# AI-DevTalk

# AI for protein folding: focus on Alphafold

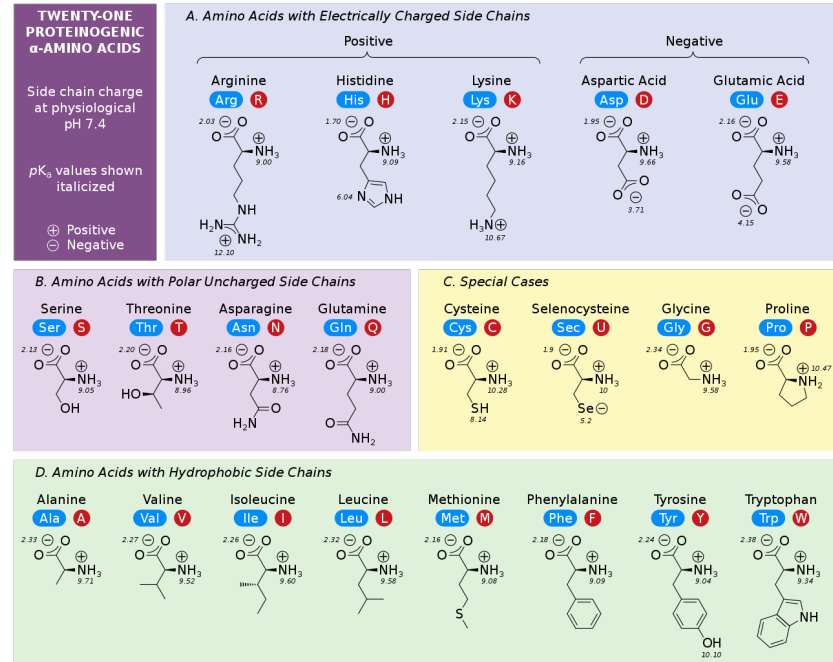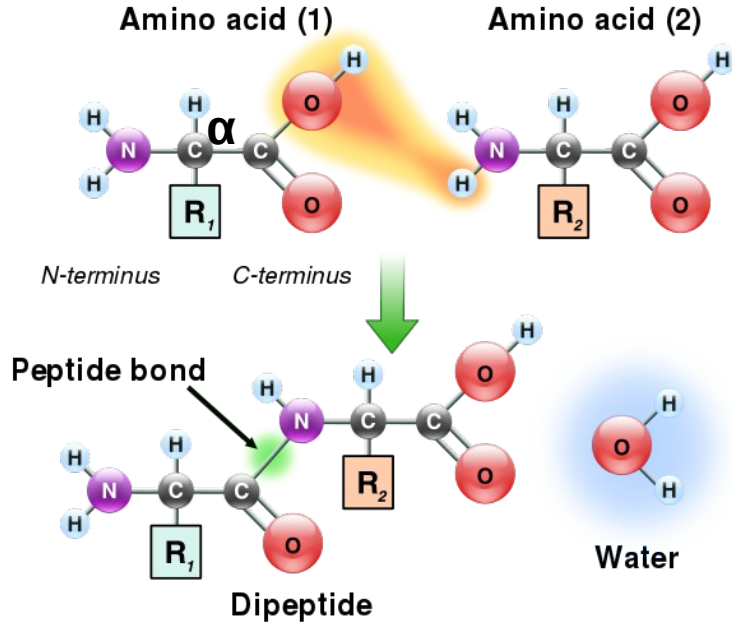Thibaut Véry – IDRIS User support

# What is a protein?

**Sequence of aminoacids: FASTA format**

```
>sp|P0DTC2|SPIKE_SARS2 Spike glycoprotein OS=Severe acute respiratory syndrome coronavirus 2 OX=2697049 GN=S PE=1 SV=1
MFVFLVLLPLVSSQCVNLTTRTQLPPAYTNSFTRGVYYPDKVFRSSVLHSTQDLFLPFFS
NVTWFHAIHVSGTNGTKRFDNPVLPFNDGVYFASTEKSNIIRGWIFGTTLDSKTQSLLIV
NNATNVVIKVCEFQFCNDPFLGVYYHKNNKSWMESEFRVYSSANNCTFEYVSQPFLMDLE
GKQGNFKNLREFVFKNIDGYFKIYSKHTPINLVRDLPQGFSALEPLVDLPIGINITRFQT
LLALHRSYLTPGDSSSGWTAGAAAYYVGYLQPRTFLLKYNENGTITDAVDCALDPLSETK
CTLKSFTVEKGIYQTSNFRVQPTESIVRFPNITNLCPFGEVFNATRFASVYAWNRKRISN
CVADYSVLYNSASFSTFKCYGVSPTKLNDLCFTNVYADSFVIRGDEVRQIAPGQTGKIAD
YNYKLPDDFTGCVIAWNSNNLDSKVGGNYNYLYRLFRKSNLKPFERDISTEIYQAGSTPC
NGVEGFNCYFPLQSYGFQPTNGVGYQPYRVVVLSFELLHAPATVCGPKKSTNLVKNKCVN
FNFNGLTGTGVLTESNKKFLPFQQFGRDIADTTDAVRDPQTLEILDITPCSFGGVSVITP
GTNTSNQVAVLYQDVNCTEVPVAIHADQLTPTWRVYSTGSNVFQTRAGCLIGAEHVNNSY
ECDIPIGAGICASYQTQTNSPRRARSVASQSIIAYTMSLGAENSVAYSNNSIAIPTNFTI
SVTTEILPVSMTKTSVDCTMYICGDSTECSNLLLQYGSFCTQLNRALTGIAVEQDKNTQE
VFAQVKQIYKTPPIKDFGGFNFSQILPDPSKPSKRSFIEDLLFNKVTLADAGFIKQYGDC
LGDIAARDLICAQKFNGLTVLPPLLTDEMIAQYTSALLAGTITSGWTFGAGAALQIPFAM
QMAYRFNGIGVTQNVLYENQKLIANQFNSAIGKIQDSLSSTASALGKLQDVVNQNAQALN
TLVKQLSSNFGAISSVLNDILSRLDKVEAEVQIDRLITGRLQSLQTYVTQQLIRAAEIRA
SANLAATKMSECVLGQSKRVDFCGKGYHLMSFPQSAPHGVVFLHVTYVPAQEKNFTTAPA
ICHDGKAHFPREGVFVSNGTHWFVTQRNFYEPQIITTDNTFVSGNCDVVIGIVNNTVYDP
LQPELDSFKEELDKYFKNHTSPDVDLGDISGINASVVNIQKEIDRLNEVAKNLNESLIDL
QELGKYEQYIKWPWYIWLGFIAGLIAIVMVTIMLCCMTSCCSCLKGCCSCGSCCKFDEDD
SEPVLKGVKLHYT
```
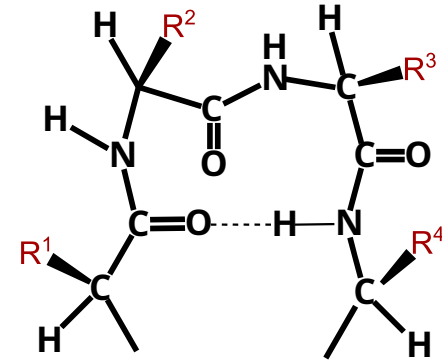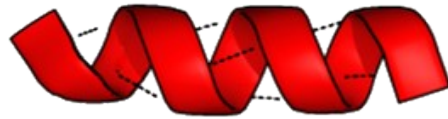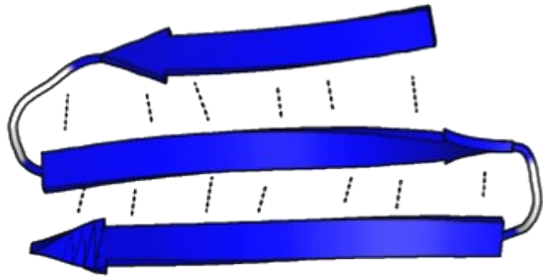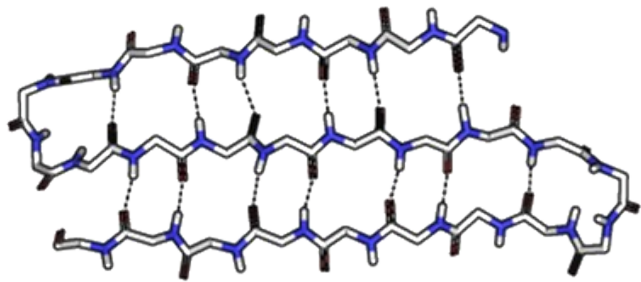
**Primary structure**: sequence of aminoacids (called residues)



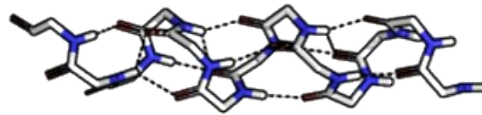Number of residues: ~40 < n < ~35000 (average ~2000)

# What is a protein?

**Secondary structure: Local structure thanks to inter-residue bonds (H-bond, ...)**



β-Sheet (3 strands)



α-helix
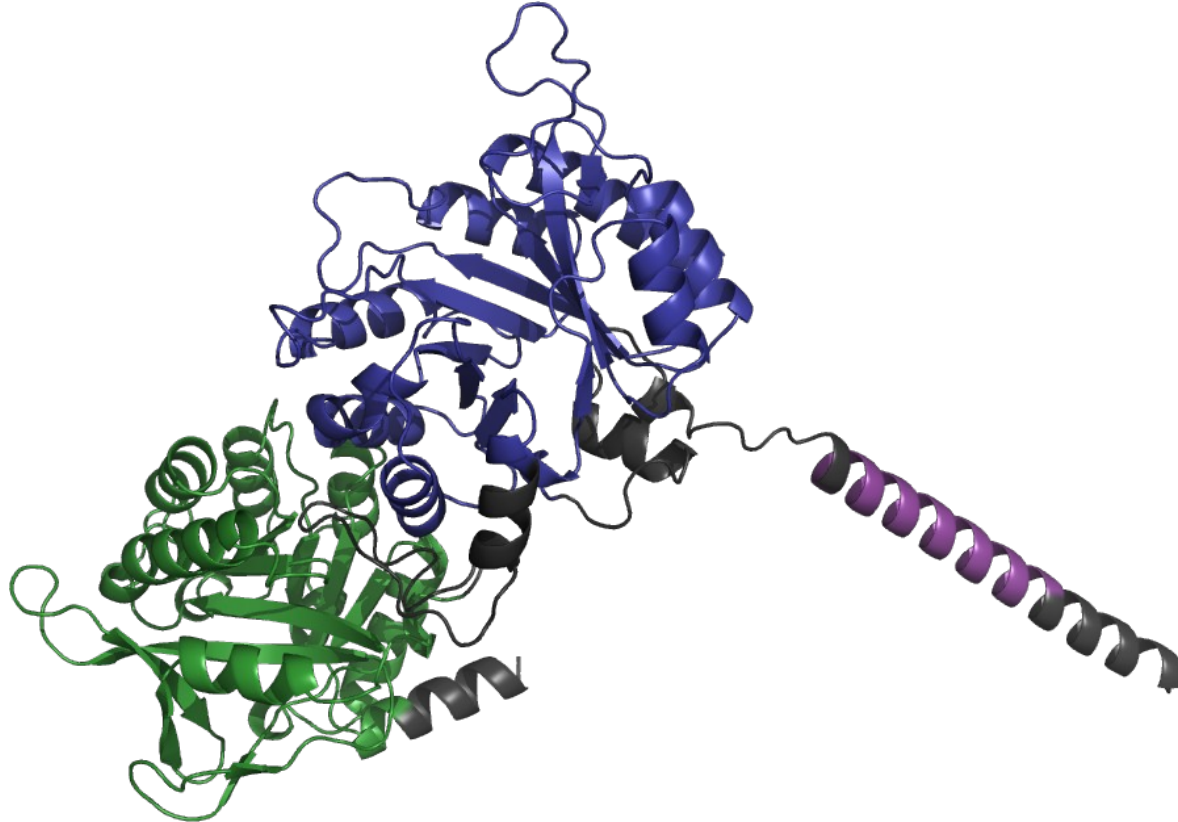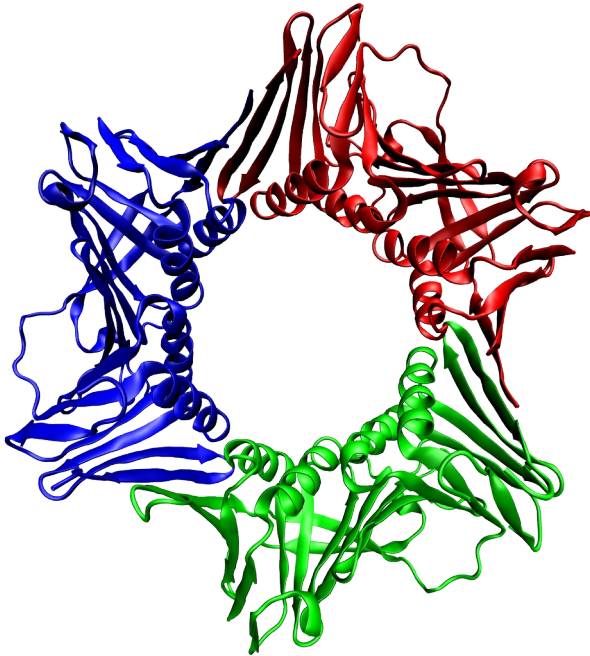


β turn: Type I



Loop

# What is a protein?

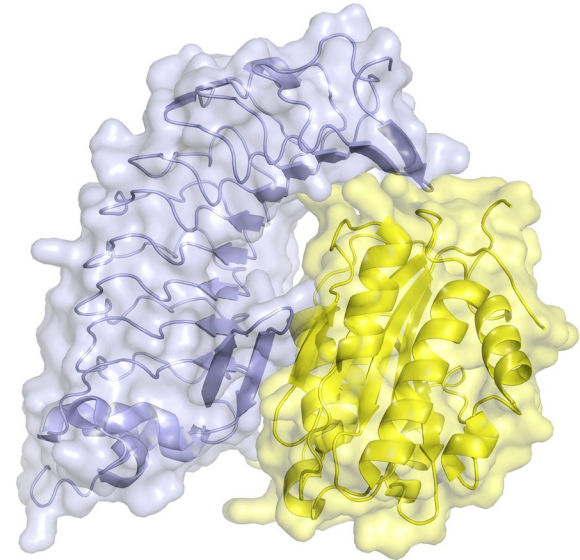**Tertiary structure**: **Global folding of the protein**

# What is a protein?

**Quaternary structure:** Assembly of several protein units

Homo-n-mer (here trimer)

Hetero-n-mer (here dimer)

# What is a protein?

**Summary**



Interaction of AA gives secondary structures

Amino acid (1)  Amino acid (2)

N-terminus  C-terminus

Peptide bond

Dipeptide

Water

β-Sheet (3 strands)  α-helix

Folding of the secondary structure

Association of monomers (same or different)
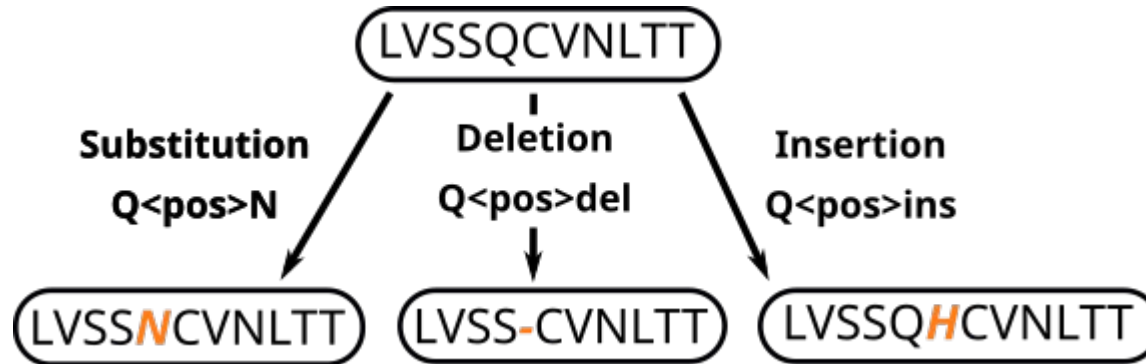
# Protein structure/function

- The structure of the protein gives its function.

- Mutations can occur in the primary sequence as long as the structure does not change too much as to break the function.

- Some amino acids are important for chemical reactions in active sites and might be difficult to replace without breaking the function.
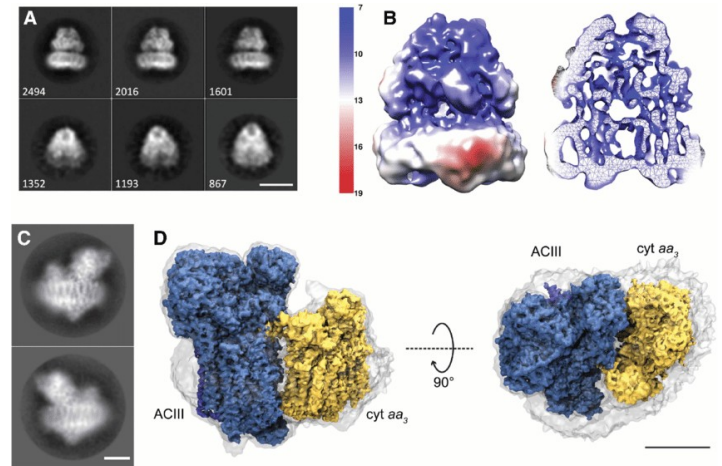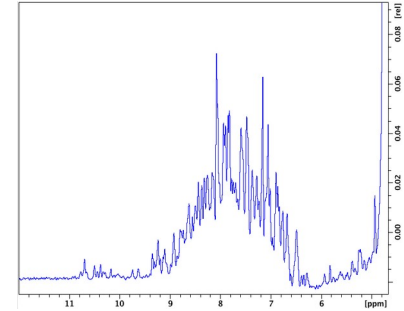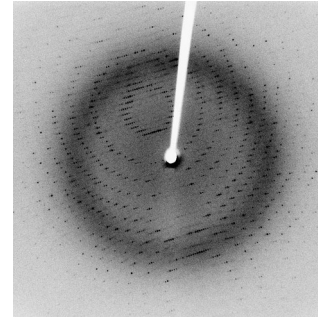
# Mutations

- Several types of mutations appear

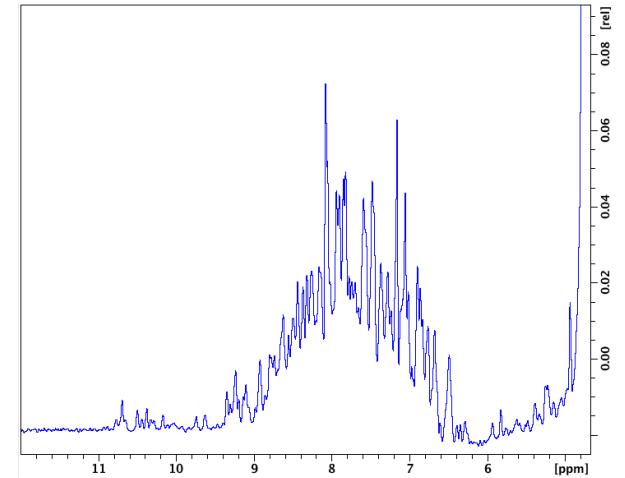# Finding protein structures: experimental methods

- Protein sequencing (no structure)

- X-ray crystallography

- NRM

- Cryo-electron microscopy

# Finding protein structures: experimental methods

- Each method has strength and drawbacks.

- It might require a long time (and money) to get the structure.

# Numerical methods

- Different groups of methods are available

    – Molecular dynamics

    – Conformational sampling

    – Comparative modeling

    – Fold recognition and threading

    – ...

# Comparative modeling

- Search homologous proteins (template): eg different species. The structure of templates are known

- Align the sequences to get information about:
  - Conserved secondary structures
  - Aminoacids that are mandatory to keep the function
  - ...

# CASP competition

- Critical Assessment of protein Structure Prediction

- Every 2 years since 1994

- Unknown protein structures resolved experimentally then compared to numerical models
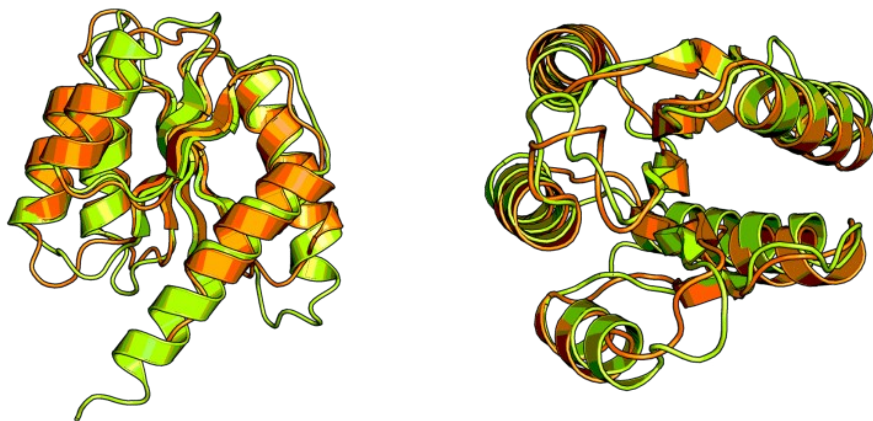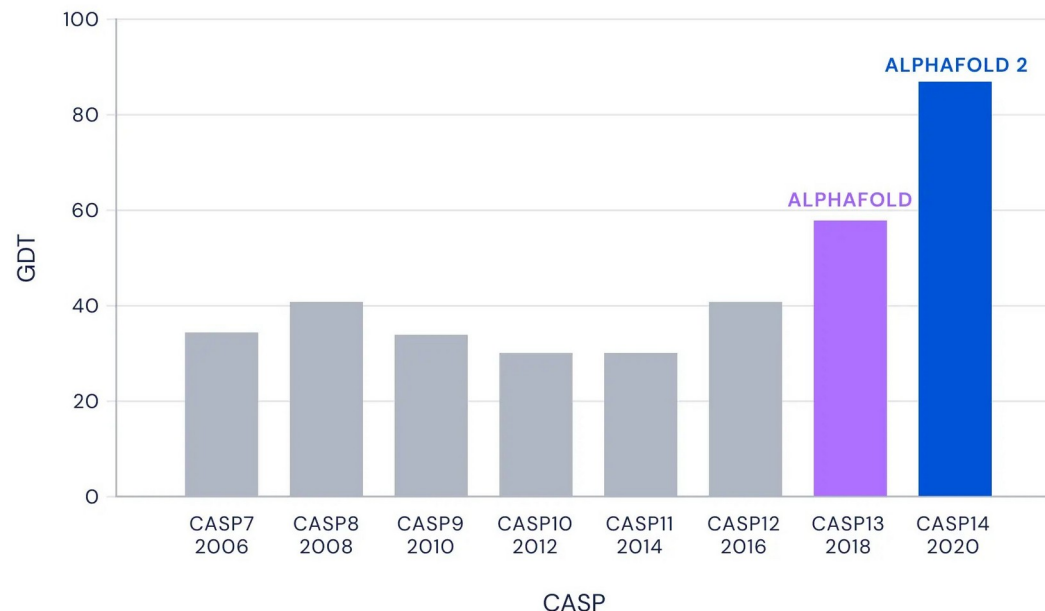
# Performance evaluation

- ## Global Distance Test (GDT)
  - ### Derived from distance of alpha carbon from target



Median Free–Modelling Accuracy

Homology model of target protein A

Experimental structure of protein homologous to protein A

# Alphafold[a] "hardware"

- Written with Tensorflow 2 + JAX

- Runs on GPU, TPU, CPU

- Depends on tools for sequence alignment
  - HH-suite (hhblits, hhsearch, ...)
  - hmmer-suite (jackhmmer)

- Part of the dataset was self-distilled[b] (noisy student)

a) Jumper et al., "Highly accurate protein structure prediction with AlphaFold," Nature, vol. 596, no. 7873, Art. no. 7873, 2021
b) Xie et al. Self-training with noisy student improves imagenet classification. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 10687–10698, 2020.

# Alphafold model

# MSA representation

- Search genetic databases for sequences
  - MGnify (metagenomics)
  - UniRef90 (protein clusters from UniProt)
  - Uniclust30 + BFD (protein clusters from various databases by Soeding lab)

- Now database with around 214M proteins → 23TB and 600M files

# Sequence Alignment

- Comparison of several related proteins (eg. different species)

```
------D-PGDF--DRNVPRICGVCGDRATGFHFNAMTCEGCKGFFRRSMKRKA--LFTCP-FNGDCRITKDNRRHCQACRLKRCVDIGMMKEFILTD
IRPQKRK-KGPAP-KMLGNELCSVCGDKASGFHYNVLSCEGCKGFFRRSVIKGA--HYICH-SGGHCPMDTYMRRKCQECRLRKCRQAGMREECVLSE
SVPGKPS-VNADE-EVGGPQICRVCGDKATGYHFNVMTCEGCKGFFRRAMKRNA--RLRCPFRKGACEITRKTRRQCQACRLRKCLESGMKKEMIMSD
EPERKRK-KGPAP-KMLGHELCRVCGDKASGFHYNVLSCEGCKGFFRRSVVRGGARRYACR-GGGTCQMDAFMRRKCQQCRLRKCKEAGMREQCVLSE
PVTKKPRMGASAG-RIKGDELCVVCGDRASGYHYNALTCEGCKGFFRRSITKNA--VYKCK-NGGNCVMDMYMRRKCQECRLRKCKEMGMLAECMYTG
QTEEKKC-KGYIPSYLDKDELCVVCGDKATGYHYRCITCEGCKGFFRRTIQKNLHPSYSCK-YEGKCVIDKVTRNQCQECRFKKCIYVGMATDLVLDD
----SPS-PPPPP---RVYKPCFVCNDKSSGYHYGVSSCEGCKGFFRRSIQKNM--VYTCH-RDKNCIINKVTRNRCQYCRLQKCFEVGMSKEAVRND
----PPS-PLPPP---RVYKPCFVCQDKSSGYHYGVSACEGCKGFFRRSIQKNM--IYTCH-RDKNCVINKVTRNRCQYCRLQKCFEVGMSKESVRND
----PPS-PPPLP---RIYKPCFVCQDKSSGYHYGVSACEGCKGFFRRSIQKNM--VYTCH-RDKNCIINKVTRNRCQYCRLQKCFEVGMSKESVRND
```
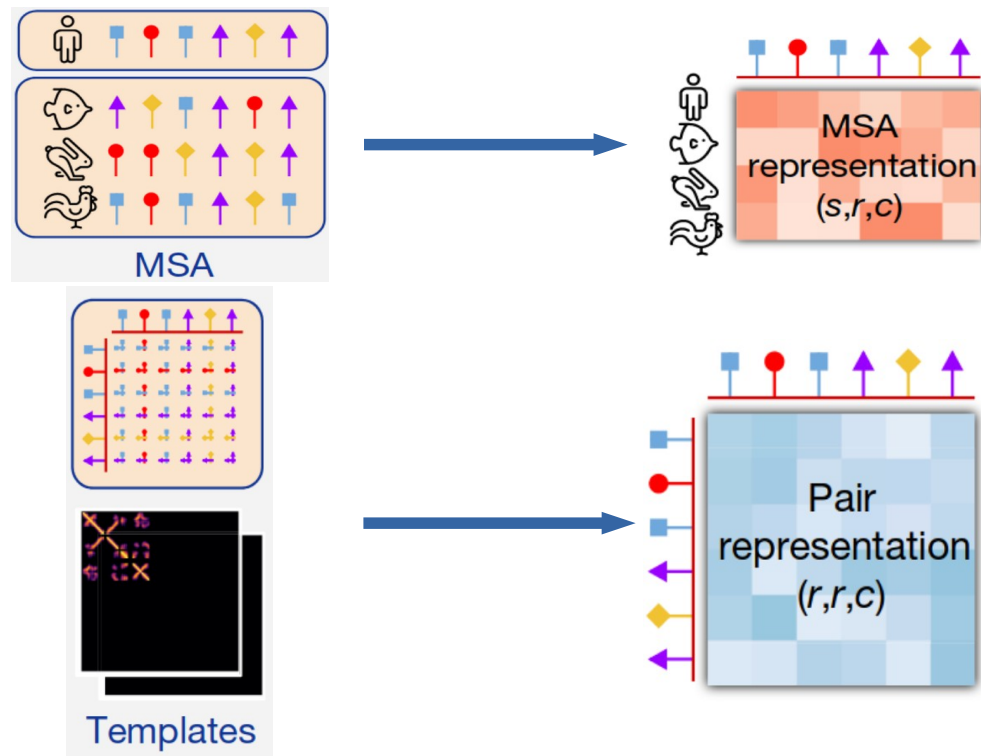
Conserved    The sequence is identical

Semi-Cons.   Some mutations are possible

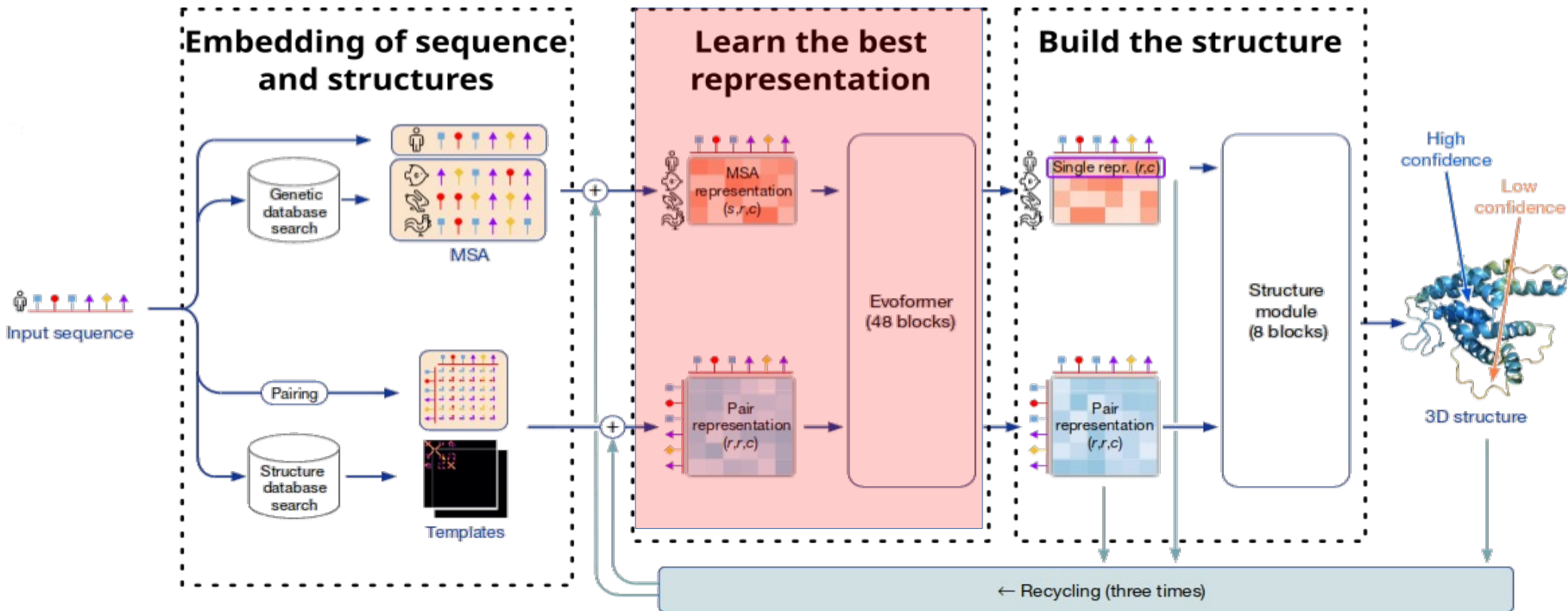**MSA: MultiSequence Alignement**

# Ingredients

- We need 2 ingredients
  - The similar sequences aligned with the input
  - Some structures close enough to serve as template
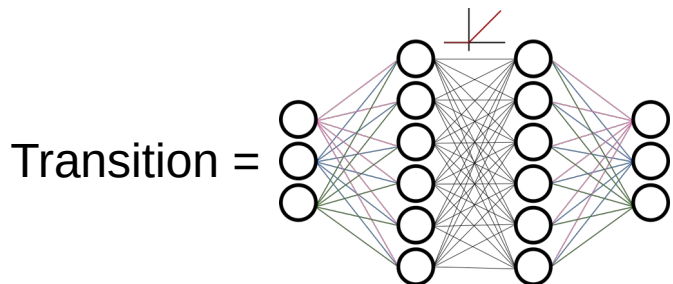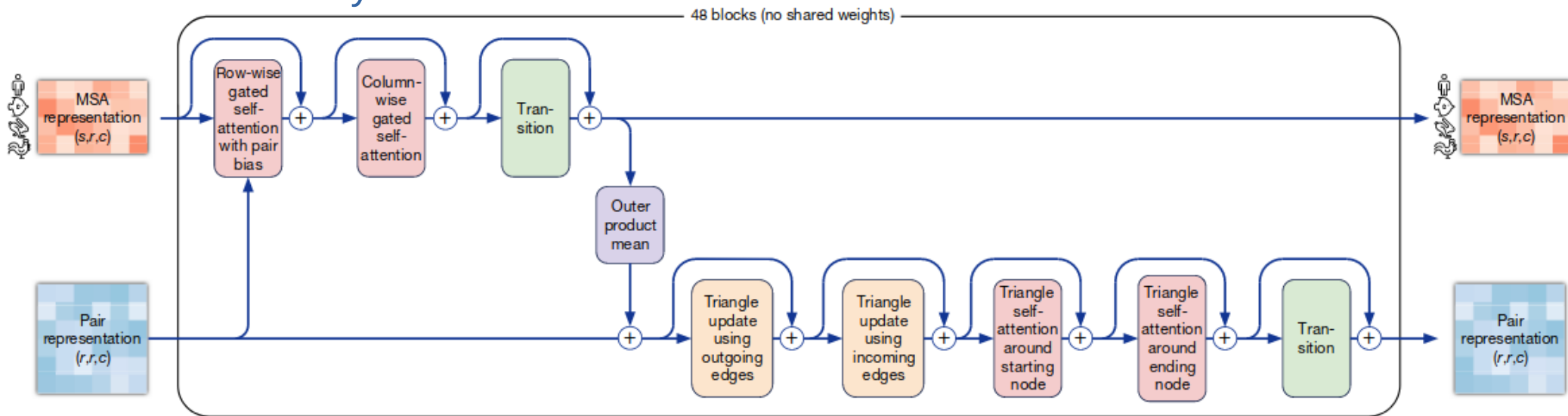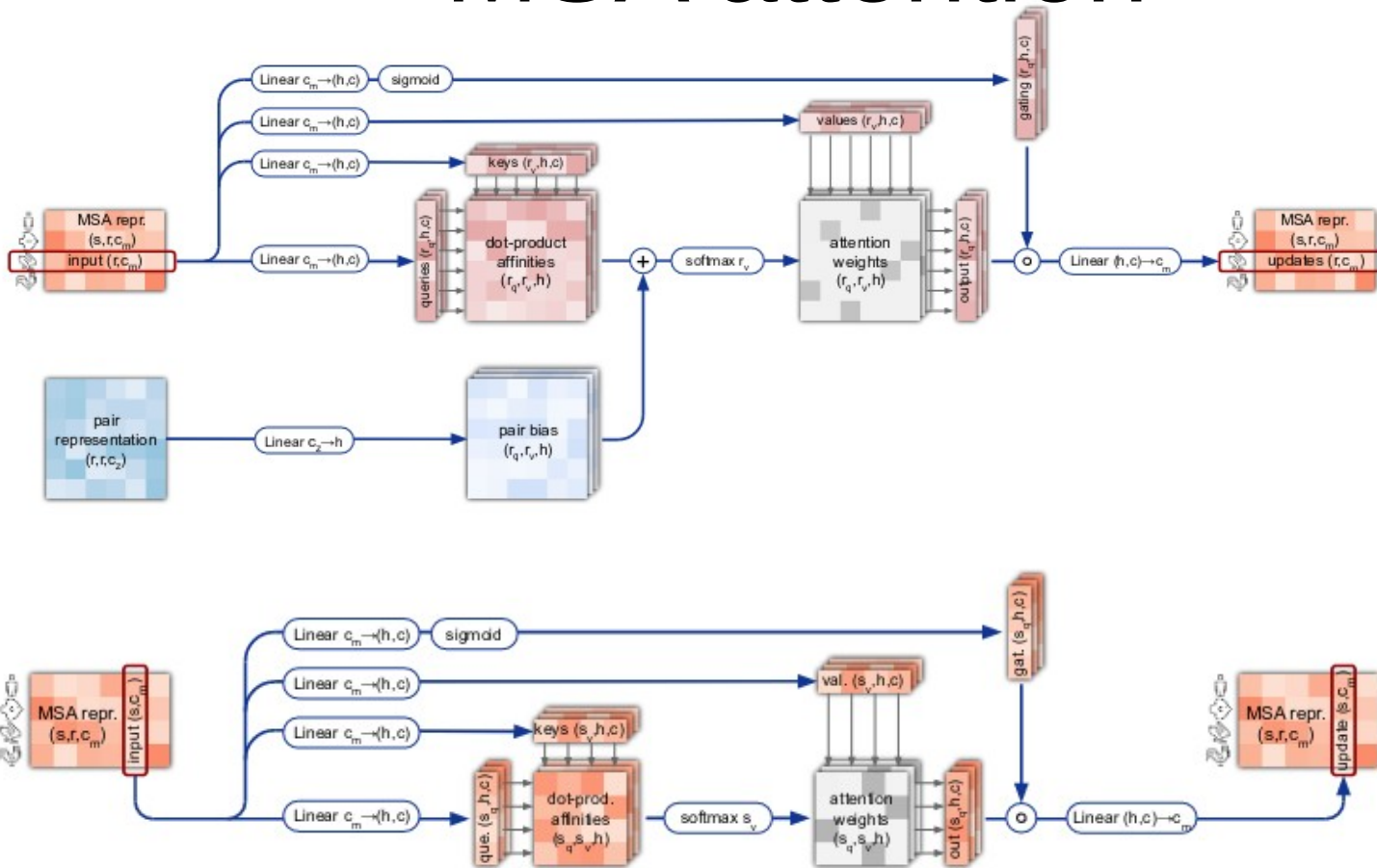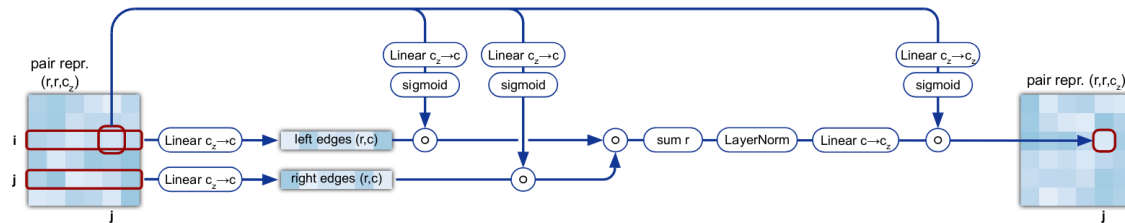- Input for the Evoformer blocks

# Alphafold model

# Evoformer

- Evolutionary Transformer



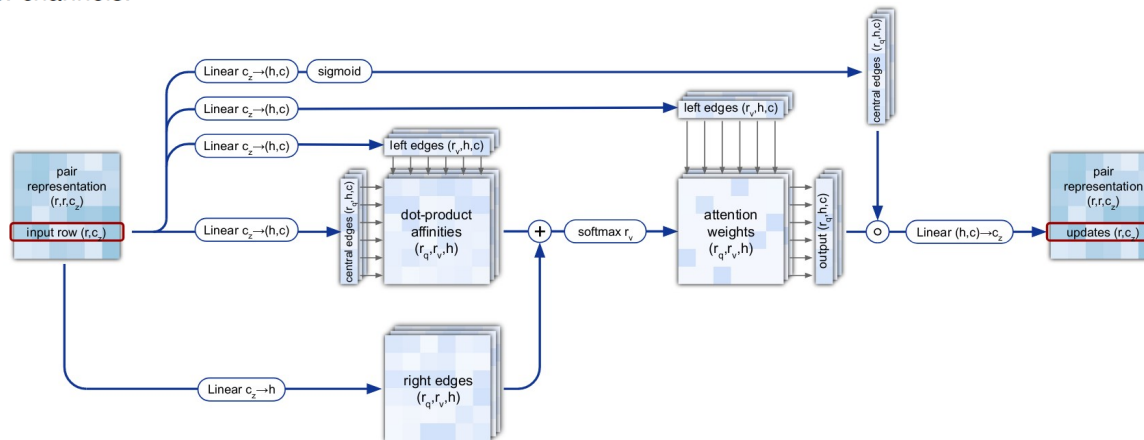Transition =

# MSA attention
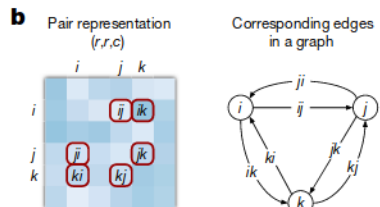
# Pair representation update

Similar for incoming edge

**Supplementary Figure 6** | Triangular multiplicative update using "outgoing" edges. Dimensions: r: residues, c: channels.
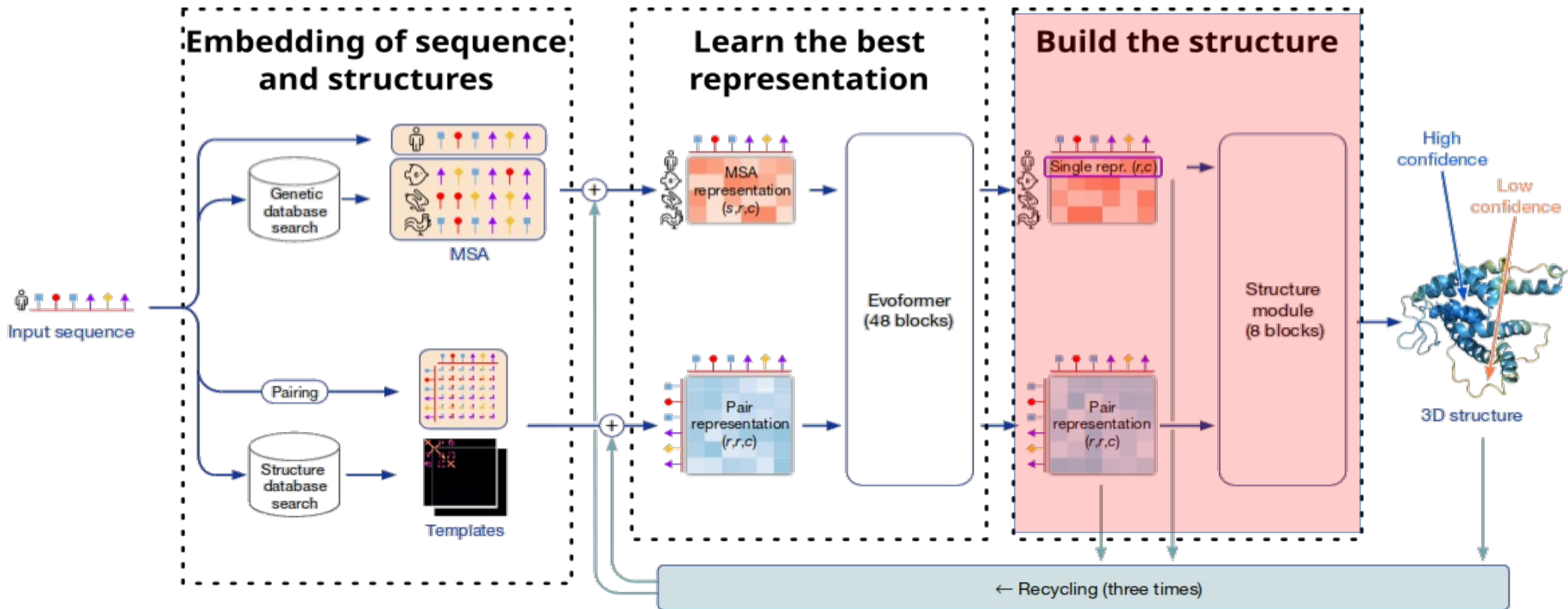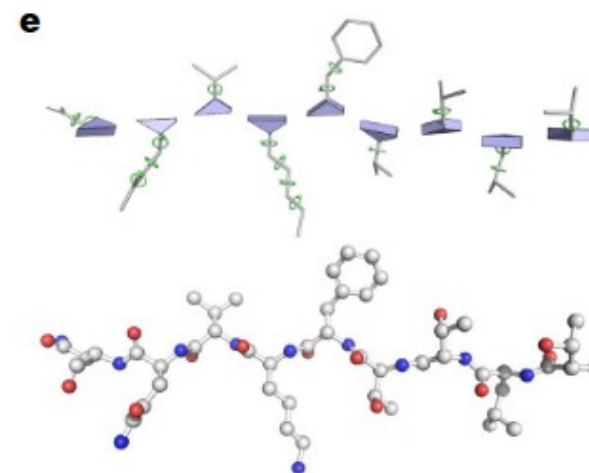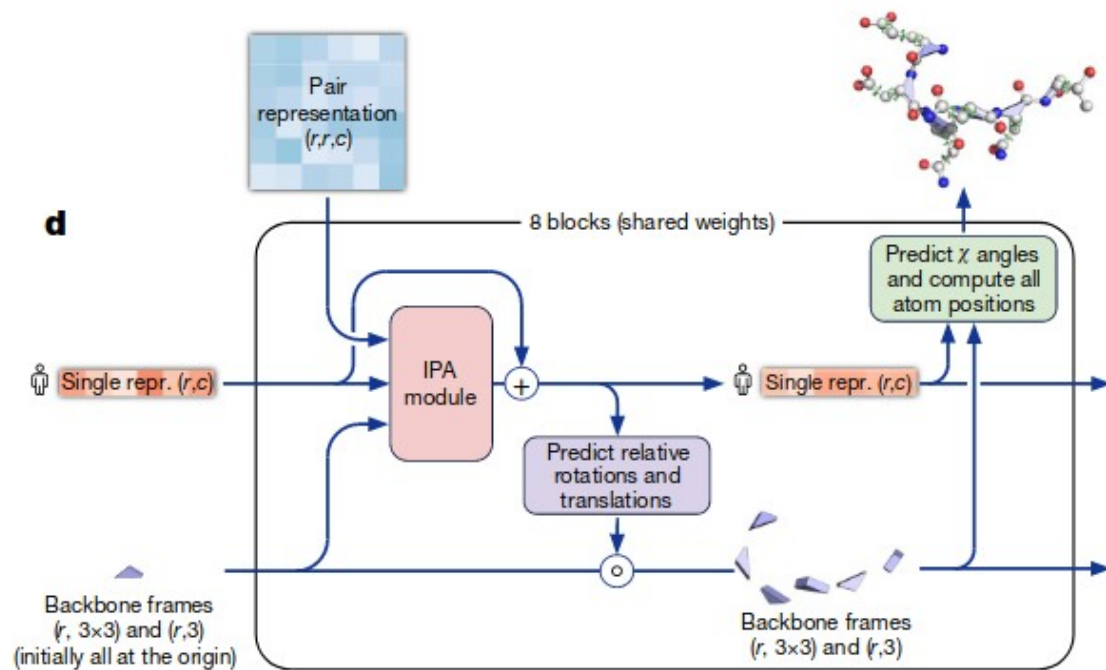


Similar for ending node

**Supplementary Figure 7** | Triangular self-attention around starting node. Dimensions: r: residues, c: channels, h: heads
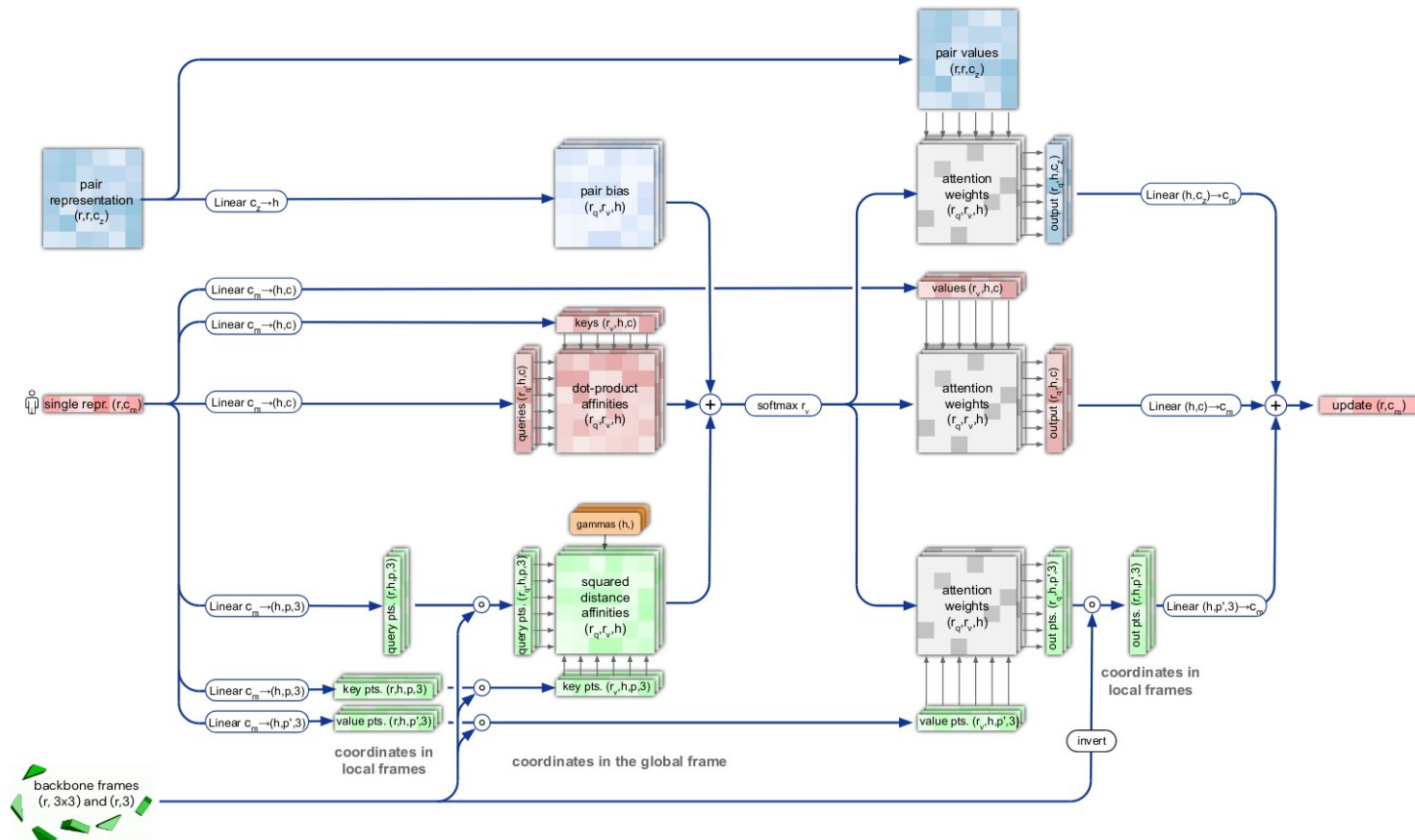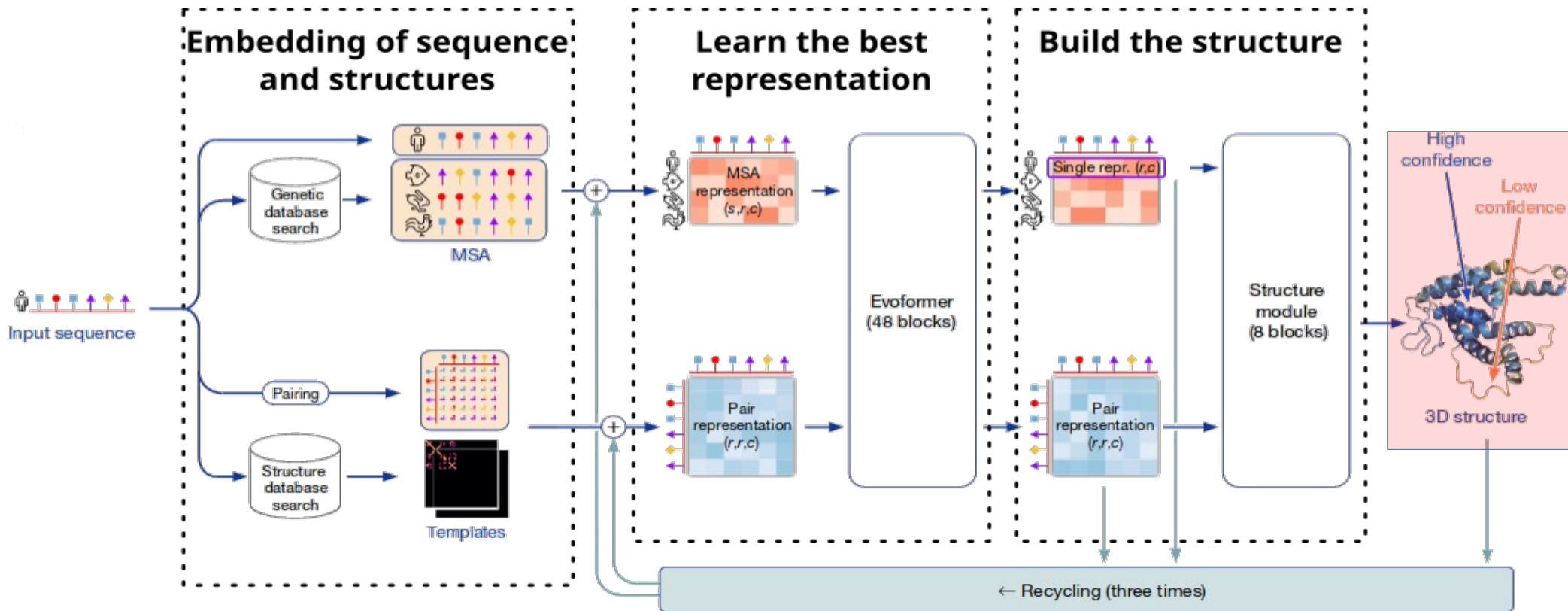
# Alphafold model

# Structure block

# Alphafold model

# The loss function

- Introduced FAPE (Frame Aligned Point Error)

**Algorithm 28** Compute the Frame aligned point error

$\textbf{def } \text{computeFAPE}(\{T_i\}, \{\vec{x}_j\}, \{T_i^{\text{true}}\}, \{\vec{x}_j^{\text{true}}\}, Z = 10\text{Å}, d_{\text{clamp}} = 10\text{Å}, \epsilon = 10^{-4}\text{Å}^2):$

$$T_i, T_i^{\text{true}} \in (\mathbb{R}^{3 \times 3}, \mathbb{R}^3)$$
$$\vec{x}_j, \vec{x}_j^{\text{true}} \in \mathbb{R}^3,$$
$$i \in \{1, ..., N_{\text{frames}}\}, j \in \{1, ..., N_{\text{atoms}}\}$$

1: $\vec{x}_{ij} = T_i^{-1} \circ \vec{x}_j$ $\qquad\qquad \vec{x}_{ij} \in \mathbb{R}^3$

2: $\vec{x}_{ij}^{\text{true}} = T_i^{\text{true}-1} \circ \vec{x}_j^{\text{true}}$ $\qquad\qquad \vec{x}_{ij}^{\text{true}} \in \mathbb{R}^3$

3: $d_{ij} = \sqrt{\|\vec{x}_{ij} - \vec{x}_{ij}^{\text{true}}\|^2 + \epsilon}$ $\qquad\qquad d_{ij} \in \mathbb{R}$

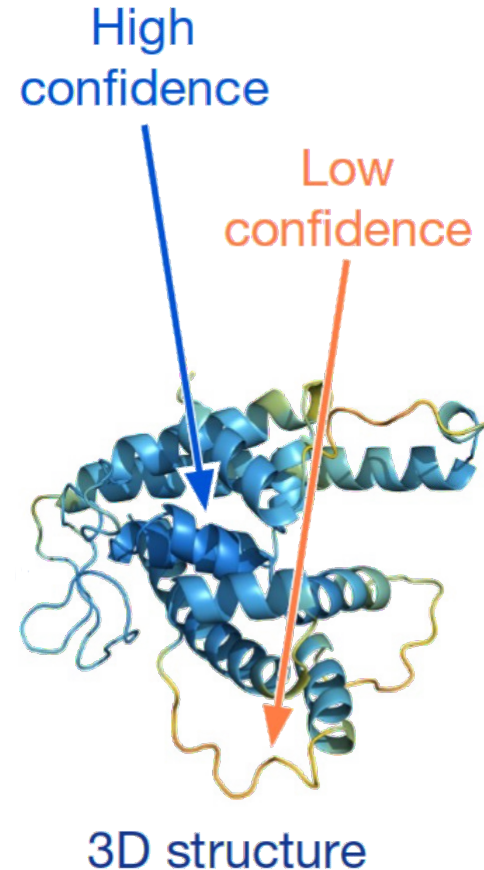4: $\mathcal{L}_{\text{FAPE}} = \frac{1}{Z} \text{mean}_{i,j}(\text{minimum}(d_{\text{clamp}}, d_{ij}))$

5: **return** $\mathcal{L}_{\text{FAPE}}$

# Inference: the structures

- At the end:
  - 3D structures
  - "confidence" score for each residue

- Refinement with parametrized physics software possible (OpenMM with Amber Force Field)



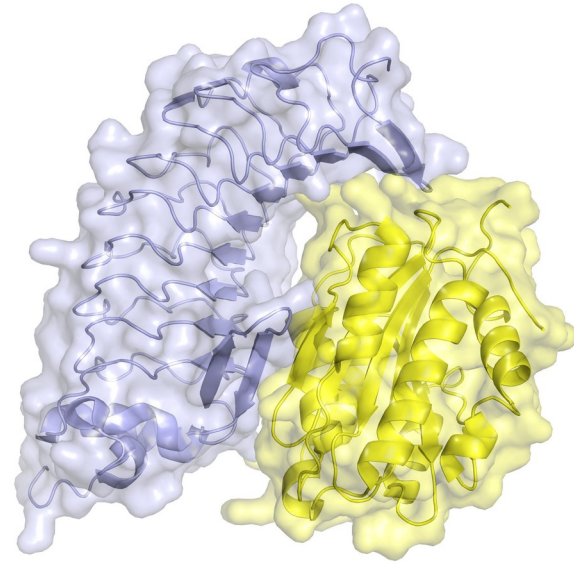High confidence

Low confidence

3D structure

# Features

- Monomer
- Multimer

# Other software

- RoseTTA

- Openfold

- Colabfold (uses alphafold model but MSA is done with Mmseqs)

- OmegaFold (pytorch port of Alphafold)

- ESMFold (based on OpenFold but no MSA needed)

- ...

# Pictures Attribution

(1) Thomas Shafee, CC BY 4.0 via Wikimedia Commons (secondary and tertiary structures of proteins)

(2) Muskid, CC BY-SA 3.0 via Wikimedia Commons (beta turn secondary structure)

(3) TungstenEinsteinium, CC BY-SA 4.0 via Wikimedia Commons (Table of amino acids)

(4) Simoncaulton, CC BY-SA 4.0 via Wikimedia Commons (Hetero dimer cloating factor)

(5) EMBL-EBI, CC BY 4.0 via Wikimedia Commons (tertiary structure FAM151A)

(6) Humphrey, W., Dalke, A. and Schulten, K., "VMD - Visual Molecular Dynamics", J. Molec. Graphics, 1996, vol. 14, pp. 33-38. (Homotrimer)

(7) VMD was developed by the Theoretical and Computational Biophysics Group in the Beckman Institute for Advanced Science and Technology at the University of Illinois at Urbana-Champaign.

(8) Jeff Dahl, CC BY-SA 3.0 via Wikimedia Commons (X-ray diffraction pattern)

(9) Loteralle, CC BY-SA 3.0 via Wikimedia Commons (NMR spectrum of calmodulin)

(10) Simon, Kailene & Pollock, Naomi & Lee, Sarah. (2018). Membrane protein nanoparticles: The shape of things to come. Biochemical Society Transactions. 46. BST20180139. 10.1042/BST20180139. (Cryoelectron microscopy of AcrB-SMALP)

(11) https://www.deepmind.com/blog/alphafold-a-solution-to-a-50-year-old-grand-challenge-in-biology (GDT CASP)

(12) https://bitesizebio.com/38005/homology-modeling-proteins/ (Thomas Warwick, homology of LytR)

(13) Opabinia regalis - Self-created from PDB ID 1A0S using PyMol, CC BY-SA 3.0 (Sucrose Porin)