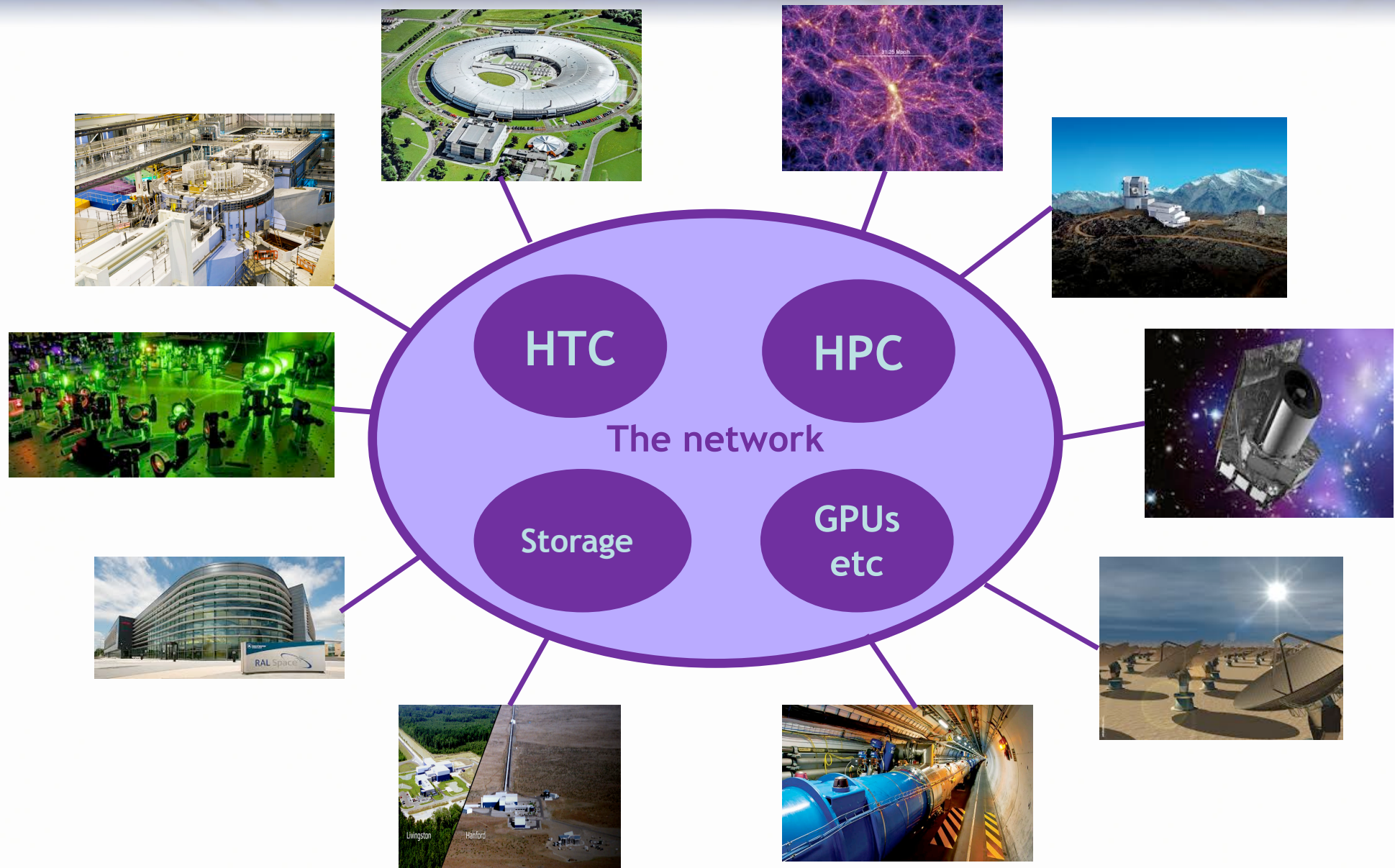
Distributed High Throughput Computing

&

Cloud computing

Pete Clarke
STFC Computing Town Meeting
Imperial College
17th Jan 2020

STFC hosts computationally leading-edge activities



STFC provides leading-edge eInfrastructure



9300 cores
50TB memory



GridPP (UK wide)
1 Million HepSpec06
70,000 cores (3-4 Pflops)
55 PB of disk storage
60 PB of Tape



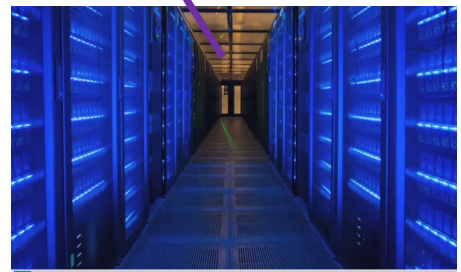
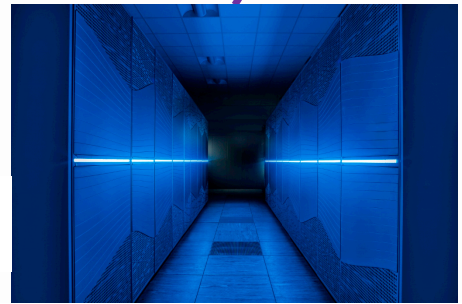
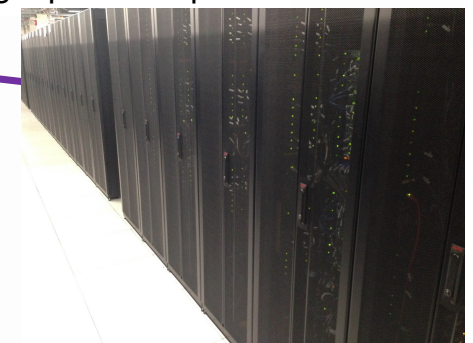
DiRAC 2.5y
4 machines
5 Pflops
HPC interconnects
Large memory

Tape store@RAL
~250 PB of Data stored



Hartree Centre
Industry focus HPC

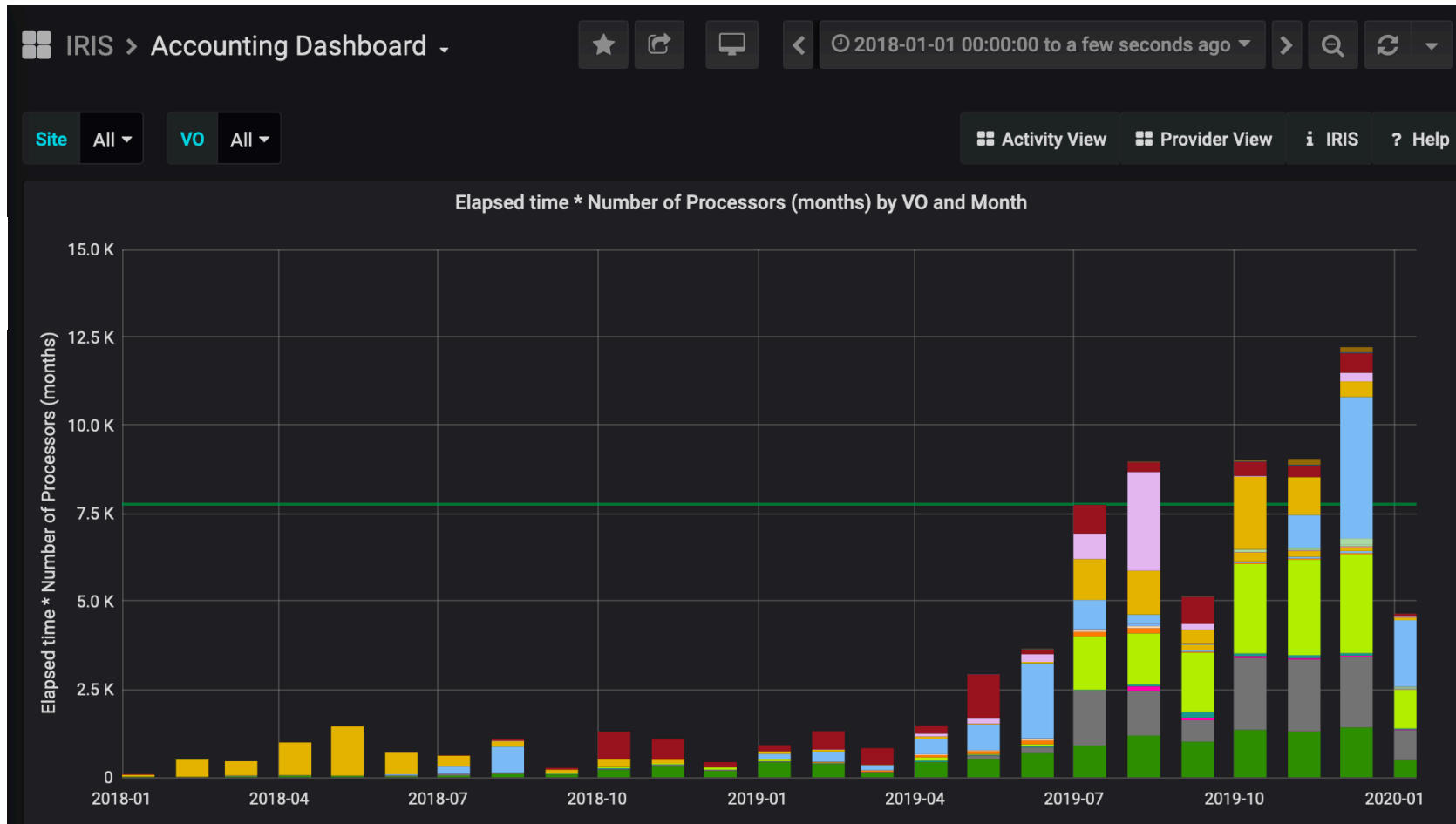
SCARF@RAL
12,000 cores
3 PB of disk storage
High speed backplane



11,500 cores
50PB storage

Recent IRIS accounting plot (excludes LHC, SCARF, ...)

Cores
in use



- ISIS
 Diamond
 CLF
 CCFE
 EUCLID
 AENEAS
 eMERLIN
 LSST
 lz
 dirac
 jintrac
 casu
 gaia
 vcycle
 gaia-dev
 gaia-prod
- gaia-test
 lsst
 lz
 skatelescope.eu
 dune
 vo.cta.in2p3.fr
 virgo
 iris.ac.uk

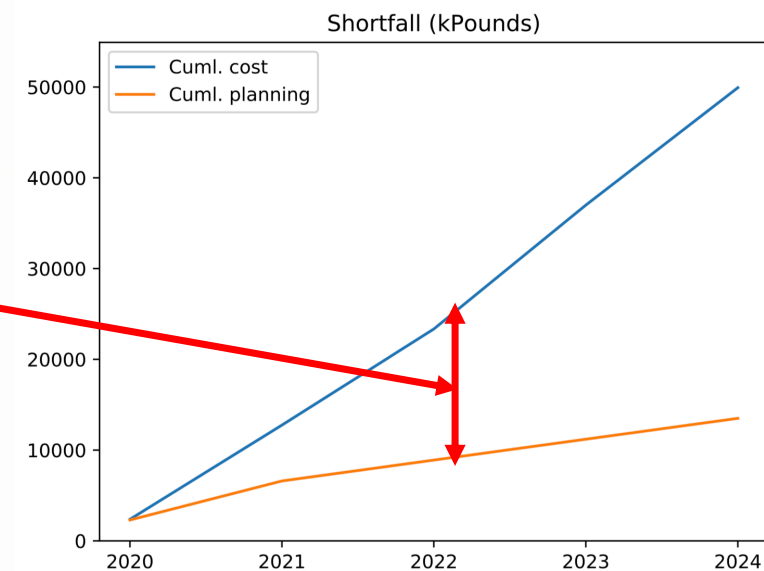
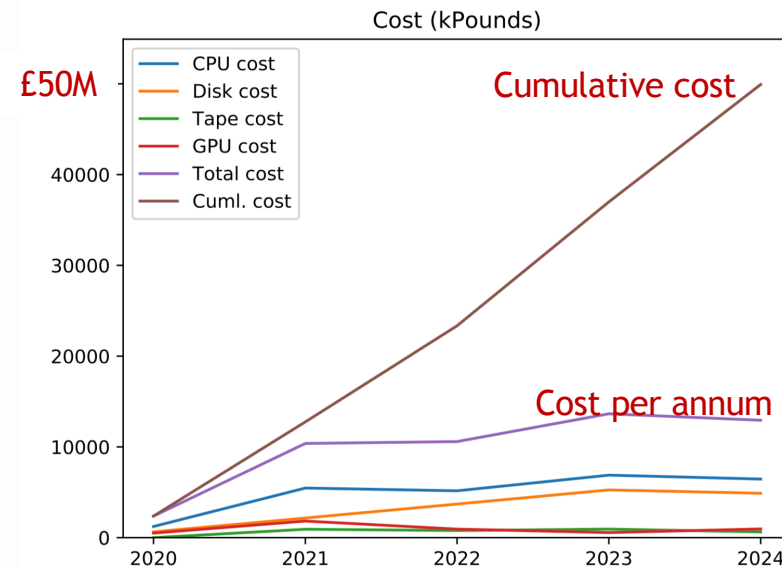
...but ... projected aggregate requirements -vs- known funding

These numbers based upon IRIS 2019
Resource Scrutiny and Allocation
round (RSAP)

Shortly to be updated for RSAP 2020
→ numbers have not gone down

There is a very large shortfall

This information has been shown
widely and is a “mater of fact”



message to UKRI, BEIS...

**There is no more important message to give to BEIS/UKRI than...
... this is a clear requirement to seriously address this UKRI wide issue**

At present the research community sees this picture:

2023



2020

**... because of this there is also a clear and present short-term problem
... structural commitments cannot be made to long term projects
.... there is a type of planning blight
.... investment “on account” is urgent in 2020**

Sustainability

A sustainability plan is needed as well

Fuel



Crew



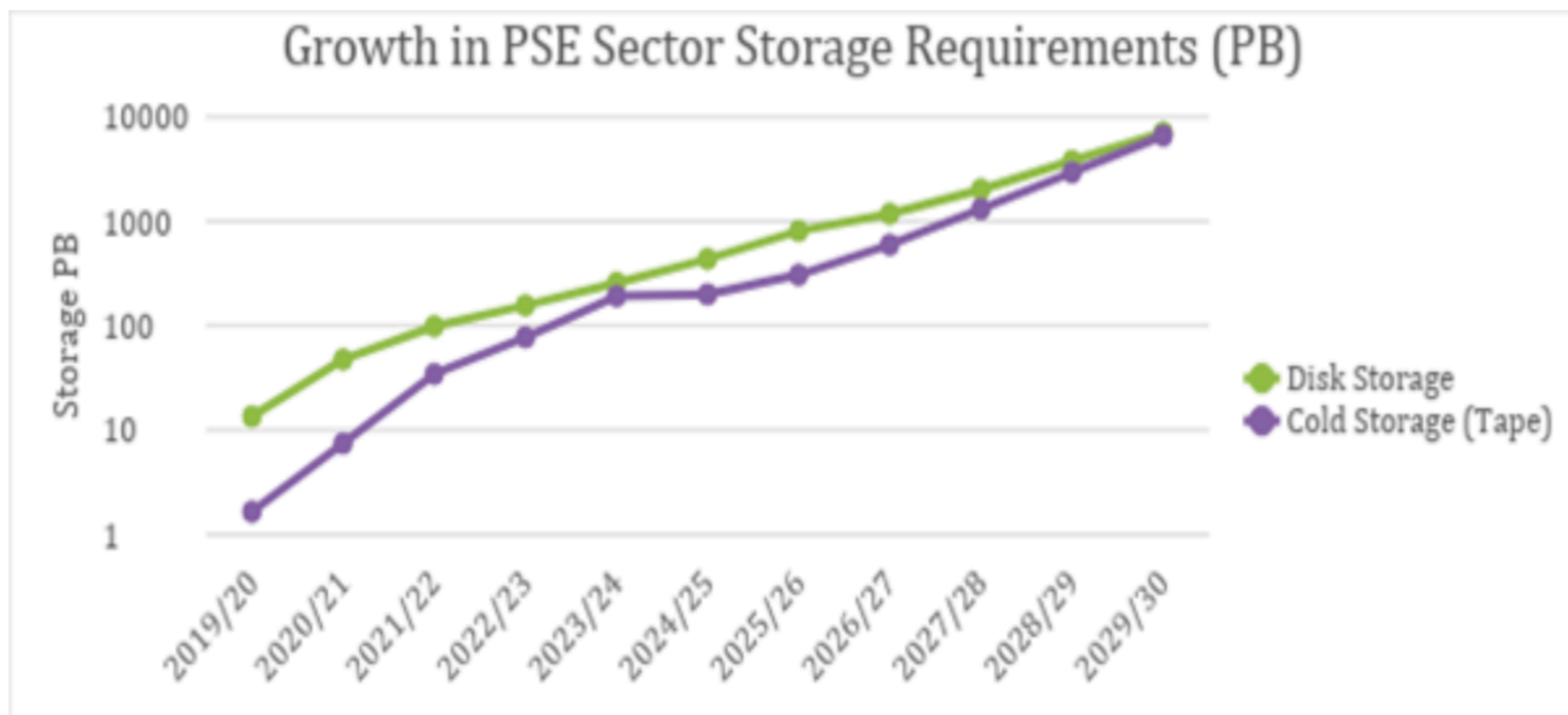
White paper:

UKRI Data Infrastructure Roadmap

Health warning: All WPS are not yet formally final.
Information presented is to be considered as draft

Compute and Data Infrastructure WP

- Figure 2: The Predicted growth in disk and cold storage for the PSE sector between 2019-2030.



UKRI Data Infrastructure WP

Introductory statements say:

First and foremost, UKRI urgently needs to restore the foundations upon which such exploitation of data can happen. This means to put in place the physical compute and storage capacity needed to host and exploit the data across UKRI. Without this all other discussions are moot. Assuming that this capacity will exist, we present recommendations for transformative modes of research which employ the linkage and analysis of heterogeneous datasets from different research sectors and other sources of data. These must be interoperable and federated to ensure data can be shared and accessed readily.

Physical RDI (computational capacity, storage capacity and networking) should be put in place to host and exploit the data. This means adequate investment in HTC and HPC facilities, both within UKRI and where appropriate using commercial cloud provision, and in the software infrastructure to manage, move, give access to and analyse the data. This should be complemented with networking connectivity via Janet. This is a prerequisite to all other recommendations in this document.

These statements apply across UKRI

UKRI Data Infrastructure WP

This means:

- Large scale HTC compute capacity (STFC HTC scale ~150k cores by 2023)
- Cloud capacity
- Large scale data storage facilities (STFC scale ~ 200 PB by 2023)
- CPU-Accelerator deployments for machine learning
- High performance databases for archives
- Visualisation
- ...and of course the software-infrastructure and people

Message of WP: scale of HTC is (on aggregate) ~ the HPC ~ several 10s Pflops

Often language is used that, inadvertently but wrongly, gives impression that HPC is the only large-scale computing needed in the UK (“unconscious bias”)

The Whitepapers

Other Points of note:

- **Software-infrastructure is as essential as physical-infrastructure**
 - UKRI AAI, Security and trust
 - Accounting
 - Virtualised customizable infrastructure (cloud, VMs, containers and federation thereof)
 - Exabyte data management, organization, transport and access

- **The UK RDI is a heterogeneous distributed infrastructure by its very nature**

- **Federation/Sharing/Interoperation is key**
 - Including Interoperation outside of UKRI -> NHS, Commerce

- **People and skills are key (crew on ship)**
 - ➔ Separate talk

Recommendations

The executive summary :

1. Research Data Infrastructure (RDI)

Investment is urgently needed now, in the period 2020-22, to put the UK on a world class footing in respect of physical infrastructure and software infrastructure, reversing the significant gap that has arisen over the last few years.

2. Research Data Exploitation and sharing

Each Sector should refine and update its RDI requirements, in terms of its own Research Data Life Cycle, such that data are supported at each stage in the life cycle and can be readily analysed, discovered, combined, reused and repurposed.

3. International Collaboration and Leadership

Coordination structures are needed commensurate with the fact that the creation and use of Resources for research data are increasingly an international activity, with major subject-specific repositories having a global reach.

4. People and Skills

Investment is needed in people needed to create, engineer and apply the advanced computing techniques to the data to extract knowledge and innovate.

Detailed recommendations : capacity

Refers to urgent need (ship sailing over edge)
Scale to solve immediate problem £100m across UKRI

	Recommendations
1	Investment is urgently needed now, in the period 2020-22, to put the UK on a world class footing in respect of physical infrastructure and software infrastructure, reversing the significant gap that has arisen over the last few years. Physical RDI (computational capacity, storage capacity and networking) should be put in place to host and exploit the data. This means adequate investment in HTC and HPC facilities, both within UKRI and where appropriate using commercial cloud provision, and in the software infrastructure to manage, move, give access to and analyse the data. This should be complemented with networking connectivity via Janet. This is a prerequisite to all other recommendations in this document.
2	The actions listed in Tables 3 and 4 for maintaining and transforming the UKRI be executed on the suggested timescales with investments of £200-300M p.a. beginning in 2020.

Refers to 2020-2027 timescale - steady state
Scale ~ physical £200M, people ~ £100M

Detailed recommendations : governance

13	A UKRI governance group be set up as soon as possible to plan for and oversee the delivery of these infrastructures and services.
14	A national coordination body be set up to co-ordinate the UK's international RDI activities. articulate the national interest, steer UK involvement in international RDI initiatives, ensure that involvement is well-informed (for example about good practice in implementing FAIR), monitor its outcomes and keep UK strategic stakeholders informed.
15	To co-ordinate the Stewardship Infrastructure (data management and curation services), a coordinating group has to be set up (10 FTEs are needed to coordinate UKRI activities). In particular, these roles will need to negotiate with the various partners and collaborators, (particularly with the implementation of FAIR principles), influencing policy and international data management activities and communicate with (i) the large research data science communities in the UK and abroad and (ii) communicate our impact to all stakeholders.



Whitepaper:

UKRI Cloud Strategy 2019

Cloud / Commercial Cloud

Distinction:

- *Cloud* = use of virtualised computing running on physical hardware. Generally OpenStack, VMs, containers, orchestration
- In the research sector this is technically a “business as usual” component of the ecosystem already. We know how to do this.
- Use of *commercial cloud* for the same workflows is beset by serious structural issues that only BEIS/UKRI can solve.
- These are addressed separately

Virtualised computing

Virtual (cloud) computing in the research sector is already commonplace for us

- CERN runs a lot of its infrastructure virtually
- RAL Cloud provides substantial resources to many users
- IC and Cambridge run IRIS cloud
- IRIS demonstrator ongoing to use commercial cloud as a back end
- HEPCloud (Fermilab, CMS)
- Many universities run clouds

Insulates science activity software stack from the underlying hardware

- Different science activities can create different environments
- Promotes sharing of infrastructure

It works perfectly fine for some workflows

- In particular for simple “worker nodes” doing bulk computing
- Currently good when there is little I/O

It remains inappropriate for other workflows

- In particular, and importantly, leading edge true HPC use
- Complex workflows in very large activities can be difficult to adapt

Virtualised computing (cloud, OpenStack)

So from technical point of view the message is

- Use of virtualized computing in STFC is well established
- STFC and its scientists have lots of expertise
- It is business as usual for those workflows for which its useful

But all of this has nothing to do with the barriers to the use of commercial cloud as a component of the ecosystem, which are:

- Policy barriers
- Cost/Business model barriers
- Inappropriate research funding instruments (not designed for this)
- Funding silos (capital | resource) mismatched

i.e. there is almost no technical impediment to the use of commercial clouds for those research workflows that it is good for, but these structural barriers are a showstopper.

[Note: Some of these barriers are diminished for University IT Departments]

Use of Commercial Cloud

The fundamental barrier is that the Cloud Vendor business model is maximally mismatched to the research model

Data hosting is the obvious case:

- Vendor models based on pay-per-access
 - This is maximally contradictory to Open Assess policies
 - Science data must be accessed arbitrarily many times freely at point of use
 - Data stored at vendor-X must be computed on with CPU at vendor-Y without cost
 -endless list of this sort of stuff.....

- Don't misunderstand: short term “free egress for x% of data” arrangement are not a long-term sustainable engagement model.

- Research needs storage costs framed as “per-TB-year”
 - There is no instrument to allow pay-per-access costs in perpetuity even (if it were desirable)
 - [Actually, there is currently no FEC consistent mechanism to pay for storage past the end of a grant anyway, but let's not go there...that's a different debate]

Use of Commercial Cloud

Similarly for CPU

- **Vendor models based on variable spot pricing**
 - This is maximally contradictory to large scale science production computing which has deterministic timing and in some cases is almost DC (continuous)
- **Of course many demonstrations have been done of “mix-n-match”**
 - E.g. HEPCloud is one: turns on use of commercial cloud (if there is a budget) when either price falls or user agrees to pay the premium rate
- **Costs are still in any case ~> 2 times greater**
- **And exactly as for data, research grant instruments are not adapted to pay these bills anyway, or to freely move money between equipment and vendor bills.**

**It feels like walking through treacle whenever trying to do any of this
Same groundhog-day discussions for years**

Serious Vendor ⇔ Research Engagement Interface

If we are to be serious about wanting to see Commercial Cloud as part of ecosystem then it has been recommended (3 years ago !) that BEIS/UKRI/JISC:

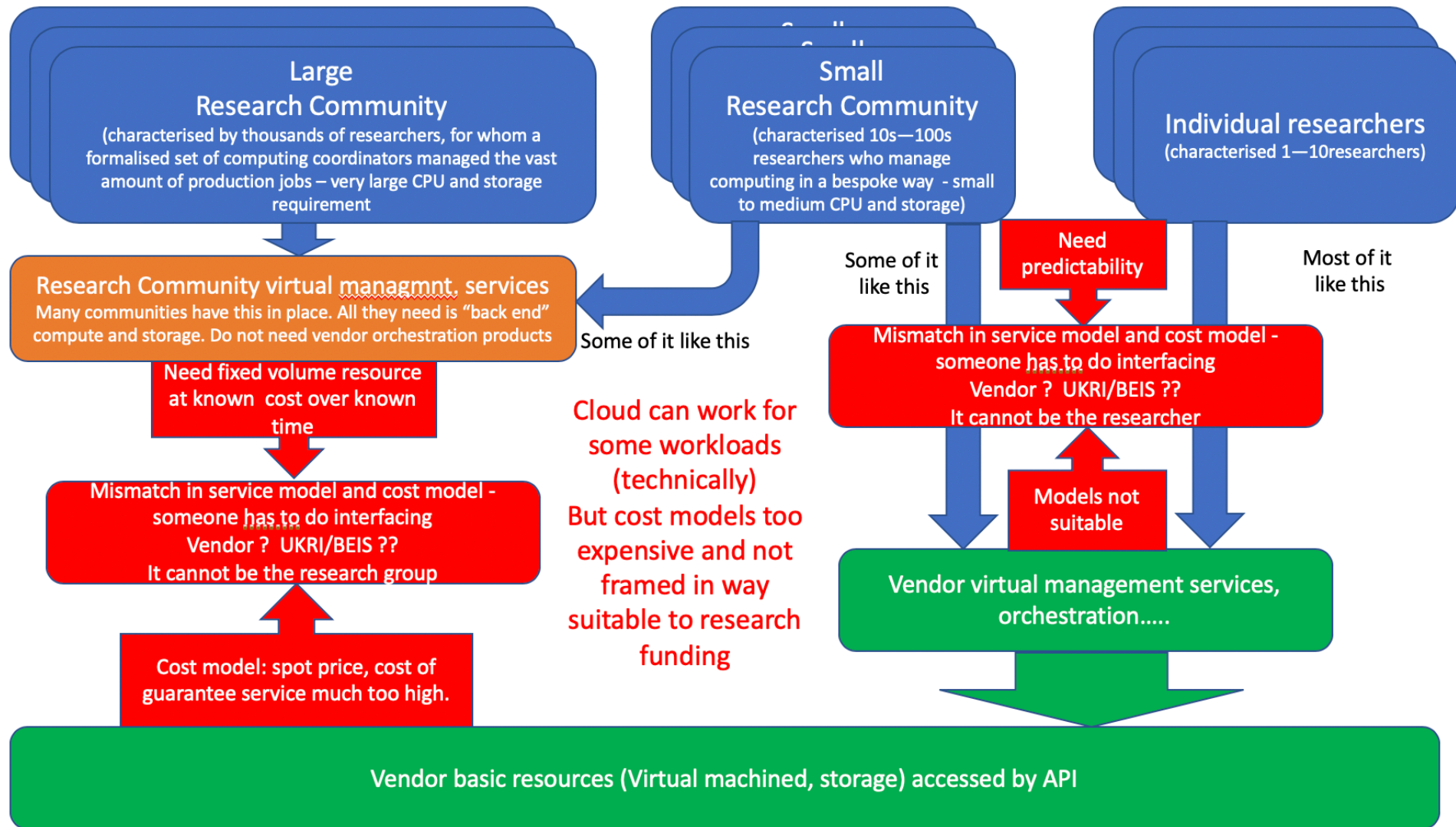
Set in place a high-level working group comprising UKRI experts and Cloud vendors to create a workable interface business model between two domains

This can only be done at BEIS/UKRI/JISC level - individual activities cannot do this

Such an exercise may, for example, conclude :

- with vendors developing a native research facing business model ? are we big enough ?
- or that that an interface layer is needed to interface the two domains
 - UKRI pays bills for per egress and spot price
 - UKRI makes this available to research on pay per CPU-hour and TB-year basis
 - UKRI has to make this consistent with UKRI's own open data policies and grant instruments
 - Agreement would have to be competitive with on-prem or it won't get used
 - Someone has to create and run this service at UKRI level (JISC ?)
- or other solutions ???

Until this is done, there is little point in continuing to have “localised” vendor ⇔ research area meetings which presume that the barrier is technical.



Summary of WP (abridged)

1. Cloud computing has brought important capability to the research community. There are many examples (technically) of using it successfully.
2. Research computing in the UK ... is diverse ... which needs to be met by diverse types of computing resources, including cloud.
3. Scientific and technical requirements must come first when selecting the most appropriate computing solutions. Cloud is not “a-priori” a solution for every research computing need. A hybrid ecosystem is appropriate
4. ... long-term *strategic* use of commercial cloud as part of a national e-infrastructure requires greater confidence around business, funding and governance models.
5. The current practice of ... large irregular capital awards is mismatched to commercial cloud charging models ... Equally, charging models from the commercial cloud providers are not adapted to the *needs* of research computing, ... making public cloud in many cases unsuited and uneconomic today.
6. Retention of staff with ResOps skills is critical

Summary

- ❑ UKRI has a clearly articulated requirement for large scale RDI comprising HTC computing, Storage and Networks.
- ❑ The scale of the physical infrastructure alone is ~ £200M per annum for UKRI
- ❑ There is a clear and present short term cut-off of funding which is blighting planning for long term science - (for STFC this is ~12M p.a.)
- ❑ The Whitepaper on RDI addresses this and notes that it is a pre-requisite to exploiting any of the data
- ❑ Operations and people are, as always, a critical part of this.
- ❑ Cloud computing technology is well understood and used widely within UKRI science
- ❑ The use of commercial cloud capacity as a systematic component of the ecosystem is blocked by structural, policy and funding instrument issues that can only be solved at UKRI/BEIS/JISC level

Backup

HPC is not a free alternative to HTC

- **HPC cycles are more expensive than HTC cycles - always in every country**
 - By construction ...literally... they cost more to construct
- **HPC use for HTC in USA is driven by policy, not economics**
 - DOE/Regions make large investments in massive overprovisioned HPC machines due to their Exascale agenda
 - Typically 100s of PF (200, exascale by 2021 has been stated)
 - Then mandate HTC users should use them - but free at point of use
 - NERSC is prime example as a “user facility”
- **UK is not equivalent to USA**
 - ARCHER is small (few PF). ARCHER-II is in construction
 - ARCHER explicitly do not support HTC use (we checked)
 - EPSRC have no role to support HEP
 - DiRAC 2.5 is struggling to support PPAN Theory+Cosmology
 - DiRAC-3 is unfunded, and if it were it would not have large underused cycles for HTC, and if it did it would not be cheaper.
 - There is no “place to go to ask for a large free HPC allocation in the UK
- **Neither ARCHER or DiRAC will be built with 20-30 Pflops of spare (unused) capacity to be made available for HTC, and even if they were it would be more expensive per cycle to do so.**
- **HPC is in any case technically difficult and labour intensive to use**
 - Always technical developments needed for internet access, CvmFS, containers, OS, batch system.....
 - Needs to be done afresh at each HPC machine (“no standard fix”)
 - So would need more staff at least initially
 - WLCG has prepared a “handshaking” document to specify what an HPC centre needs to provide to be “usable

The costs are indicative and it should be noted that they are subject to audit to establish the existing baseline. This baseline should be available by Q4 2019.

Table 4: Action Items and indicative costs

Item	Time Period	Resource	Indicative Annual Costs £M p.a.	Responsibility
1. Lights on Investment in Physical RDI	2020-2022	Equipment & Operations	150-250	UKRI & Sector/Large Projects/Facilities
2. Setting up Governance Structures	2020-2026	10 FTE	1	UKRI
3. Review of Current capabilities	2019	10 FTE	1	UKRI
4-8 Transforming e-Infrastructure: Software	2020-2022	240 FTE	24	UKRI
9. Develop data storage and re-use capability within JISC	2020-2022	Equipment & Operations	10	UKRI, JISC
10. Setting up Data Stewardship Infrastructure at the national facility level	2020-2022	500 FTE 50 FTE (Fellows)	55	UKRI, Research Councils
11. Physical RDI Dependency work	2020-2027	Equipment & 20 FTEs	8	UKRI, JISC
12. Physical Infrastructure Hardware: sustained	2022-2027	Equipment & Operations	150-220	UKRI
13. Physical Infrastructure Software: Sustained	2022-2027	200 FTE	24	UKRI and Research Councils
14. Data Stewardship Infrastructure at the national facility level: Sustained.	2022-2027	500 FTE 50 FTE (Fellows)	55	UKRI, Research Councils

Capital Funding matrix (hardware) - responsibilities for sectors and status

	SCD	DiRAC	GridPP	IRIS
Facilities	Only 20% of requirement is funded annually			Makes up some of shortfall until 2021
PPAN Theory		DiRAC-3 is as yet unfunded. STFC/BEIS funds DiRAC-2.5y to keep lights on		
PPAN HEP			Only partially, even if GridPP6 hardware request fully funded.	Makes up some or all of of shortfall until 2021
Other PPAN (Astro, PA, nuclear)				Provides some of requirements until 2021