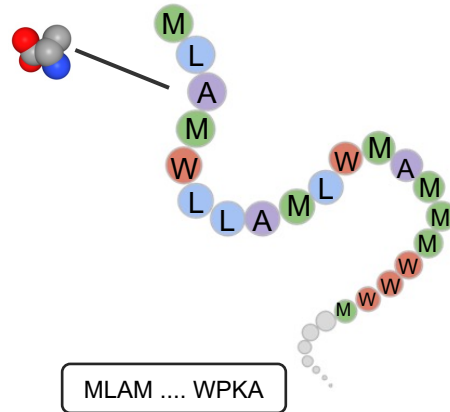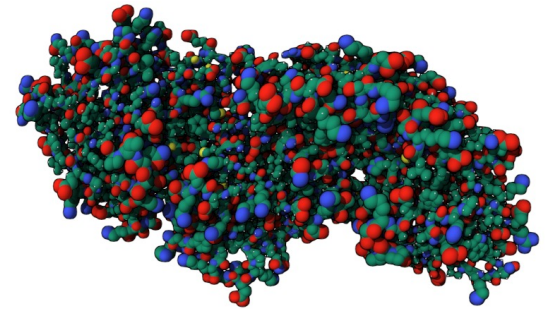# Determining protein 3D structures

# Proteins

- Proteins consist of amino acid sequences that fold into 3D structures
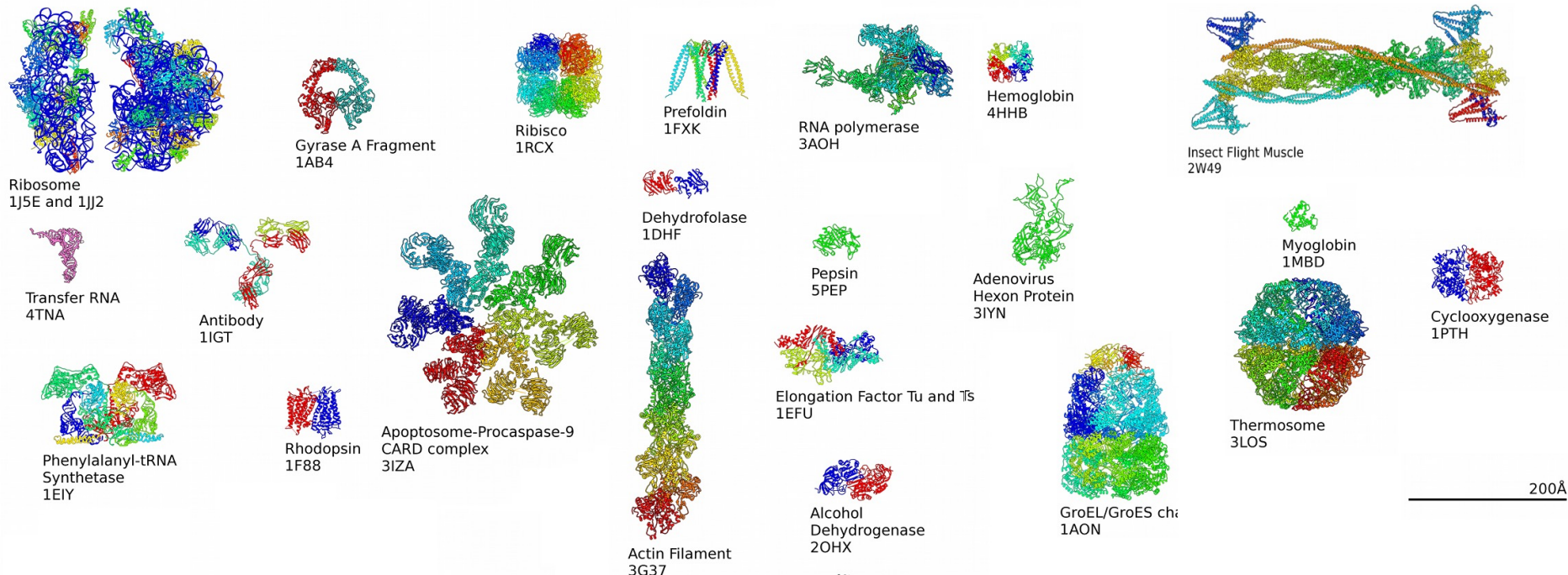
20 different amino acids; each amino acid is a small molecule

MLAM .... WPKA

Folding into
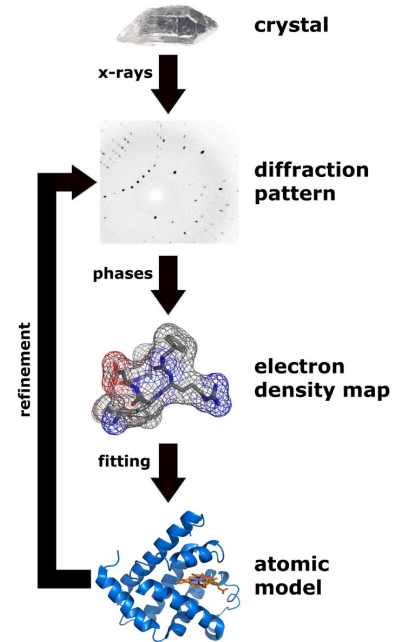
3D structure

# Protein structure space



Ribosome
1J5E and 1JJ2

Gyrase A Fragment
1AB4

Ribisco
1RCX

Prefoldin
1FXK

RNA polymerase
3AOH

Hemoglobin
4HHB

Insect Flight Muscle
2W49

Dehydrofolase
1DHF

Transfer RNA
4TNA

Antibody
1IGT

Pepsin
5PEP

Adenovirus
Hexon Protein
3IYN

Myoglobin
1MBD

Cyclooxygenase
1PTH

Phenylalanyl-tRNA
Synthetase
1EIY

Rhodopsin
1F88

Apoptosome-Procaspase-9
CARD complex
3IZA

Elongation Factor Tu and Ts
1EFU

Thermosome
3LOS

Actin Filament
3G37

Alcohol
Dehydrogenase
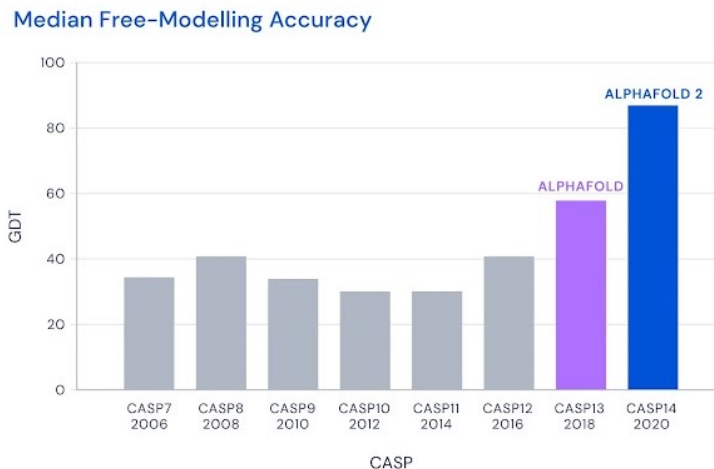2OHX

GroEL/GroES cha
1AON

200Å

# Experimentally determining 3D structures

- Nuclear Magnetic Resonance (NMR) Spectroscopy
- Cryogenic Electron Microscopy (Cryo-EM)
- Electron Crystallography
- X-ray Crystallography
  - 1. Protein Purification
  - 2. Protein Crystallization
  - 3. X-ray Diffraction
  - 4. Electron Density Maps
  - 5. Atom Model
  - 6. Refinements
- For novel proteins, this process typically takes several months or even up to over a year
- Protein Data Bank (PDB): stores ~180,000 experimentally determined protein 3D structures (~600M proteins are sequenced)



crystal

x-rays

diffraction pattern

phases

refinement

electron density map

fitting

atomic model

hhu.de

# Predicting protein 3D structures

- Predicting protein 3D structure from protein sequence is desirable
  - Was not possible (with high accuracy) until 2020:

**Median Free-Modelling Accuracy**



- Deep Learning Models that can predict protein 3D structures
  - AlphaFold 2
  - ESMFold
  - RoseTTAFold

hhu.de

# AlphaFold2 - Input



Multiple Sequence Alignment (MSA)

embedding

$\in (N_{seq}, N_{res}, c_m)$

co-evolution

Structure templates

embedding

$(i, j)$

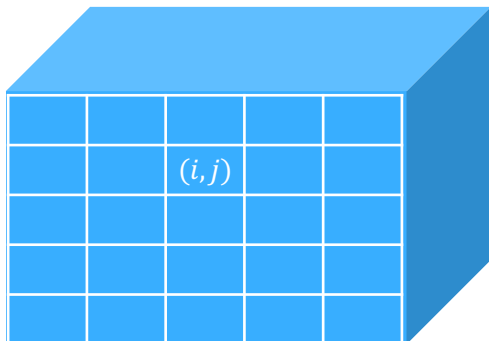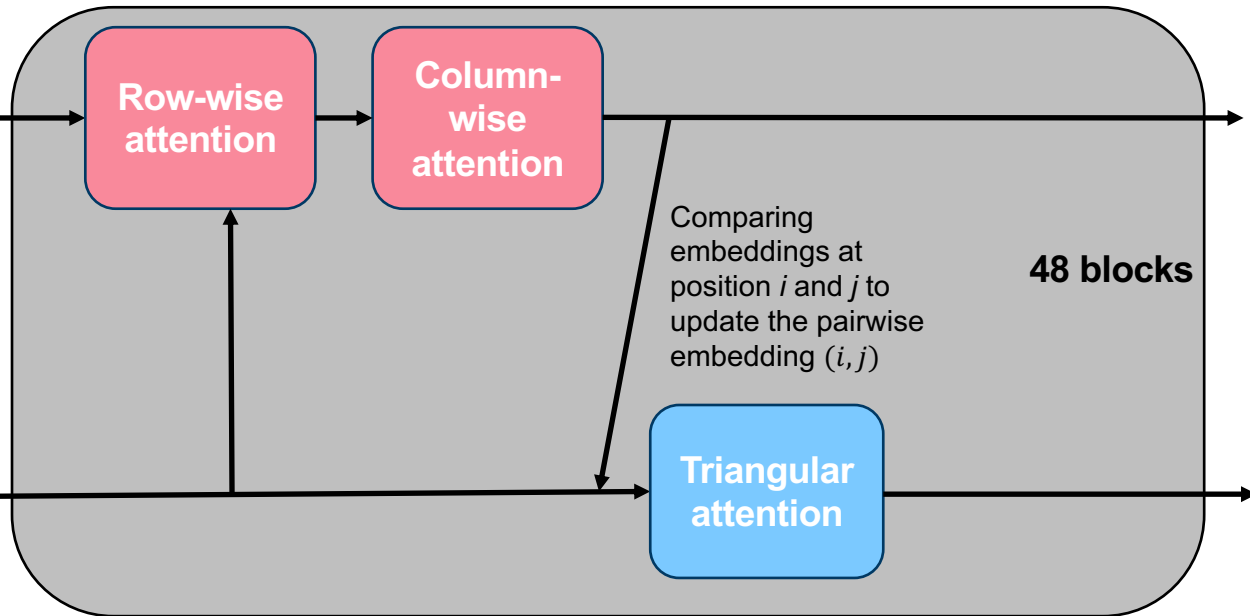$\in (N_{res}, N_{res}, c_r)$

MSA

MLAM .... WPKA

# AlphaFold2 - Evoformer

**MSA**



$\in (N_{seq}, N_{res}, c_m)$

**Pair representations**

$(i,j)$

$\in (N_{res}, N_{res}, c_r)$

**Evoformer**

**Row-wise attention**

**Column-wise attention**

Comparing embeddings at position $i$ and $j$ to update the pairwise embedding $(i,j)$
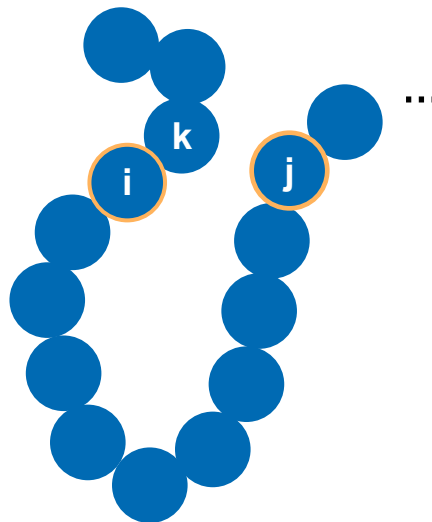
**48 blocks**

**Triangular attention**
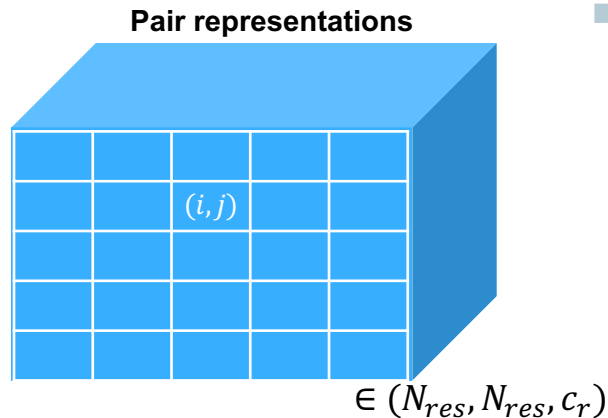
- Triangular attention: Node rep. $(i,j)$ is updated using all edges between $i$ & $k$ and $j$ & $k$

# AlphaFold2 – Triangular Attention

**Pair representations**



$\in (N_{res}, N_{res}, c_r)$

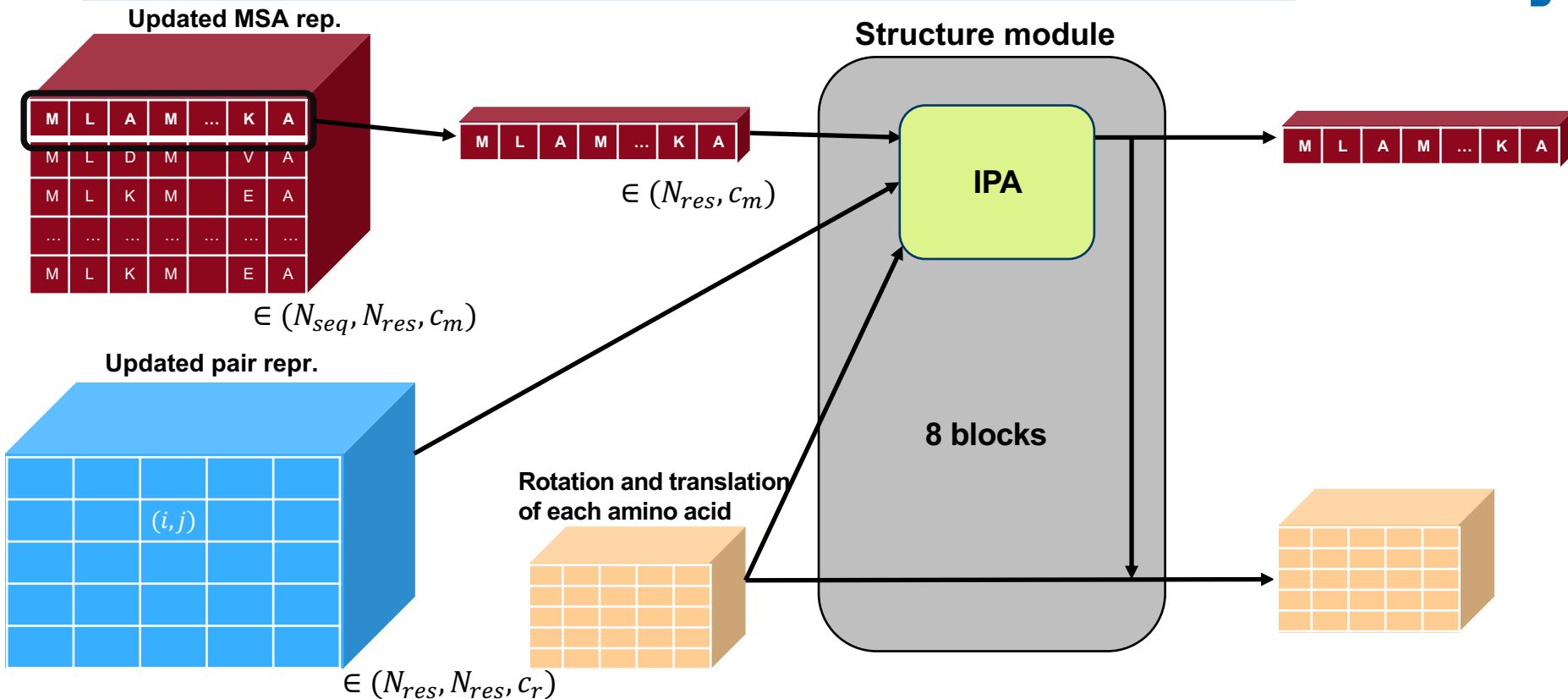- Triangular attention: Node rep. $(i,j)$ is updated using all edges between $i$ & $k$ and $j$ & $k$

  - Integrating Euclidean geometry into the network:
    Knowing distance between $i$ & $k$ and between $j$ & $k$ put a strong constraint of distance between $i$ & $j$

# AlphaFold2 – Structure module



**Updated MSA rep.**

$\in (N_{seq}, N_{res}, c_m)$

$\in (N_{res}, c_m)$

**Structure module**

**IPA**

**8 blocks**

**Updated pair repr.**

$(i, j)$

$\in (N_{res}, N_{res}, c_r)$
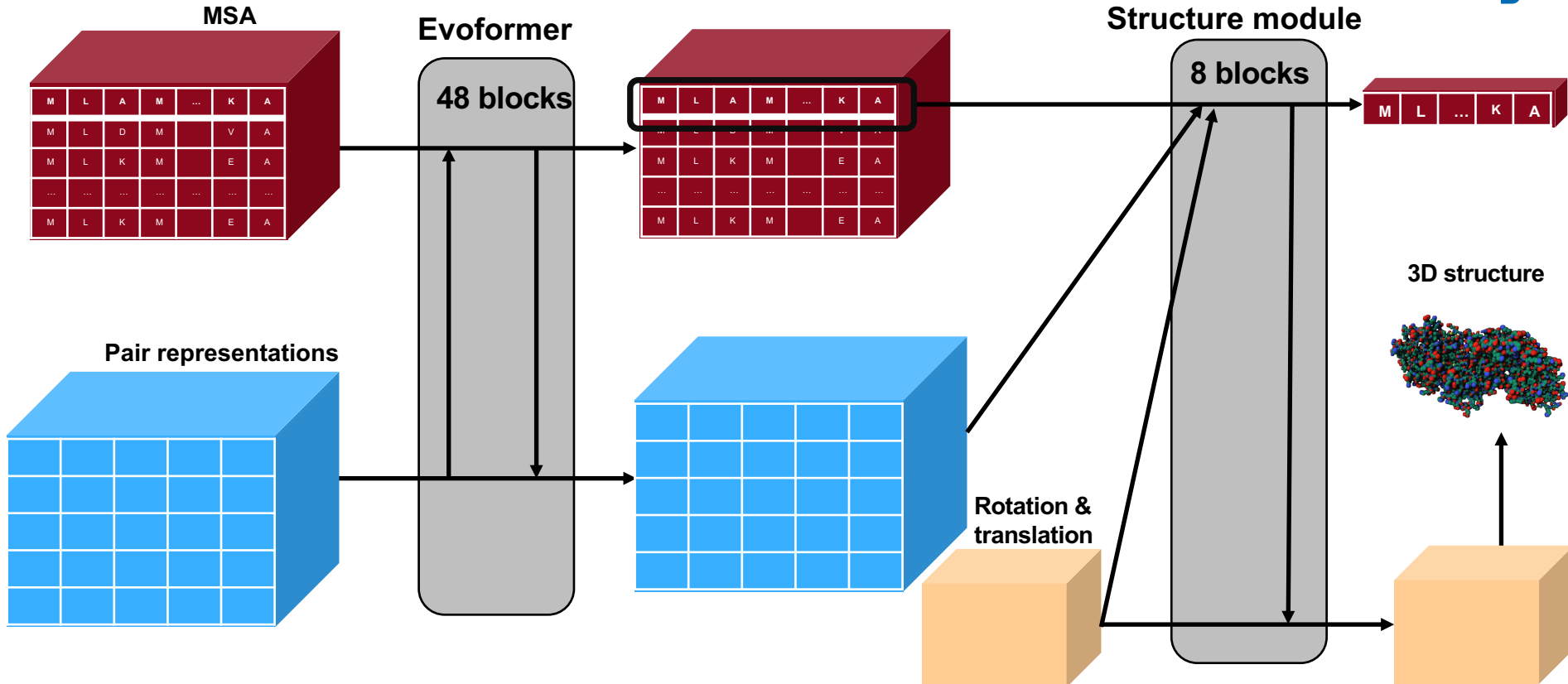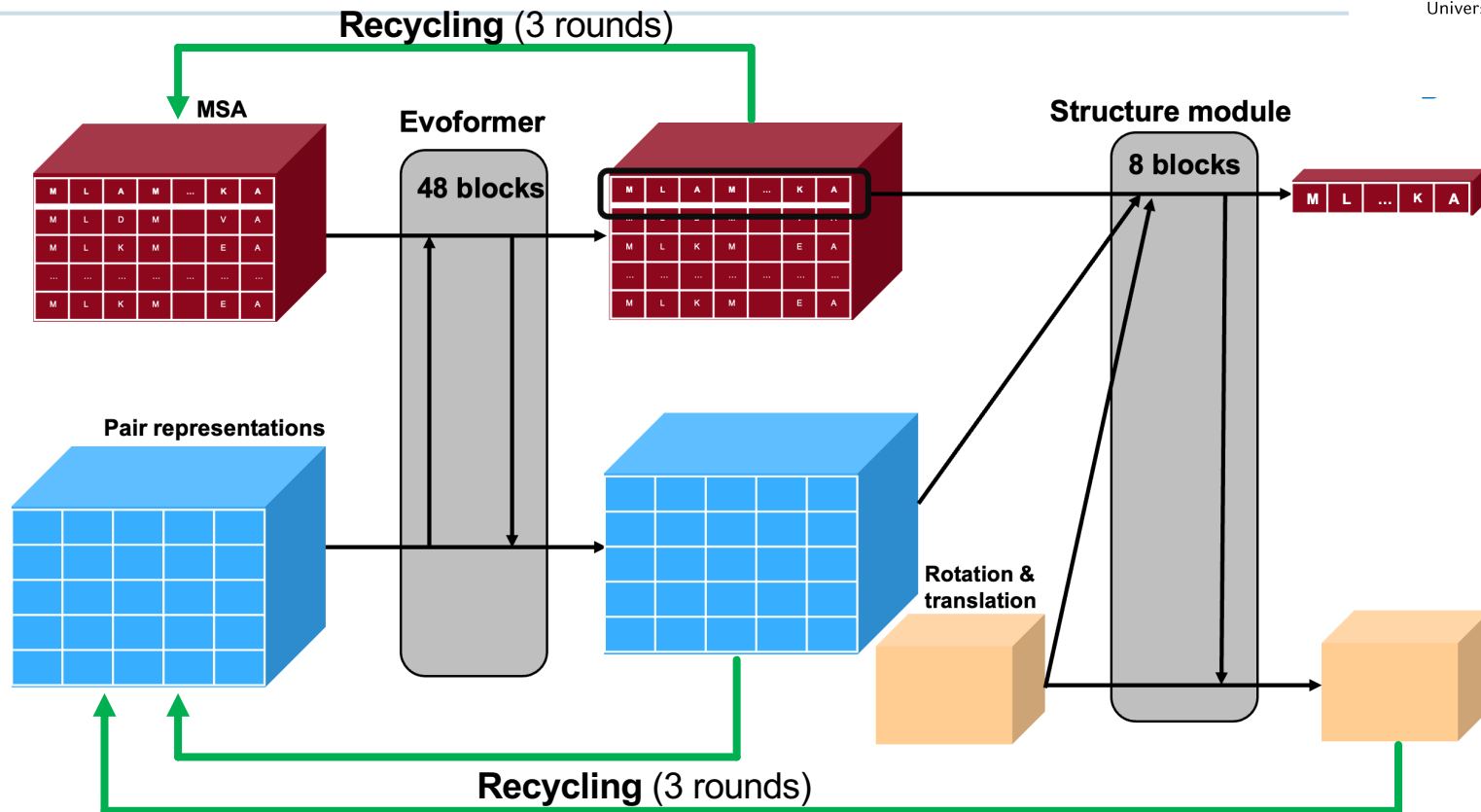
**Rotation and translation of each amino acid**

# AlphaFold2 – Full Model

# AlphaFold2 – Recycling

# AlphaFold2 – Training

- Loss function:
  - Structural loss named FAPE: Similar to RMSD (root mean squared deviation) of atomic positions
  - Auxiliary losses:
    - Distogram loss: Comparing the pairwise distances between amino acids to ground truth
    - MSA masking: Some tokens of the MSA are masked out and are predicted during training
    - …
- Training data:
  - ~120k experimentally determined protein 3D structures from the PDB
  - After first run of training:
    - Predicting structures of ~350k proteins with yet unknown structures
    - New training on experimental and predicted 3D structures

# AlphaFold2 – Capabilities and limitations

- Capabilities
  - Single protein chains
  - Protein multimers
  - Protein-protein complexes
  - AlphaFold2 can often predict protein structure, even if there are no known related protein structures, if a sequence has many related sequences
- Limitations
  - Struggles to predict the structures of protein with few closely related protein sequences
  - Cannot predict effect of point mutations
  - AlphaFold2 does not capture such conformational changes
  - Cannot predict interactions with other molecules, e.g., small molecules