# Introduction to Module & Bioinformatics

Tim Hubbard @timjph

King's College London, King's Health Partners

Genomics England

Bioinformatics, Interpretation and Data Quality in Genome Analysis

MSc in Genomics Medicine

15th February 2016

# 5 days of Bioinformatics

- Day 1 - Introduction to variant analysis using NGS data and quality control
- Day 2 - Introduction to variant calling and Annotation
- Day 3 - Variant Annotation and Interpretation
- Day 4 - Researching links between genotype to clinical phenotype
- Day 5 – Additional annotation and genomic analyses

# Wider overview of Bioinformatics

- …the application of computer technology to the management of biological information. Computers are used to gather, store, analyze and integrate biological information

- **Pipelines** to process experimental data
- **Repositories** to store archive
- **Resources** to organise, present, interogate
- **Analysis** to make discoveries
- **Algorithms** to make predictions

# Biological Information

- **DNA** (copy of genetic material in every cell)
  - makes
- **RNA** (transcripts of expressed genes)
  - makes
- **Protein** (translation of coding regions of genes)
  - Linear sequences folds into a 3D structure
- Within **Cells**
  - Proteins interact with each other, metabolites, DNA, RNA
- Within **Organisms**
  - Cells divide, develop, interact

- **Genetics differences**
  - Change an Organism's behavior

# Biological Information – 1980s

- **DNA** (copy of genetic material in every cell)
  - makes
- **RNA** (transcripts of expressed genes)
  - makes
- **Protein** (translation of coding regions of genes)
  - Linear sequences folds into a 3D structure
- Within **Cells**
  - Proteins interact with each other, metabolites, DNA, RNA
- Within **Organisms**
  - Cells divide, develop, interact

- **Genetics differences**
  - Change an Organism's behavior

# Bioinformatics & Structural Biology

- **Pipelines:**
  - Xray, NMR data processing packages

- **Repositories:**
  - PDB (Protein Data Bank)

- **Resources:**
  - e.g. SCOP (Structural Classification of Proteins), CATH (Class, Architecture, Topology, Homology) databases

- **Analysis:**
  - e.g. molecular mechanisms

- **Algorithms:**
  - e.g. for Protein Structure Prediction

# CATH database

# Protein 3D structures: TIM barrels

Example Domain     2vxnA00 [PDB]        Example Domain     2vwsA00 [PDB]

# **Algorithms**: Protein structure Prediction

- Why is there a need for a prediction algorithm?
  - Vastly more protein sequences than structures
  - Xray/NMR slow, difficult, expensive
  - Ideally algorithm would predict consequence of mutation
- Strategies
  - Show sequence of unknown structure related to sequence of known structure and 'model' from it
  - Predict directly using knowledge of atomic interactions (physics) and simulation (*ab initio*)
  - Predict many models from 'lego' components from database of known structures
- History
  - 1970s–1994 Delusion
  - 1994 First CASP Competition; Steady progress since, but still not 'solved'

# Algorithms: Protein structure Prediction
# http://predictioncenter.org/
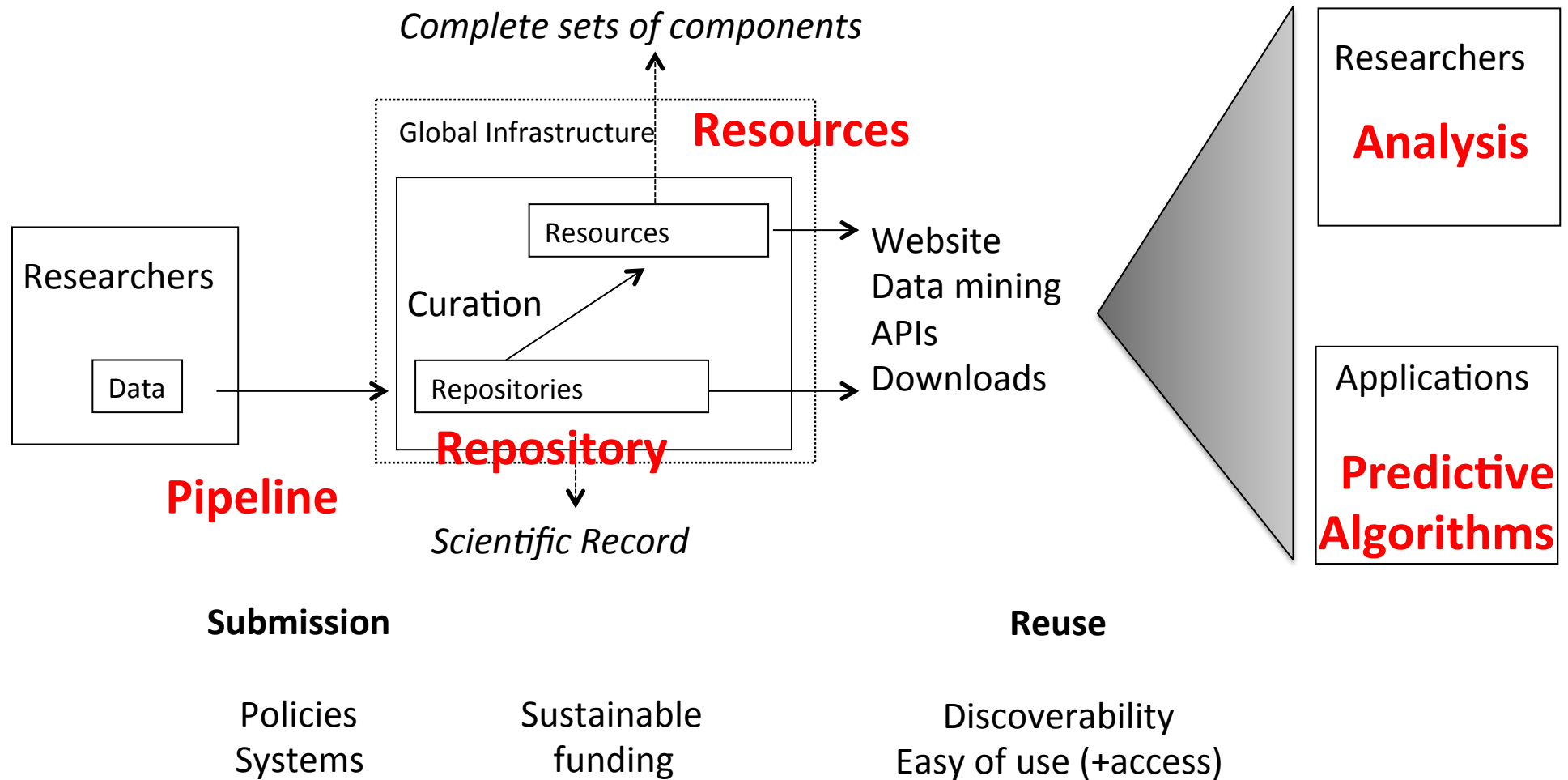
# CASP progress: 1994-2010

# Bioinformatics & Genomics

- **Pipelines:**
  - assemble whole genome from fragments of DNA

- **Repositories:**
  - Genbank/ENA/DDBJ

- **Resources:**
  - e.g. Ensembl

- **Analysis**:
  - e.g. genes, evolution, conservation etc.

- **Algorithms:**
  - e.g. gene prediction

# **Algorithms**: Gene prediction

- Why is there a need for a prediction algorithm?
  - Vastly more genome sequences than annotated genomes
  - Ideally algorithm would predict consequence of mutation
- Strategies
  - Collect transcriptome data; map back onto genome to annotate
  - Look for conserved regions between related genomes
  - Predict directly using knowledge of regulatory regions and simulation (*ab initio*)
- History
  - Ab initio algorithms performed well on single gene regions, but once large regions of genome sequenced, showed to perform poorly for vertebrates (~1998)
  - Transcriptomics data increasing easy to collect and algorithms relying on it because gold standard (Ensembl) + Curation for better accuracy (GENCODE)
  - Can still find new genes missing from transcriptome collection
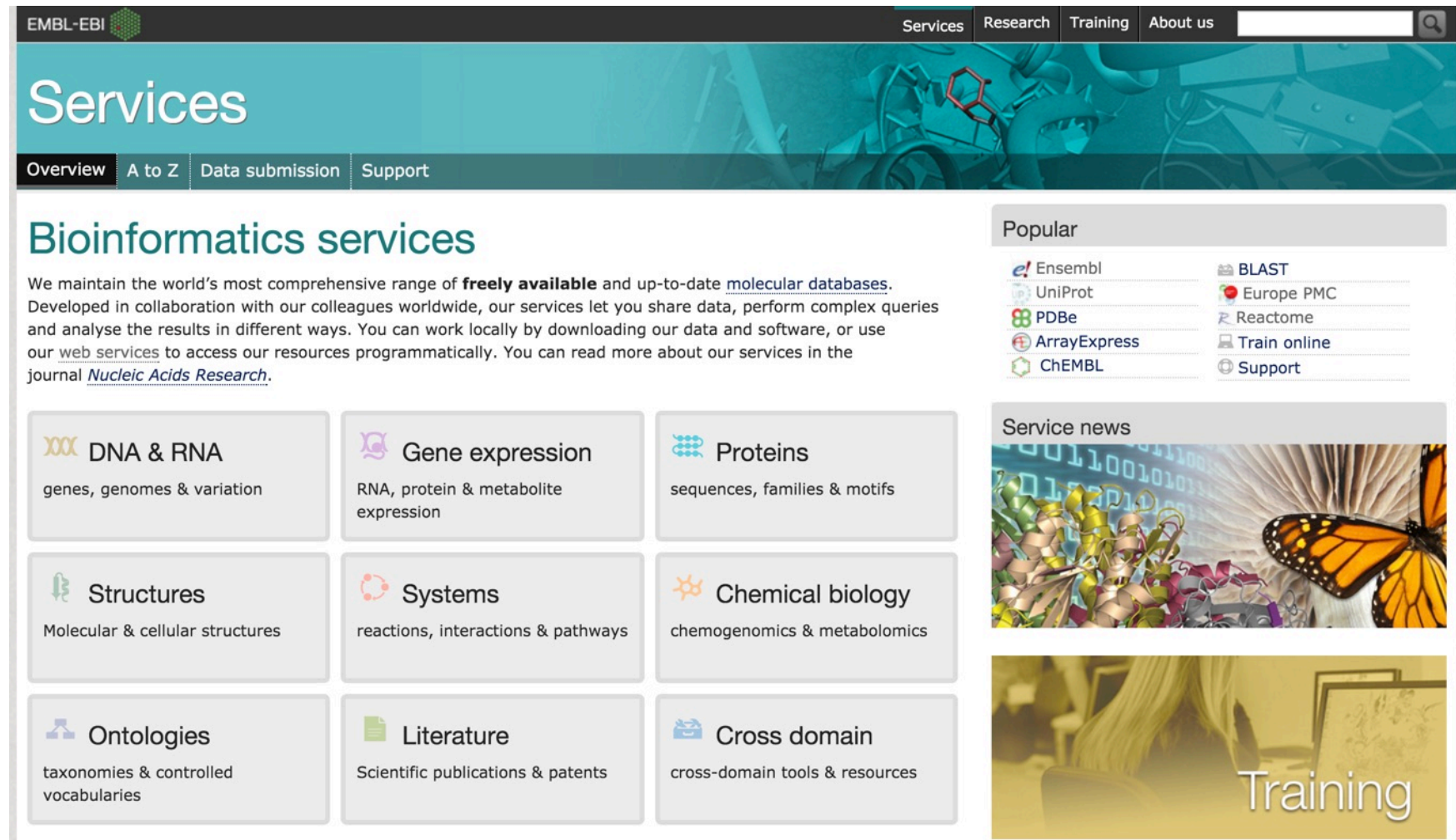
# Data, Databases & Bioinformatics



*Complete sets of components*

Global Infrastructure

**Resources**

Researchers

**Analysis**

Researchers

Data

Curation

Resources

Website
Data mining
APIs
Downloads

Repositories

**Repository**

**Pipeline**

*Scientific Record*

Applications

**Predictive
Algorithms**

**Submission**

**Reuse**

Policies
Systems

Sustainable
funding

Discoverability
Easy of use (+access)

# Organisation, Scale, History

- Centralised repositories for raw data
  - one data type, one repository
  - mandatory submission linked to publication
- Infrastructure to organise raw data for access
  - human genome presented to user as whole chromosomes instead of thousands of fragments
- Curated databases of biological objects
  - supported by evidence from raw data repositories

- First repositories >40 years old
- 1,000s of full time staff supporting infrastructure distributed worldwide

# Repositories in EU: EBI
# http://www.ebi.ac.uk/services

# Repositories in USA: NCBI
# http://www.ncbi.nlm.nih.gov/

# Scale up in EU: ELIXIR
# https://www.elixir-europe.org/

# Scale up in USA: BD2K
# https://datascience.nih.gov/bd2k

# Biological research is a grand project

- Build complete models of biological systems



- Future application to human medicine
  - Disease redefinition
  - Improved drug development
  - Personalised medicine

# Bioinformatics & Clinical Data

- **Pipelines:**
  - extract, clean, anonymise data from Electronic Health Records (EHRs)
- **Repositories:** NONE (privacy)
- **Resources:**
  - e.g. HSCIC (Health and Social Care Information Center), CPRD (Clinical Practice Research Datalink), PHE (Public Health England), NIHR HIC (Health Informatics Collaborative), GEL (Genomics England)
- **Algorithms:**
  - e.g. Text extraction

# Bioinformatics & Genomic Medicine

- **Pipelines:**
  - align individual genome to reference and call *variants*

- **Repositories:**
  - e.g. dbSNP (variants), [But not for clinical – privacy]

- **Resources:**
  - e.g. Decipher, Ensembl

- **Analysis**:
  - e.g. New targets for drug development

- **Algorithms:**
  - e.g. Prediction of disease causing variants

# Bioinformatics algorithm assessments

- Protein Structure
  - CASP – Critical Assessment of Structure Prediction (since 1994, CASP11 in 2014)
- Gene prediction
  - GASP, RGASP – Gene prediction and RNAseq assessments
- Variant effect prediction
  - CLARITY Challenge – 2012
    - http://genes.childrenshospital.org/
  - CAGI – 2010, 2011, 2013
    - https://genomeinterpretation.org/

# Biological Information – in theory

- **DNA** (copy of genetic material in every cell)
  ↓ *predict*
- **RNA** (transcripts of expressed genes)
  ↓ *predict*
- **Protein** (translation of coding regions of genes)
  ↓ *predict*
- Interactions that make up **Cell**
  ↓ *predict*
- Interactions that make **Organisms**

- **Genetics differences**
  ↓ *predict consequences for disease*

# Biological Information – in reality

- **DNA** (copy of genetic material in every cell)
  ↓ *predict*

- **RNA** (transcripts of expressed genes) ← *Expression, other omics data*
  ↓ *predict*

- **Protein** (translation of coding regions of genes)
  ↓ *predict*

- Interactions that make up **Cell** ← *iPS cells, organoids*
  ↓ *predict*

  ← *Tests, sensors*

- Interactions that make **Organisms**
  ← *Observation, self reporting*

- **Genetics differences**
  ↓ *predict consequences for disease*

High throughput data collection cost effective way to extend understanding

# From Genome Wide Association Studies (GWAS) to Whole Genome Analysis (WGA)



Genotypes

Exomes

Low Coverage Genomes

…to Whole Genome Analysis of Individual Genomes

Clinical Coverage Genome