



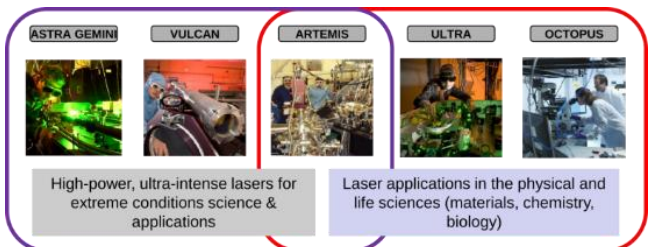
Science and  
Technology  
Facilities Council

# Data challenges at the CLF

Dan Rolfe  
CLF Octopus Imaging Facility  
STFC Rutherford Appleton Laboratory  
UK

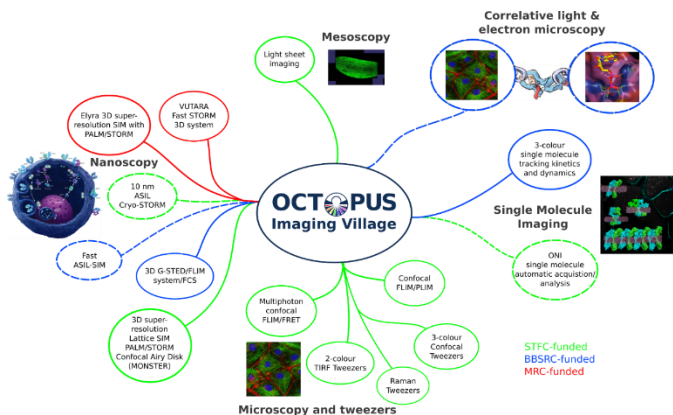


# Central Laser Facility



- Inputs are many
- Output is scientific data and papers
- Data management & exploitation is a growing challenge in CLF as it moves to big data regime

## OCTOPUS



- Cluster of microscopes and lasers and expert end-to-end multidisciplinary support
- Many different workflows

## Extreme Photonics Applications Centre (EPAC)

- £81.2M centre for applications of laser-driven sources in industry, medicine, security etc.
- New type of accelerator – plasma accelerator – driven by powerful lasers, producing x-ray and high energy particle pulses with unique properties
- **Will operate at 10Hz (upgradable to potentially 100Hz)**





# Data challenges at the CLF

- Several different facilities
  - a challenge and an opportunity
- Wide variety of
  - domains & applications, users, platforms, open/closed tools, techniques, workflows, granularities, throughputs, funding models/policies, rapidly evolving
- Deliver latest techniques at first opportunity
- Complex workflows with rapidly increasing data & compute volumes
- Live data and real time analysis
- Data acquisition, instrument control & diagnostics
- Analyses
  - image segmentation, feature detection, tracking, registration, correlation, parameter estimation, diagnostics, tomography, interpretation, simulation...
- Automation
  - to scale up throughput, make techniques practical & accessible
- Post-experiment data access, compute and analysis support
- Varied storage and compute solutions with common authentication & data ownership model. Flexible, agile, cross-platform...





# Data challenges at the CLF

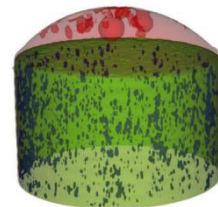
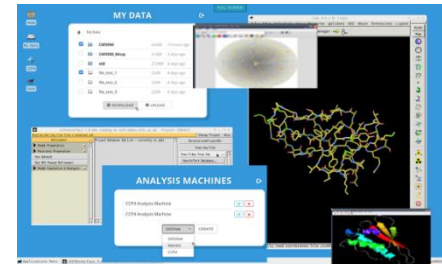
- Aligning approaches and working together – in which direction?
  - Across organisation?
  - Across technique?
  - Across application?
  - Multidisciplinary workflows and applications don't simplify this.
- Effort needed to exploit common tools & infrastructure can scale with number of scientific workflows.
- Archiving, sharing & curation
  - User responsibility vs facility responsibility
- Not all users will have expertise and/or compute/data resources to handle facility data or specialist tools
- Volumes of data to move, especially if external cloud solutions to be used.
- Metadata/datamodel different for each technique & application.
  - Essential for high throughput/automated workflows and meaningful ongoing exploitation/sharing/open data.
- Balance between using shared approaches and solutions or agile/optimal solutions for specific workflows/techniques/projects.



# Data challenges at the CLF

Many approaches in use and development

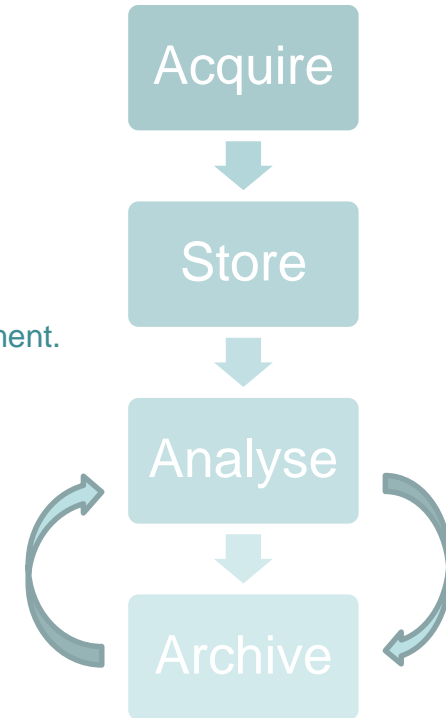
- **Great expertise on Campus to collaborate on solutions and increasing initiatives to support this**
- Analysis approach is workflow-specific e.g. in Octopus
  - Host data and compute on site for data intensive workflows
  - Users take data and analyse themselves for less demanding workflows
- DAaaS
  - Remote desktop via web with configured software/tools and access to data e.g. CT processing (EPAC) or single molecule analysis in OCTOPUS.
- STFC Cloud – onsite supported compute cloud
- SCD CephFS – scalable live disk storage
  - Octopus from 100 to 800TB in two years
- SCARF
  - STFC HPC cluster for facilities
- IRIS (eInfrastructure for Research & Innovation for STFC)
  - STFC community sharing & coordinating digital research infrastructure
  - Funded 25 GPU nodes and 54 CPU nodes for Octopus SLURM cluster in Cloud
- Ada Lovelace Centre (ALC)
  - Focused on improve data science and data exploitation in the facilities by funding projects and cross-disciplinary approaches
  - Funded development of OCTOPUS data infrastructure
- SCD SciML group working with us to deploy AI solutions in many areas
- EPICS for instrument control developments
- Uncoupling things for maximum flexibility



Science and  
Technology  
Facilities Council

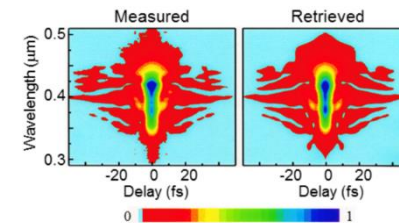
# Data Management in EPAC

- Many diagnostics and detectors producing data at high repetition rates
  - Current estimates ~ 5 GB/s peak – 10Hz, if we move all data through system
- Data transfer may not be very easy
  - Data volume/experiment could be significant; will need centralised storage and management.
  - Annual data volume could be 1-2 PB at the start (but could be theoretically 10's of PB, if we chose to save everything). Storing all the data would be very expensive
  - Depends on operational modes –we will need to be choosy



# EPAC: Data types

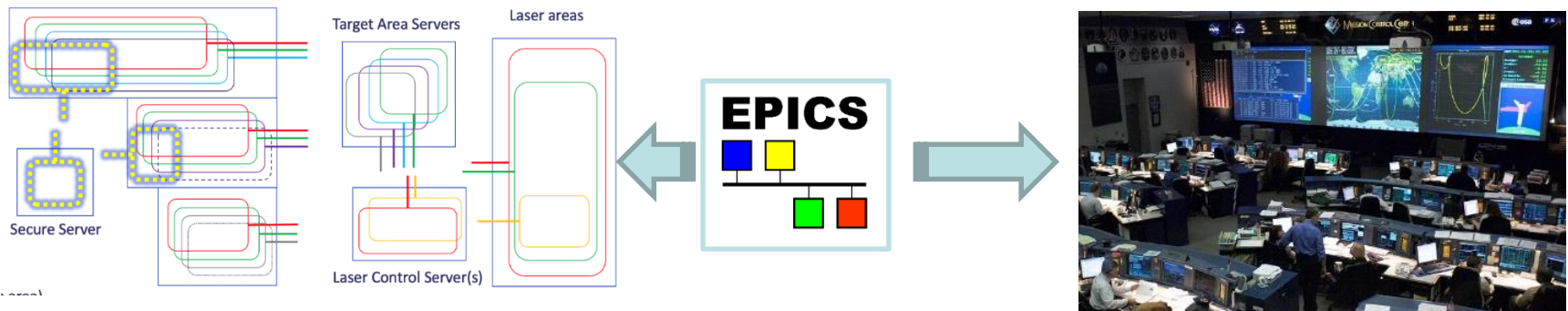
- Live data
  - ✓ Continuous running data streams – facility monitoring
  - ✓ Only some raw data needs to be archived
  - ✓ Needs continuous display
  - ✓ Most of it need not go through the data management system
  
- Dip-stick data
  - ✓ May need processing– eg. pulse length
  - ✓ Dipstick – at different rates: every second/minute/10 mins
  - ✓ Need metadata
  
- High rep-rate data
  - ✓ Most of Experimental data, some with processing, eg. CT
  - ✓ Need to archive both raw data, metadata and processed data
  - ✓ This will decide the architecture suitability



Machine Data & Facility Data may need different treatment

# Data Capture Concepts

- **Not ALL data** need to go through the data management system
- Some will just need to be displayed throughout (Live Data)
- Still need **network infrastructure** coping with peak data rates – EPAC will have single mode Fibre connections for this



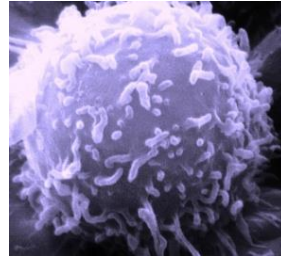
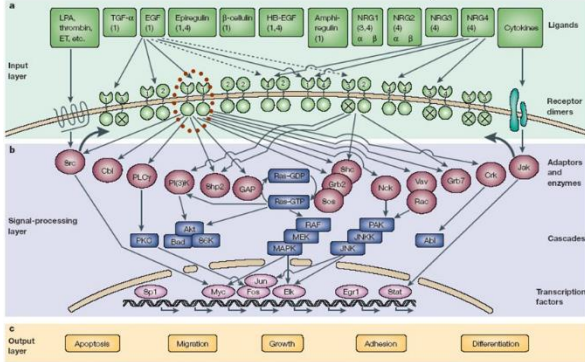
A significant fraction of the data, however, will go through the data management system



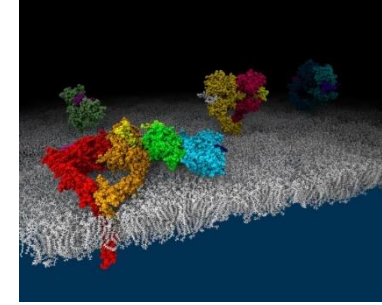
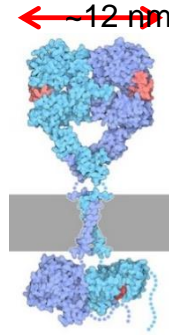
# OCTOPUS FLImP example

## ErbB Family Receptors and the Signaling Pathways

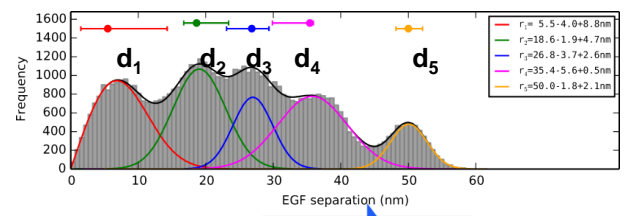
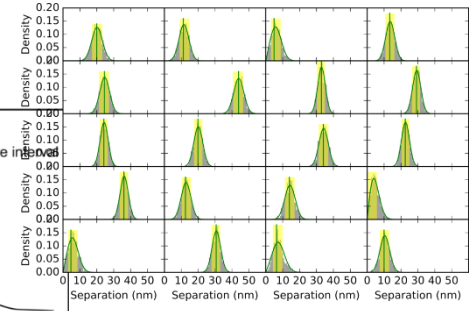
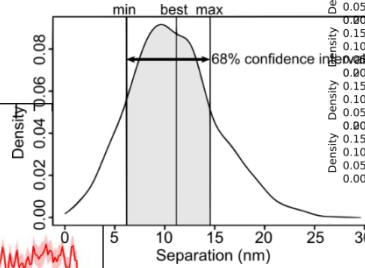
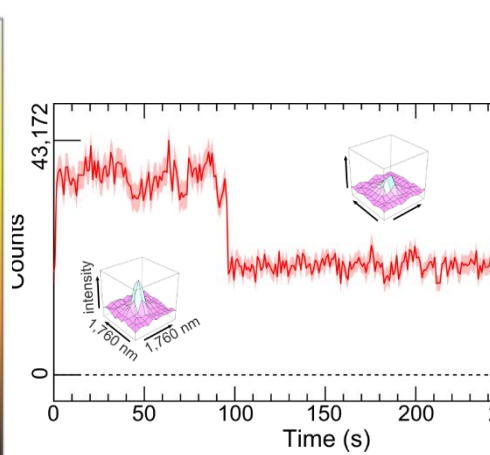
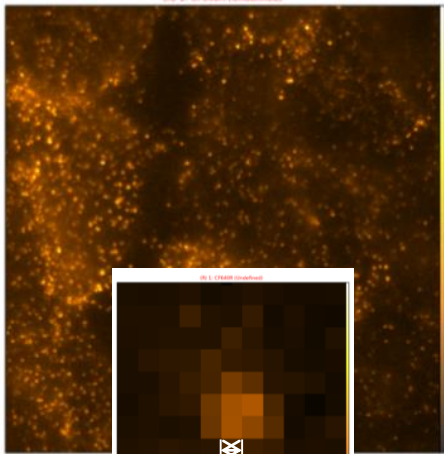
Yarden and Slivkowsky, nature reviews, 2001



CANCER CELL



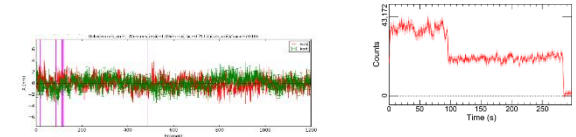
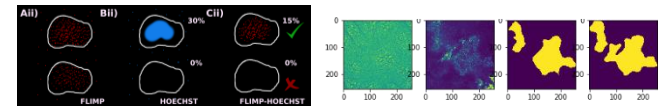
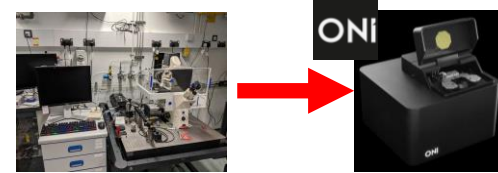
## CHO cells expressing ~ 10<sup>5</sup> EGFR/cell



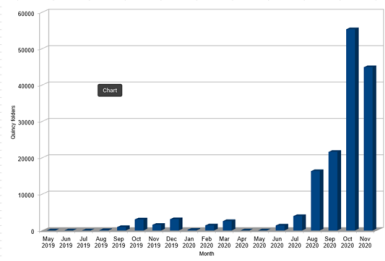


# FLImP example

- Acquisition
  - ONI Nanoimager + custom Python tools for autofocus, ROI identification
- Analysis
  - Automate drift correction & track selection
  - Pipeline coordinating multistage workflow
- Infrastructure enabling this
  - CephFS
  - SCARF, STFC Cloud + SLURM + our tools + Singularity
- Speed up achieved ~ 50x
  - 1000s of cluster jobs and ~TB data per day
  - Now limited by raw data archiving solution
- Funded
  - external grants (BBSRC), internal (STFC CLASP), facility development, IRIS, ALC
- Room for improvement in all areas but able to offer FLImP to facility users next year



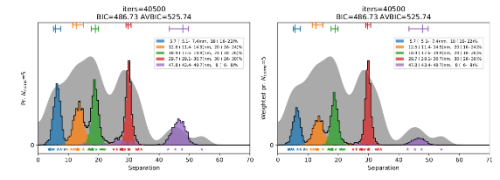
010019	0
010019	1
010019	0
010019	100
010019	847
010019	1004
011119	1048
011219	3104
012119	141
013220	1385
014220	2975
015420	0
016520	1336
017620	3005
018820	18317
019920	23811
011020	50314
011120	44883



Our @Imaging Nanoimager has been a great platform for automation of our @STFC @Oxford microscopical imaging method for cancer research and diagnosis.

• ONI @Imaging-22 Oct 2020

• We're bringing microscopy out of the darkroom! With remote operation on the Nanoimager, you can acquire and analyse your data anytime, anywhere! 📱 📊 Check out @STFC @STFC making #SuperResolutionMicroscopy look like a walk in the park! #Microscopy #Science



# I have a dream...



- Variety of separate, uncoupled data storage and compute systems supporting a common model/choice of tools for authentication and access to resources
- Cross platform, open APIs
- Fast, scalable disk storage for live access (~PBs)
- Archive
- CPU and GPU compute clusters, clouds, job management systems
- Remote desktop via web to VMs
- Containerisation to package software & reduce porting effort
- Minimal common data model to link with per-workflow existing and bespoke databases
- **Allow facility scientists and users to focus on science**