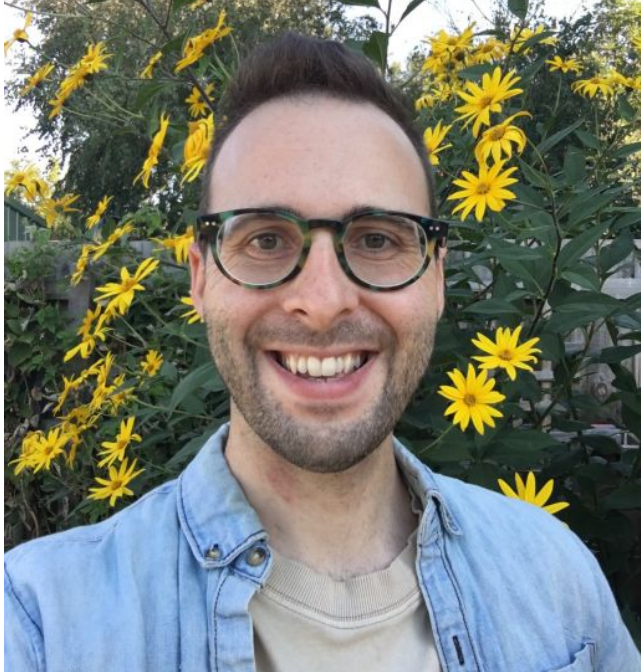# COMP90014

Algorithms for Bioinformatics

Week 6A - Genomic Features & Regions

# Assignment 2 Released Tomorrow

(Wednesday 30th @ midnight)

# Upcoming Guest Speakers!



*Adam Taranto*

This Thursday! (31st August)



*Ryan Wick*

Next Tuesday!  (5th September)

# Genomic Features & Regions

Introduction

Leveraging Data

Prokka: Annotation Pipeline

Training Based Approaches

- Markov Models
- Deep Learning: AlphaFold 2.0

# Introduction

# Introduction

Genome assembly just the beginning!

Want to understand genomic features / regions

# Introduction

Genome assembly just the beginning!
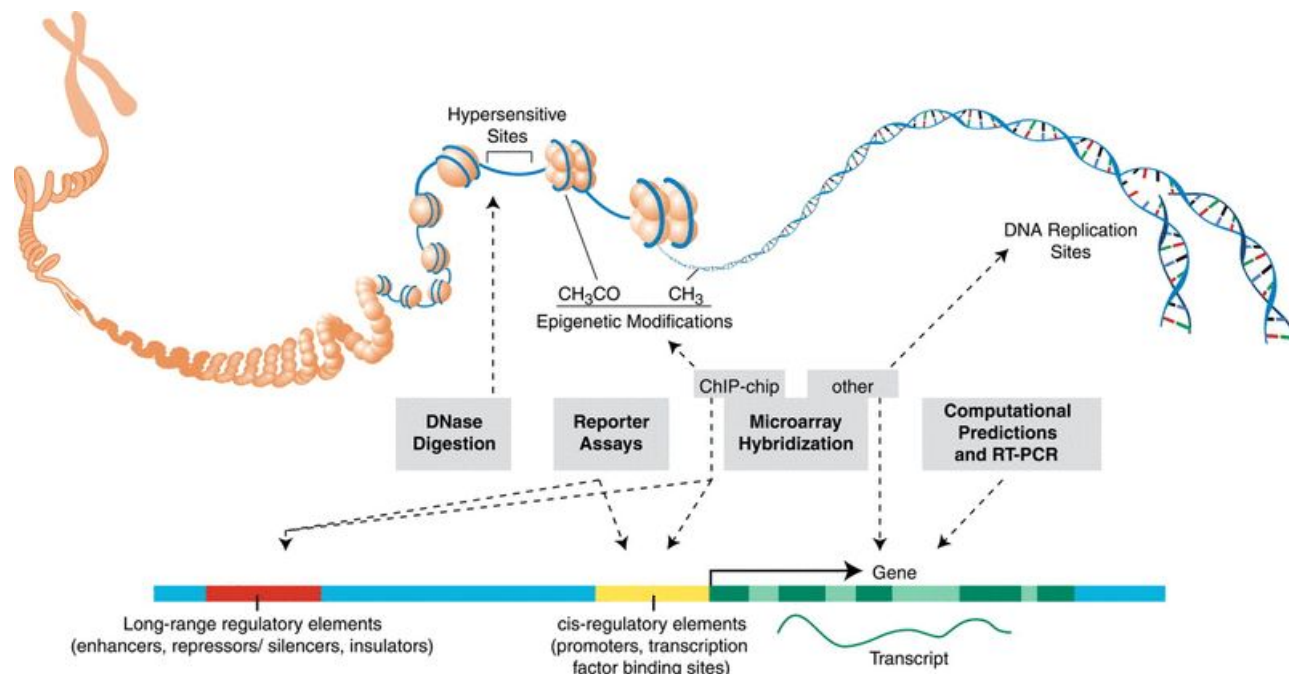
Want to understand genomic features / regions

Genome annotation

"Identifying and labeling the relevant features on a genome sequence"

Coding genes and predicted AA products
(most important)

Promoter / Enhancers, Non-coding RNAs
Signal Peptides, etc
(also important)



The ENCODE (ENCyclopedia Of DNA Elements) Project. DOI:10.1126/science.1105136

# Genomic Features & Regions

Too broad to cover implementations

- Repeat masking
- Identifying promotors / enhancers
- Identifying signal peptides (SPs)
- Identifying genes
- Gene annotation
- Assembly QC: BUSCO
- Identifying structural variation
- Identifying sequence variation
- Characterisation

Will instead cover

- Some main ideas
- Leveraging Data
- Prokka: Annotation Pipeline (main ideas)
- Training Based Approaches
- Markov Models
- Deep Learning: AlphaFold 2.0

# Leveraging Data

# Leveraging Data

Over last 50 years, accumulated lots of data!

This data can be helpful for tasks we wish to perform.

- Some tasks purely isolated
- Some tasks partially leverage existing data
- Some tasks fully use existing data

# Leveraging Data

Over last 50 years, accumulated lots of data!

This data can be helpful for tasks we wish to perform.

- Some tasks purely isolated
- Some tasks partially leverage existing data
- Some tasks fully use existing data

Leveraging data is a spectrum.

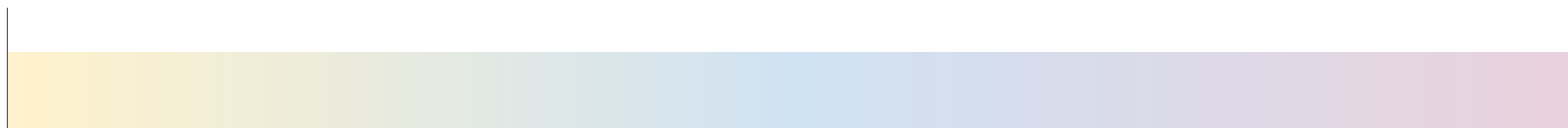For a given task, often have different software / algorithms which sit at different points along the spectrum.

These perform better / worse for a particular set of input data. Why?

" De Novo"                    " Homology Based"

" De Novo"                                          " Homology Based"

De Novo                                                        Homology Based

# De Novo

Existing data not used.

Existing data:

- May not be applicable, or
- May not help

- Indexing & Alignment
- Genome Assembly (De Novo)
- Transcriptome Assembly (De Novo)
- Evolutionary Trees
- Dimensionality Reduction
- Sequence / Structural Variation

" De Novo"

```
Indexing & Alignment
Genome Assembly (De Novo)
Transcriptome Assembly (De Novo)
Evolutionary Trees
Dimensionality Reduction
Sequence / Structural Variation      /
```

De Novo

Homology Based

# Homology Based

Existing data used from large sequence databases.

Far end of spectrum:
- Existing data fully answers our question.

Depending on task, may require strong evidence:
- Large set of corroborating data
- Experimentally validated data

For sequences: Homology

- Characterisation (MetaPhLan)
- Gene Annotation (BLAST)
- Identifying Homologues
- Conservation

De Novo

Homology Based

# Training Based

Existing data can be used as model.

Combines existing data with new data.

Statistical approaches (often machine learning)

Large group.

- Repeat Masking
- Signal Peptide Identification
- Gene Identification
- Gene Annotation
- Variant Effect Prediction

```
Characterisation (MetaPhlAn)          MetaPhlAn
Gene Annotation (BLAST)               BLAST
Identifying Homologues
Conservation
```

De Novo                    Training Based                    Homology Based

# Leveraging Data

| Task | De Novo | Training Based | Homology Based |
|---|:---:|:---:|:---:|
| Genome QC | ✖ | ✖ | ✔ |
| Repeat Masking | ✔ | ✔ | ✖ |
| Gene Identification | ✖ | ✔ | ✖ |
| Gene Annotation | | | |
| Promoter / Enhancer Identification | ✖ | ✔ | ✔ |
| Signal Peptide Identification | ✖ | ✔ | ✔ |
| Sequence / Structural Variation | ✔ | ✔ | ✖ |
| Variant Effect Prediction | ✔ | ✔ | ✖ |
| Characterisation | ✖ | ✖ | ✔ |

# Leveraging Data

| Task | De Novo | Training Based | Homology Based |
| --- | :---: | :---: | :---: |
| Genome QC | ✖ | ✖ | ✔ |
| Repeat Masking | ✔ | ✔ | ✖ |
| Gene Identification | ✖ | ✔ | ✖ |
| Gene Annotation | ✖ | | |
| Promoter / Enhancer Identification | ✖ | ✔ | ✔ |
| Signal Peptide Identification | ✖ | ✔ | ✔ |
| Sequence / Structural Variation | ✔ | ✔ | ✖ |
| Variant Effect Prediction | ✔ | ✔ | ✖ |
| Characterisation | ✖ | ✖ | ✔ |

# Leveraging Data

| Task | De Novo | Training Based | Homology Based |
|------|:-------:|:--------------:|:--------------:|
| Genome QC | ✖ | ✖ | ✔ |
| Repeat Masking | ✔ | ✔ | ✖ |
| Gene Identification | ✖ | ✔ | ✖ |
| Gene Annotation | ✖ | ✔ | |
| Promoter / Enhancer Identification | ✖ | ✔ | ✔ |
| Signal Peptide Identification | ✖ | ✔ | ✔ |
| Sequence / Structural Variation | ✔ | ✔ | ✖ |
| Variant Effect Prediction | ✔ | ✔ | ✖ |
| Characterisation | ✖ | ✖ | ✔ |

# Leveraging Data

| Task | De Novo | Training Based | Homology Based |
|------|:-------:|:--------------:|:--------------:|
| Genome QC | ✖ | ✖ | ✔ |
| Repeat Masking | ✔ | ✔ | ✖ |
| Gene Identification | ✖ | ✔ | ✖ |
| Gene Annotation | ✖ | ✔ | ✔ |
| Promoter / Enhancer Identification | ✖ | ✔ | ✔ |
| Signal Peptide Identification | ✖ | ✔ | ✔ |
| Sequence / Structural Variation | ✔ | ✔ | ✖ |
| Variant Effect Prediction | ✔ | ✔ | ✖ |
| Characterisation | ✖ | ✖ | ✔ |

# Leveraging Data

| Task | De Novo | Training Based | Homology Based |
|---|---|---|---|
| Genome QC | ✖ | ✖ | ✔ |
| Repeat Masking | ✔ | ✔ | ✖ |
| Gene Identification | ✖ | ✔ | ✖ |
| Gene Annotation | ✖ | ✔ | ✔ |
| Promoter / Enhancer Identification | ✖ | ✔ | ✔ |
| Signal Peptide Identification | ✖ | ✔ | ✔ |
| Sequence / Structural Variation | ✔ | ✔ | ✖ |
| Variant Effect Prediction | ✔ | ✔ | ✖ |
| Characterisation | ✖ | ✖ | ✔ |

Particular software may perform better / worse for a particular set of input data. Why?

# Prokka: Annotation Pipeline

# Prokka

Prokka     pipeline

FASTA     DNA     contigs/ scaffolds

/

## Prokka: Torsten Seemann

Prokaryote genome annotation pipeline.
Coordinates existing tools in a pipeline.

**Table 1.** Feature prediction tools used by Prokka

| Tool (reference) | Features predicted |
| --- | --- |
| Prodigal ( Hyatt 2010 ) | Coding sequence (CDS) |
| RNAmmer ( Lagesen et al. , 2007 ) | Ribosomal RNA genes (rRNA) |
| Aragorn ( Laslett and Canback, 2004 ) | Transfer RNA genes |
| SignalP ( Petersen et al. , 2011 ) | Signal leader peptides |
| Infernal ( Kolbe and Eddy, 2011 ) | Non-coding RNA |

## Input

DNA sequences in Fasta format (contigs / scaffolds)
Ideal: 1 sequence per chromosome / plasmid, no gaps
Real: Fragmented, gaps

## Output:

Annotated features with start / stop positions
Need to both identify, and annotate these features.

# Prokka

Prokka's use of data: It's all about trust.

# Prokka

Prokka's use of data: It's all about trust.

Hierarchical approach (coding sequence annotation):

1. Optional user-provide set of annotated proteins.
   Expected to be curated datasets - high confidence

# Prokka

Prokka's use of data: It's all about trust.

Hierarchical approach (coding sequence annotation):

1. Optional user-provide set of annotated proteins.
   Expected to be curated datasets - high confidence

2. All bacterial proteins in UniProt that have real protein or transcript evidence
   Experimentally validated - high confidence
   Typically covers >50% of the core genes in most genomes

# Prokka

Prokka's use of data: It's all about trust.

Hierarchical approach (coding sequence annotation):

1. Optional user-provide set of annotated proteins.
   Expected to be curated datasets - high confidence

2. All bacterial proteins in UniProt that have real protein or transcript evidence
   Experimentally validated - high confidence
   Typically covers >50% of the core genes in most genomes

3. All proteins from finished bacterial genomes in RefSeq for a specified genus
   Proteins this organism is reasonably likely to possess

# Prokka

Prokka's use of data: It's all about trust.

Hierarchical approach (coding sequence annotation):

UniProt

50%

RefSeq

Hidden Markov Model, HMM

"    /    "

1. Optional user-provide set of annotated proteins.
   Expected to be curated datasets - high confidence

   Prokka

2. All bacterial proteins in UniProt that have real protein or transcript evidence
   Experimentally validated - high confidence
   Typically covers >50% of the core genes in most genomes

3. All proteins from finished bacterial genomes in RefSeq for a specified genus
   Proteins this organism is reasonably likely to possess

4. A series of hidden Markov model profile databases
   "Putative / predicted" proteins

# Training Based Approaches

# Training Based Approaches

Homology approaches

Rely on relevant & high confidence data.

What if your organism is too different from those in databases?

# Training Based Approaches

Homology approaches

Rely on relevant & high confidence data.

What if your organism is too different from those in databases?

Training based approaches

If no direct comparisons can be made, can turn to generalised features / patterns.

Underlying structure.

```
                Homology approaches



            Training based approaches
                                /
    "       "    Underlying structure
```



**Bicep Curls**

LICENCED UNDER CC BY-SA 4.0

# Training Based Approaches



A mostly complete chart of
# Neural Networks
©2019 Fjodor van Veen & Stefan Leijnen    asimovinstitute.org

# Training Based Approaches

Common in bioinformatics (many flavors)
Latent diffusion: Stable Diffusion, Midjourney



A mostly complete chart of
## Neural Networks
©2019 Fjodor van Veen & Stefan Leijnen    asimovinstitute.org

Dota2: OpenAI Five

Nanopore Basecalling

Generate Examples for Image Datasets

# Markov Models

# Markov Models

Widely used in bioinformatics.

Visible data (eg a sequence) modelled as a system with states.

Will explore these flavors:

- Markov Chain
- Interpolated Markov Models
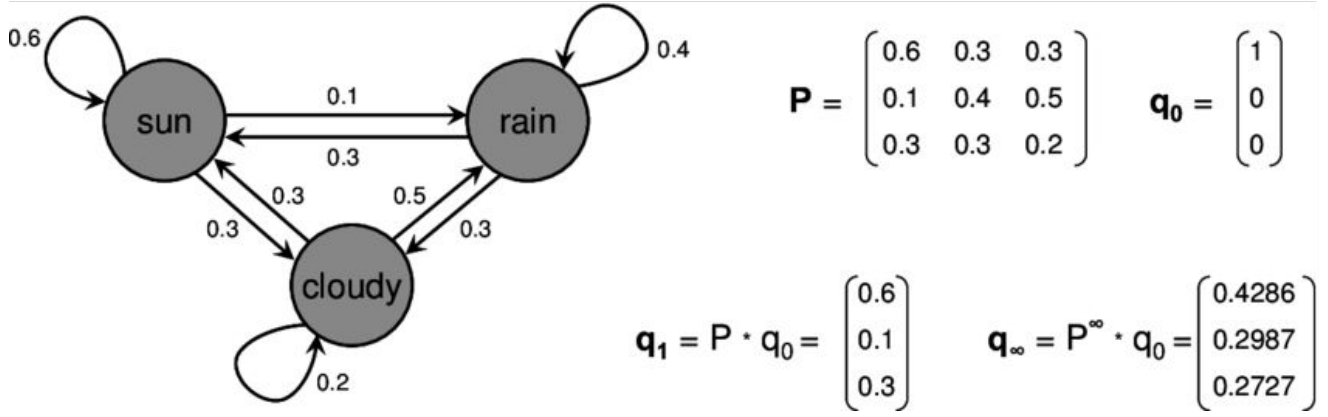- Hidden Markov Models
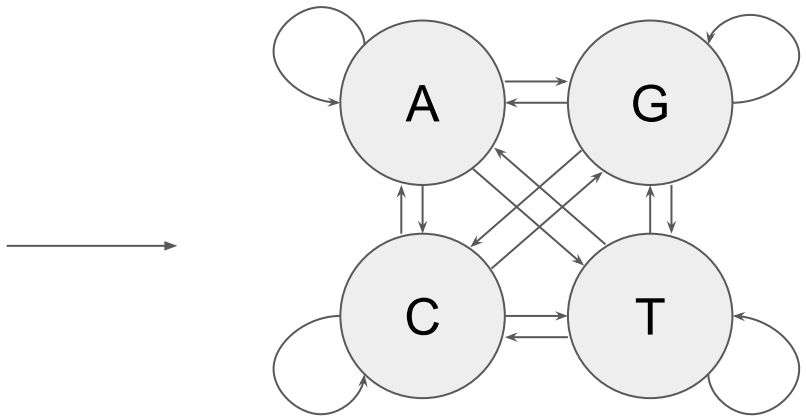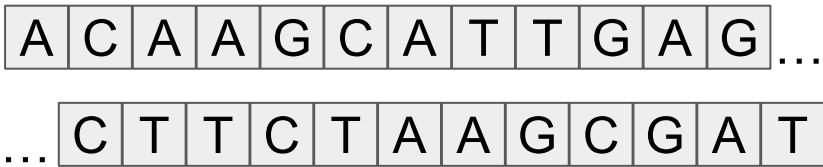
# Markov Chain

Simplest example.

Stochastically model of a series of events or 'states'

Current state only depends on previous state (memorylessness)

Transition matrix models the transition probabilities from one state to another.



$$P = \begin{bmatrix} 0.6 & 0.3 & 0.3 \\ 0.1 & 0.4 & 0.5 \\ 0.3 & 0.3 & 0.2 \end{bmatrix} \qquad q_0 = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$$

$$q_1 = P \cdot q_0 = \begin{bmatrix} 0.6 \\ 0.1 \\ 0.3 \end{bmatrix} \qquad q_\infty = P^\infty \cdot q_0 = \begin{bmatrix} 0.4286 \\ 0.2987 \\ 0.2727 \end{bmatrix}$$

Arnold, Ruedi. (2023). Interactive Learning Environments for Mathematical Topics.

memorylessness

transition matrix

0.1                    0.3

P

P * q0    q
       q0

q1

q0

P

# Markov Chain

Simplest example.

Stochastically model of a series of events or 'states'

Current state only depends on previous state (memorylessness)

Transition matrix models the transition probabilities from one state to another.



$$P = \begin{bmatrix} 0.6 & 0.3 & 0.3 \\ 0.1 & 0.4 & 0.5 \\ 0.3 & 0.3 & 0.2 \end{bmatrix} \quad q_0 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

$$q_1 = P \cdot q_0 = \begin{bmatrix} 0.6 \\ 0.1 \\ 0.3 \end{bmatrix} \quad q_\infty = P^\infty \cdot q_0 = \begin{bmatrix} 0.4286 \\ 0.2987 \\ 0.2727 \end{bmatrix}$$

Arnold, Ruedi. (2023). Interactive Learning Environments for Mathematical Topics.



| A | C | A | A | G | C | A | T | T | G | A | G | … |

… | C | T | T | C | T | A | A | G | C | G | A | T |

# How does this apply to gene prediction?

Goal: Predict a genomic feature.

Underlying structure of coding regions is different to non-coding regions.

# How does this apply to gene prediction?

coding regions

DNA          non-coding regions

Goal: Predict a genomic feature.

codons

Underlying structure of coding regions is different to non-coding regions.

Encoding a protein puts different pressures on the sequence.

Certain amino acids are more common
Leucine used more commonly than Cysteine

Codon wobble
Alanine has 4 codons, Methionine has 1 codon.

Codon preference
For a given species, usually one codon is used more often than others

# How does this apply to gene prediction?

1. Generate two models:
   - One from known coding sequence
   - One from known non-coding sequence

# How does this apply to gene prediction?

1.  Generate two models:
    -   One from known coding sequence
    -   One from known non-coding sequence

2.  Identify open reading frames (ORFs)
    -   Start / stop codons in the same frame

```
1. ATG CAA TGG GGA AAT GTT ACC AGG TCC GAA CTT ATT GAG GTA AGA CAG ATT TAA
2. A TGC AAT GGG GAA ATG TTA CCA GGT CCG AAC TTA TTG AGG TAA GAC AGA TTT AA
3. AT GCA ATG GGG AAA TGT TAC CAG GTC CGA ACT TAT TGA GGT AAG ACA GAT TTA A
```

All 3 possible ORFS for +ve strand, indicated by start & stop codon

# How does this apply to gene prediction?

1. Generate two models:
   - One from known coding sequence
   - One from known non-coding sequence

2. Identify open reading frames (ORFs)
   - Start / stop codons in the same frame

3. For each ORF, score the sequence by:
   - Calculating the probability that it was generated by the coding model.
   - Calculating the probability that it was generated by the non-coding model.

```
1. ATG CAA TGG GGA AAT GTT ACC AGG TCC GAA CTT ATT GAG GTA AGA CAG ATT TAA
2. A TGC AAT GGG GAA ATG TTA CCA GGT CCG AAC TTA TTG AGG TAA GAC AGA TTT AA
3. AT GCA ATG GGG AAA TGT TAC CAG GTC CGA ACT TAT TGA GGT AAG ACA GAT TTA A
```

All 3 possible ORFS for +ve strand, indicated by start & stop codon

# How does this apply to gene prediction?

1. Generate two models:
   - One from known coding sequence
   - One from known non-coding sequence

2. Identify open reading frames (ORFs)
   - Start / stop codons in the same frame

3. For each ORF, score the sequence by:
   - Calculating the probability that it was generated by the coding model.
   - Calculating the probability that it was generated by the non-coding model.

4. Higher probability for coding model indicates possible gene
   - Confidence can be measured using some form of ratio.

```
1. ATG CAA TGG GGA AAT GTT ACC AGG TCC GAA CTT ATT GAG GTA AGA CAG ATT TAA
2. A TGC AAT GGG GAA ATG TTA CCA GGT CCG AAC TTA TTG AGG TAA GAC AGA TTT AA
3. AT GCA ATG GGG AAA TGT TAC CAG GTC CGA ACT TAT TGA GGT AAG ACA GAT TTA A
```

All 3 possible ORFS for +ve strand, indicated by start & stop codon

# Interpolated Markov Model (IMM)

Previously: 1st order model

States: 'A', 'G', 'T', 'C'

# Interpolated Markov Model (IMM)

Previously: 1st order model

  States: 'A', 'G', 'T', 'C'

Would be nice if we could include more "history"

  Eg. Predict the next word:

  "... what"  {"do", "is", "the", "to"}

# Interpolated Markov Model (IMM)

Previously: 1st order model

  States: 'A', 'G', 'T', 'C'

Would be nice if we could include more "history"

  Eg. Predict the next word:

  "... what"  {"do", "is", "the", "to"}

# Interpolated Markov Model (IMM)

Previously: 1st order model

   States: 'A', 'G', 'T', 'C'

Would be nice if we could include more "history"

   Eg. Predict the next word:

   "... what"  {"do", "is", "the", "to"}

   "... I don't know what"  {"do", "is", "the", "to"}

# Interpolated Markov Model (IMM)

Previously: 1st order model

   States: 'A', 'G', 'T', 'C'

Would be nice if we could include more "history"

   Eg. Predict the next word:

   "... what"  {"do", "is", "the", "to"}

   "... I don't know what"  {"do", "is", "the", "to"}

# Interpolated Markov Model (IMM)

Previously: 1st order model

   States: 'A', 'G', 'T', 'C'

Would be nice if we could include more "history"

   Eg. Predict the next word:

   "... what"  {"do", "is", "the", "to"}

   "... I don't know what"  {"do", "is", "the", "to"}

4th order model

# Interpolated Markov Model (IMM)

Issue: not enough data!

Number of parameters (transition probabilities)
we must estimate grows exponentially with
order

$n^{th}$ order: $4^{n+1}$

$1^{st}$ order: $4^2 = 16$

$4^{th}$ order: $4^5 = 1024$

# Interpolated Markov Model (IMM)

**Issue:** not enough data!

Number of parameters (transition probabilities) we must estimate grows exponentially with order

  $n^{th}$ order: $4^{n+1}$

  $1^{st}$ order: $4^2 = 16$
  $4^{th}$ order: $4^5 = 1024$

Imagining we have 200k bases of sequence data to estimate parameters:

  For $2^{nd}$ order Markov chain, expect to see each state ~12k times

  For $8^{th}$ order Markov chain, expect to see each state ~3 times

# Interpolated Markov Model (IMM)

Issue: not enough data!

Number of parameters (transition probabilities) we must estimate grows exponentially with order

$n^{th}$ order: $4^{n+1}$

$1^{st}$ order: $4^2 = 16$

$4^{th}$ order: $4^5 = 1024$

Imagining we have 200k bases of sequence data to estimate parameters:

For $2^{nd}$ order Markov chain, expect to see each state ~12k times

For $8^{th}$ order Markov chain, expect to see each state ~3 times

## Interpolated Markov Models (eg in GLIMMA)

- Build $w$ orders (individual chains)
- For each base position in sequence to assess, combine probabilities from each chain
- Weight by how much evidence was witnessed for that state for that particular chain

# Hidden Markov Model (HMM)

Markov Chain incorporating hidden 'states' which affect witnessed output
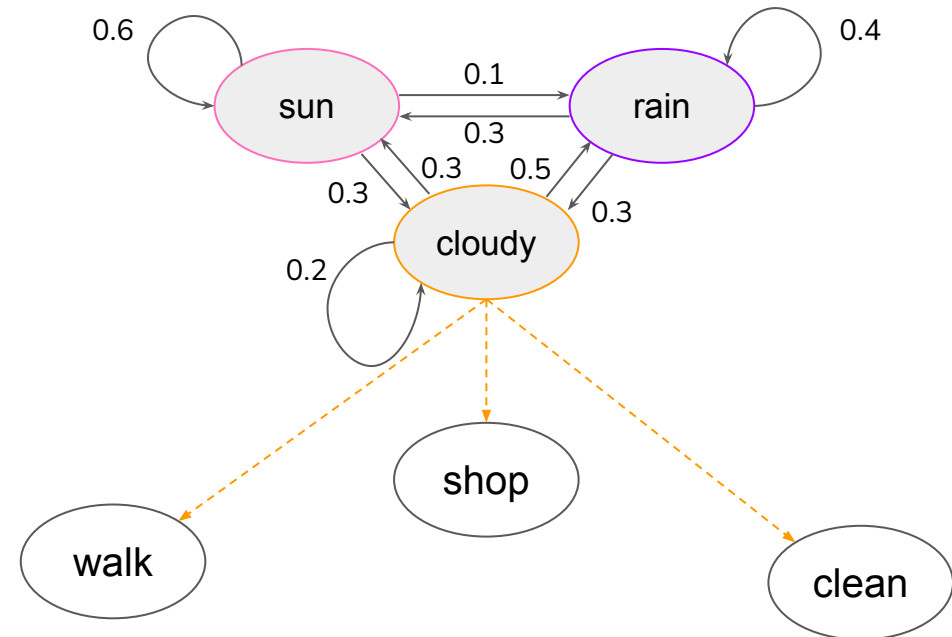
# Hidden Markov Model (HMM)

Markov Chain incorporating hidden 'states' which affect witnessed output

## Markov Model



Arnold, Ruedi. (2023). Interactive Learning
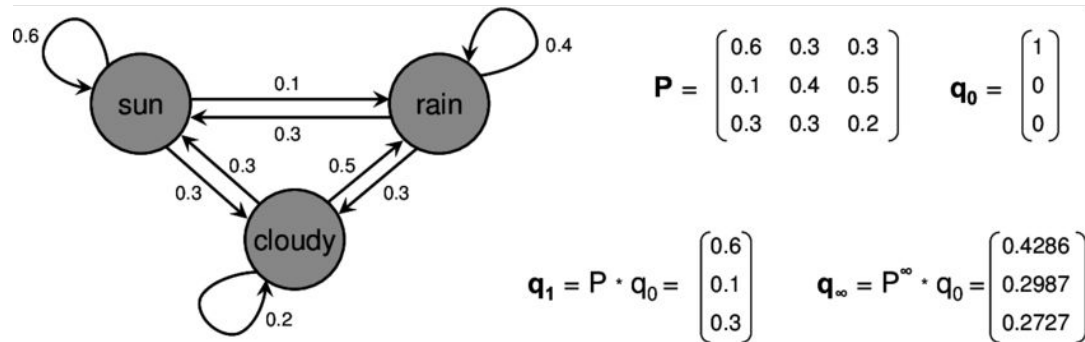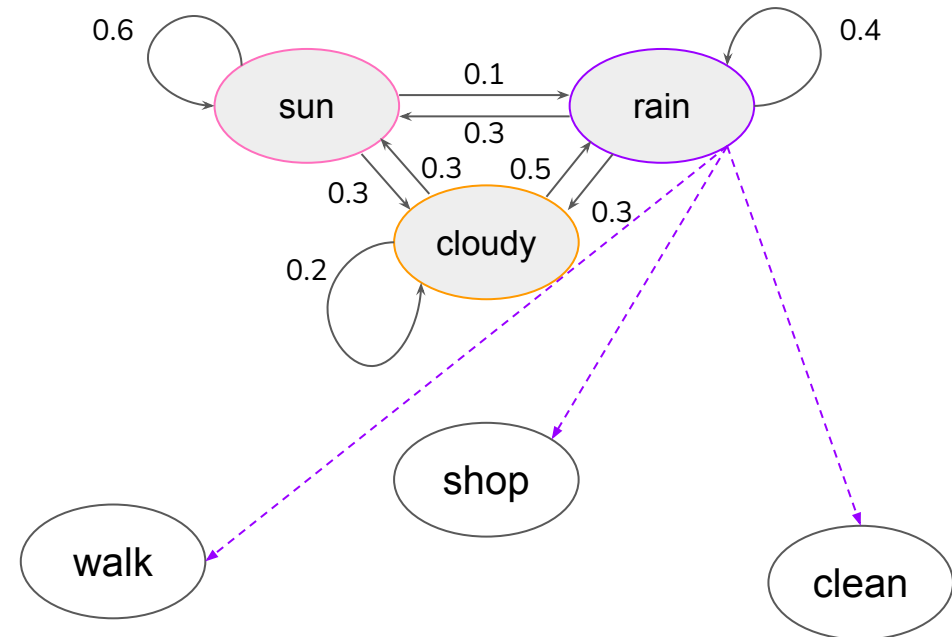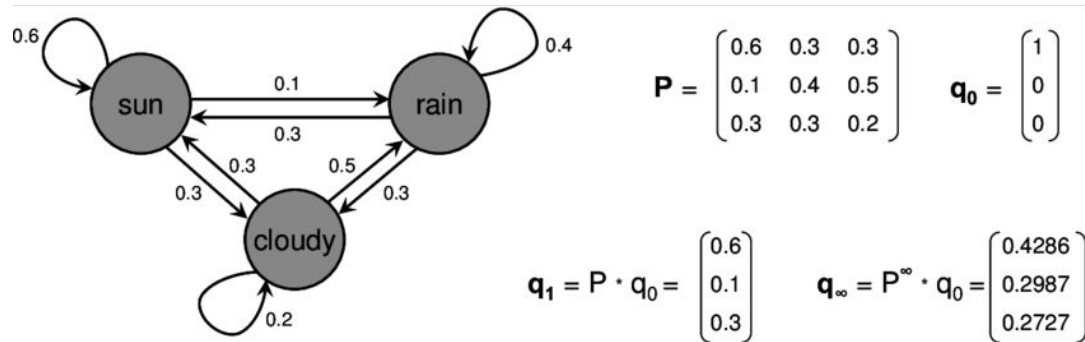Environments for Mathematical Topics.
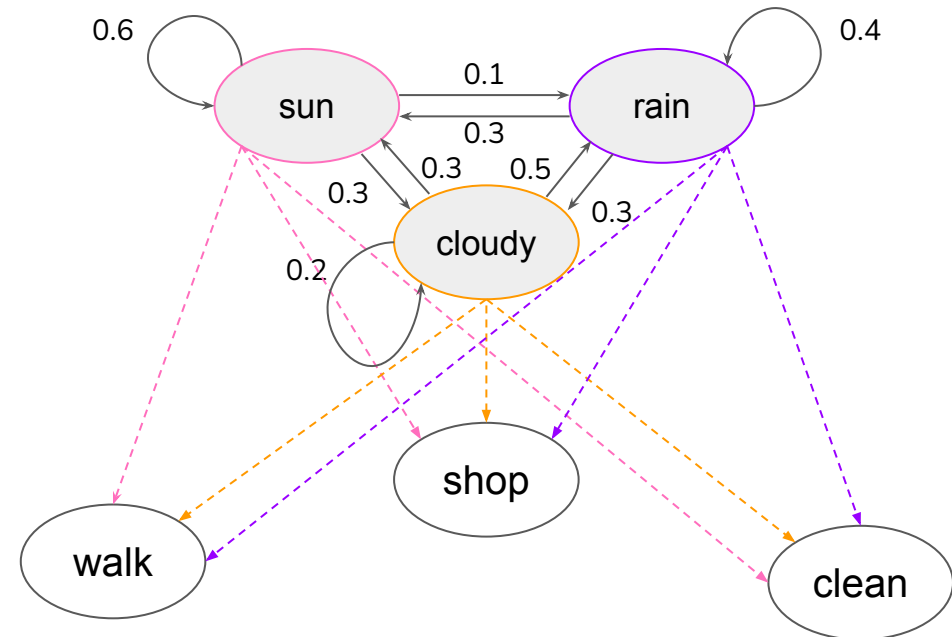
# Hidden Markov Model (HMM)

Markov Chain incorporating hidden 'states' which affect witnessed output

## Markov Model



Arnold, Ruedi. (2023). Interactive Learning
Environments for Mathematical Topics.

## Hidden Markov Model

# Hidden Markov Model (HMM)

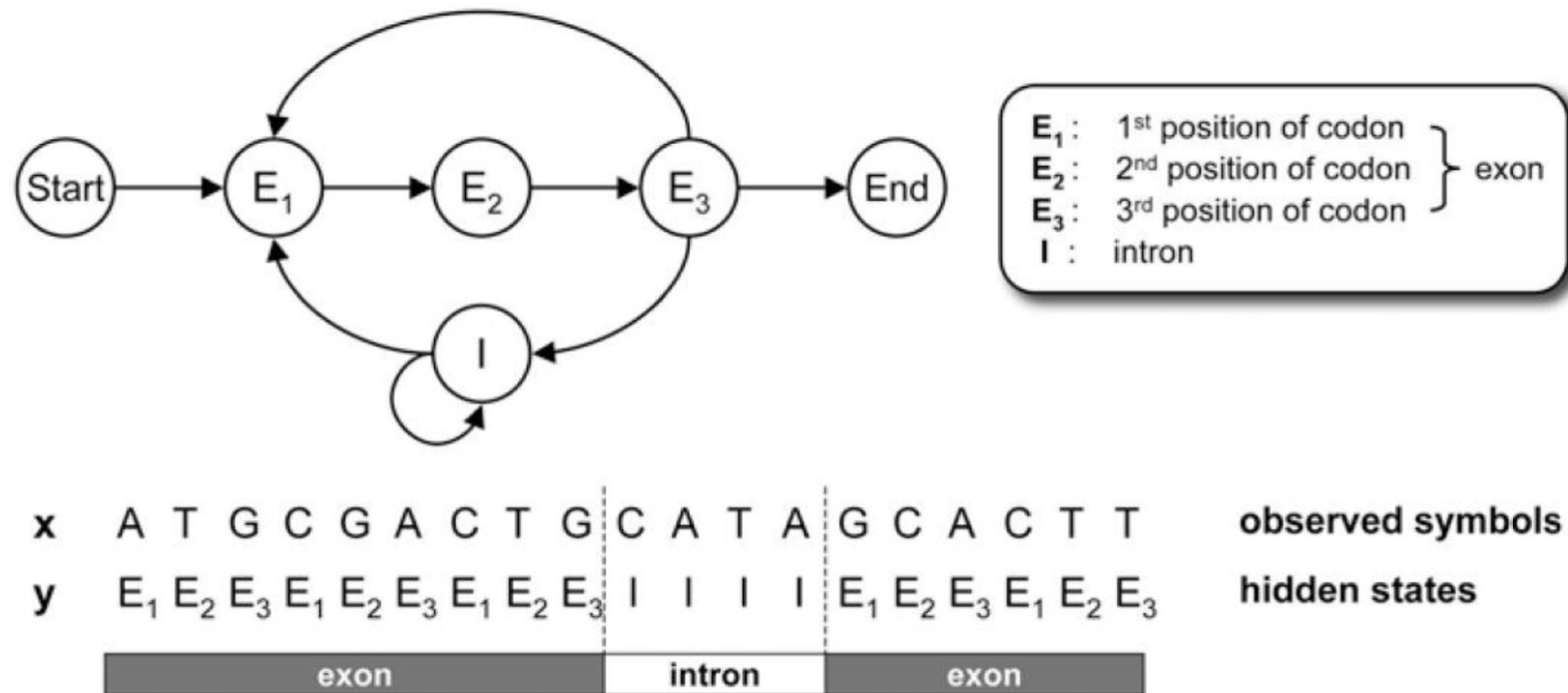Markov Chain incorporating hidden 'states' which affect witnessed output

### Markov Model



Arnold, Ruedi. (2023). Interactive Learning
Environments for Mathematical Topics.

### Hidden Markov Model



Emission
probabilities
for state "sun"

# Hidden Markov Model (HMM)

Markov Chain incorporating hidden 'states' which affect witnessed output

## Markov Model



Arnold, Ruedi. (2023). Interactive Learning Environments for Mathematical Topics.
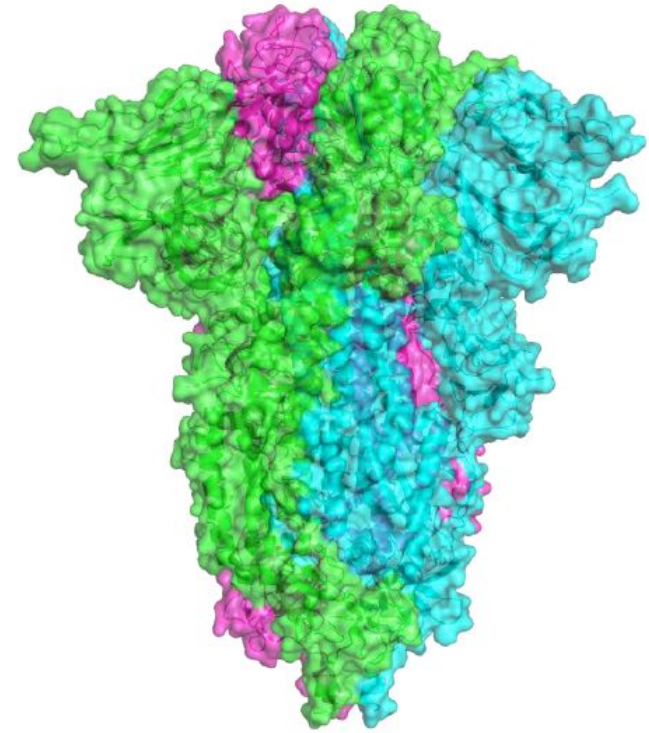
## Hidden Markov Model

# Hidden Markov Model (HMM)

Markov Chain incorporating hidden 'states' which affect witnessed output

### Markov Model



Arnold, Ruedi. (2023). Interactive Learning Environments for Mathematical Topics.

### Hidden Markov Model

# Hidden Markov Model (HMM)

Markov Chain incorporating hidden 'states' which affect witnessed output

## Markov Model



Arnold, Ruedi. (2023). Interactive Learning Environments for Mathematical Topics.

## Hidden Markov Model

# Hidden Markov Model (HMM)



Yoon, BJ 2009. doi: 10.2174/138920209789177575

# Deep Learning: AlphaFold 2.0

# AlphaFold 2.0

Given a sequence of amino acids, predicts 3D fold.



**"Spike" protein (7CAB)**

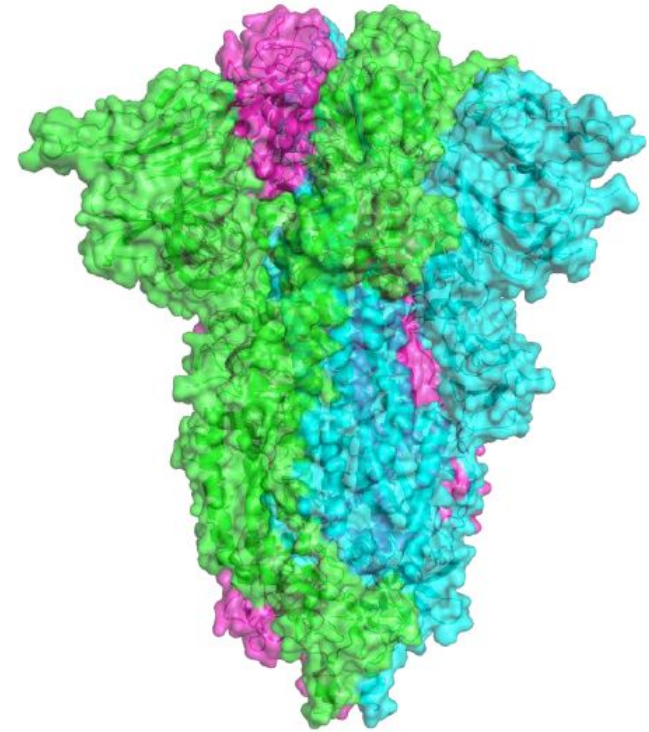Mediates viral entry into the host cell

# AlphaFold 2.0

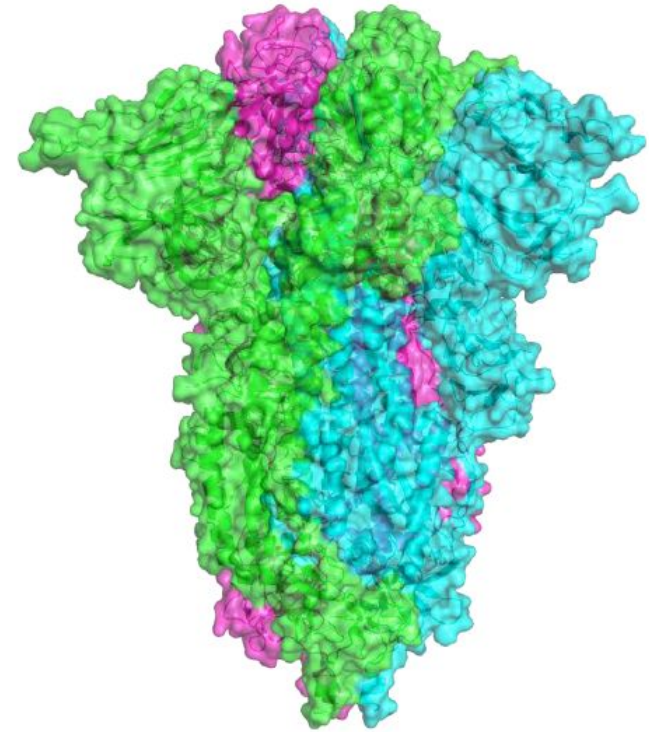Given a sequence of amino acids, predicts 3D fold.

Can be used for gene annotation:

For situations where DNA / AA identity to known proteins devolves below recognisable threshold

Conservation: DNA (least) -> AA -> Fold (most)



"Spike" protein (7CAB)

Mediates viral entry into the host cell

# AlphaFold 2.0

Given a sequence of amino acids, predicts 3D fold.

Can be used for gene annotation:

For situations where DNA / AA identity to known proteins
devolves below recognisable threshold

Conservation: DNA (least) -> AA -> Fold (most)

1. Identify possible genes (maybe HMM / IMM)
2. Fold each candidate
3. Compare each fold to structures on PDB



"Spike" protein (7CAB)

Mediates viral entry into the
host cell

# AlphaFold 2.0

Given a sequence of amino acids, predicts 3D fold.

Can be used for gene annotation:

For situations where DNA / AA identity to known proteins devolves below recognisable threshold

Conservation: DNA (least) -> AA -> Fold (most)

1. Identify possible genes (maybe HMM / IMM)
2. Fold each candidate
3. Compare each fold to structures on PDB

Machine learning (training & inference), but then uses homology (RMSD) to known proteins



"Spike" protein (7CAB)

Mediates viral entry into the host cell

# Protein Folding - Levinthal's Paradox (1969)

- The number of possible conformation available to a protein is astronomically large
  - 100 Amino acid protein
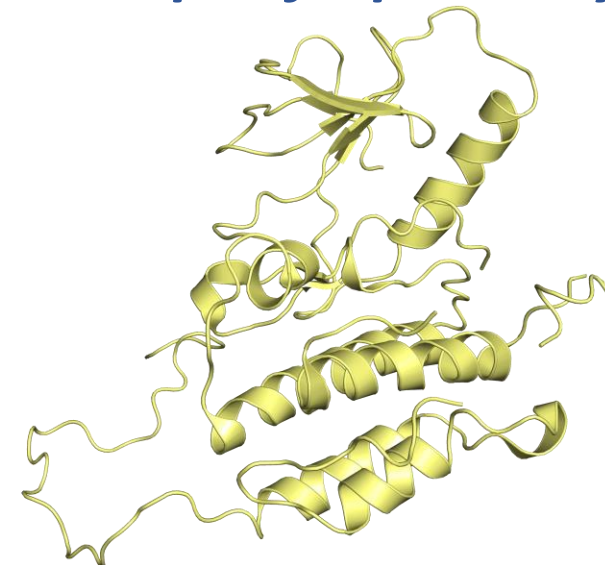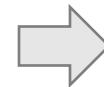  - 3 conformations per amino acid

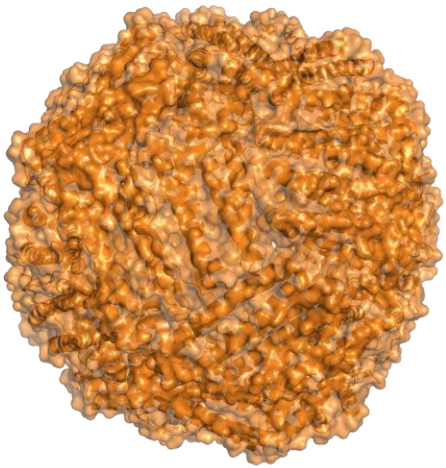$3^{100}$ conformations ⟹ $10^{12}$ conformations/second

2 x $10^{28}$ years

*Protein folding is not random and must have a specific pathway*

>P15056|BRAF_HUMAN
MAALSGGGGGGGAEPGQALFNGDMEPEAGAGAGAAASSAADPAIPEEVWNIKQMIKLTQEH
IEALLDKFGGEHNPPSIYLEAYEEYTSKLDALQQREQQLLESLGNGTDFSVSSSASMDTV
TSSSSSSLSVLPSSLSVFQNPTDVARSNPKSPQKPIVRVFLPNKQRTVVPARCGVTVRDS
LKKALMMRGLIPECCAVYRIQDGEKKPIGWDTDISWLTGEELHVEVLENVPLTTHNFVRK
TFFTLAFCDFCRKLLFQGFRCQTCGYKFHQRCSTEVPLMCVNYDQLDLLFVSKFFEHHPI
PQEEASLAETALTSGSSPSAPASDSIGPQILTSPSPSKSIPIPQPFRPADEDHRNQFGQR
DRSSSAPNVHINTIEPVNIDDLIRDQGFRGDGGSTTGLSATPPASLPGSLTNVKALQKSP
GPQRERKSSSSSEDRNRMKTLGRRDSSDDWEIPDGQITVGQRIGSGSFGTVYKGKWHGDV
AVKMLNVTAPTPQQLQAFKNEVGVLRKTRHVNILLFMGYSTKPQLAIVTQWCEGSSLYHH
LHIIETKFEMIKLIDIARQTAQGMDYLHAKSIIHRDLKSNNIFLHEDLTVKIGDFGLATV
KSRWSGSHQFEQLSGSILWMAPEVIRMQDKNPYSFQSDVYAFGIVLYELMTGQLPYSNIN
NRDQIIFMVGRGYLSPDLSKVRSNCPKAMKRLMAECLKKKRDERPLFPQILASIELLARS
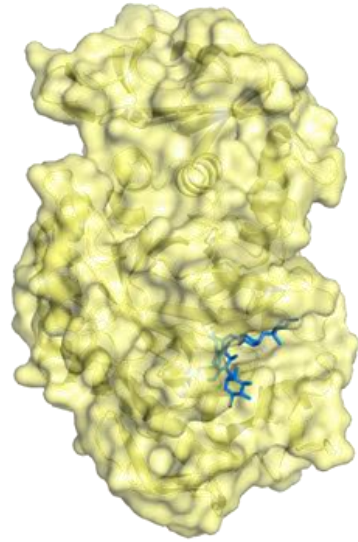LPKIHRSASEPSLNRAGFQTEDFSLYACASPKTPIQAGGYGAFPVH
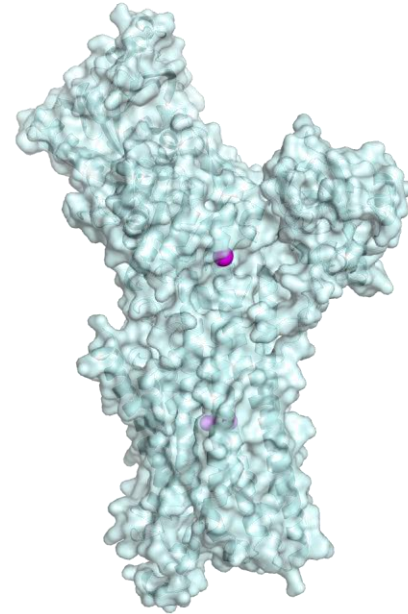


5

# Proteins: Structure and Function

**Ferritin (1FHA)**
Forms a hollow shell that
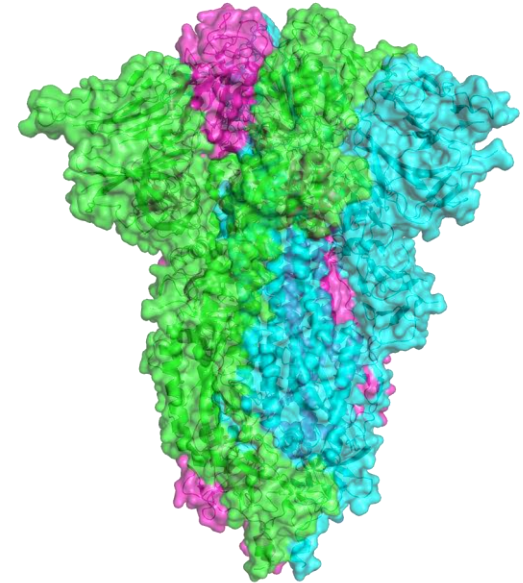stores iron from our food

**Alpha-amylase (1PPI)**
An enzyme with a catalytic site
that begins the breakdown of
carbohydrates in our saliva

**Calcium Pump (1SU4)**
Moves ions across cell
membrane

**"Spike" protein  (7CAB)**
Mediate viral entry in
the host cell

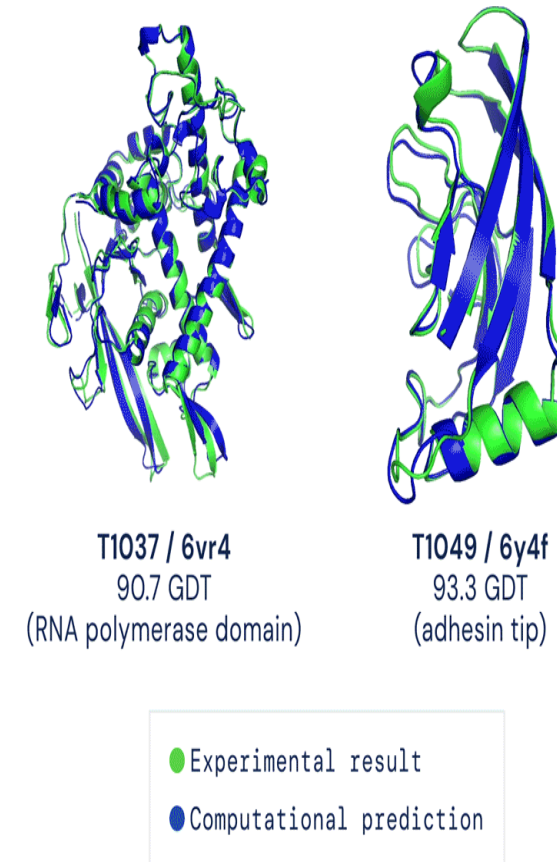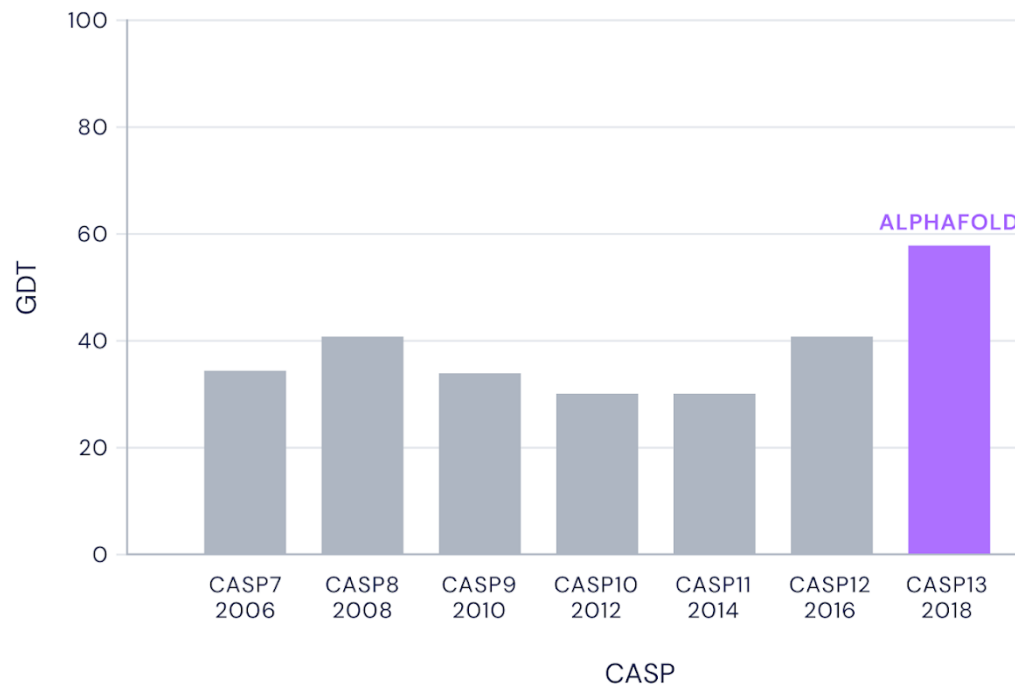# Protein Data Bank

- 3D data for biological macromomolecules
  - Proteins
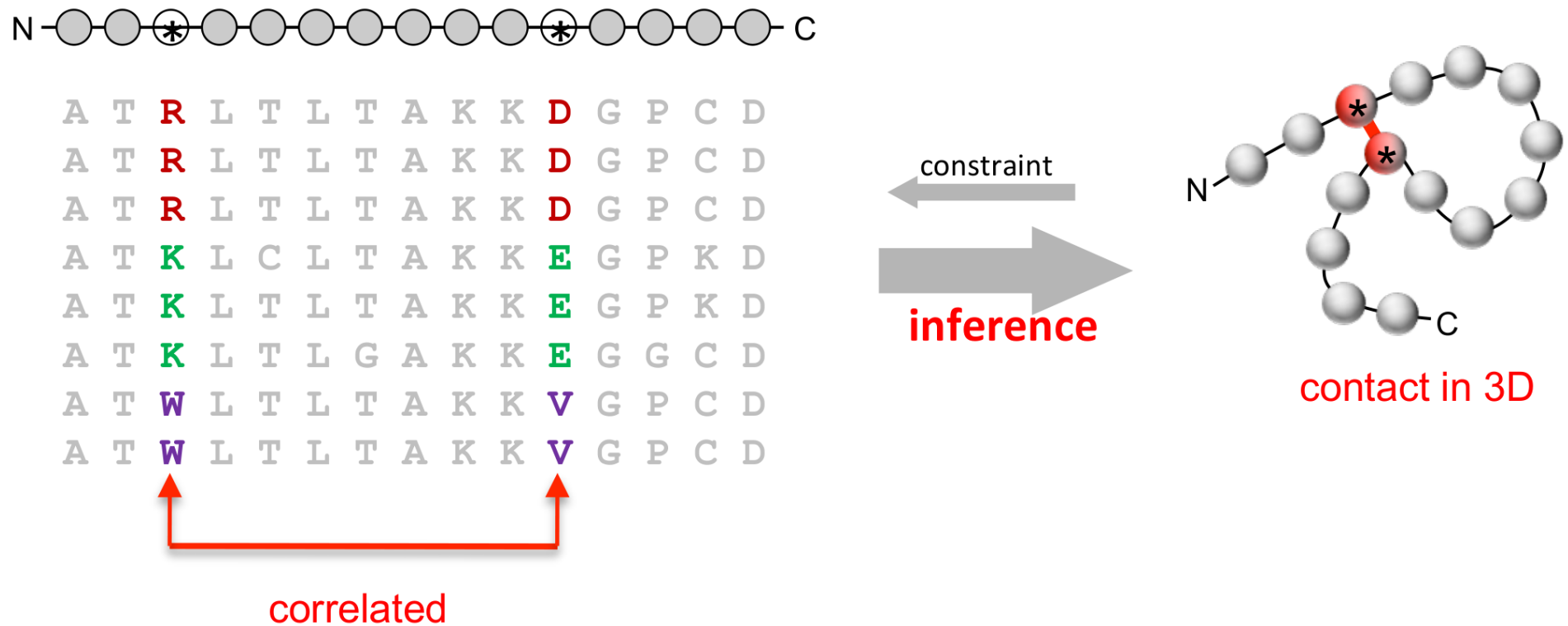  - Nucleic Acids (DNA – RNA)
  - Organelles
  - Viruses

| Molecular Type | X-ray | NMR | EM | Multiple methods | Neutron | Other | Total |
|---|---|---|---|---|---|---|---|
| Protein (only) | 134134 | 11512 | 4263 | 162 | 67 | 32 | 150170 |
| Other | 8122 | 92 | 551 | 6 | 0 | 4 | 8775 |
| Protein/NA | 7104 | 269 | 1517 | 3 | 0 | 0 | 8893 |
| Nucleic acid (only) | 2103 | 1309 | 54 | 6 | 2 | 1 | 3475 |
| Total | 151463 | 13182 | 6385 | 177 | 69 | 37 | 171313 |

# Critical Assessment of Structure Prediction 2018 (CASP13)



Median Free-Modelling Accuracy

T1037 / 6vr4
90.7 GDT
(RNA polymerase domain)

T1049 / 6y4f
93.3 GDT
(adhesin tip)

- Experimental result
- Computational prediction

Senior, Andrew W., et al. "Improved protein structure prediction using potentials from deep learning." Nature 577.7792 (2020): 706-710.

# Coevolution from Multiple Sequence Alignments (MSA)



contact in 3D

correlated

Marks, Debora S., et al. "Protein 3D structure computed from evolutionary sequence variation." *PloS one* 6.12 (2011): e28766.

# Critical Assessment of Structure Prediction 2018 (CASP13) - AlphaFold1



Senior, Andrew W., et al. "Improved protein structure prediction using potentials from deep learning." Nature 577.7792 (2020): 706-710.

# AlphaFold2



**NEWS** | 30 November 2020

## 'It will change everything': DeepMind's AI makes gigantic leap in solving protein structures

Google's deep-learning program for determining the 3D shapes of proteins stands to transform biology, say scientists.

The New York Times

## London A.I. Lab Claims Breakthrough That Could Accelerate Drug Discovery

Researchers at DeepMind say they have solved "the protein folding problem," a task that has bedeviled scientists for more than 50 years.

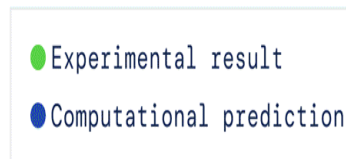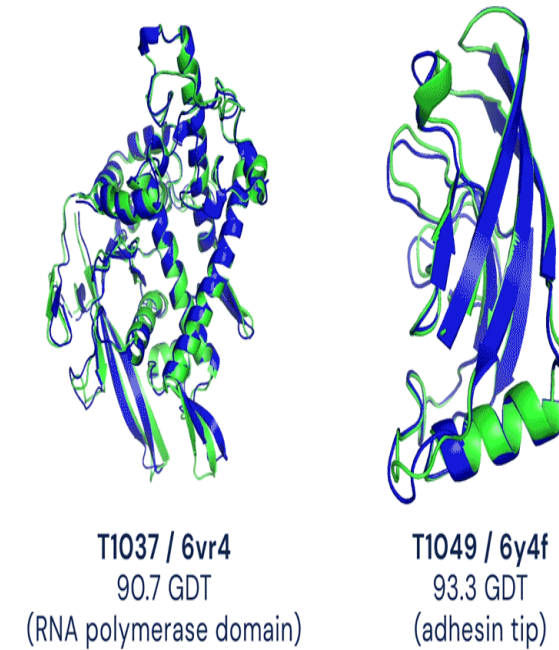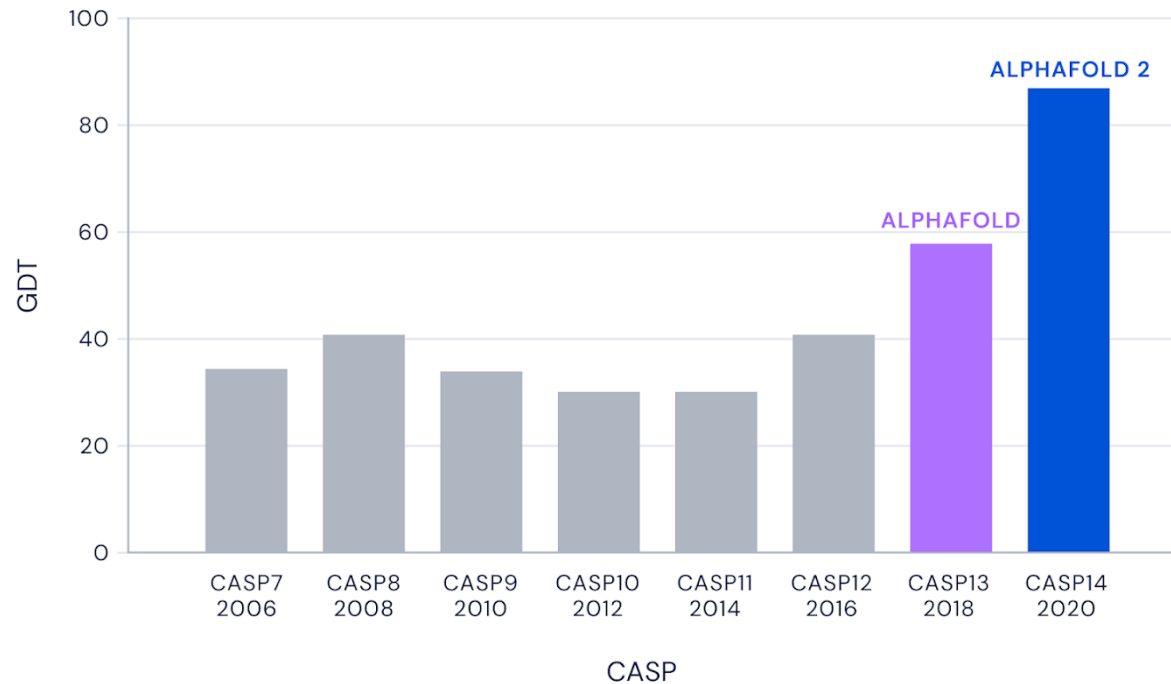**BLOG POST**
RESEARCH

30 NOV 2020

## AlphaFold: a solution to a 50-year-old grand challenge in biology

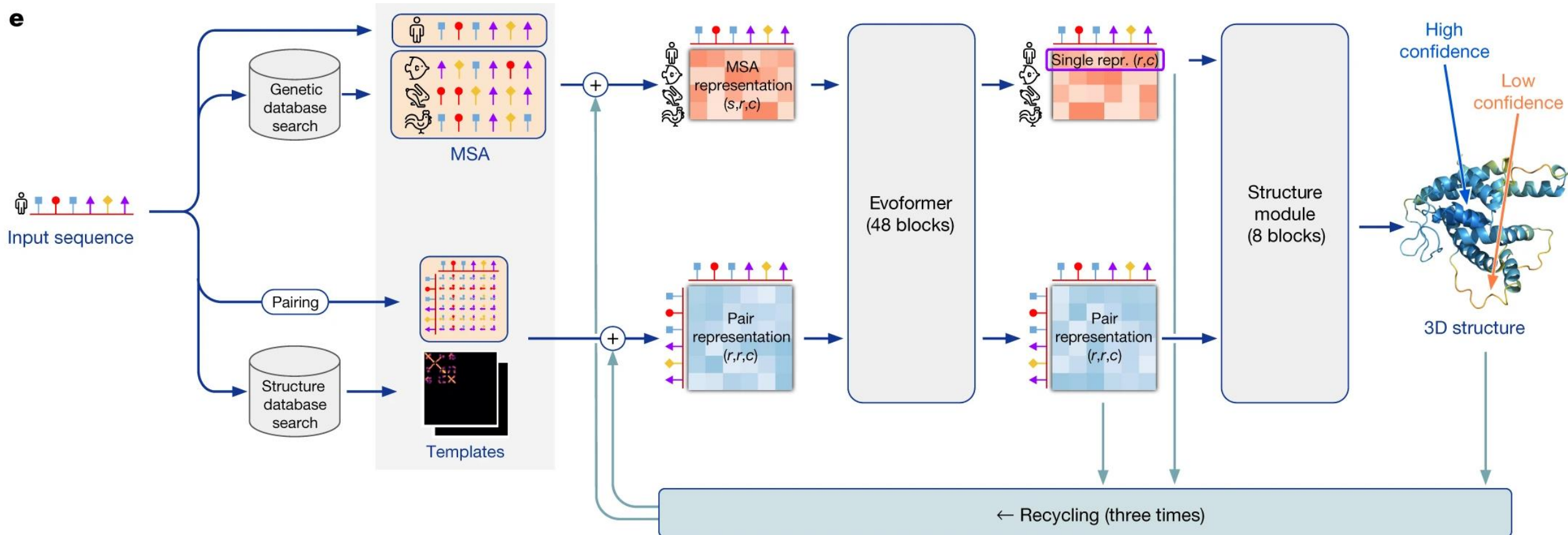## Great expectations – the potential impacts of AlphaFold DB

A discussion of the applications that AlphaFold DB may enable and the possible impact of the resource on science and society

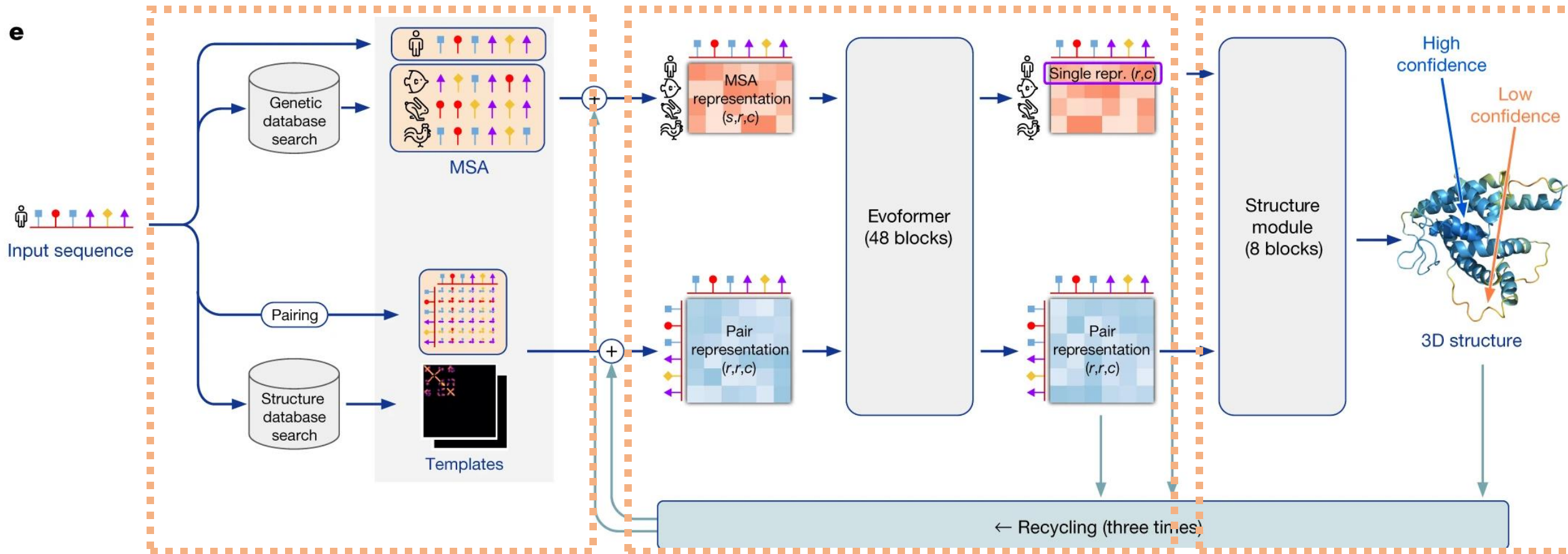# Critical Assessment of Structure Prediction 2020 (CASP14) - AlphaFold2
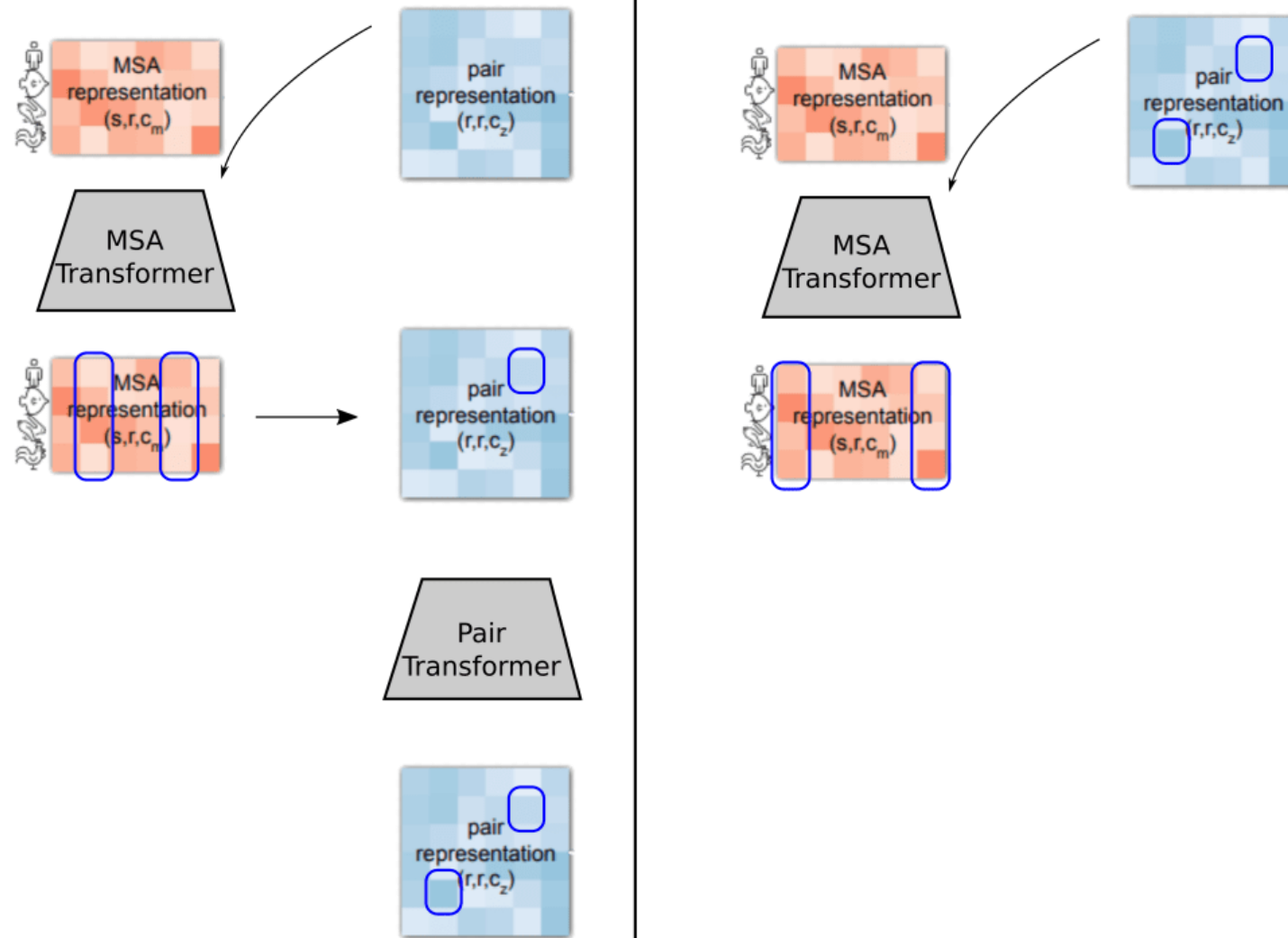


**Median Free-Modelling Accuracy**

ALPHAFOLD 2

ALPHAFOLD

GDT

CASP7 2006 | CASP8 2008 | CASP9 2010 | CASP10 2012 | CASP11 2014 | CASP12 2016 | CASP13 2018 | CASP14 2020

CASP



**T1037 / 6vr4**
90.7 GDT
(RNA polymerase domain)

**T1049 / 6y4f**
93.3 GDT
(adhesin tip)

● Experimental result
● Computational prediction

# Critical Assessment of Structure Prediction 2020 (CASP14) - AlphaFold2



Jumper, John, et al. "Highly accurate protein structure prediction with AlphaFold." Nature 596.7873 (2021): 583-589.

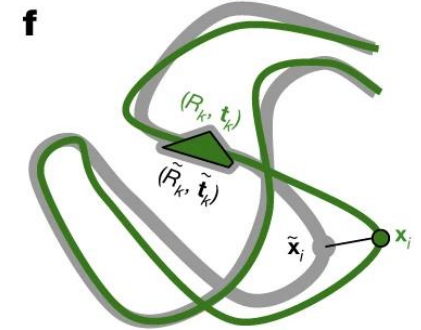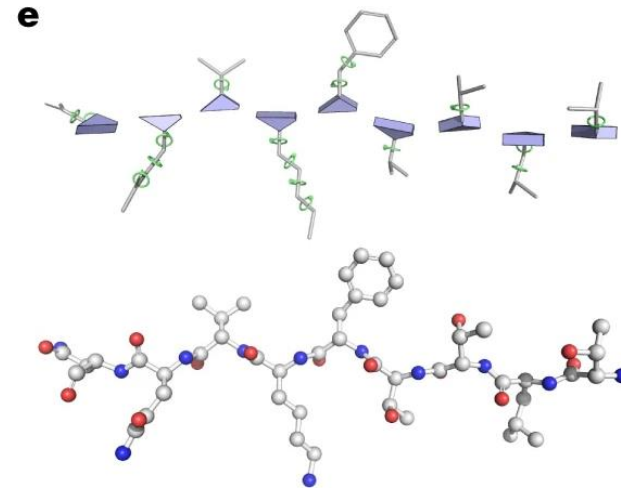# Critical Assessment of Structure Prediction 2020 (CASP14) - AlphaFold2



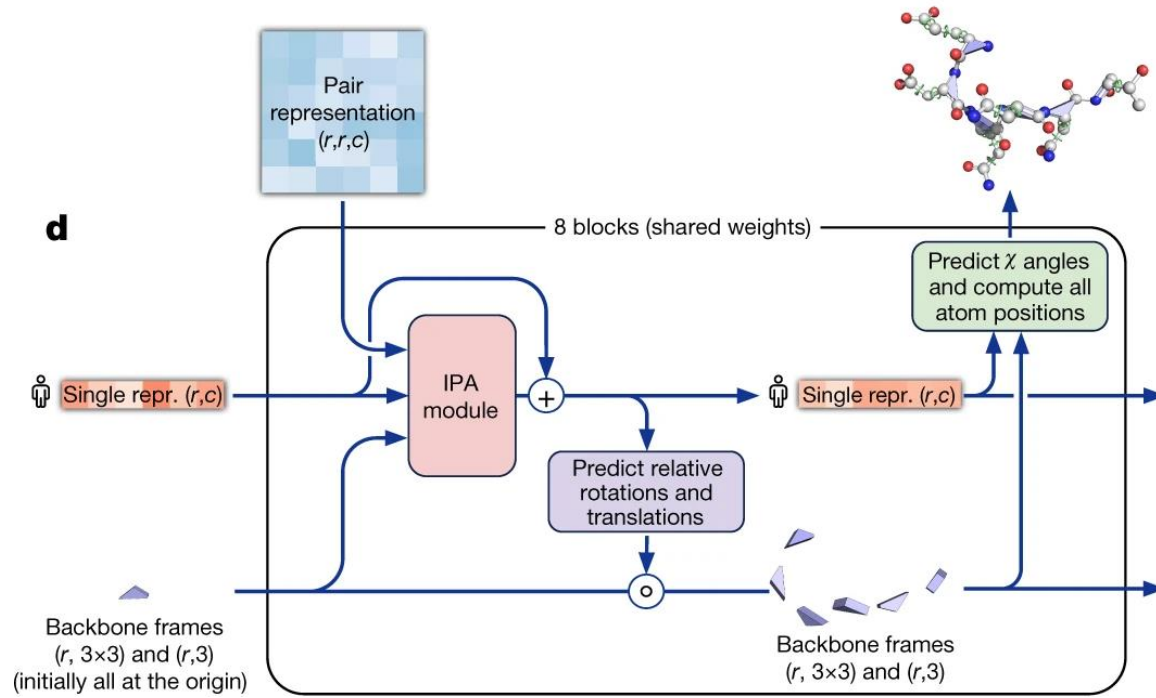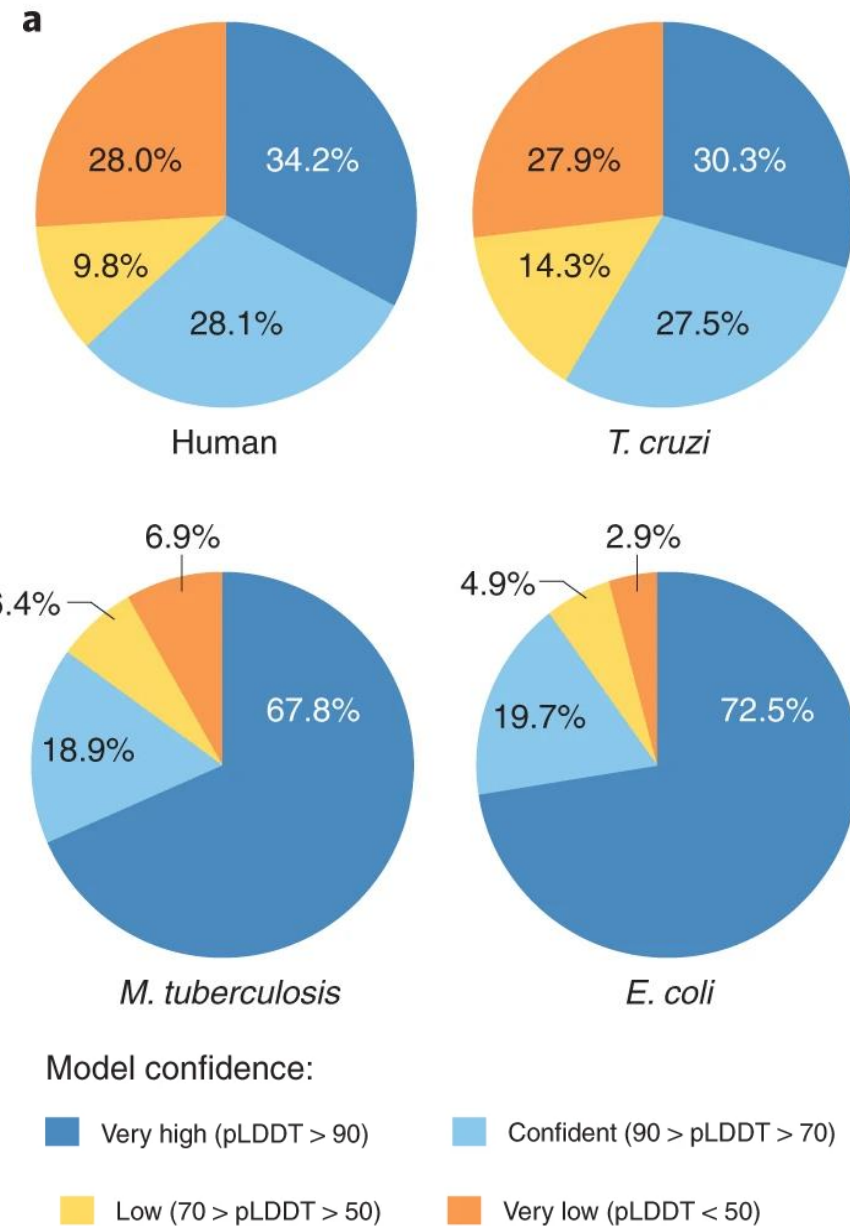Jumper, John, et al. "Highly accurate protein structure prediction with AlphaFold." Nature 596.7873 (2021): 583-589.

# AF2 - "Evolutionary Transformer" (Evoformer)

Jumper, John, et al. "Highly accurate protein structure prediction with AlphaFold." Nature 596.7873 (2021): 583-589.

# AF2 - Structure Module

Jumper, John, et al. "Highly accurate protein structure prediction with AlphaFold." Nature 596.7873 (2021): 583-589.

# AF2 – Coverage and Quality of predictions



Thornton, J. M., Laskowski, R. A. & Borkakoti, N. AlphaFold heralds a data-driven revolution in biology and medicine. Nat Med 27, 1666–1669 (2021). 21

# AF2 – Most recent updates



**ALPHAFOLD MANIA**
The number of research papers and preprints citing the AlphaFold2 AI software has shot up since its source code was released in July 2021*.

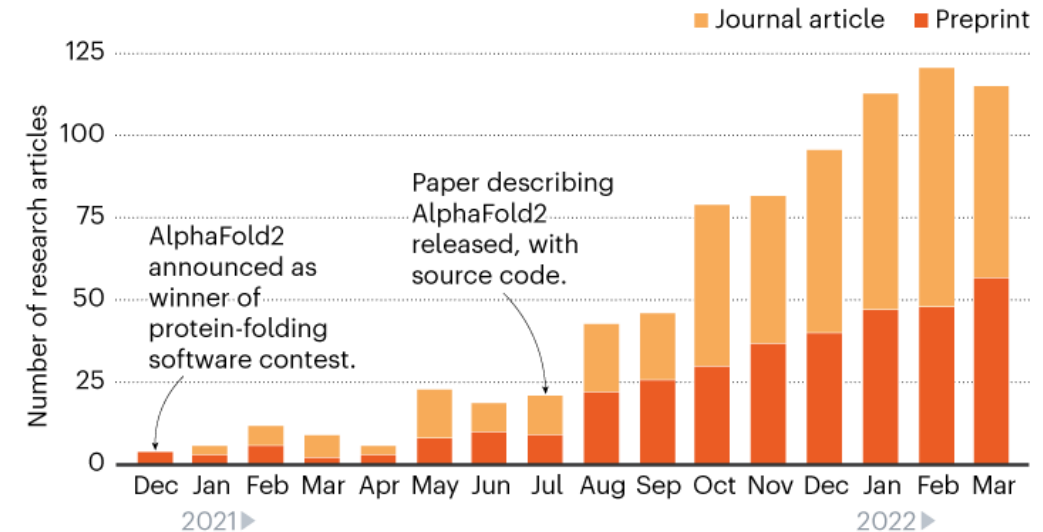*Nature analysis using Dimensions database; removing duplicate preprints and papers/R. Van Noorden, E. Callaway.

©nature

- AF2-Multimer
  - Prediction of protein-protein interactions
- Release of over structures for over 200k different proteins
  - Estimated human proteome ~25k
- Limitations
  - Min of 16 and max of 2,700 amino acids
    - Model smaller overlapping fragments
    - Alignment of predicted structures
    - "Glue" fragments on overlapping regions

# Thank you!

*Adam Taranto this Thursday! (31st August)*

*Ryan Wick next Tuesday!  (5th September)*

*Assignment 2 released tomorrow! (Wednesday 30th @ midnight)*

**Today:** Genomic Features & Regions

**Next time:** Genomic Intervals (actual algorithms)