# The Illustrated AlphaFold

Elana Simon, Jake Silberg

# Background and Existing Work

# Protein Folding Challenge

- Proteins are linear chains of amino acids
- Amino acids interact with each other and naturally folds into lowest-energy 3D structure
- This structure determines protein's function
- Applications: drug design, disease research, protein engineering

# Existing Work - AlphaFold 1 (2018):

- DeepMind's first attempt at the protein folding problem
- Built on previous work using evolutionary data to predict protein structure
- Used deep learning to predict distances between amino acid and construct protein structures.
- Significantly outperformed other methods at the time, but still not at experimental accuracy levels
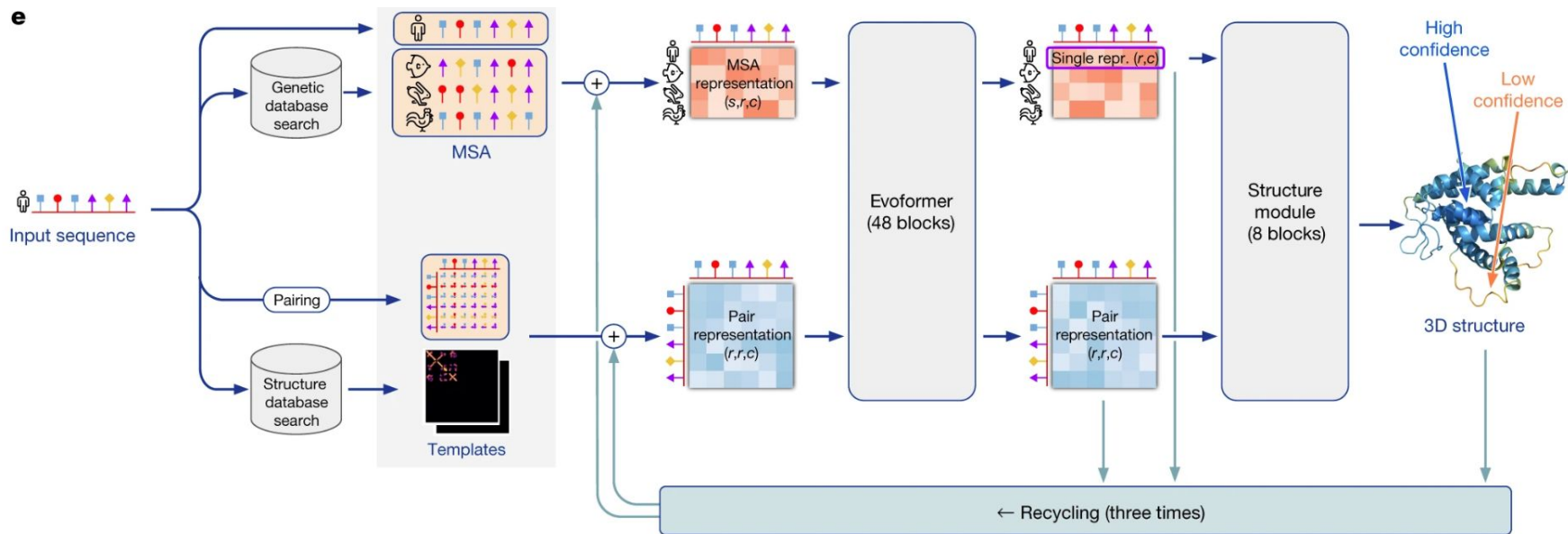
# Existing Work - AlphaFold 2 (2020):

**Machine Learning Framework:**

- Utilized an attention-based neural network architecture called Evoformer
- Incorporated a multiple sequence alignment (MSA) and pair representation

**Performance:**

- Used pattern recognition for structure prediction, with minimal physics-based refinement
- Achieved breakthrough accuracy in protein structure prediction
- Predicted structures for nearly all cataloged proteins known to science

# Existing Work - AlphaFold 2 (2020):

# Molecular Structure Prediction

**Limitation of AlphaFold 2**:

- Focused on single protein structures
- Input: Amino acid sequences only
- Limited to standard amino acids

**Real Biology is Complex:**

- Proteins rarely work alone
- Drug design needs protein-ligand interactions
- DNA/RNA interactions crucial for gene regulation

# AlphaFold 3 (2024):

**Function:**

- Expanded beyond proteins to predict structures and interactions of various biomolecules
- Can model proteins, DNA, RNA, ligands, and their interactions
- Predicts chemical modifications that control cell functioning

**Machine Learning Framework:**

- Enhanced version of the Evoformer module
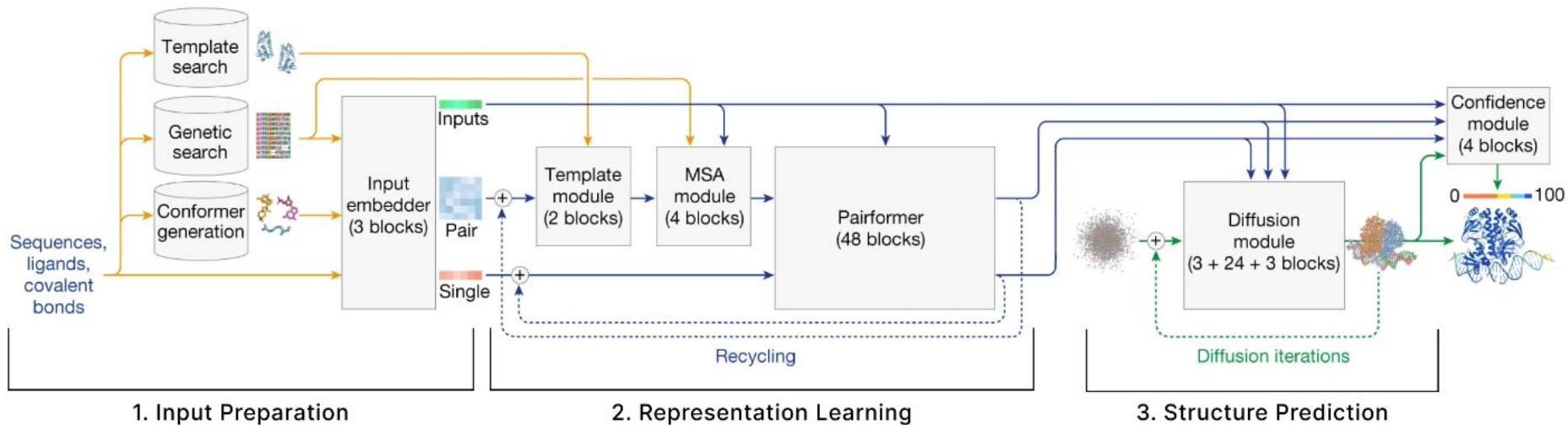- Introduced a diffusion network for structure assembly

**Performance:**

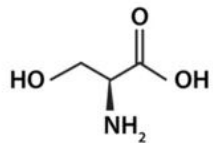- Improved accuracy in predicting molecular interactions by at least 50% compared to existing methods
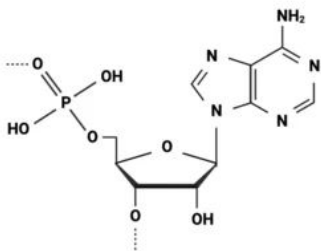
# Key Idea & Designs

# Architecture Overview

# Input Preparation - Two key representations
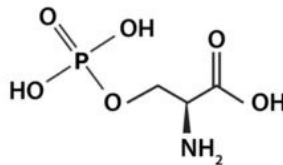


**Standard Amino Acid** (serine)
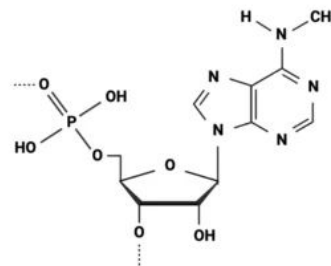7 atoms*
1 tokens

**Standard Nucleotide** (adenosine)
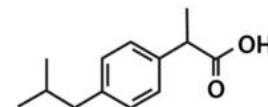23 atoms*
1 token

**Modified Amino Acid** (phosphoserine)
11 atoms*
11 tokens

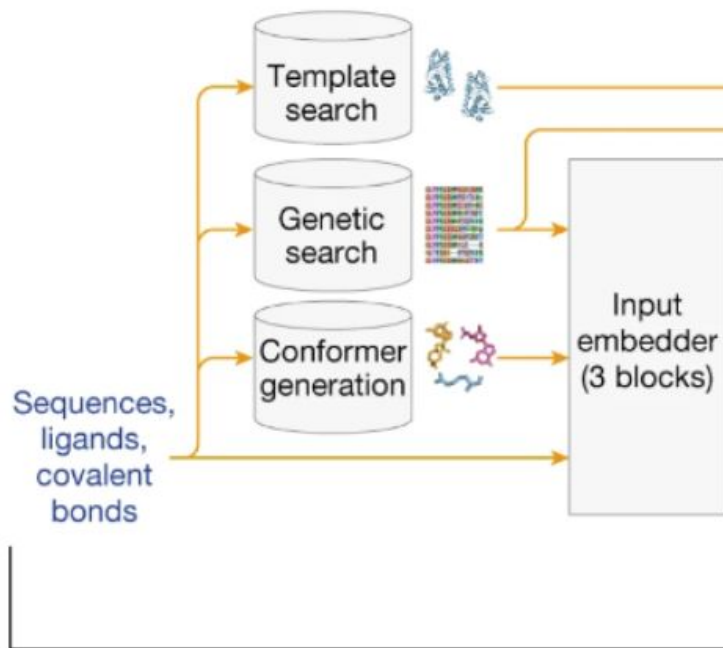**Modified Nucleotide** (methyladenosine)
24 atoms*
24 tokens

**Ligand** (ibuprofen)
15 atoms*
15 tokens

*atoms=heavy atoms

# Input Preparation



1. Input Preparation

**Genetic Search:** finds similar molecular sequence across different species, get Multiple Sequence Alignment (**MSA**) (N_MSA*N_token*C_m)

**Template Search:** finds similar known structures, get template matrix **t** (N_token*N_token*C_t*N_templates)

**Conformers** : local chemical arrangements (initial atom positions, chemical properties), get original atom-level representation matrix **c** (C_atom*N_atom)
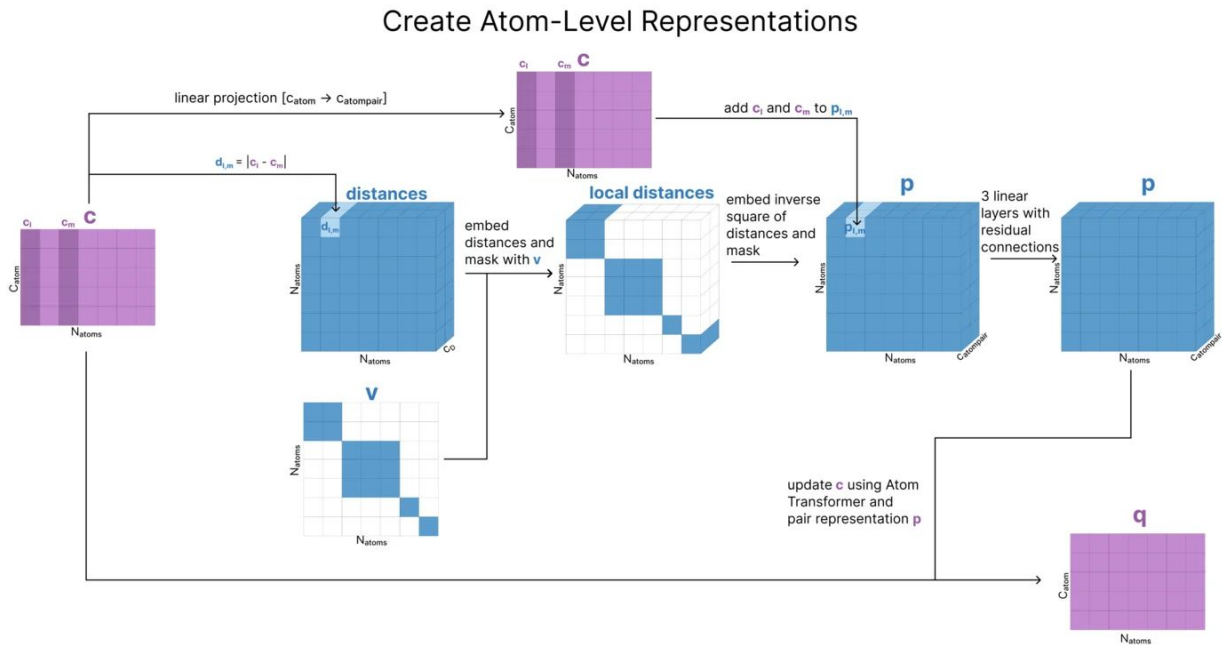
# Create Atom-Level Representations

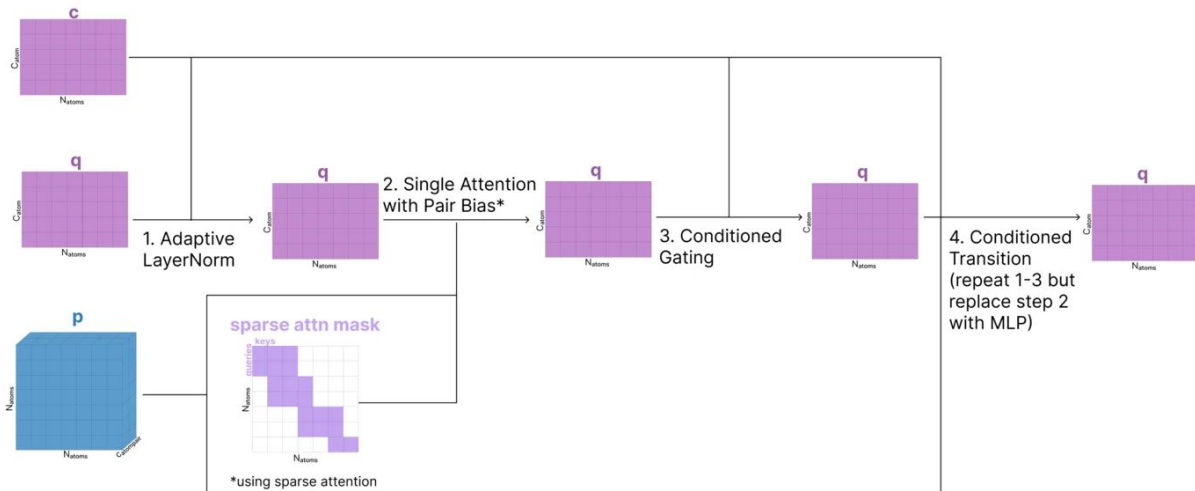**Single Representation (c → q):**

**Matrix c**: stores atom properties

**Matrix q**: copy of c that will be updated during processing

**Pair Representation (p)**: Stores distances between atoms within each token



Create Atom-Level Representations

# Update Atom-Level Representations (Atom Transformer)



Overview of Atom Transformer

**Adaptive LayerNorm**: uses input c to dynamically generate normalization parameters for q

**Attention with Pair Bias:** enhances standard self-attention by using pair representation as bias

**Conditioned Gating:** Controls information flow using gates generated from original atom representation c

**Conditioned Transition**: Modified MLP layer using SwiGLU activation, sandwiched between AdaNorm and gating, both conditioned on c

# Aggregate Atom-Level → Token-Level

**Token-Level Single Representation (s):** atoms (q) → projection → averaging → feature addition → projection → s_init (s)

**Token-Level Pair Representation (z):** s_init (s) → projection (c_token → c_z ) → z_ij = si, sj→ add r.p.e → add bonds → z_init
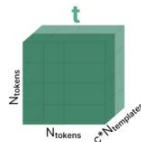
Set aside the **atom-level** representations (**c, q, p**) and focus on updating our token-level representations **s** and **z** in the next section (with the help of m and t).

Information about related sequences and their structures

Multiple Sequence Alignment

**m**
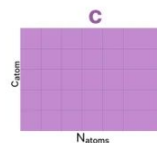
$N_{MSA}$   $N_{tokens}$   $C_m$

Structure Templates

**t**

$N_{tokens}$   $N_{tokens}$   $C*N_{templates}$

Information about all the atoms ("single")

Original Atom-Level Single Representation

**c**

$C_{atom}$   $N_{atoms}$

Updated Atom-Level Single Representation

**q**

$C_{atom}$   $N_{atoms}$

Token-Level Single Representation

**s**

$C_{token}$   $N_{tokens}$

Information about all the pairs of atoms ("pair")

Atom-Level Pair Representation

**p**

$N_{atoms}$   $N_{atoms}$   $C_{atompair}$

Token-Level Pair Representation

**z**

$N_{tokens}$   $N_{tokens}$   $C_z$

# Representation Learning

**Template module**: updates **z** using the structure templates **t**

**MSA module:** first updates the MSA using input token level single representation (**q**)
- adds **MSA** to update token-level pair representation (**z**) using Outer Product Mean
- updates the **MSA** based on **z** with a simplified version of self attention with pair bias

# Pairformer Module



Updates s and z with geometry-inspired (triangle) attention

# Why Look at Triangles

**Geometric Principle**

- Based on triangle inequality: sum of any two sides
  > third side
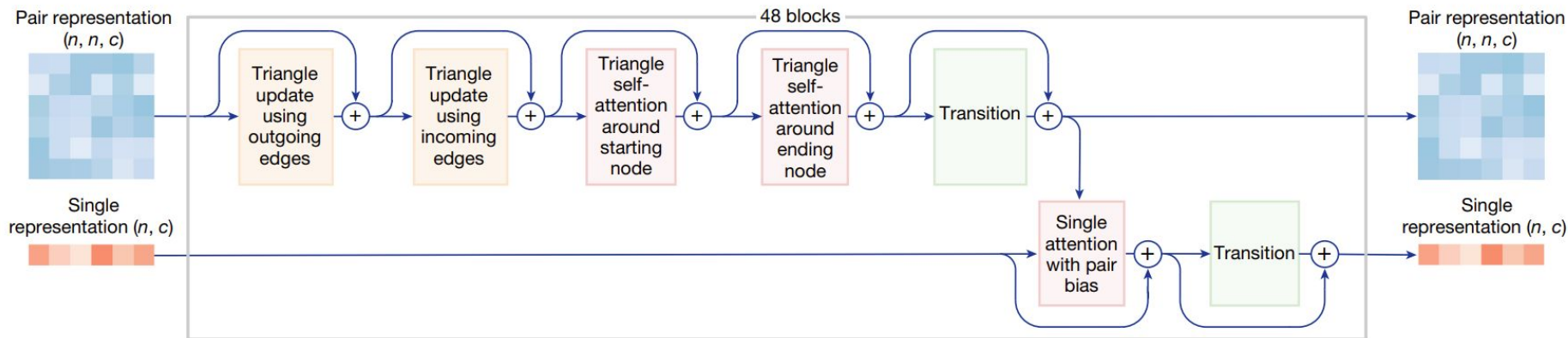- Helps constrain predictions using geometric
  relationships

**Implementation in AF3**

- Each pair relationship ($z\_ij$) is updated using all
  possible third points (k)

**Directional Relationships**

- Two types of paths considered: outgoing edges and
  incoming edges



"outgoing edges"

"incoming edges"

# Triangle Updates - Outgoing



Triangle Update (Outgoing)

Multiply corresponding elements:
- Take row i from a
- Take row j from b
- Element-wise multiplication
- Sum over k dimension

Update zij using zik and zjk

"outgoing edges"

# Triangle Updates - Incoming



Triangle Update (Incoming)

Multiply corresponding elements:
- Take column i from a
- Take column j from b
- Element-wise multiplication
- Sum over k dimension

Update zij using zki and zkj

"incoming edges"

# Triangle Attention (Starting Node)

## Triangle Attention (Starting Node)

(operations not part of standard multi-head self-attention are **bolded**)

For position zij:
**Query**: from zij
**Keys**: from zik (all k in row i)
**Bias**: from zjk

"outgoing edges"

# Triangle Attention (Ending Node)



## Triangle Attention (Ending Node)

(operations not part of standard multi-head self-attention are **bolded**)

For position zij:
**Query**: from zij
**Keys**: from zki (all k in column i)
**Bias**: from zkj

attention along column i with bias from column j, with each row in a column indexed by k

"incoming edges"

# Single Attention with Pair Bias



Attention with pair-bias (Token-Level)

demonstrating the process for updating one token at a time (token s)
operations not part of standard multi-head self-attention are **bolded**

# Structure Prediction via Diffusion

**Training Phase**:

- Start: Real atomic coordinates (xt=0)
- Process: Add noise gradually → xt=T
- Learn: Predict noise added at each step
- Loss: Compare predicted vs actual noise

**Inference Phase:**

- Start: Random coordinates (xt=T)
- Process: Iteratively remove predicted noise
- End: Final denoised structure (x0)

# Conditional Diffusion

**Conditional Diffusion**: final generation matches the information represented by dataset and conditioning input.

**AF3 Implementation:**

- **data**: a matrix x with the x,y,z coordinates of all the atoms.
- **Training phase**:  add noise to matrix x and predict the noises while having input condition
- **Inference phase**:
  - starting with random coordinates
  - first randomly rotate and translate our entire predicted complex.
  - then add a small amount of noise to the coordinates to encourage more heterogeneous generations.
  - Finally, we predict a de-noising step using the Diffusion Module.

# Diffusion Module



Diffusion Module Overview

**1. Prepare token-level conditioning tensors**

**1a.** Prepare **token-level** pair conditioning tensor (**z**): Combine z_trunk with the relative positional encodings through projection and transitions

**1b.** Prepare token-level single conditioning tensor (**s**): Merge s_inputs and s_trunk, add timestep information through Fourier embedding

# Diffusion Module



Diffusion Module Overview

**2. Prepare atom-level tensors, apply atom-level attention, and aggregate back to token-level**

**2a**. **2b.** create atom-level conditioning tensors (**q, p**), based on the current token-level representations (**s, z**)

**2c,2d,2e,2f** use the atom's current coordinates (**x**) by the variance of the data to update (**q**). Finally, we update (**q**) with the Atom Transformer using (**p**), and aggregate back to tokens level (**a**)

# Diffusion Module



Diffusion Module Overview

## 3. Apply attention at the token-level

**3a**. apply attention to update token-level representation (**a**) of the atom coordinates and sequence information, which mirrors the Atom Transformer at input preparation but for tokens.

# Diffusion Module



Diffusion Module Overview

**1a** Prepare token-level pair conditioning tensor

**2a.** Create atom-level pair conditioning tensor

**1b** Prepare token-level single conditioning tensor

**2b** Create atom-level single conditioning tensor

**2c** Create "dimensionless" version of x by scaling to have unit variance

**2d** Update with information from coords

**2e** Update with atom attention

**3a** Token-level attention biased by conditioning tensors (s,z)

**2f** Aggregate to token-level

**4a** Expand to atom-level and update with atom attention

**4b** Predict de-noising update

**4c** Apply update

**4. Apply attention at the atom-level to predict atom-level noise updates**

**4a**. use our updated **a** to update **q** using the Atom Transformer, by broadcasting our **a** first.

**4b**. **4c.** predict de-noising update by maps this atom-level representation q back to R3 and apply update to x.

# Evaluation

# Evaluation



**c**

**Ligands PoseBusters set**

**Nucleic acids**
- AF3
- RoseTTAFold2NA
- Alchemy_RNA2 (has human input)

**Covalent modifications**

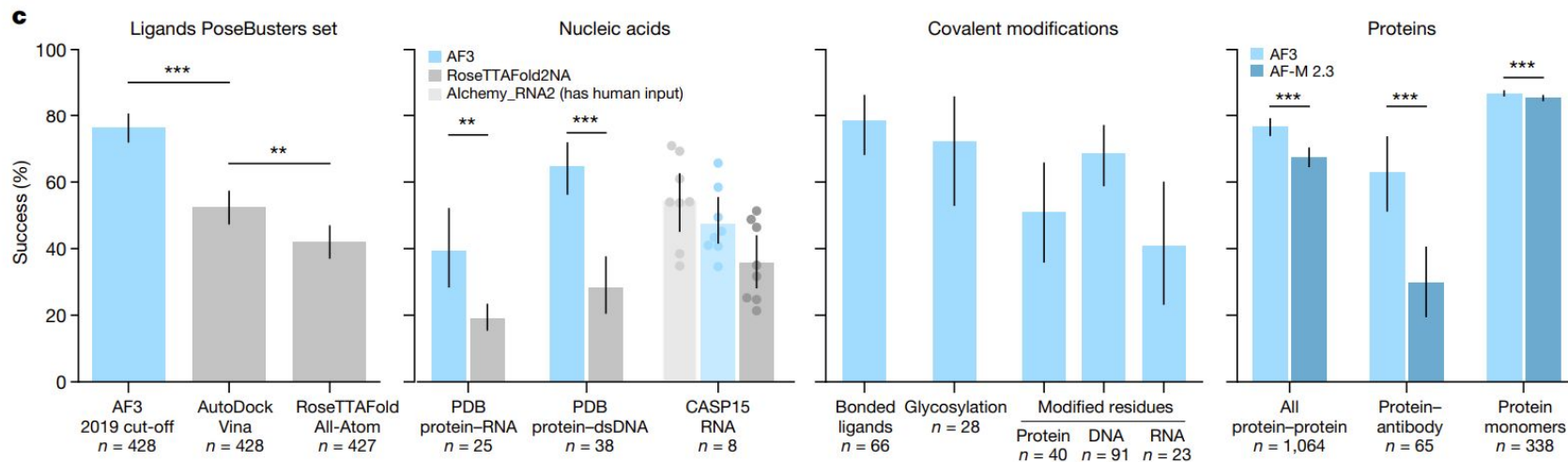**Proteins**
- AF3
- AF-M 2.3

Success (%)

Ligands PoseBusters set:
- AF3 2019 cut-off, *n* = 428
- AutoDock Vina, *n* = 428
- RoseTTAFold All-Atom, *n* = 427

Nucleic acids:
- PDB protein–RNA, *n* = 25
- PDB protein–dsDNA, *n* = 38
- CASP15 RNA, *n* = 8

Covalent modifications:
- Bonded ligands, *n* = 66
- Glycosylation, *n* = 28
- Modified residues: Protein *n* = 40, DNA *n* = 91, RNA *n* = 23

Proteins:
- All protein–protein, *n* = 1,064
- Protein–antibody, *n* = 65
- Protein monomers, *n* = 338

# Thoughts

# ML Musings

**AlphaFold as Retrieval-Augmented Generation:**

Include retrievals from the training set at inference time, by utilizing an MSA and template search

Large Language Models routinely use Retrieval Augmented Generation systems such as a traditional web search at inference time to orient the model toward relevant information

**Pair-Bias Attention**

Attention where the queries, keys, and values all originate from the same source (like in self-attention), but there is a bias term added to the attention map from another source

This particular type of cross-biasing is not seen to be used in other fields