Exploring the Infrastructure of EMBL's European Bioinformatics Institute (EMBL-EBI)

Sarah Butcher Software Development and Operations Team Leader sarahb@ebi.ac.uk



What is EMBL-EBI?

- Europe's home for biological data services, research and training
- A trusted data provider for the life sciences
- Part of the European Molecular Biology Laboratory, an intergovernmental research organisation
- Home of the ELIXIR Technical hub



Deliver excellent research



Deliver scientific services



Train the next generation of scientists Engage with industry



Coordinate bioinformatics in Europe



OUR MISSION

North Pacific Ocean

Pacific

To provide freely available data and bioinformatics services to all facets of the scientific community in ways that promote scientific progress

Argentina



Big Data, Big Demand

62 million requests to EMBL-EBI websites every day

Requests from 24 million unique IP addresses

30 000 participants to EMBL-EBI Training events

307 petabytes

of raw storage in our data centres



Infrastructure = Hardware + Software



Maintaining the Platform

- Transforming Data \rightarrow Information \rightarrow Knowledge
 - Over 40+ data resources collecting, processing and releasing data
 - Correlation and cross-linking between different data resources
 - Frequently undertaken in international collaborations
- EMBL-EBI employs over 500 technical staff (+ research & admin staff)
 - Huge range: Software Engineers, DevOps, Sys Admins, ...
 - People key to maintaining the pipelines and services for our data resources
- Need to leverage a wide range of IT resources for big data analysis



Example Data Flow: Molecular Archival Resources





Data Centre & Public Cloud Infrastructure - 2021





Hardware

- Challenge is to scale the storage to match the compute
 - Large-scale analysis on a compute cluster can overload a storage system
 - Support a very mixed bursty workload
- Operate a variety of technologies
 - Storage: Flash, POSIX Disk, Object Store, Tape
 - Compute: Clusters with interconnects (& without)
 - Virtualisation: VMware Hypervisors, Delphix Database virtualisation
 - Clouds: Local OpenStack and partnering with public cloud providers
- Public Clouds
 - Provide opportunities for scaling out, bursty workloads, global services, etc.
 - Data governance, security and performance access a concern



Hardware Landscape

- Much of the access to datasets generates both high bandwidth (large data transfers) and high numbers of IO operations (e.g. file creations, deletions) placing substantial load on storage systems
- Workload distributed by spreading datasets and services across storage systems
- 300+ Petabytes of raw storage holding around 25 billion files and objects

- Raw storage
 - Object Store: 103 PB
 - NAS: 81 PB
 - HPC Storage: 27 PB
 - Tape: 55 PB

- Analysis capacity
 - Compute: 49,700 job slots (LSF)
 - Cloud (EMBASSY): 7,00 vCPUs
 - Virtual (Vmware): 16,600 vCPUs
 (>4600 VMs in use)



Strategic Growth

- Data doubling time slowing from every 12 to every 18 months over last 5 years
 - Major data repository (55 PB) still seeing an average of 1.5PB/month growth
 - Need the capacity to scale rack space, electrical power and cooling
- We make data available for long-term use: Open Data \rightarrow Open Science
 - Data made FAIR (Findable, Accessible, Interoperable & Reusable)
 - Freely available due to the investment from member (28) & associate (2) states
 - Additional grant funding from UK, EC, USA and others
- Given the growth of data networking a key enabler
 - Use UK's Network Research & Education Network (JANET)
 - Peers with international networks through GEANT



EMBL-EBI Embassy Cloud

Available for EMBL groups and collaborators



EMBASSY 🕥 cloud

Designed for Bioinformatics

Embassy has been developed with EMBL-EBI's research and service teams to provide an infrastructure tuned for bioinformatics workloads

Secure



Embassy provides a secure collaborative workspace with devolved administration



http://www.embassycloud.org/

FIRE – Online access to EMBL-EBI Data

- The FIRE system is a wrapper to an underlying object store (HGST) and a tape replica archive
- Data deposited into FIRE can be referenced via a persistent path within its single coherent namespace
- Used by some of the largest EMBL-EBI data resources for genomic and bio-imaging data
- Data are uploaded from external submitters through centralised file transfer services, processed internally and stored for public access
- Some sensitive datasets are controlled access
- Once public, data remain online for fast access





FIRE Architecture





FIRE Data are Active Data



Monthly Aggregated Data Ingest (TB)



Current size 55 PB (and 55PB tape replica)

- Traffic hard to predict accurately as most comes from external communities
- Also which datasets are 'hot'
- Max daily ingest 124TB (Jan 2021)
- Max monthly ingest 1.9PB (Jan 2021)
- Daily read average 136TB
- Max. daily read 170 TB
 - ~9 million daily actions, mostly GET i.e. read data
- 5PB of data going in and out in last 30 days (not internal HGST traffic or replication to tape)



IaaS Services

- 'Search as a Service'
 - EBI Search Apache Lucene Core since 2008
 - indexing 3.9 billion 'documents' for all EMBL-EBI's data resources
 - API •
 - Consumed by various core data services
- 'Tools as a Service'
 - JDispatcher bioinformatics Analysis tools
 - ~100 distinct sequence analysis tools
 - ~1.5 million jobs/day
 - APIs consumed by various core data services
 - OpenAPI, FAIR, OpenAIR and CWL compliant



Browse by type

XXX DNA & RNA	Gene Expression	₩ Proteins
Î ş Structures	C) Systems	्रद्ध Chemical biology
A.	Literature	Cross domain

Programmatic access

EMBL-EBI web services allow you to query our large biological data resources programmatically, so that you can develop data analysis pipelines or integrate public data with your own applications. The Web Services technology we use are built on open standards to ensure client and server software from various sources

Browse EMBL-EBI web services

Latest APIs paper has over 1500 citations since 2019 most during 2020 – 2021

Nucleic Acids Res 2019 Jul 2;47(W1):W636-W641. doi: 10.1093/nar/gkz268.



https://www.ebi.ac.uk/services/all

Provisioning Infrastructure for Web Public Services

- Web VMs (~2000)
 - public service-related workloads
 - public services web applications
 - team collaboration services Confluence, Jira
 - web services monitoring -Nagios, ganglia, Elastic Stack
 - web security e.g. WASF OWASP
- Web K8s (~40 clusters)
 - web deployment pipelines
 - 40 CI/CD services (GitLab)
 - web stats framework
 - ~80 million requests/day (2020)
 - ~37 million unique IP addresses (2020)

Containerisation of web services

docker

GitLab

- Giving the web developer more control
 - Remove software management barriers
 - Accelerate product cycle
 - Predictable environment
 - Isolation and segregation
 - Portability
 - Web security
 - Self-service

But with more freedom comes more responsibility



Infrastructure for Two Types of Web Services

- Centrally Managed (All EBI Services)
 - VM/K8s managed at the admin level
 - Visible on the internet
 - Keep root privs, Traffic management
 - Updates:
 - OS updates
 - Software level
 - Web Security (Web Application Security Framework)
 - SAST/DAST-style scans and mitigation tools

- User managed (not ebi.ac.uk services)
 - Users manage VMs and K8s containers and PODS e.g. via GitLab pipelines
 - U2TU
 - OS and SW level updates
 - Security (may use WASF)



'ResOps' Training Courses

- Bespoke course for EMBL-EBI teams developed by EBI Cloud Consultants
 - http://bit.ly/resops-2019
- Covering cloud-native technologies, tools, architecture and design
 - Starting with basics, through to fairly advanced topics
 - Docker, Gitlab, Kubernetes, with hands-on exercises
 - All material online, suitable for self-paced learning
- Course given 6 times @ EBI from 2019H2, to a total of 140 people
 - Also externally e.g. EOSC-Life WP1 in Berlin, EMBL in 2020



Two Classes of Web Services

- In-premise (now)
 - Typically well established services
 - Not designed for Cloud
 - Will be centrally supported
 - Traffic is managed via central Traffic Managers

- Cloud hosted (near future)
 - Nascent services
 - Designed for Cloud
 - Teams need to provide their own support
 - May be via EBI Cloud consultants
 - Traffic managed outside of central services in Cloud (e.g. Route 53)



Hybrid Cloud Use Cases





External Collaborations

- International Infrastructures
 - ELIXIR Compute Platform and related Implementation Studies
 - European Open Science Cloud (EOSC)
 - Academic Providers: EOSCpilot, EOSC-Hub, EOSC-Life
 - Commercial Providers: Helix Nebula (OCRE & Archiver)
 - EIROForum IT
 - GA4GH: Active within the Cloud Work Stream
- International Applications
 - Human Cell Atlas
 - Open Targets
 - Human Data Research UK
 - BioExcel II









Global Alliance for Genomics & Health

HUMAN



EMBL-EB

Summary

- Big Data is a driver across everything we do
- Meeting the 'big data' infrastructure challenges
 - Services for repeatability and resilience
 - Integration with cloud providers
 - Supporting services, training, research and administration
- Innovating to deliver new infrastructure
 - Working with application communities and service providers
 - Building a software development capability to support new infrastructure
- See a very cloud-native infrastructure-based future
 - Across internal and external infrastructure

