# Highly diverse computing (including the long tail)

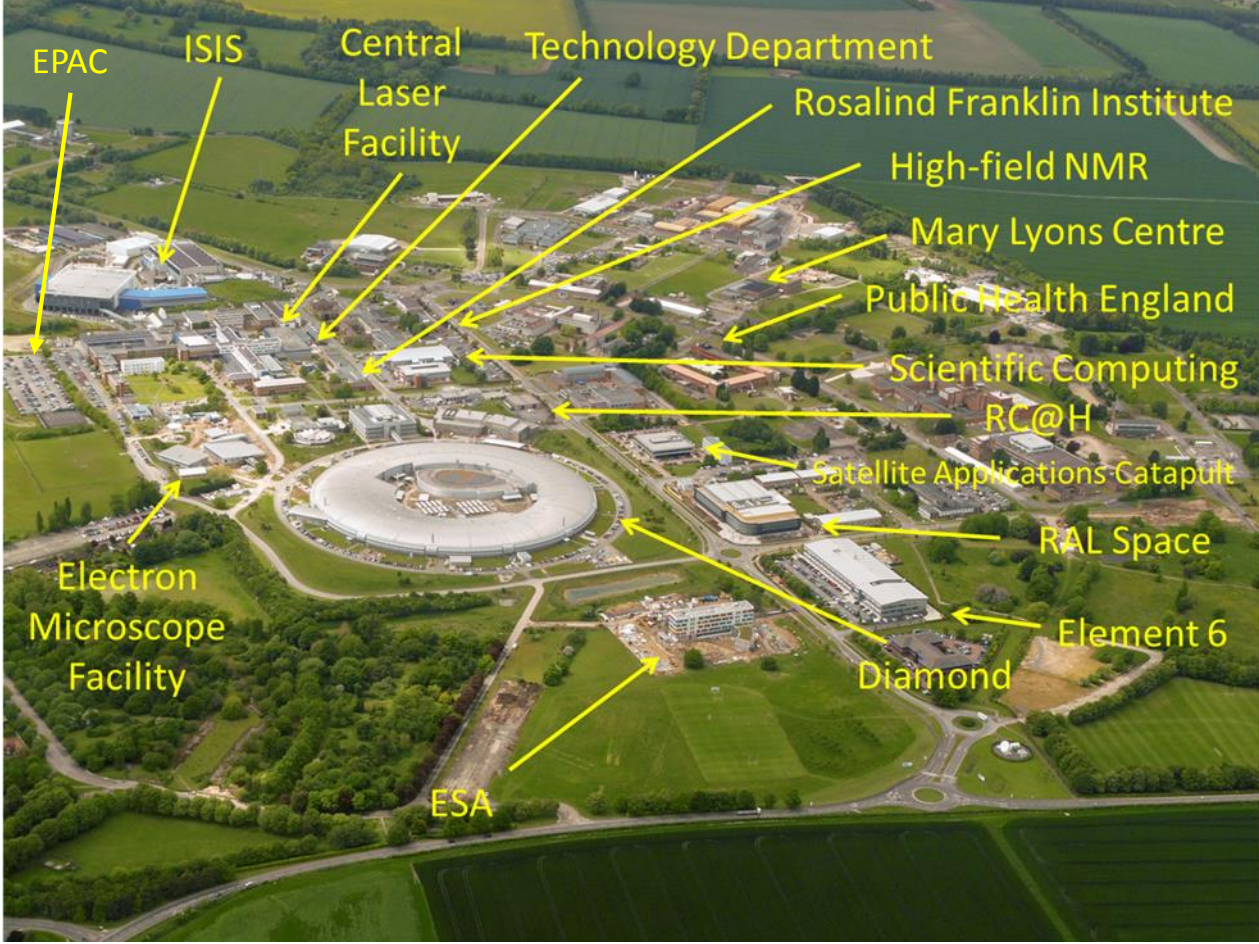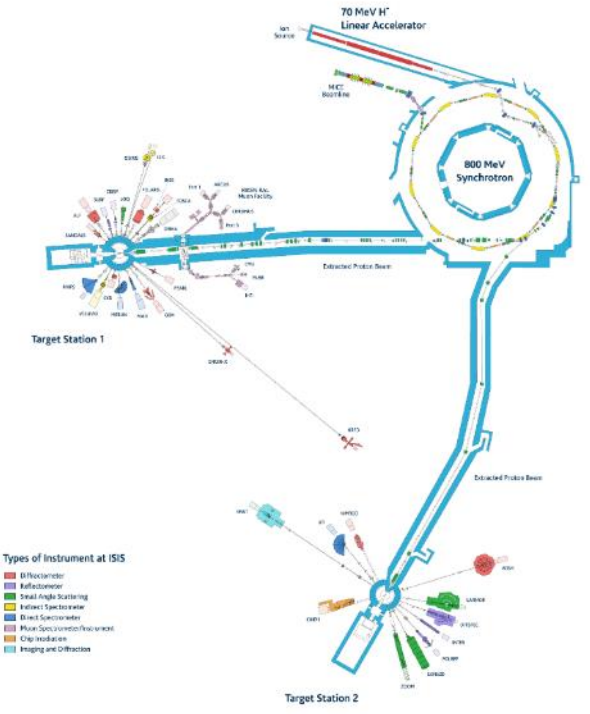## or

## Computing for STFC Facilities

# What are the STFC Facilities?

# What are the STFC Facilities?
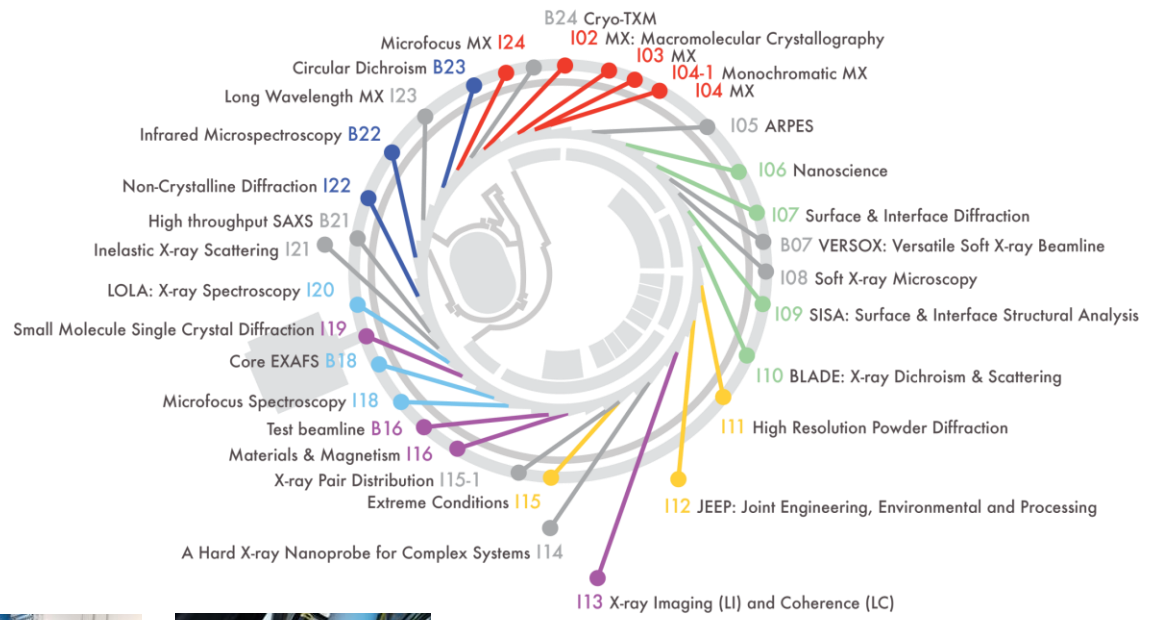
# Who are the users?



- Diverse …

- Thousands of users – national and international, on-site and remote

- Biologists to engineers

- Experts to non-experts

- 24/7 operation, multiple simultaneous user groups, individual experiments from minutes to weeks

- Continuous development, small and large (eBIC/ePSIC, EPAC, Diamond Lattice upgrade …)

# Typical User Data Flow

# Data

- Volumes
- Fast/slow access
- Metadata
- Preservation, curation
- Policy





Total data archived from Diamond&eBIC



Raw Data (GB) per Year

# Diverse demands for computing: Examples

**Macromolecular Crystallography**
- High throughput beamlines requiring near real-time feedback on data collected
- Typical collecting 300 to 900 data sets per day,
- Data range: 5TB to 35TB per day RAW data, 2.5TB to 17TB per day processed data

**Imaging Beamline**
- 10 to 200 datasets per day
- Data range: 0.5TB to 20TB per day (technique dependant and detector dependant)

**Electron Microscopy (eBIC – Life Sciences)**
- 1 dataset collected per 48 hours
- Ranges: 4TB to 17TB per day

**Electron Microscopy (ePSIC – Physical Sciences)**
- 20 to 300 datasets per day
- Data range few GB's to 100's GB per day

diamond

# Computing infrastructure

- Hardware – at facility or elsewhere (e.g. SCD Scarf)
- Fast/slow access
- Cloud (e.g. iDAaaS)
- Support for modelling/theory?
- Networking (including WiFi, Eduroam)

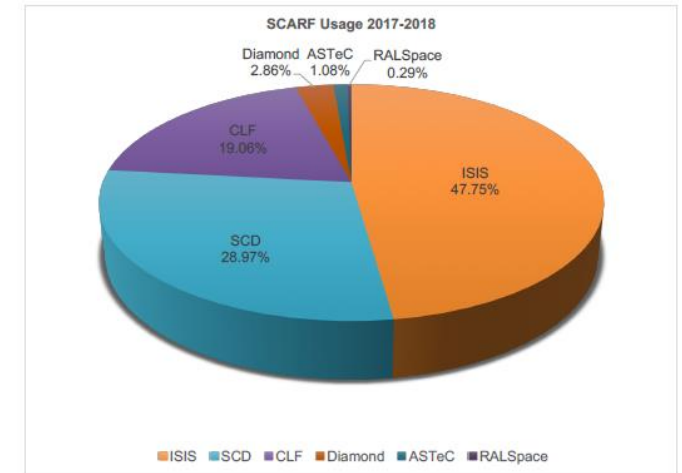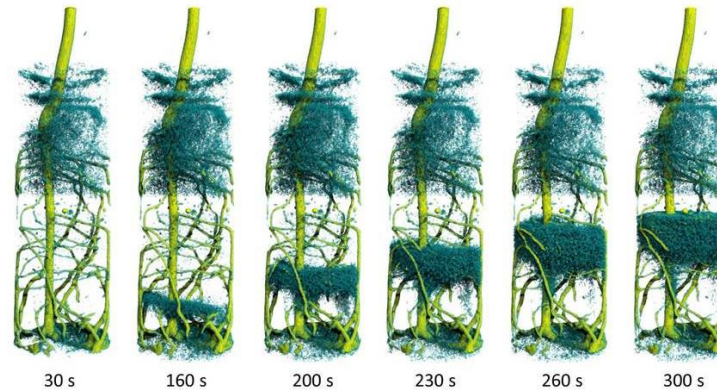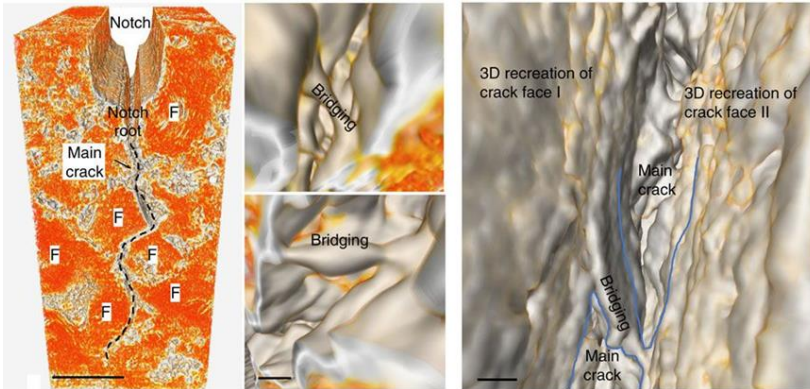# Diverse demands for computing: Infrastructure

Increasing data rates and data volume require ever increasing computational infrastructure to be deployed to support the demand

- Diamond currently operates and consumes:
- Two onsite data centres
- ~6000 X86 cores
- ~200GPUs
- ~16PB of high performance storage (GPFS) for data acquisition
- Utilises offsite IRIS infrastructure for a few pioneering projects and exploring other 'cloud' services particularly for post-visit analysis
- Retains experimental data on-disk for 30 days
  - High data rate experiments are impacted most by this retention period in managing their workflow for data analysis. After 30 days data has to be restored from archive for further processing.
  - Some low data rate beamlines in the 'long tail' can be afforded the luxury of keeping data online 'indefinitely'.
  - Overall data management has not kept up with the growth in data volumes. Poor data management is an increasing risk to data and timely analysis of data than just lacking computational resource to process data.

diamond

# Diamond Future Cloud

- **Offline / Post Processing** – enabling users to process data sets beyond the std. 40 day period '**easily**'
  - Need to transfer data from DLS to service (driven by DLS, DLS requirement to move offsite certain types of data processing)
  - Make 'Diamond Software stack' easily available e.g. via Docker/singularity

- **Burst computing** – 'top-up' backend infrastructure for Diamond
  - Transparent to end users
  - Only works for defined, relatively small (depends on target site) datasets.
  - Time – dependant analysis

- **Hosting** analysis services based around **Jupyter Hub** notebooks – user driven/user customisable
  - Need to transfer data from DLS to service (on-demand user driven service?)
  - Scalable infrastructure need

# Software development



- Full experiment lifecycle: Preparation, Experiment, Reduction, Analysis, Visualisation, Modelling

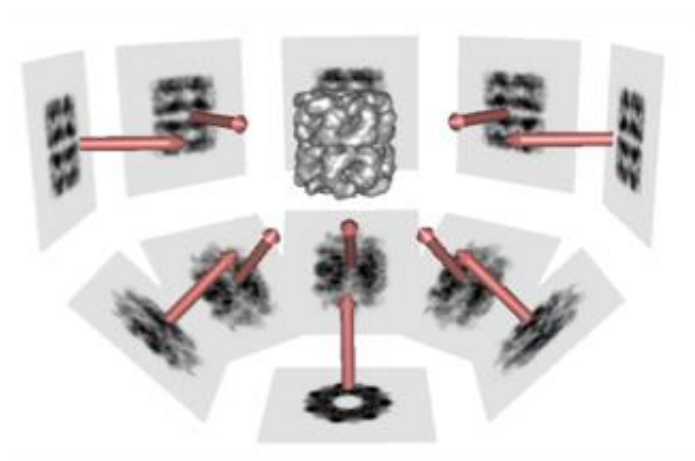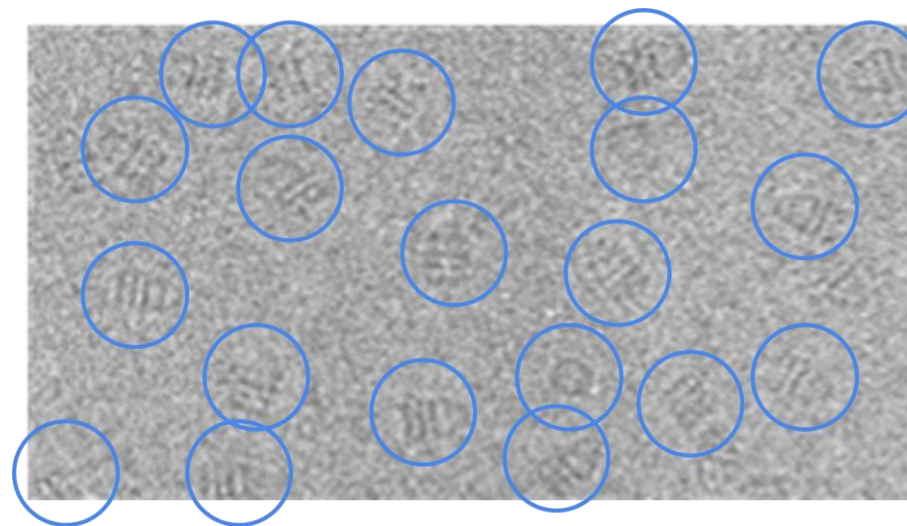- Transition from 'lone amateur' to 'professional team'' e.g. Mantid

- Collaborative Computational Projects

# Machine Learning



- 'Early days' for facilities
- Imaging was an early adopter
- Exploring many other possibilities
- Data quality/relevance screening
- 'Downstream' databases?

# User expectations

- More … and more … and more …

# Skills

- Recruitment, training and retention …

# Source

## Instruments

### Sample environment

Data treatment