

# AI, Data and their impact on Life Sciences research

Pistoia Meeting  
London  
April 2024



**Janet Thornton**

Director Emeritus



In the last 50 years, we have been living through  
**SEVERAL** revolutions.

DNA sequencing: Imaging methods: Computational power: AI methods

The impact on The Life Sciences is already huge



What did the invention of cars mean ~1900 for

- a horse-carriage driver
- a farrier; a blacksmith
- a bakery shop owner
- the city council of Cambridge/Heidelberg/Rome/...
- seaside resorts
- urban planners and developers
- the army
- ...?

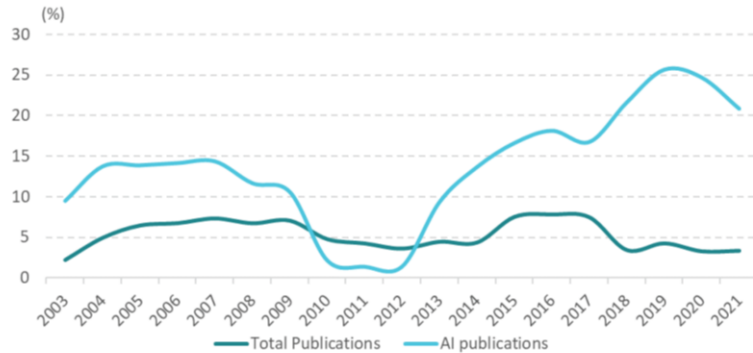


# What AI “thinks” how AI will change laboratories



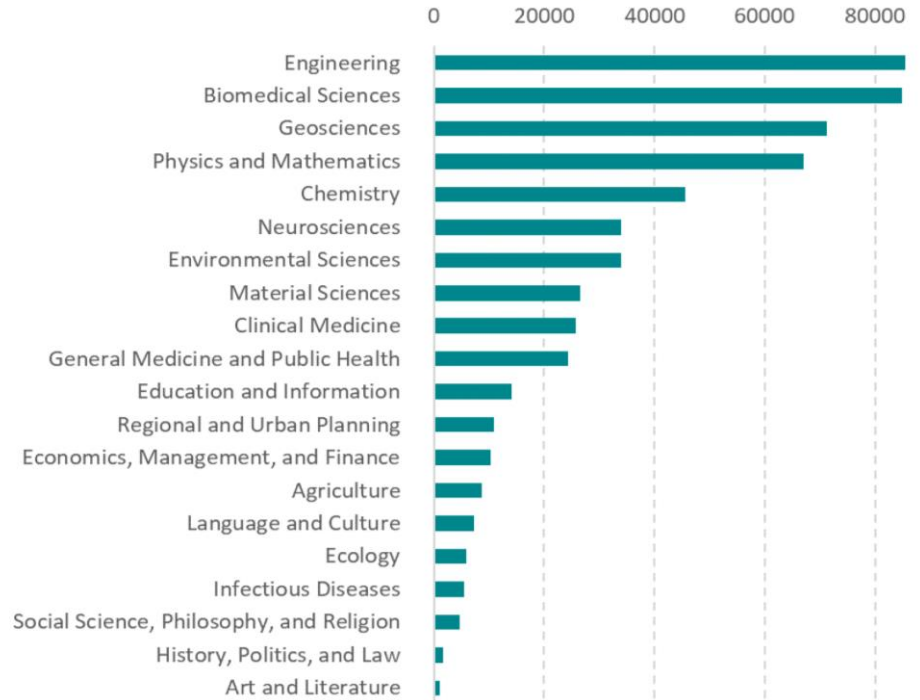
# Use of AI is growing rapidly in all fields of Research

## Growth in Scientific Activity – 3year average



Source: European Commission, DG Research & Innovation, calculations based on Web of Science data. Annual growth calculated as a 3-year rolling average.

## Number of AI publications (2017-2021) per Scientific Domain



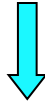
# Overview

- AI & AlphaFold
  - A solution to a long-standing challenge
- Biological Data & Infrastructure
  - EMBL-EBI
  - Community Actions
- Impact on Life Sciences Research
  - Impact in Biology
  - Impact in Medical Sciences

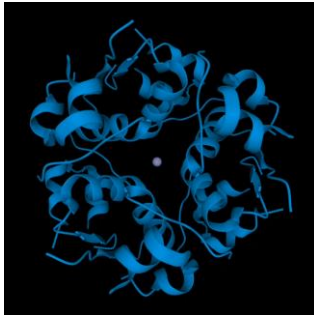
# The Challenge: Protein Structure Prediction

Protein Sequence

AGLYFE.....



3D structure



The unique 3D structure of a protein is determined by its sequence.

Can we predict structure (ie coordinates) from sequence?

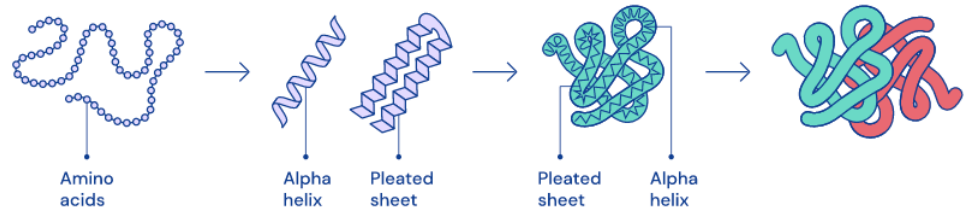
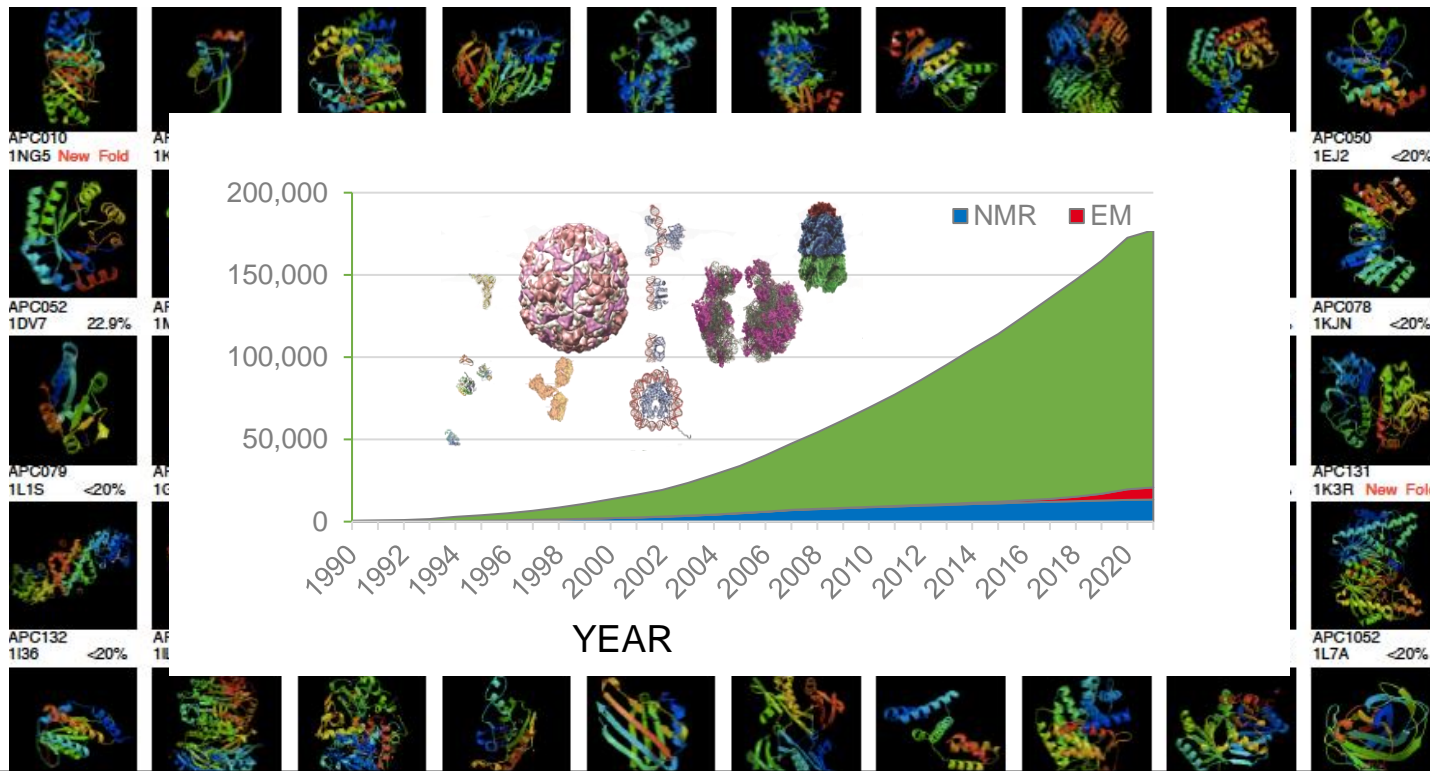
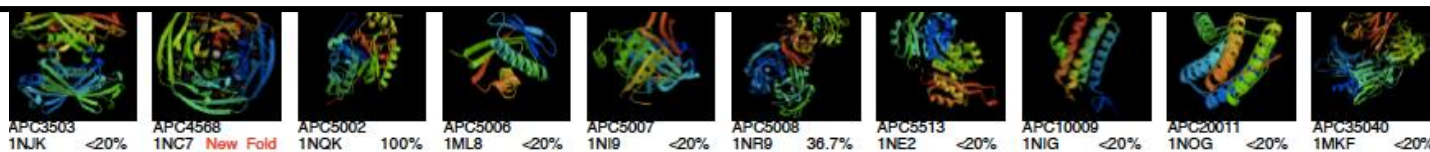


FIGURE 1: COMPLEX 3D SHAPES EMERGE FROM A STRING OF AMINO ACIDS.

# More & More Structures

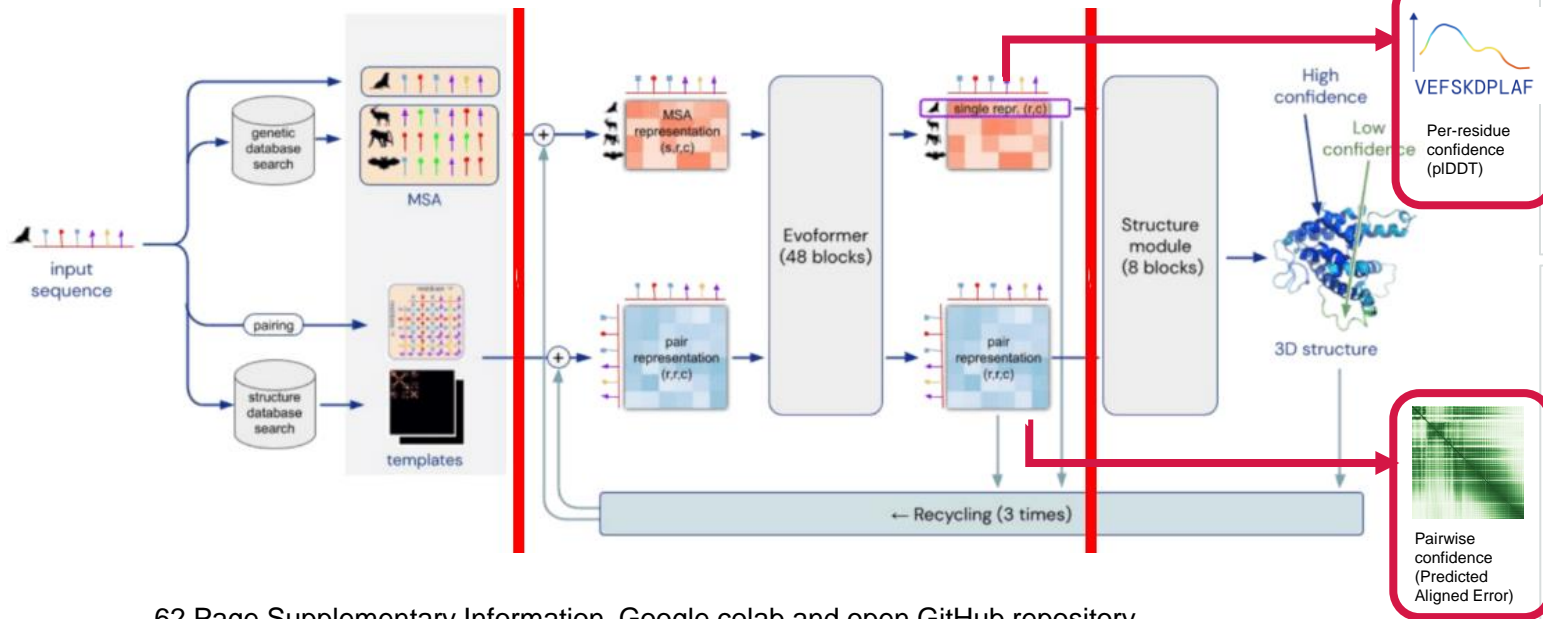


Only 0.027% of UniProt accessions are directly mapped to the PDB (55,000+)





# DeepMind's AlphaFold – core method



62 Page Supplementary Information, Google colab and open GitHub repository  
See Jumper et al. 2021 (especially the SI) for details

This is the first time a serious scientific problem has been solved by AI. - John Moutl

**Article**  
**Highly accurate protein structure prediction with AlphaFold**

John Jumper<sup>1,2\*</sup>, Richard Evans<sup>1</sup>, Alexander Pritzel<sup>1</sup>, Tom Green<sup>1</sup>, Michael Figurno<sup>1</sup>, Oriol Ronneberg<sup>1</sup>, Kathryn Tunyasubunrui<sup>1</sup>, Ryan Bates<sup>1</sup>, Michael Kohler<sup>1</sup>, Andrew Senior<sup>1</sup>, David Dalrymple<sup>1</sup>, Andrew Dalrymple<sup>1</sup>, Michael Senior<sup>1</sup>, Edward Taylor<sup>1</sup>, Iain Barr<sup>1</sup>, Andrew Brayn<sup>1</sup>, Saeed Adhikari<sup>1</sup>, Jeremy Chitambar<sup>1</sup>, Adam Britten<sup>1</sup>, Michael Bruff<sup>1</sup>, Alex P. Senior<sup>1</sup>, Edward G. Simonsian<sup>1</sup>, Boris B. Serdyuk<sup>1</sup>, David Reid<sup>1</sup>, Adam Senior<sup>1</sup>, Richard Woodcock<sup>1</sup>, Andrei S. Zhukov<sup>1</sup>, Sargur Sundaresan<sup>1</sup>, Mark Allwardt<sup>1</sup>, David Greig<sup>1</sup>, Drago Srebric<sup>1</sup>, Andrew Senior<sup>1</sup>, Andrei Roșca<sup>1</sup>, Andrew W. Senior<sup>1</sup>

**Abstract**  
Recent advances in AI, and in particular deep learning, have led to the development of highly accurate methods for predicting the structure of proteins. We have developed a deep learning-based method for protein structure prediction, AlphaFold, which achieves state-of-the-art performance on a standard benchmark. AlphaFold is a deep learning-based method for protein structure prediction, which achieves state-of-the-art performance on a standard benchmark. AlphaFold is a deep learning-based method for protein structure prediction, which achieves state-of-the-art performance on a standard benchmark.

**Introduction**  
The development of computational methods to predict the structure of proteins from their amino acid sequences is a long-standing goal in computational biology. The development of deep learning-based methods for protein structure prediction, such as AlphaFold, has led to significant advances in this field. AlphaFold is a deep learning-based method for protein structure prediction, which achieves state-of-the-art performance on a standard benchmark.

**Methods**  
AlphaFold is a deep learning-based method for protein structure prediction, which achieves state-of-the-art performance on a standard benchmark. AlphaFold is a deep learning-based method for protein structure prediction, which achieves state-of-the-art performance on a standard benchmark.

**Results**  
AlphaFold achieves state-of-the-art performance on a standard benchmark. AlphaFold is a deep learning-based method for protein structure prediction, which achieves state-of-the-art performance on a standard benchmark.

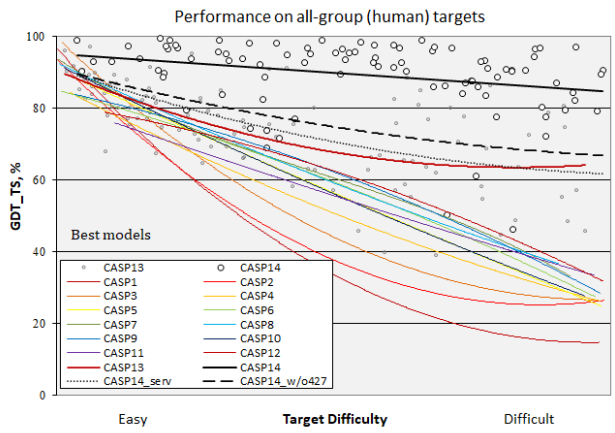
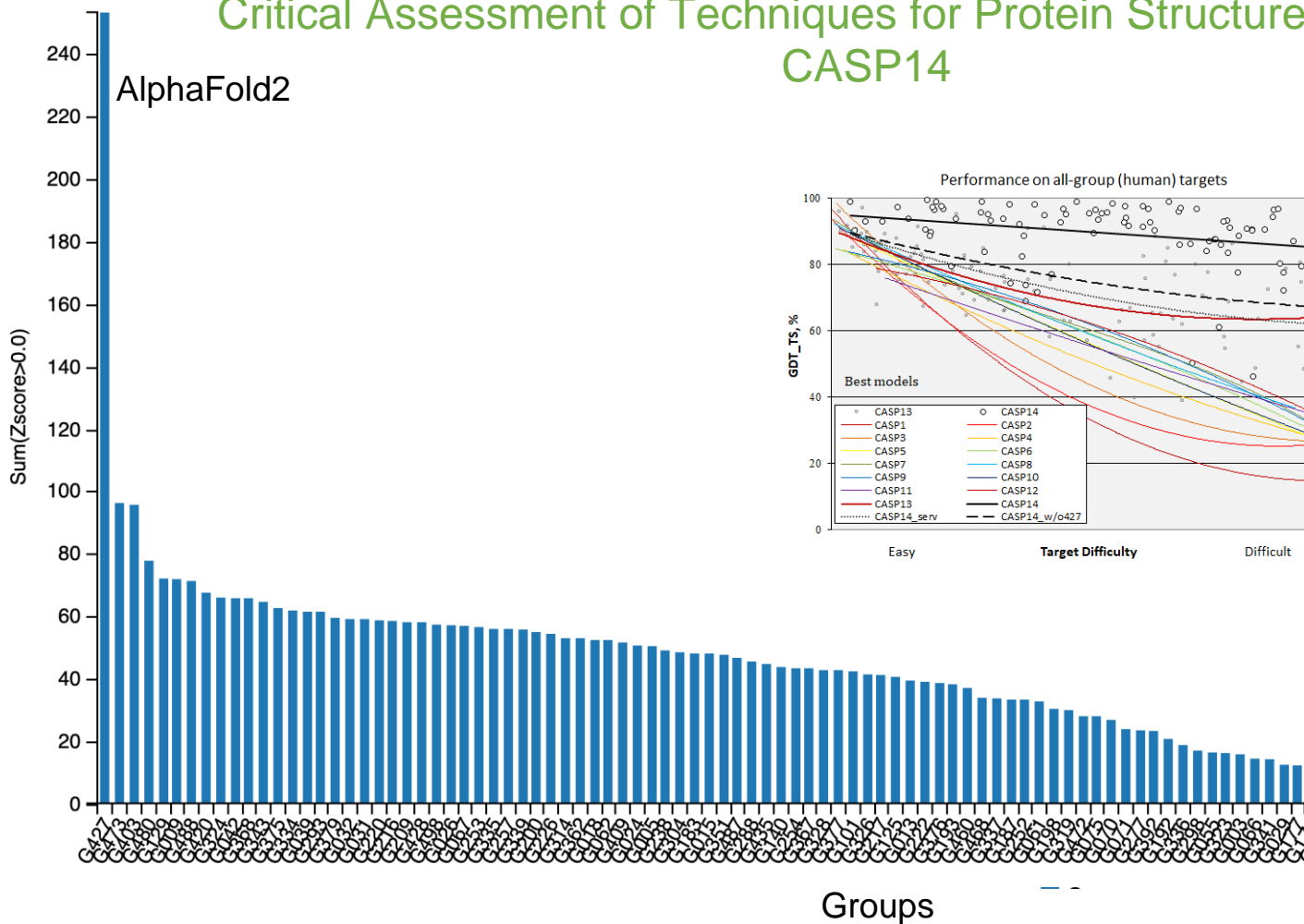
**Discussion**  
AlphaFold is a deep learning-based method for protein structure prediction, which achieves state-of-the-art performance on a standard benchmark. AlphaFold is a deep learning-based method for protein structure prediction, which achieves state-of-the-art performance on a standard benchmark.

**Supplementary Information**  
Supplementary Information for: Highly accurate protein structure prediction with AlphaFold

**Contents**

- 1.1. Supplementary Methods
- 1.2. Overview
- 1.3. Pairing
- 1.4. Genetic search
- 1.5. Template search
- 1.6. Training data
- 1.7. MSA token definition
- 1.8. Residue-embedding
- 1.9. Self-supervised distillation
- 1.10. Self-supervised distillation
- 1.11. Self-supervised distillation
- 1.12. Self-supervised distillation
- 1.13. Self-supervised distillation
- 1.14. Self-supervised distillation
- 1.15. Self-supervised distillation
- 1.16. Self-supervised distillation
- 1.17. Self-supervised distillation
- 1.18. Self-supervised distillation
- 1.19. Self-supervised distillation
- 1.20. Self-supervised distillation
- 1.21. Self-supervised distillation
- 1.22. Self-supervised distillation
- 1.23. Self-supervised distillation
- 1.24. Self-supervised distillation
- 1.25. Self-supervised distillation
- 1.26. Self-supervised distillation
- 1.27. Self-supervised distillation
- 1.28. Self-supervised distillation
- 1.29. Self-supervised distillation
- 1.30. Self-supervised distillation
- 1.31. Self-supervised distillation
- 1.32. Self-supervised distillation
- 1.33. Self-supervised distillation
- 1.34. Self-supervised distillation
- 1.35. Self-supervised distillation
- 1.36. Self-supervised distillation
- 1.37. Self-supervised distillation
- 1.38. Self-supervised distillation
- 1.39. Self-supervised distillation
- 1.40. Self-supervised distillation
- 1.41. Self-supervised distillation
- 1.42. Self-supervised distillation
- 1.43. Self-supervised distillation
- 1.44. Self-supervised distillation
- 1.45. Self-supervised distillation
- 1.46. Self-supervised distillation
- 1.47. Self-supervised distillation
- 1.48. Self-supervised distillation
- 1.49. Self-supervised distillation
- 1.50. Self-supervised distillation
- 1.51. Self-supervised distillation
- 1.52. Self-supervised distillation
- 1.53. Self-supervised distillation
- 1.54. Self-supervised distillation
- 1.55. Self-supervised distillation
- 1.56. Self-supervised distillation
- 1.57. Self-supervised distillation
- 1.58. Self-supervised distillation
- 1.59. Self-supervised distillation
- 1.60. Self-supervised distillation
- 1.61. Self-supervised distillation
- 1.62. Self-supervised distillation
- 1.63. Self-supervised distillation
- 1.64. Self-supervised distillation
- 1.65. Self-supervised distillation
- 1.66. Self-supervised distillation
- 1.67. Self-supervised distillation
- 1.68. Self-supervised distillation
- 1.69. Self-supervised distillation
- 1.70. Self-supervised distillation
- 1.71. Self-supervised distillation
- 1.72. Self-supervised distillation
- 1.73. Self-supervised distillation
- 1.74. Self-supervised distillation
- 1.75. Self-supervised distillation
- 1.76. Self-supervised distillation
- 1.77. Self-supervised distillation
- 1.78. Self-supervised distillation
- 1.79. Self-supervised distillation
- 1.80. Self-supervised distillation
- 1.81. Self-supervised distillation
- 1.82. Self-supervised distillation
- 1.83. Self-supervised distillation
- 1.84. Self-supervised distillation
- 1.85. Self-supervised distillation
- 1.86. Self-supervised distillation
- 1.87. Self-supervised distillation
- 1.88. Self-supervised distillation
- 1.89. Self-supervised distillation
- 1.90. Self-supervised distillation
- 1.91. Self-supervised distillation
- 1.92. Self-supervised distillation
- 1.93. Self-supervised distillation
- 1.94. Self-supervised distillation
- 1.95. Self-supervised distillation
- 1.96. Self-supervised distillation
- 1.97. Self-supervised distillation
- 1.98. Self-supervised distillation
- 1.99. Self-supervised distillation
- 2.00. Self-supervised distillation

# Critical Assessment of Techniques for Protein Structure Prediction CASP14



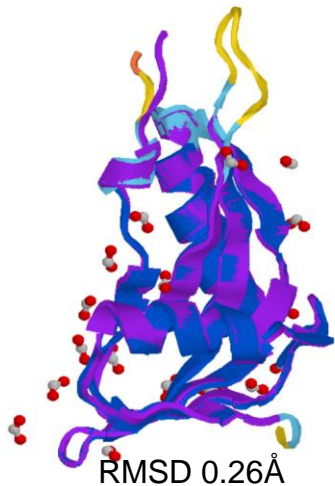
**John Moutl**

*Institute for Bioscience and  
Biotechnology Research  
(IBBR), University of Maryland,  
Washington, US*



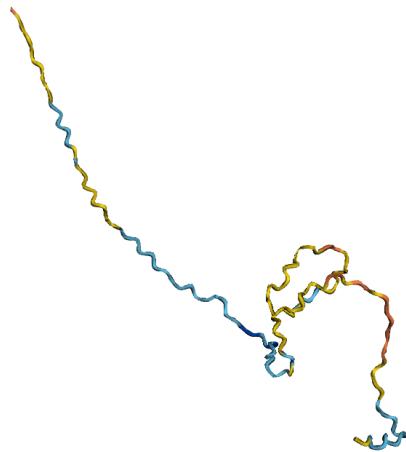
# What do the AlphaFold models look like?

Q9NRX4 –  
14 kDa phosphohistidine phosphatase



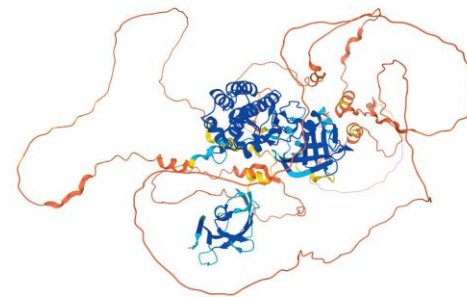
Average confidence score: 90.1  
76.8% of residues score above 90

O60927  
E3 ubiquitin-protein ligase PPP1R11



Average confidence score: 66.4  
3.1% of residues score above 90

Q99558  
MAP3 Kinase K14

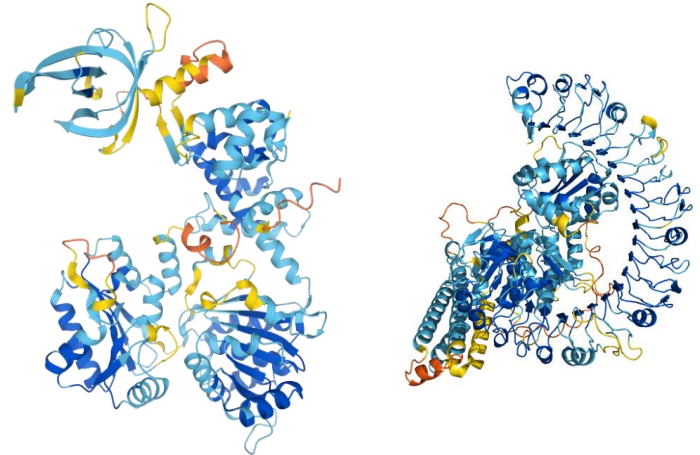
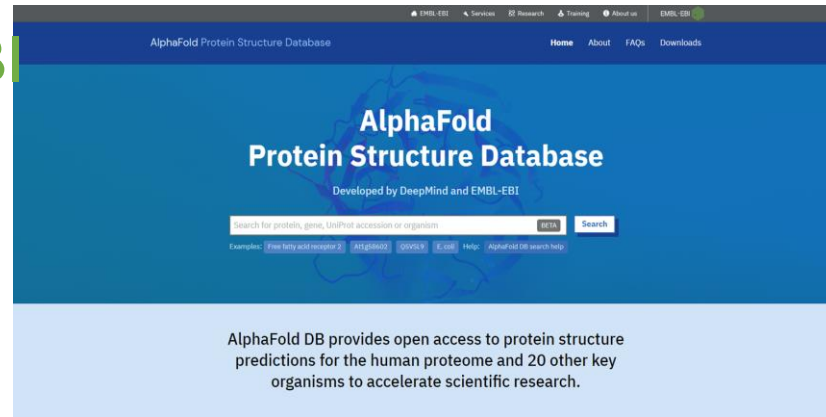


3 domains well predicted,  
but connections unknown

Confidence Scores are essential

# AlphaFold Database @ EMBL-EBI

- Launched on 22nd July 2021
  - Collaboration between EMBL and DeepMind
  - Open access (CC-BY-4.0 license)
  - Data available for bulk download via FTP (<ftp.ebi.ac.uk/pub/database/alphafold>)
- Structure for every known protein in UniProt database are available or can be modelled
- 3-D structures for virtually all (98.5%) of the human proteome
- C.f. Only 17% human proteome in PDB



# Protein Structure Prediction – The perfect problem for AI

- Long standing important problem
  - All life is based on proteins;
  - Sequencing DNA is rapid and cheap; Sequencing proteins is much more difficult and expensive;
  - Structure determines biological function; Structure is critical for drug design
- Lots of data – 170K Protein structure entries in PDB
  - Data are curated, clean, accurate and freely available (sequences and structures)
  - Problem is well defined: - given a protein sequence, predict the 3D structure (coordinates)
  - Good metrics to measure success
- Evolution
  - Most proteins belong to a relatively small number of families
  - Conservation of protein structure during evolution and in different species facilitates prediction
- Many helpful scientific and technical developments in last 30 years
- CASP – Comparative Assessment of Structure Prediction: independent test of quality of prediction.

# AlphaFold Structures – Fundamental limitations

AlphaFold does not predict protein folding pathway

AlphaFold not trained or validated for predicting the effect of mutations

The original AlphaFold2 did not consider multimers

Ligands are not included

-

Low reliability regions have poor stereochemistry

Predictions may (or may not) lead to hypotheses about protein function – any hypotheses have to be tested by further experimentation

Access to good data were critical for the success of AlphaFold!

- Biological Data & Infrastructure
  - EMBL-EBI
  - Community Efforts

# The data challenges in the life sciences



Data growth



Diversity of  
data



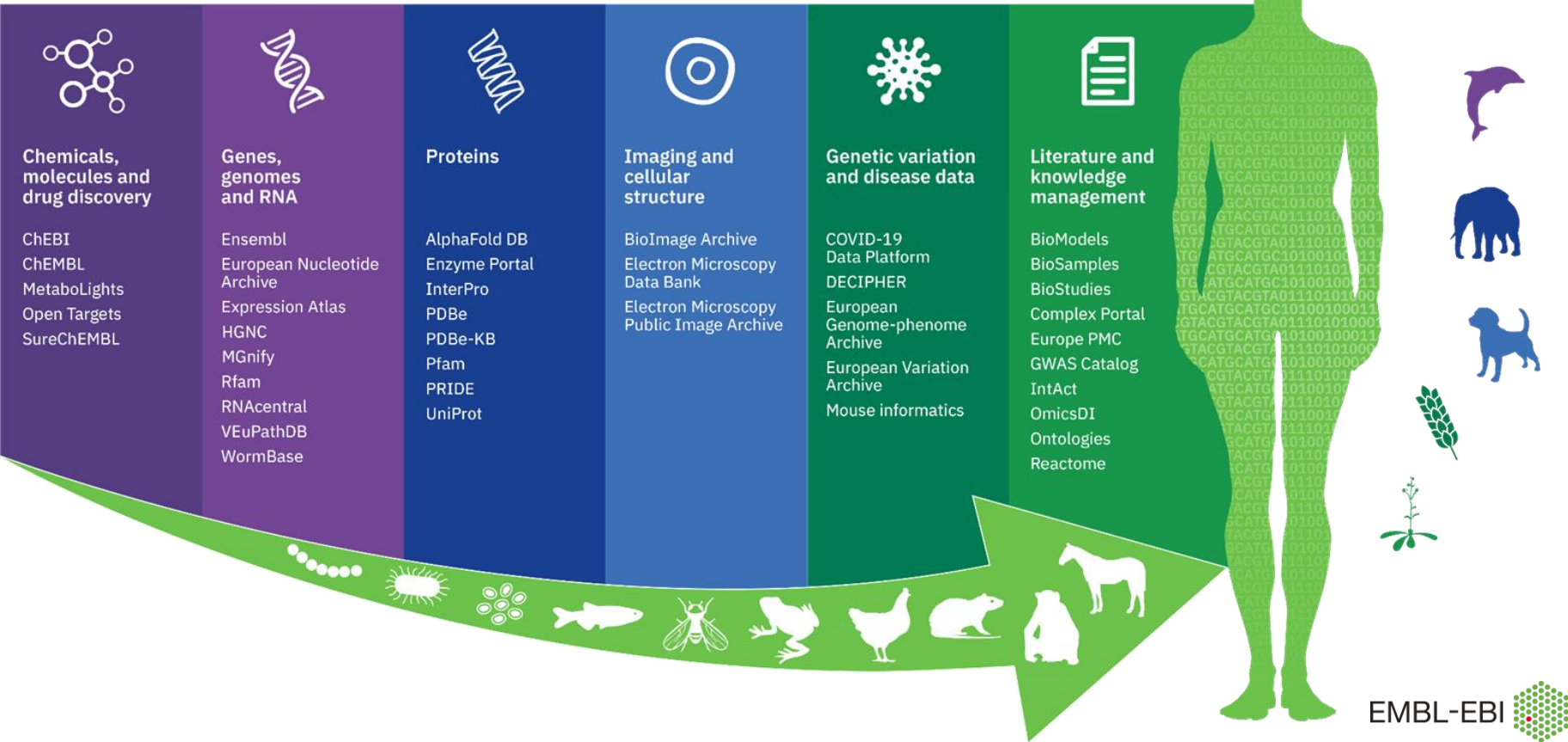
Geographic  
spread



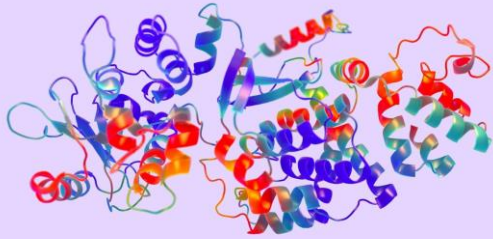
Sensitive  
data



# Data resources at EMBL-EBI



## Impact of EMBL-EBI data resources



**Over 40**  
open data resources

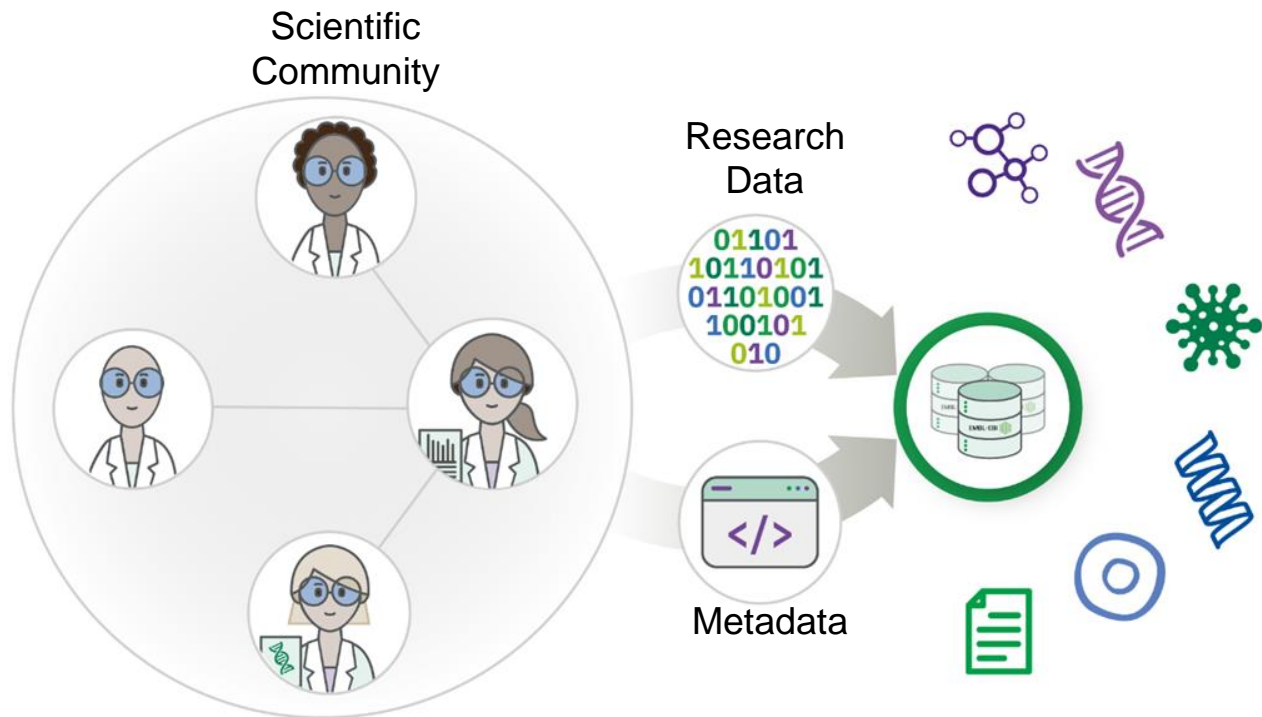


**Over 100 million**  
requests to  
our websites on an  
average day



Our data resources are  
referenced in scientific papers  
200 times every day

# AI and data science at EBI



## EMBL-EBI data resources:

- Curation
- Annotation
- Classification
- Enrichment
- Analysis

Our 40+ data resources span genetics, genomics, proteins, chemistry, literature data and more.

High quality  
OPEN data

Scientific  
Community

Research  
Data

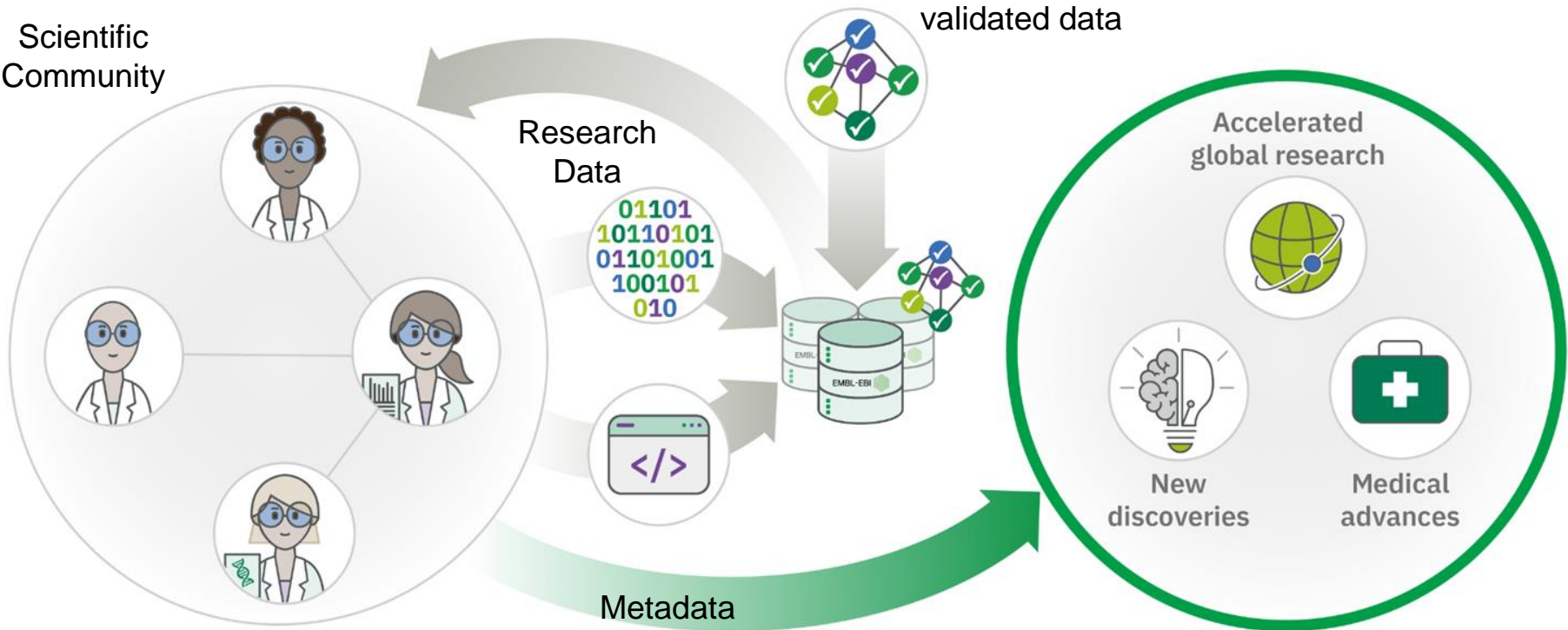
AI Predicted and  
validated data

Metadata

Accelerated  
global research

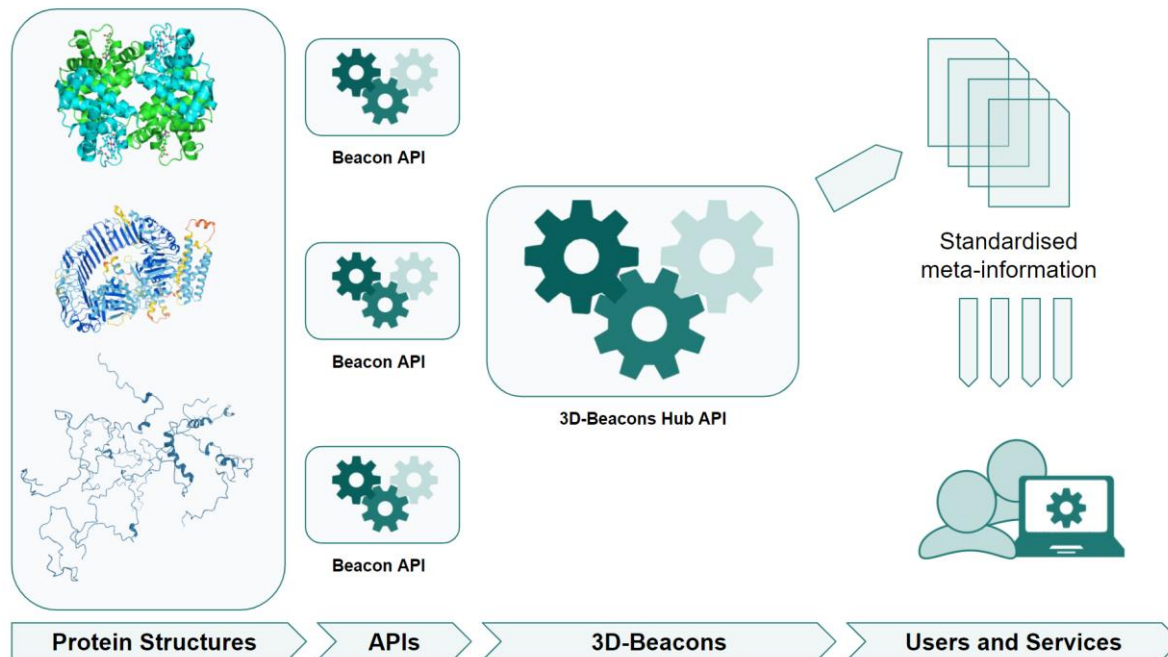
New  
discoveries

Medical  
advances



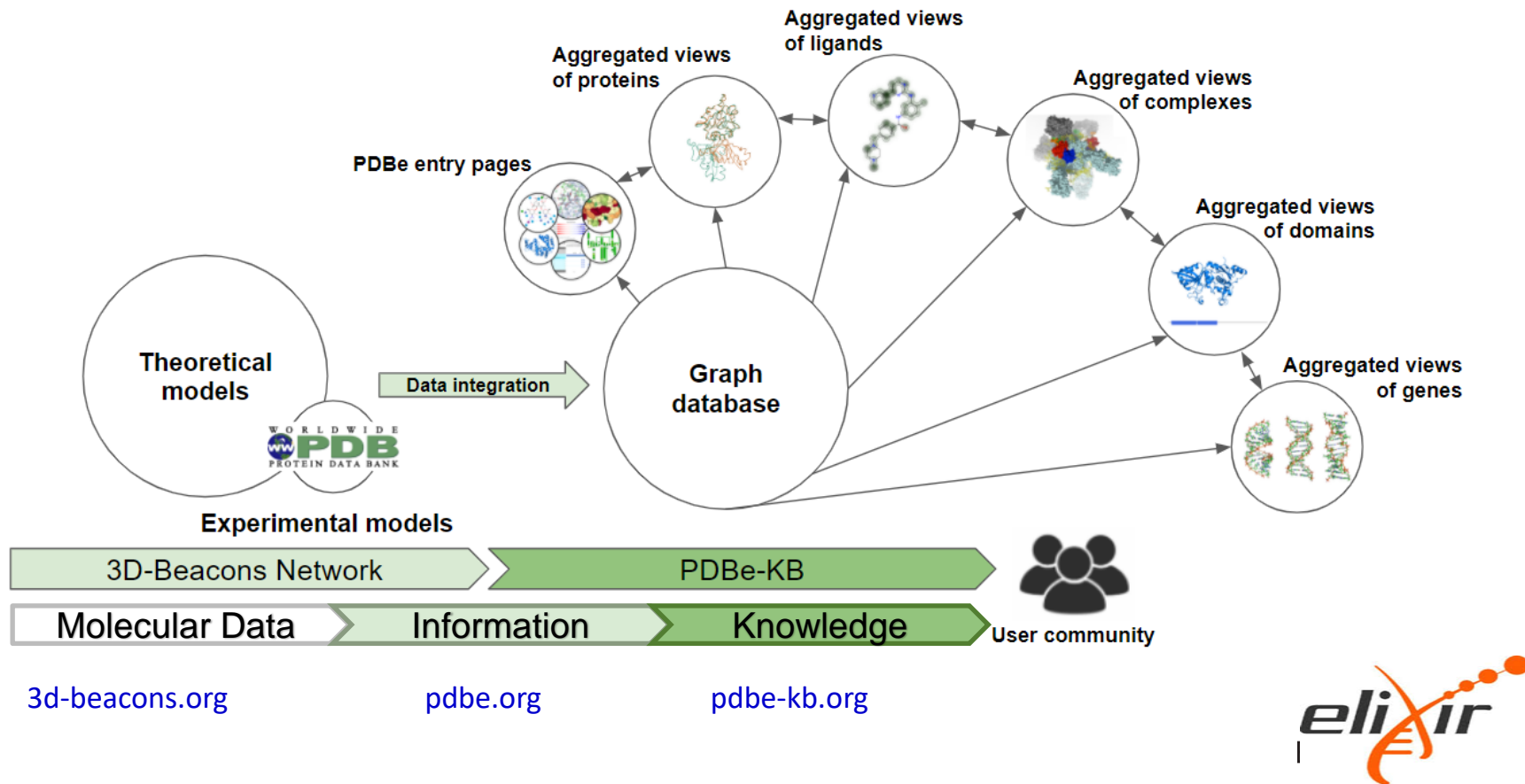
# Community Efforts: 3D-Beacons – an infrastructure to make all protein structure models available, including AlphaFold

- Create a common API specification to access both experimental and predicted structures and annotations
- Provide quality assessment data for all structure models



# Infrastructure for Structural Biology

## The PDB Knowledge Base - from data to knowledge PDBe-KB



# Reaching out to Industry: EMBL-EBI Industry Programme today

EMBL's mission to disseminate cutting-edge technologies to industry.

- Explore emerging areas for R&D
- Network of experts
- Solve big data challenges
- Developing data standards
- Efficiency savings
  
- Quarterly strategy meetings
- Knowledge-exchange workshops

## Pharmaceutical and diagnostic



## Agrofood and consumer goods



AI is dominating computational approaches in the Life Sciences

## EMBL-EBI Industry Programme Workshops: 2024 schedule

1. Chemical spaces and synthesis automation informatics | UK : March 13-14 - [EMBL-EBI](#)
2. Recent developments in large language models for biology | US : April 9-10; [BMS, Boston](#)
3. **Advances in machine learning for protein design** | UK : May 21-22 - [EMBL-EBI](#)
4. **Data-driven research I.T. exchange** | UK : June 12-13 - [EMBL-EBI](#)
5. **Data Literacy in R&D** | online June 2024 & In-person Oct 9-10 [EMBL-EBI](#)
6. Cell-cell communication analysis | US : June 25-26 - [BMS, Boston](#)
7. Preclinical immunogenicity assessment of biologics | UK : Sept 11-12 - [EMBL-EBI](#)
8. **Advances in machine learning for protein design II** | US Edition: Sept 18-19 - [Sanofi, Boston](#)
9. Genetic Biomarkers in pharmaceutical development | US : October 29-28 - [EMD Serono, Boston](#)
10. Genetic and Genomic risk stratification in clinical trial design | US : Nov 13-14 - [Genentech, CA](#)
11. **Multi-omics concepts and applications in biotechnology research** | UK : Nov 20-21 - [EMBL-EBI](#)

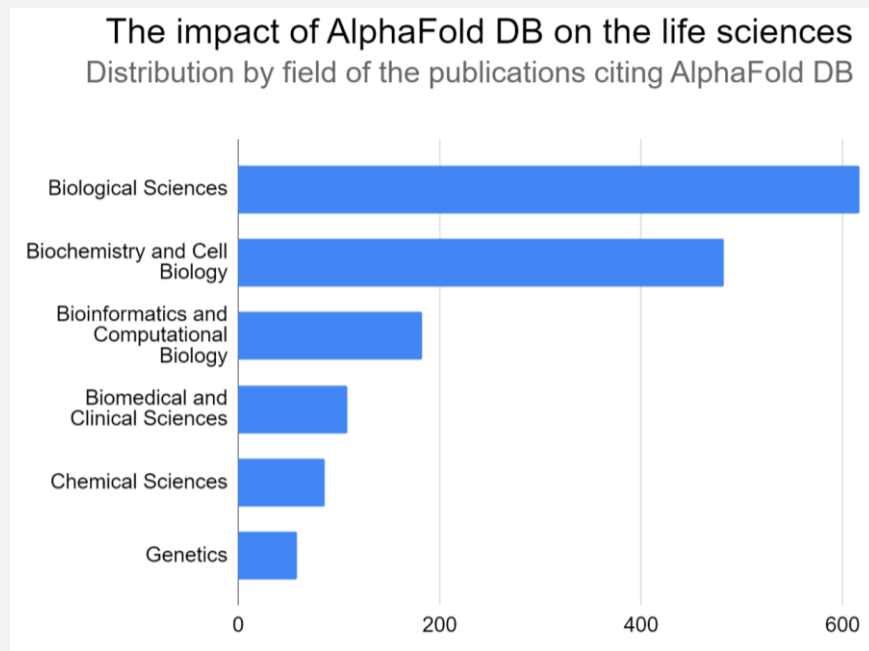
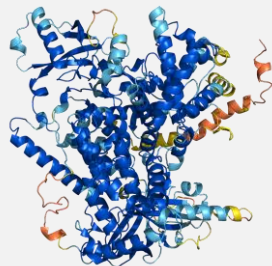


# AI & Impact on Life Sciences Research

- Impact in Structural Biology
- Impact in Medical Sciences

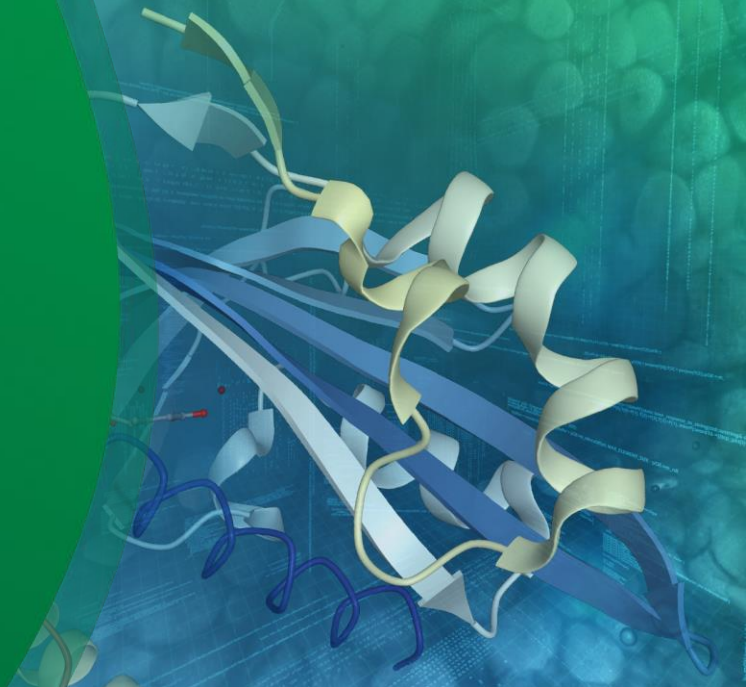
# AlphaFold Database Usage

- >1.8 M unique users from every Continent
- >1B million API requests
- Users have downloaded 1.2M files
- >23 K bulk archives
- AF2 papers cited > 18,000 times



# Enabling discoveries using macromolecular structure data

Impact of structure prediction using AI methods



 **PDBe**  
Protein Data Bank in Europe

 **PDBe-KB**  
Protein Data Bank in Europe - Knowledge Base

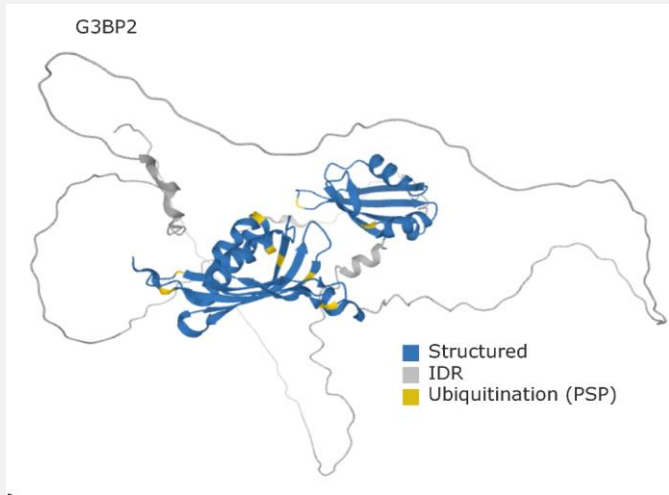
 **AlphaFold**  
Protein Structure Database

[pdbe.org](http://pdbe.org) | [pdbe-kb.org](http://pdbe-kb.org) | [alphafold.ebi.ac.uk](http://alphafold.ebi.ac.uk) | [3d-beacons.org](http://3d-beacons.org)

EMBL-EBI 

**Sameer Velankar**

# Providing structural context



- Integrating large scale **PTM information** and predicted models
  - Phosphorylation, ubiquitination and acetylation
- **Chemical cross-linking studies** to identify and validate macromolecular complexes
- Combining **experimental data (cross-linking, co-fractionation)** and **structure prediction methods** to get insights into functional complexes and interactions

## The structural context of posttranslational modifications at a proteome-wide scale

Isabell Bludau, Sander Willems, Wen-Feng Zeng, Maximilian T. Strauss, Fynn M. Hansen, Maria C. Tanzer, Ozge Karayel, Brenda A. Schulman, Matthias Mann 

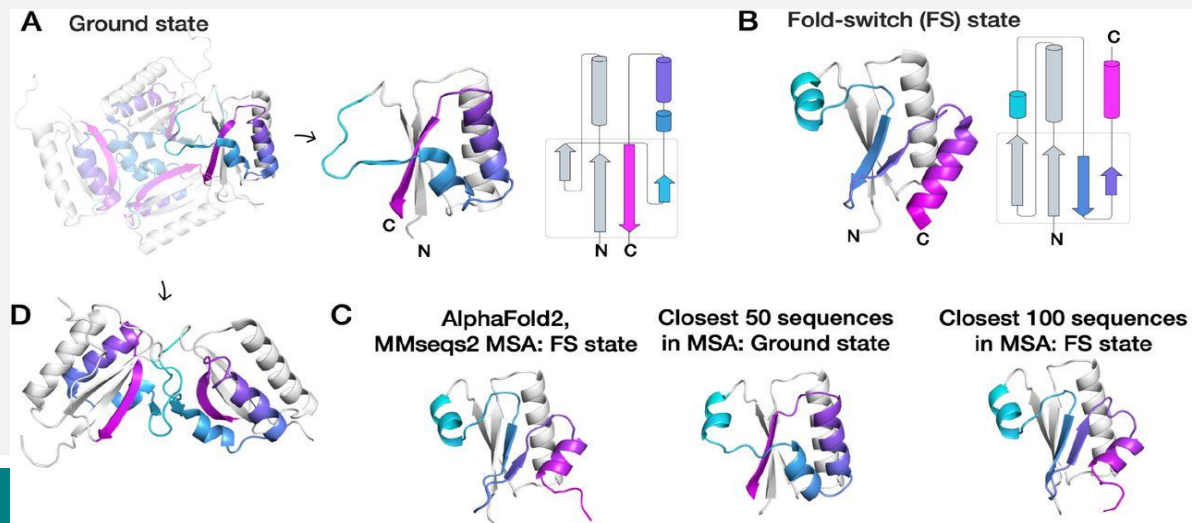
Published: May 16, 2022 • <https://doi.org/10.1371/journal.pbio.3001636>

# Predicting multiple conformational states

## Prediction of multiple conformational states by combining sequence clustering with AlphaFold2

 Hannah K. Wayment-Steele,  Sergey Ovchinnikov,  Lucy Colwell,  Dorothee Kern

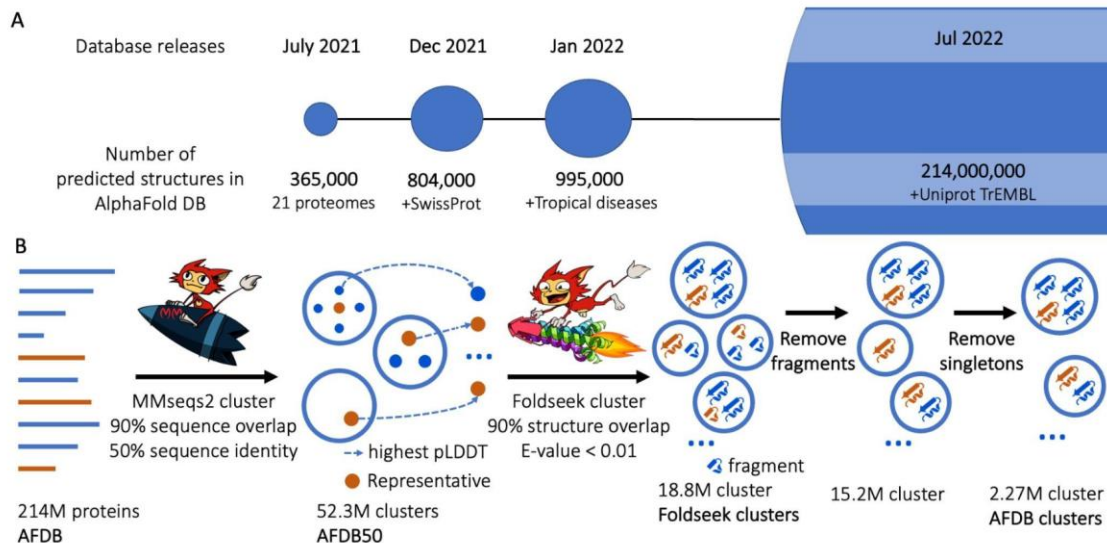
**doi:** <https://doi.org/10.1101/2022.10.17.512570>



Using only the closest 50 sequences by sequence distance returned from the MSA returns the ground state, but the closest 100 returns the fold-switch state.

# Comparing and clustering structures at scale

- 2.27M non-singleton structural clusters.
  - >1.1M (50% of AFDB clusters) found to be, at least partially, similar to previously known structures in the PDB
  - 4% seem species specific
- 31% lack annotations representing likely novel structures
  - Tend to have few representatives covering only 4% of all proteins



Clustering predicted structures at the scale of the known protein universe

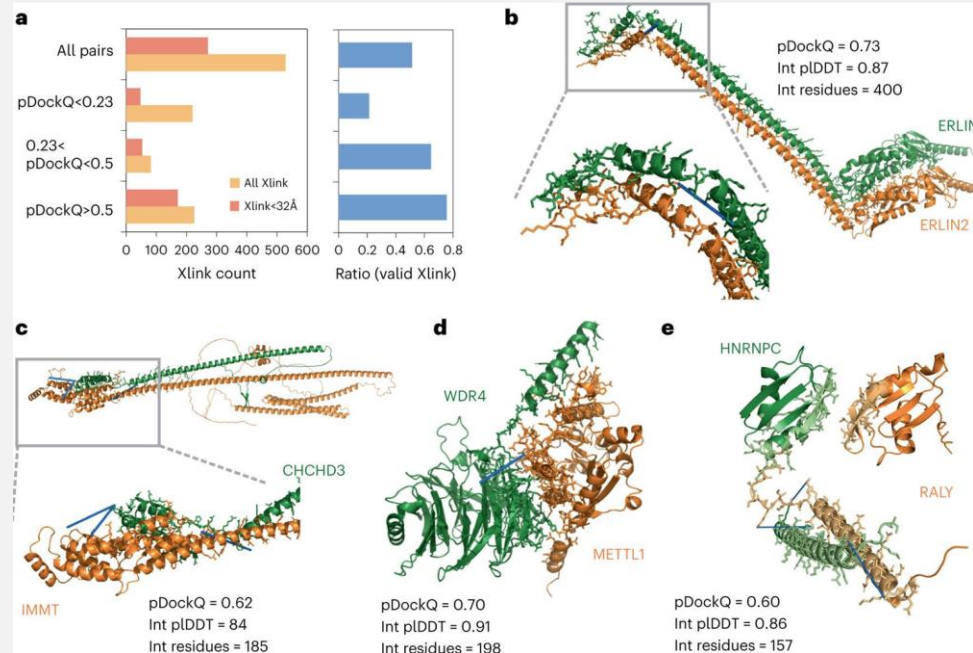
Inigo Barrio Hernandez<sup>1\*</sup>, Jingli Ye<sup>2\*</sup>, Jürgen Jänes<sup>3</sup>, Milot Mirdita<sup>2</sup>, Cameron L.M. Gilchrist<sup>4</sup>, Tanita Wein<sup>4</sup>, Mihaly Varadi<sup>1</sup>, Sameer Velankar<sup>1</sup>, Pedro Beltrao<sup>2,3,5,6</sup>, Martin Steinegger<sup>2,6,7,8,9</sup>

What is hidden in the darkness?  
Deep-learning assisted large-scale protein family curation  
uncovers novel protein families and folds

Janani Durairaj<sup>1,2</sup>, Andrew Waterhouse<sup>1,2</sup>, Toomas Mets<sup>3</sup>, Tetiana Brodiazhenko<sup>5</sup>, Minhhal Abdullah<sup>3</sup>, Mehmet Akdel<sup>4</sup>, Antonina Andreeva<sup>6</sup>, Alex Bateman<sup>5</sup>, Vasili Hauryluk<sup>3,6,7</sup>, Tanel Tenson<sup>3</sup>, Torsten Schwede<sup>1,2</sup>, Joana Pereira<sup>1,2</sup>

# Large scale predictions of protein interactions

- Predicted the structures of **65,000** pairs of *human* proteins, obtained **3,000** high-confidence pairs
- **1,400** of the newly predicted high-confidence complexes lacked homology to existing structures
- Cross-linking validation:
  - Experimental evidence confirmed 171 out of 246 (**70%**) complexes predicted with high confidence



Burke, D. F. et al., Towards a structurally resolved human protein interaction network. Nat. Struct. Mol. Biol. (2023).

# More and more protein sequences available: Metagenomics is tapping into the 99% of the unknown microbes



Case study - harnessing microbes for bioremediation



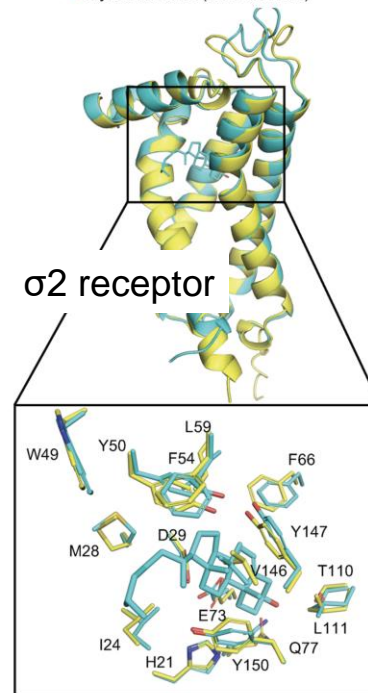
# Specific Applications in progress now

- Drug discovery
- Tackle neglected diseases
- Understand antibiotic resistance
- Fight plastic pollution
- Increase health of honey bees
- Understand how ice forms

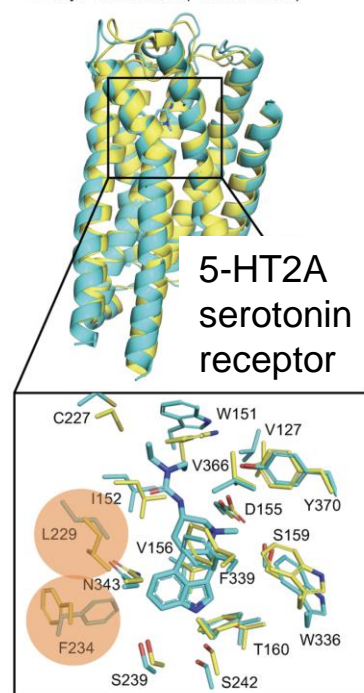
# Are AF2 models good for drug discovery?

- Retrospective simulation of structure-based ligand discovery may underestimate the ability of AF2 structures to template new ligand discovery prospectively
- The prospective large library docking campaigns against the AF2 models were no less effective than those against experimental structures
- The hit rates were high and were not significantly different between the modelled and experimental structures.
- Against the 5-HT<sub>2A</sub> serotonin receptor, the AF2 structure led, if anything, to more potent and selective compounds

**$\sigma_2$  receptor**  
■ AlphaFold2 structure (RMSD: 0.5Å)  
■ Crystal structure (PDB ID: 7MFI)



**5HT<sub>2A</sub> receptor**  
■ AlphaFold2 structure (RMSD: 1.6Å)  
■ Cryo-EM structure (PDB ID: 8UWL)



**AlphaFold2 structures template ligand discovery**

# New normal in structural biology

## Improved AlphaFold modeling with implicit experimental information

Thomas C. Terwilliger<sup>1,2\*</sup>, Billy K. Poon<sup>3</sup>, Pavel V. Afonine<sup>3</sup>, Christopher J. Schlicksup<sup>3</sup>, Tristan I. Croll<sup>5</sup>, Claudia Millán<sup>5</sup>, Jane. S. Richardson<sup>6</sup>, Randy J. Read<sup>5</sup> and Paul D. Adams<sup>3,5</sup>

## Accelerating crystal structure determination with iterative AlphaFold prediction

Thomas C. Terwilliger, Pavel V Afonine, Dorothee Liebschner, Tristan I Croll, A J McCoy, Robert D Oeffner, Christopher J Williams, Billy K Poon, Jane S Richardson, Randy J Read, Paul D. Adams

doi: <https://doi.org/10.1101/2022.11.18.517112>

- Predicted model rebuilt in one cycle is used as a template for prediction in the next cycle.
- 87% of the recent 215 structures yielded a model with at least 50% of C $\alpha$  atoms matching those in the deposited models within 2Å.

## AlphaFold predictions: great hypotheses but no match for experiment

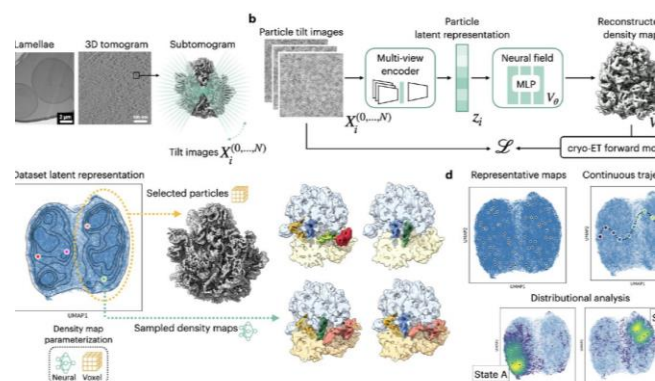
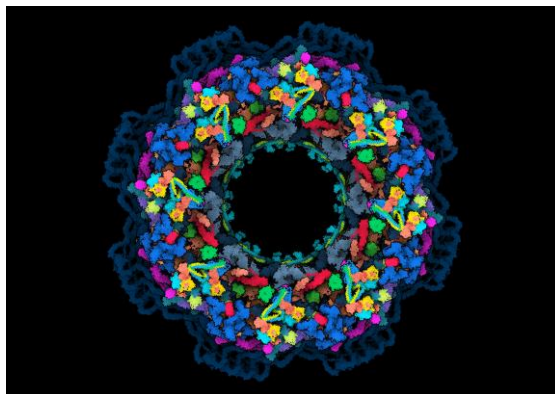
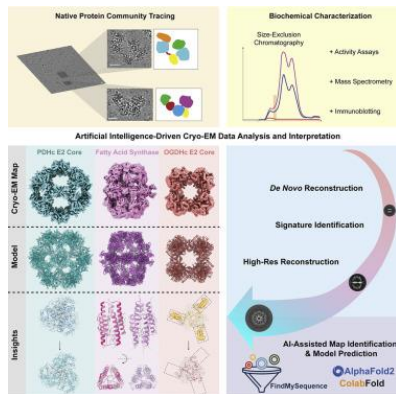
Thomas C. Terwilliger, Dorothy L Liebschner, Tristan Croll, Christopher J Williams, A J McCoy, Billy K Poon, Pavel Afonine, Robert D Oeffner, Jane Shelby Richardson, Randy J Read, Paul D. Adams

doi: <https://doi.org/10.1101/2022.11.21.517405>

AlphaFold confidence (pLDDT)	Median prediction error (Å)	Percentage with error over 2 Å
>90	0.6	10
80 - 90	1.1	22
70 - 80	1.5	33
<70	3.5	77

- Goal of structural biology is to shed light on e.g. biological function or mechanism of biological macromolecules/complexes, not “stamp collecting”
- Many structural studies can now start with a model as the null-hypothesis
- Design experiments based on models
  - Construct engineering to express individual domains without the “spaghetti”
  - Design and determine structures of mutants
  - Determine structures of complexes
  - Fragment/compound screens

# A(nother) Golden Age for Structural Biology: From Structural Inventories to Processes



High-resolution cryo-EM of cell extracts volume-based *de novo* identification and AI-assisted atomic modeling

Skalidis, Ioannis, et al. "Cryo-EM and artificial intelligence visualize endogenous protein community members." *Structure* 30.4 (2022): 575-589.

Total mass: 120 Mda  
Old model: 35 Mda  
**New model: 70 MDA**

**Artificial intelligence reveals nuclear pore complexity**

Shyamal Mosalaganti, Agnieszka Obarska-Kosinska, Marc Siggel, Beata Turonova, Christian E. Zimmerli, Katarzyna Buczak, Florian H. Schmidt, Erica Margiotta, Marie-Therese Macknull, Wim Hagen, Gerhard Hummer, Martin Beck, Jan Kosinski

doi: <https://doi.org/10.1101/2021.10.26.465776>

Visualisation of structural landscape and dynamics

**Deep reconstructing generative networks for visualizing dynamic biomolecules inside cells**

Ramya Rangan, Sagar Khavnekar, Adam Lerer, Jake Johnston, Ron Kelley, Martin Obr, Abhay Kotecha, Ellen D. Zhong

doi: <https://doi.org/10.1101/2023.08.18.553799>

Electron microscopy/AI-assisted modelling

Integrative/hybrid methods

Structure dynamics in their biological context (visual proteomics)

# Visual Proteomics : A Proof of Concept with *Chlamydomonas reinhardtii*

**PLANT:** The Chloroplast

**Thylakoids:** Light Harvesting

**Pyrenoid:** Carbon Fixation

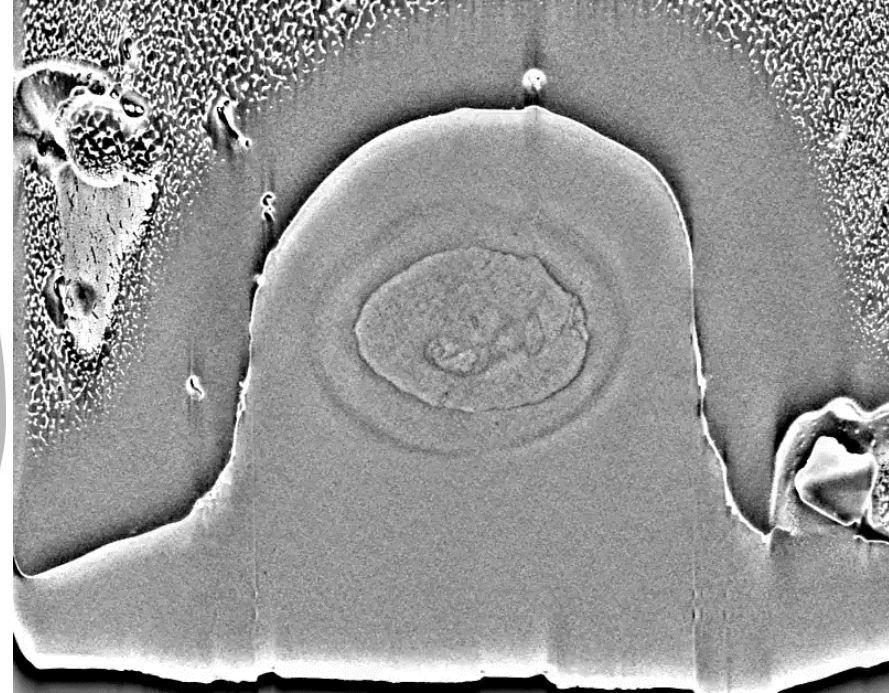
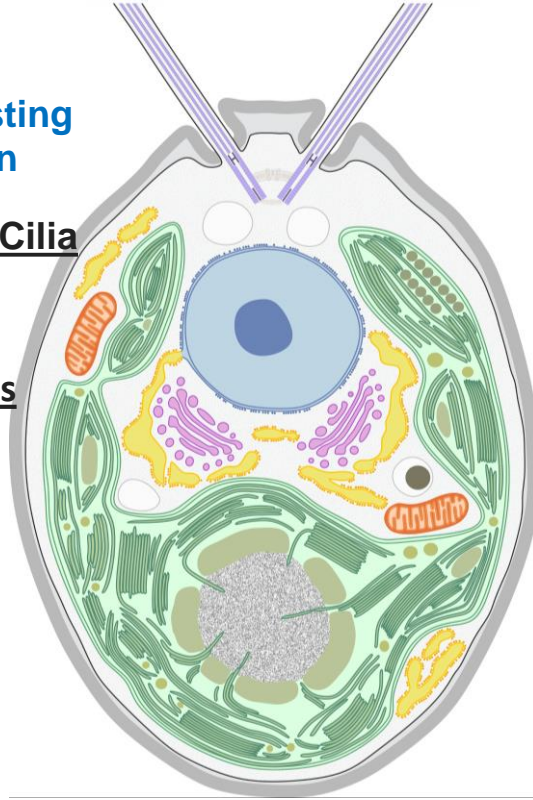
**ANIMAL:** Centrioles and Cilia

"Core" Eukaryotic Organelles

**Nucleus**

**Endoplasmic Reticulum &  
Golgi Apparatus**

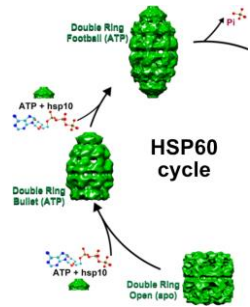
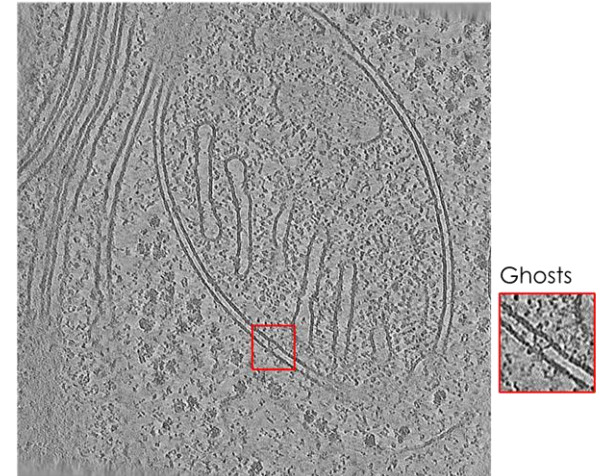
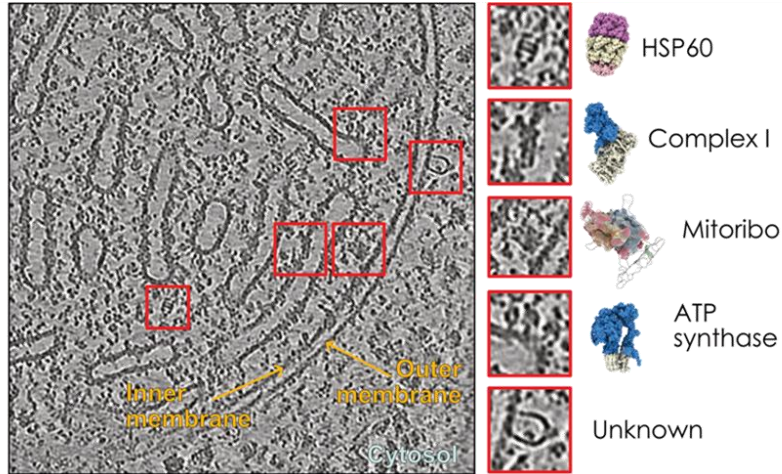
**Mitochondria**



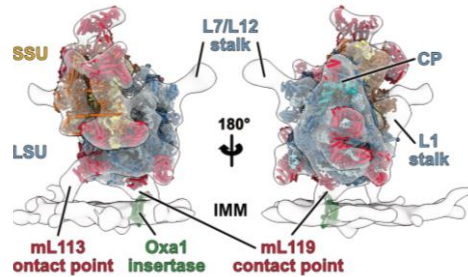
cryo slice and view of a single *Chlamy* cell

# In situ structural cell biology of mitochondria

Need better identification of known and rare targets, *need better classification*



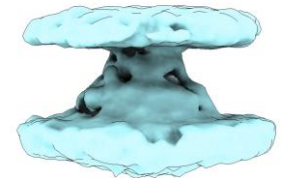
HSP60



Mitoribosomes



ATP Synthase



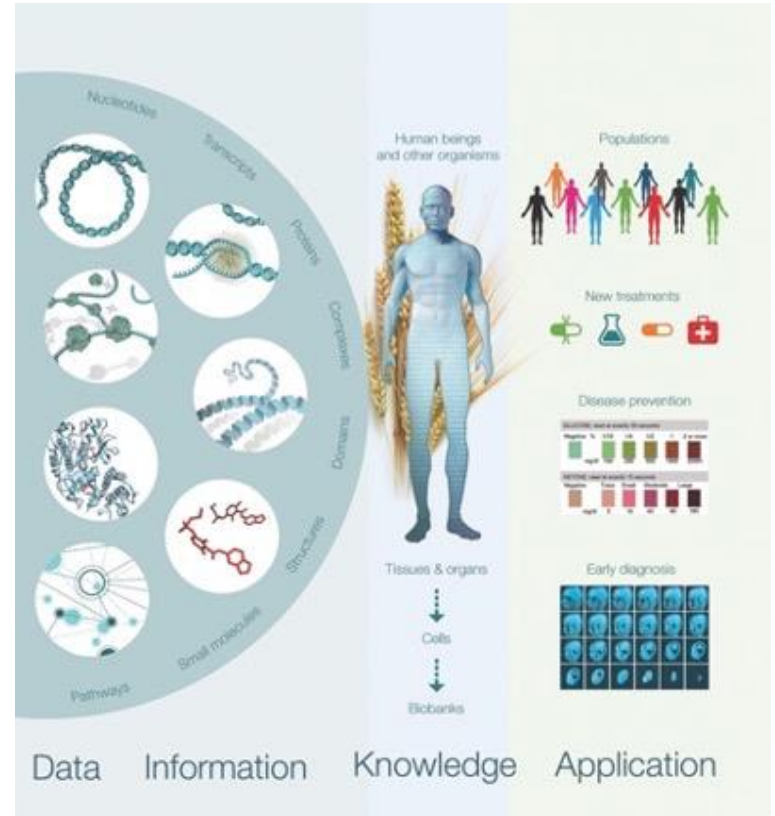
Something ?

# AI/Data in the Medical Sciences

In many ways, medicine is an ideal arena for the application of AI:

- Lots of data
- Many expert clinicians
- Highly regulated profession
- Current approach is mainly based on previous experience
- Imaging is critically important

BUT



# AI & Medicine: Diagnosis & Prognosis

- Long standing important problem – ✓
- Lots of data – ✓
  - Data are curated, clean, accurate and freely available ✗
  - Problem is well defined: ✗
  - Good metrics to measure success - ?? Maybe in extremis
- Evolution
  - Many diseases are genetic – BUT impact of environment – more complicated
  - Different species often behave differently eg mouse & human
- Many helpful technical developments in last 30 years - SOME, but underlying causes of diseases often unknown
- CASP ✗

Major challenges:- DATA quality: lack of standards: country specific: data confidentiality: fragmentation of medical 'disciplines'; Ethics: Implementation etc etc;

Scale is enormous!

**BEGIN WITH DATA**



# AI & Data for Medicine



Mission is to accelerate trustworthy  
data use to enable discoveries that  
improve people's lives

25/04/2024

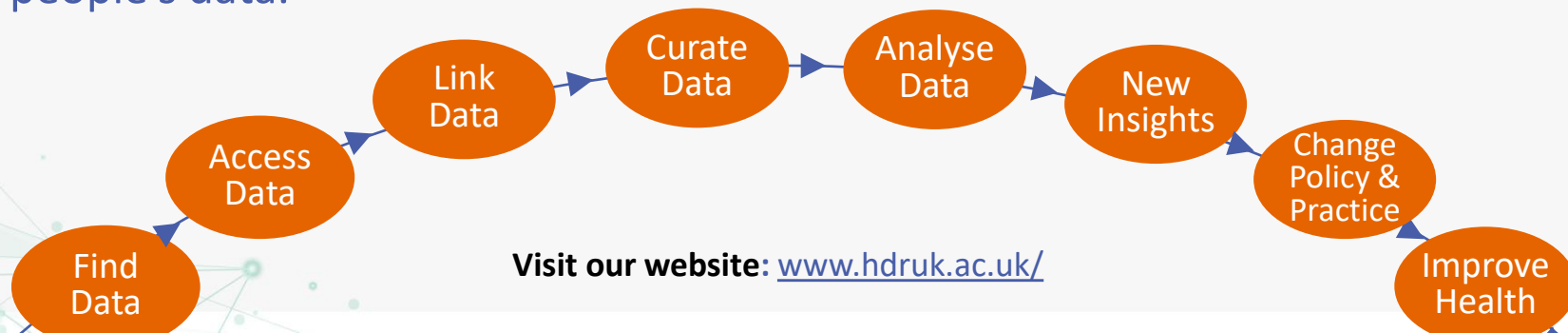
**Professor Dame Janet Thornton**

Member of HDR UK's Board of Trustees



# Why is HDR UK needed?

- In principle the NHS has cradle-to-grave records for 65 million people
- Access to this data for researchers is often a lengthy, fragmented process. It can take years to find, access and 'wrangle' data from different places into a state that researchers can use for analysis
- HDR UK is working with multiple partners to make it easier for researchers to find, access and work with the data they need to generate new discoveries, using streamlined systems that are designed to ensure the privacy and security of people's data.



Visit our website: [www.hdr.uk.ac.uk/](http://www.hdr.uk.ac.uk/)

# RETFound: Using retinal imaging to predict and detect disease

- Detects eye diseases and predicts risk of health conditions
  - e.g., Parkinson's disease, heart attack and strokes
- One of the first AI foundation models in healthcare and the first in ophthalmology
- Pre-trained on 1.6m retinal images
- Adaptable and validated in multiple disease detection tasks
- Freely available for use on [GitHub](#)



“If the UK can combine high quality clinical data from the NHS, with top computer science expertise from its universities, it has the true potential to be a world leader in AI-enabled healthcare. We believe that our work provides a template for how this can be done.”

- Professor Pearse Keane

Learn more: [Eye-scanning AI predicts and detects disease in world-first, study suggests - HDR UK](#)

# Foresight: Predicting future health of patients

- Predicts medical events, aids decision making and informs clinical research
- Trained using NHS ECRs from over 811,000 people in the UK
- Precision rates of 68% to 88%
- Clinically validated

Predictions were 93% relevant



“One of the main advantages of Foresight is that it can easily scale to more patients, hospitals or disorders with minimal or no modifications, and the more data it receives, the better it gets.”

- Zeljko Kraljevic

**Learn more:** [New AI tool may offer insights into patients' future health - HDR UK](#)

# The importance of high-quality, representative data for AI models

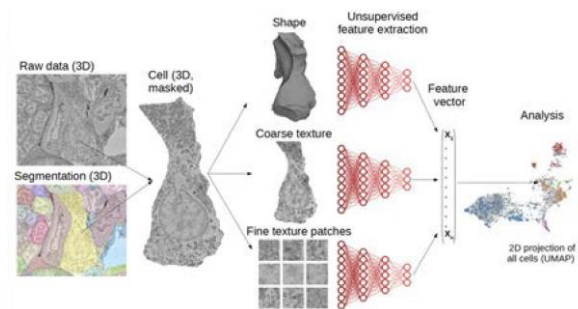
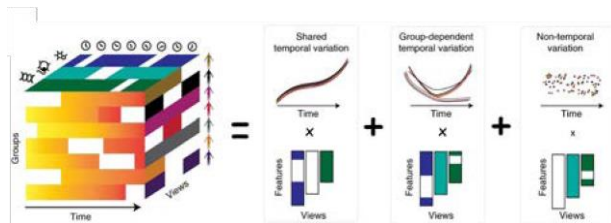
- Analysed ethnicity data from GP and hospital records of over 61m people
- Highlighted that ethnicity information was missing for almost 1 in 10 patients
- First phase of a three-phase project aiming to reduce bias in AI health prediction models



“Because AI-based healthcare technology depends on the data that is fed into it, a lack of representative data can lead to biased models that ultimately produce incorrect health assessments. Better data from real-world settings, such as the data we have collected, can lead to better technology and ultimately better health for all.”

- Professor Sara Khalid

# Today: Opening up the world's Biodata for AI



## Macromolecular Structures

Protein design & data generation informed by AI

## Multi-modal integration

of multiple modalities across space and time (omics, imaging, text...)

## Imaging

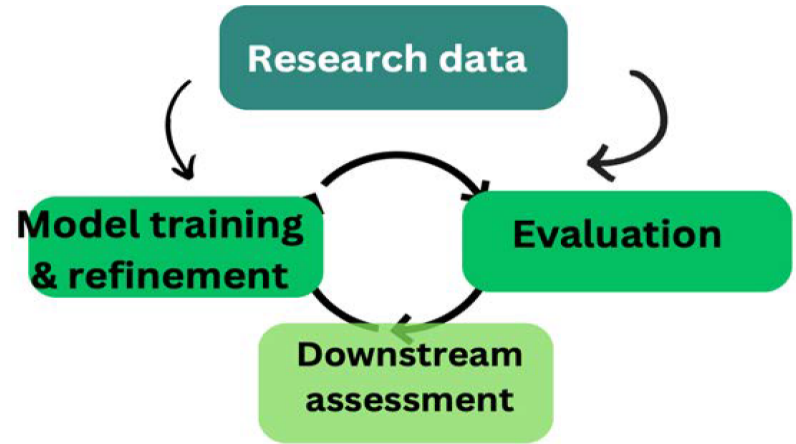
AI-driven cellular imaging models across scales

# The Future?



# Ingredients of (today's) AI advances in biology

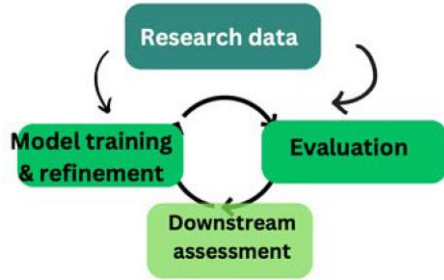
- Well-defined question & metrics
- Computable, high-quality data
- Creative AI developers
- Sufficient compute capacity
- Continuous model benchmarking
- Community assessment of model outputs in practice



=> Future needs the same for Medical Research



# The Evolution of AI Paradigms



## Task-specific AI

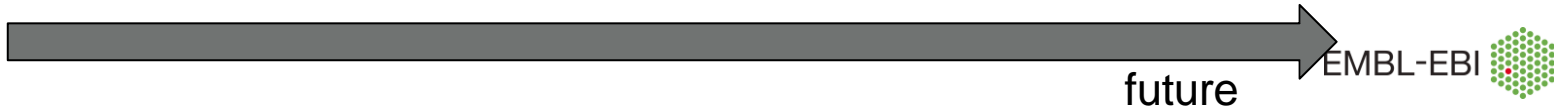
Often building on ground truth labels. Incompasses AlphaFold and many EMBL AI tools today.

## Foundation models

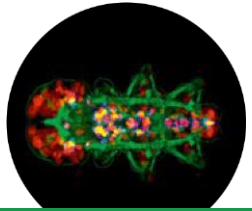
Task agnostic AI. Trained on large volumes of data, typically label free.

## Lab-in-the-loop AI

Tie data generation and technology to AI.



# AI spans all aspects of any organisation (eg. EMBL)



## Research

- Application of AI to scientific discovery
- Lab-in-the-loop AI
- Rigorous, explainable, causal AI modelling & theory



## Services

- Enhance EMBL data resources through AI
- Deliver AI-ready data resources



## European coordination

- Coordination of European AI networks
- Engage in European AI policy



## Technology transfer

- Initiate new pre-competitive partnerships in BioAI
- Set examples for impactful collaborations with industry



## Training & ethics

- Enhance training using AI
- Interdisciplinary training: AI & life sciences
- Guidelines and best practices



## Operations

- Information management
- Building maintenance
- Space usage
- Controlling & legal

# Biology Research and Medical Research

- It has taken the last 30 years to develop the data infrastructure for **biological research**, characterized by:
  - Data Sharing – culture and deposition tools
  - FAIR data
  - Agreed standards
  - A public data infrastructure – many different types of data
  - Developments of many new methods

Huge progress in understand- though still a long way to going of biological processes

- **Medical Research** is just starting on this journey
  - Data sharing remains challenging – in part because of patient confidentiality but also culture
  - Scale is at least ten times bigger
  - Data are even more siloed than biological data – by country/even by hospital
  - Data infrastructure just now beginning to develop – TREs, data wrangling etc
  - Potential commercial gains
  - Systems biology

BUT – this journey is underway and I have no doubt that in the next 10-20 years, progress will be enormous and affect every one of us both scientifically and personally

# Acknowledgements

**Sameer Velankar &  
Colleagues at EMBL**

**Scientists at DeepMind**

**PDBe & AlphaFold  
Database@EMBL-EBI**

**Health Data Research UK**

**Oliver Stegle (EMBL-HD)**

**My group**



Sameer Velankar

