# Highly Accurate Macromolecular Structure Modeling by Deep Learning

**Xiao Wang**

# What is Macromolecular Structure?
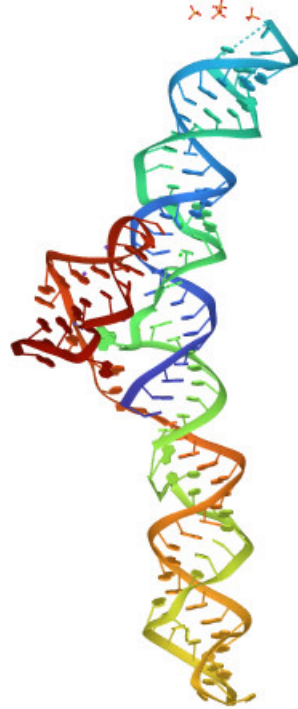
- Large molecules: DNA, RNA, protein

**DNA**

**RNA**

**Protein**

**Complex**



**Gene Information**

**Gene Regulation
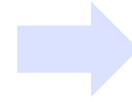Protein Synthesis**

**Cellular Functions**
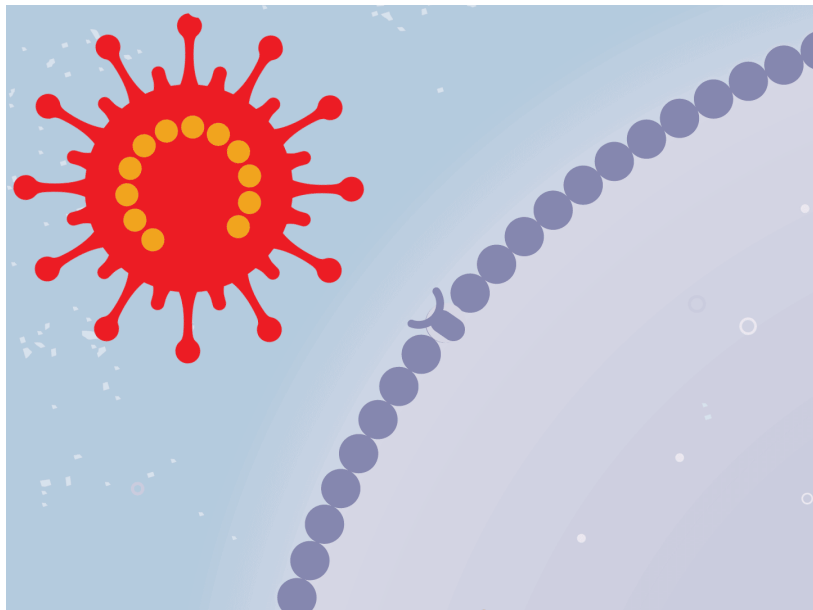
**Gene Regulation
Protein Synthesis**

**Background**

# Why we study Macromolecular Structure?

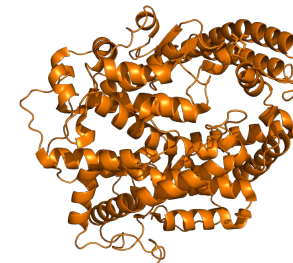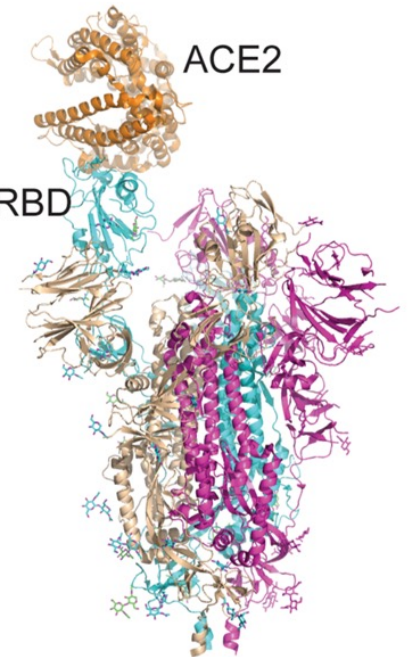Structure → Function

SARS-CoV-2 enters human cells



Spike protein

Closed RBD

ACE2

Opened RBD

ACE2 protein

Huang, Yuan, et al. "Structural and functional properties of SARS-CoV-2 spike protein: potential antivirus drug development for COVID-19." *Acta Pharmacologica Sinica* 41.9 (2020): 1141-1149.
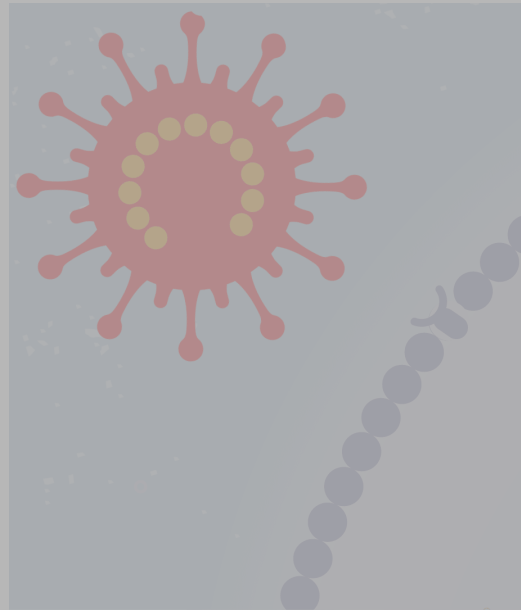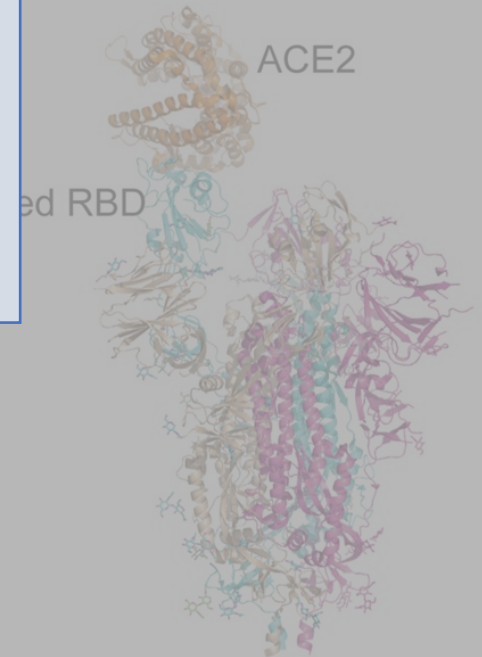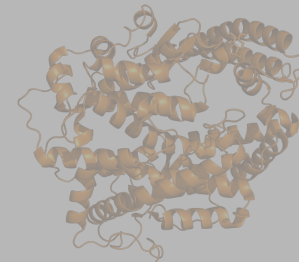
# Why we study Macromolecular Structure?

Structure → Function

SARS-CoV-2 enters human cells

Spike protein

Structure is important!
**3D coordinates** of atoms are important!

ACE2

ed RBD

ACE2 protein

Huang, Yuan, et al. "Structural and functional properties of SARS-CoV-2 spike protein: potential antivirus drug development for COVID-19." *Acta Pharmacologica Sinica* 41.9 (2020): 1141-1149.

# How to Determine Macromolecular Structure?

**X-ray Crystallography**

**Nuclear Magnetic Resonance (NMR)**

**Cryo-Electron Microscopy(cryo-EM)**



(Photo by Charles Christoffer, 2022)
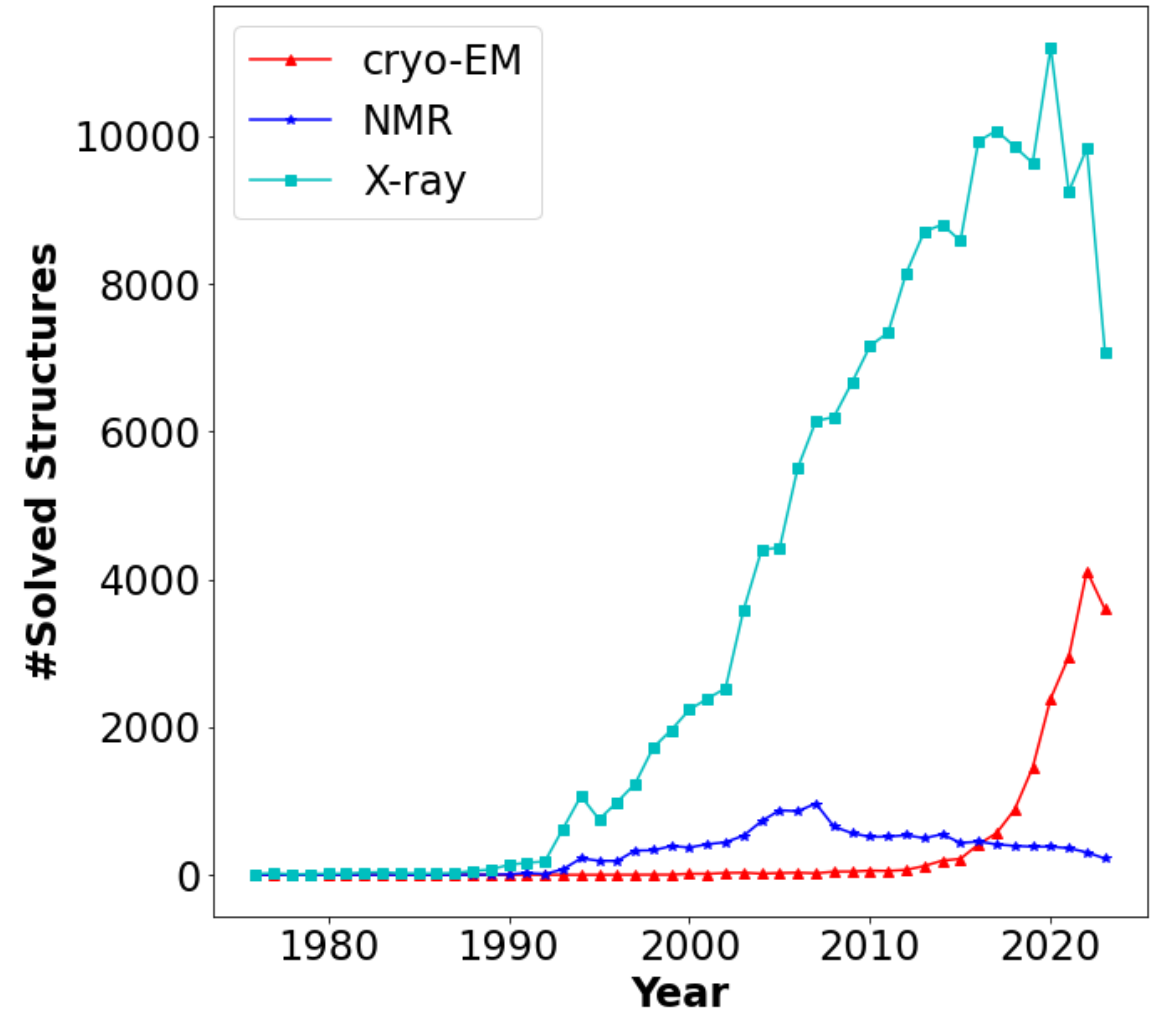
X-Ray Crystallography Facility
Florida State University

Bruker 800MHz NMR
BRWN @ Purdue

Krios G4 Cryo-EM
HOCK @ Purdue
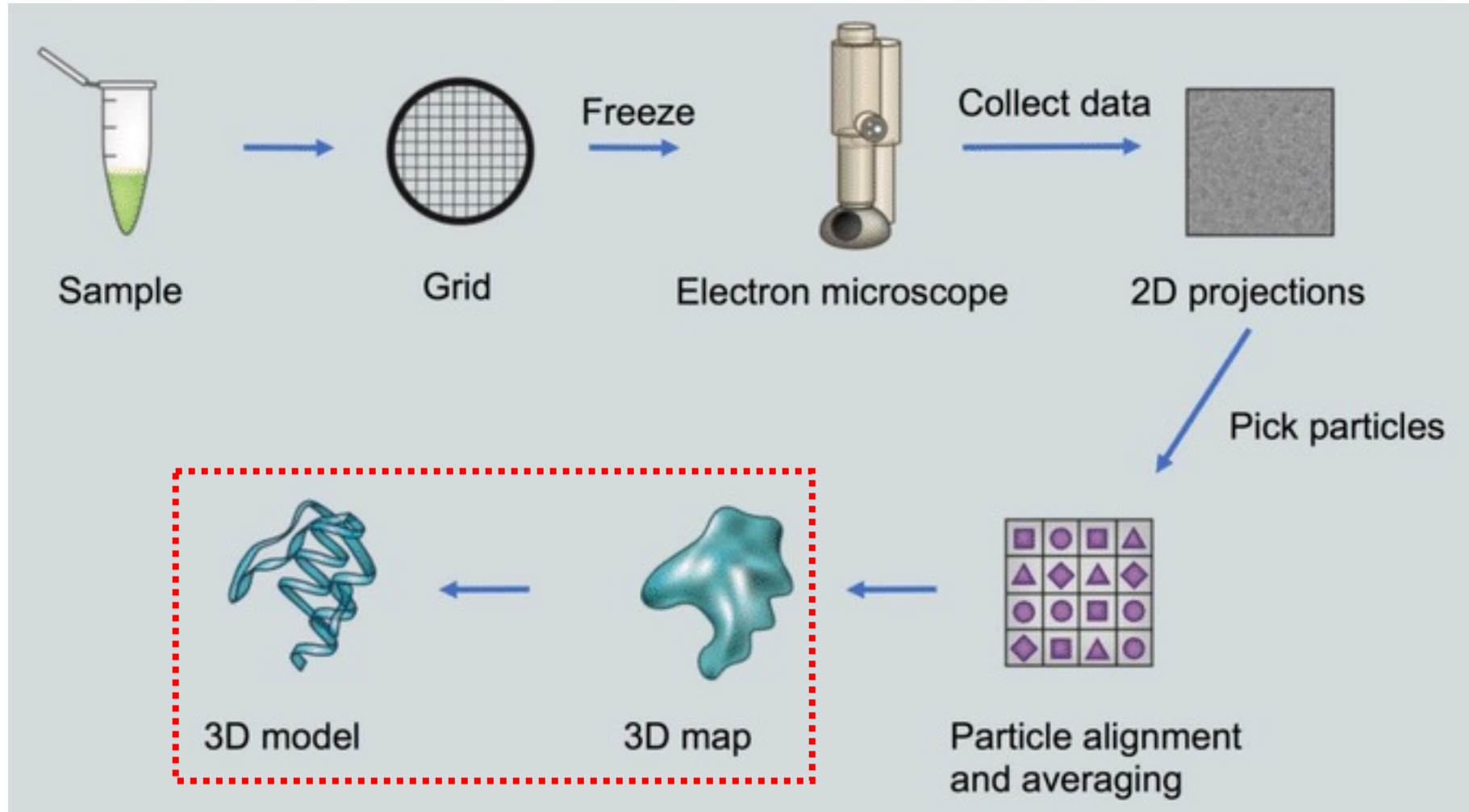
# Statistics – Structure Determination Methods

- Number of solved structures per year by different methods.

- Cryo-EM become popular!

- Advantage of cryo-EM:
  - No need to be crystallized.
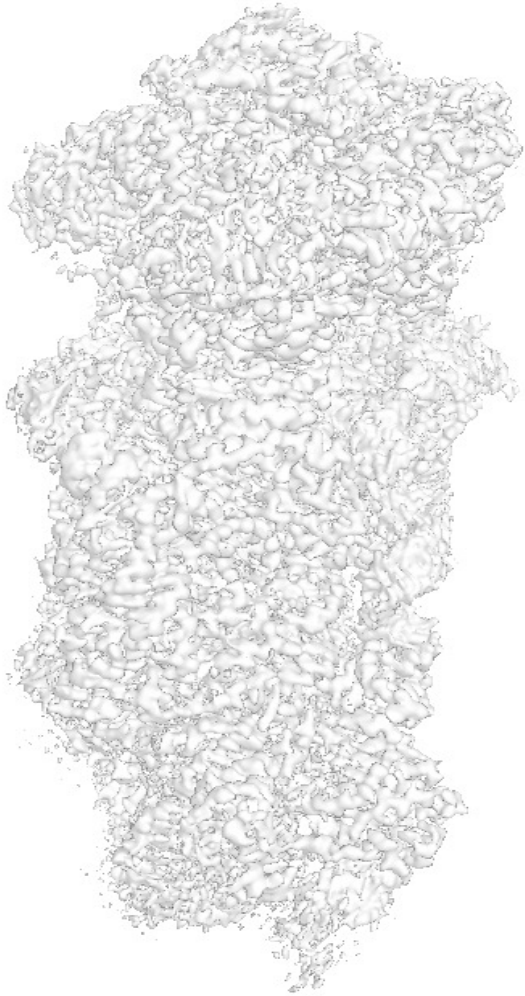  - Can determine large macromolecules.
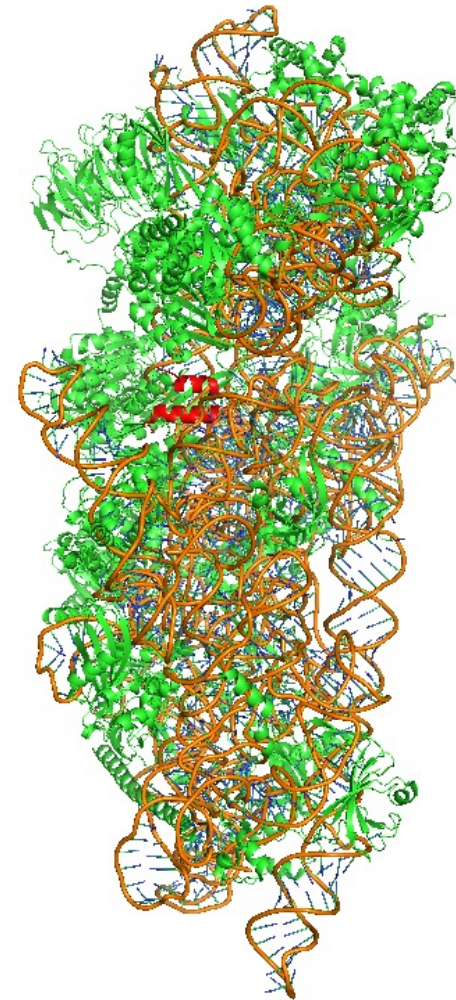
# Cryo-EM Structure Determination Pipeline

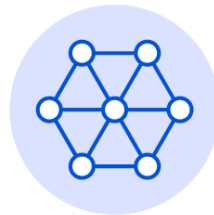**Background**

# AI for Macromolecular Structure Modeling

## Cryo-EM Map

## Structure



**Structure Modeling**

Xiao Wang, Genki Terashi, & Daisuke Kihara. "Cryo-READ: De novo structure modeling for nucleic acids in cryo-EM maps using deep learning". Nature Methods. (2023)
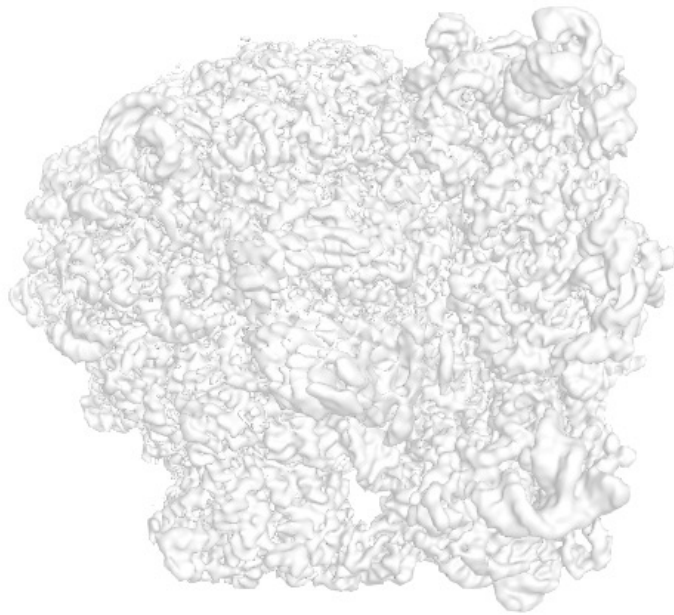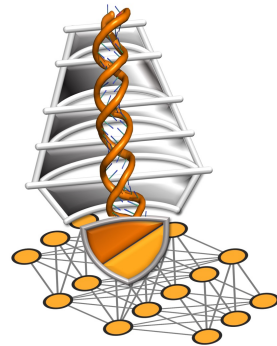
**Overview**

# Method 1: CryoREAD for DNA/RNA structure modeling

Cryo-EM Map

**Structure Modeling**

DNA/RNA
Structure



**CryoREAD**

Xiao Wang, Genki Terashi, & Daisuke Kihara. "Cryo-READ: De novo structure modeling for nucleic acids in cryo-EM maps using deep learning". Nature Methods. (2023)

**Overview**

# Previous Structure Modeling Methods

Cryo-EM Map

Structure
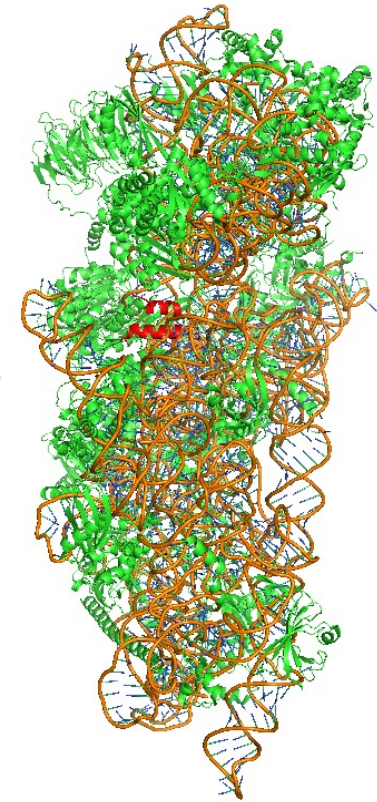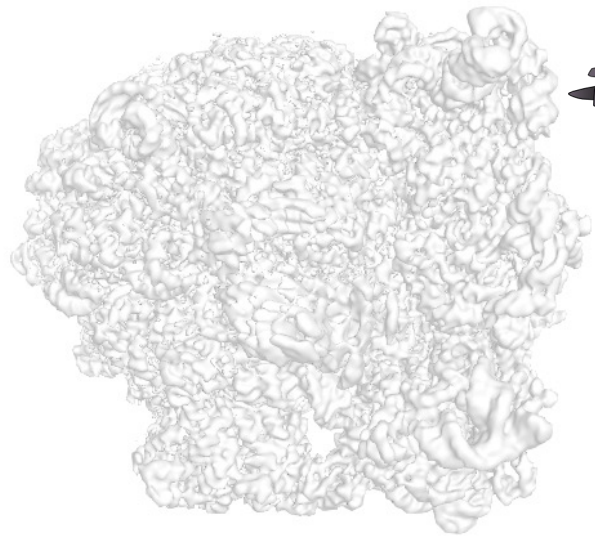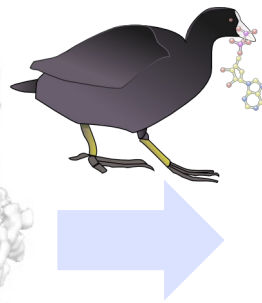
Manual Modeling

Computation Tools

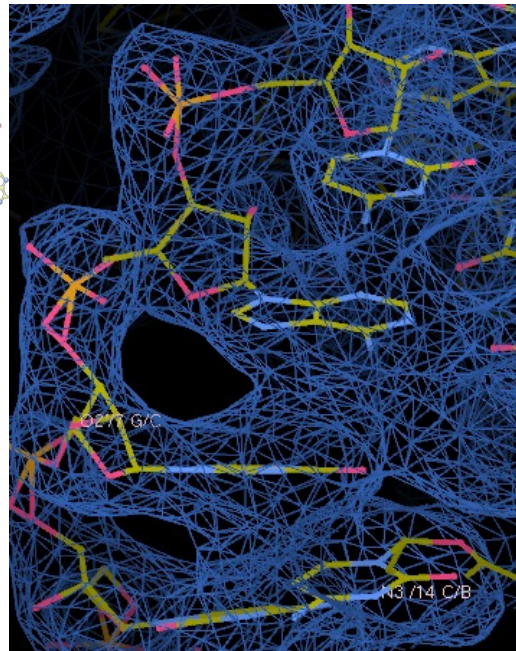# Manual Modeling: Interactive Structure Modeling

Cryo-EM Map

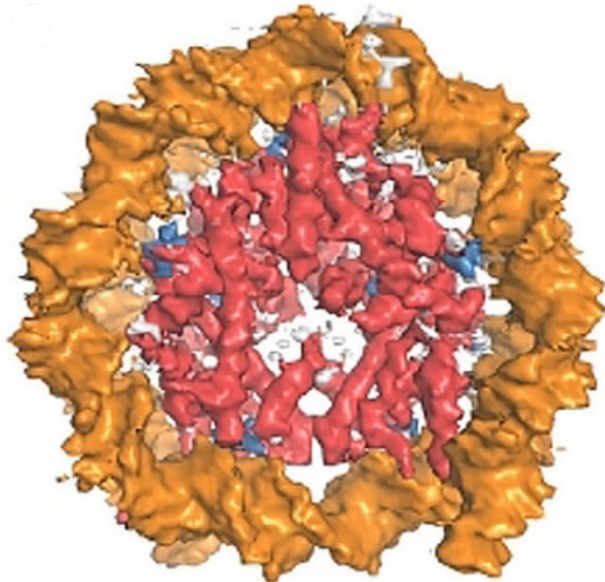Detailed
Visualization

Coot

**Limitations:**
➢ Time-consuming
➢ Local low resolution
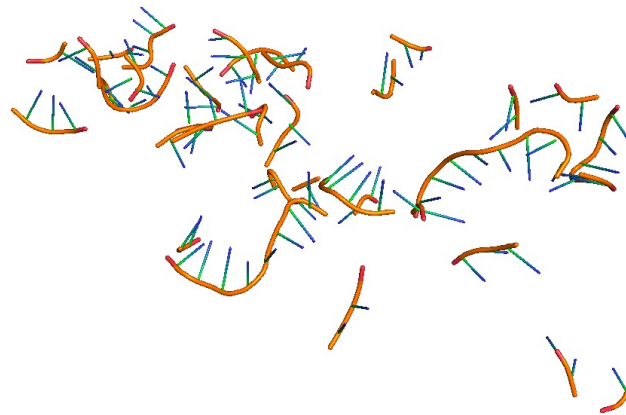➢ Human errors

# Previous computational modeling

- Haruspex: Only structure detection
- Phenix:

1. Focused on protein
2. DNA/RNA atomic model is not accurate.
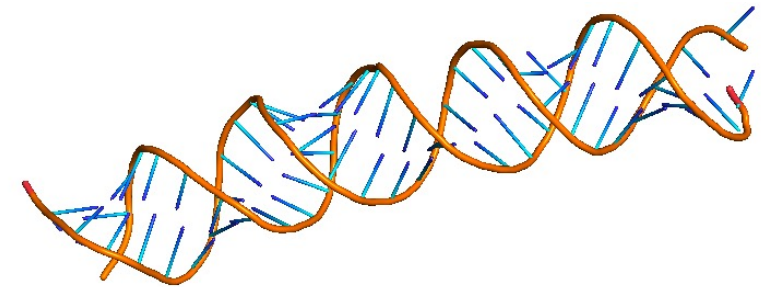
**No good DNA/RNA computational tools!**
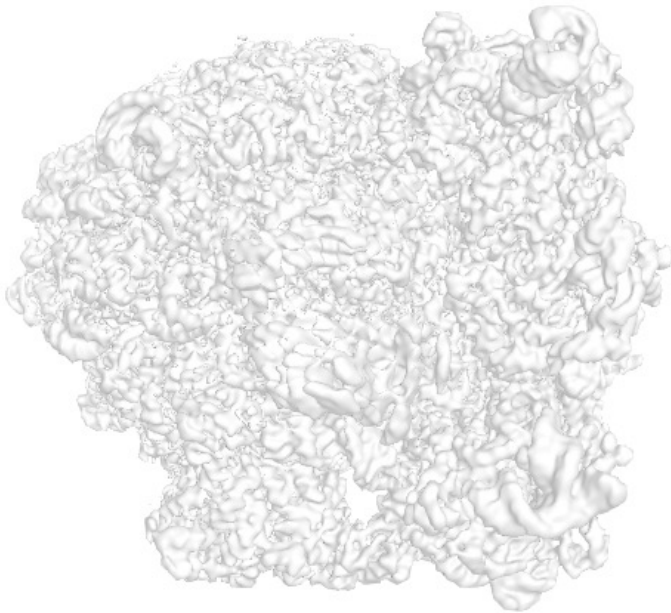
HaruSpex

Phenix

Native Structure (Ground Truth)

Wang Xiao, et al. "CryoREAD: De novo structure modeling for nucleic acids in cryo-EM maps using deep learning". Nature Methods. (2023)
Mostosi, P., Schindelin, H., Kollmannsberger, P., & Thorn, A. (2020). Haruspex: a neural network for the automatic identification of oligonucleotides and protein secondary structure in cryo-electron microscopy maps. Angewandte Chemie International Edition, 59(35), 14788-14795.
Terwilliger, T. C., Adams, P. D., Afonine, P. V., & Sobolev, O. V. (2018). A fully automatic method yielding initial models from high-resolution cryo-electron microscopy maps. Nature methods, 15(11), 905-908.

**Motivation**

# CryoREAD: *De Novo* DNA/RNA Structure Modeling
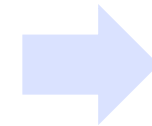
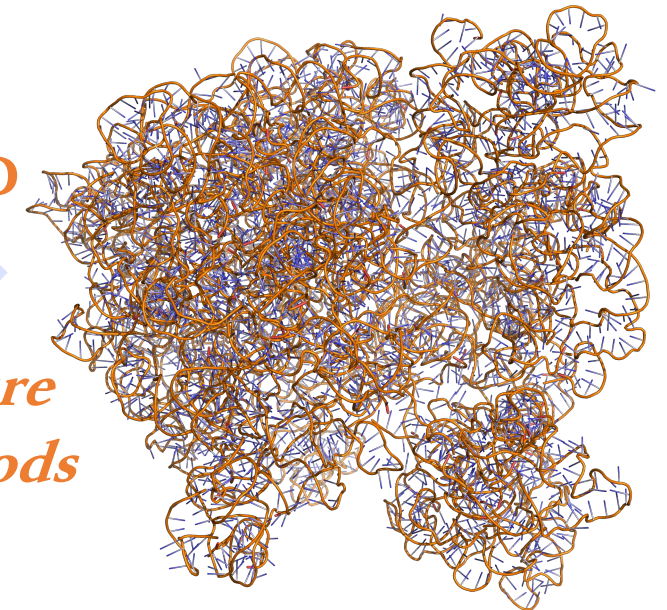Cryo-EM Map

Sequence (Optional)

DNA/RNA Structure

```
>chain A
UACCUGGUUGAUCCUGCCAGU
AGCAUAUGCUUGUCUCAAAGA
UUAAGCCAUGCAUGUCUAAGU
ACGCACGGCCGGUACAGUGAA
ACUGCGAAUGGCUCAUUAAAU
CAGUUAUGGUUCCUUUGGUCG
>chain B
UAACUGUGGUAAUUCUAGAGC
UAAUACAUGCCGACGGGCGCU
GACCCCCUUCGCGGGGGGGAU
GCGUGCAUUUAUCAGAUCAAA
ACCAACCCGGUCAGCCCCUC
```
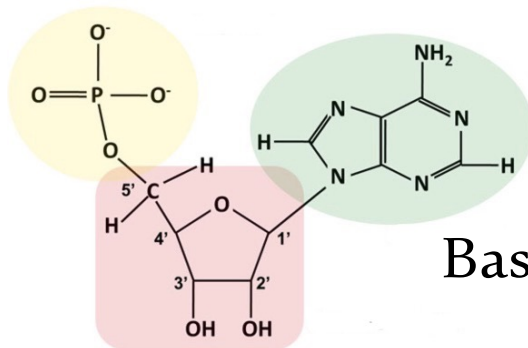
Cryo READ

*Nature Methods*

Wang Xiao, et al. "Cryo-READ: De novo structure modeling for nucleic acids in cryo-EM maps using deep learning". Nature Methods. (2023)

**Overview**

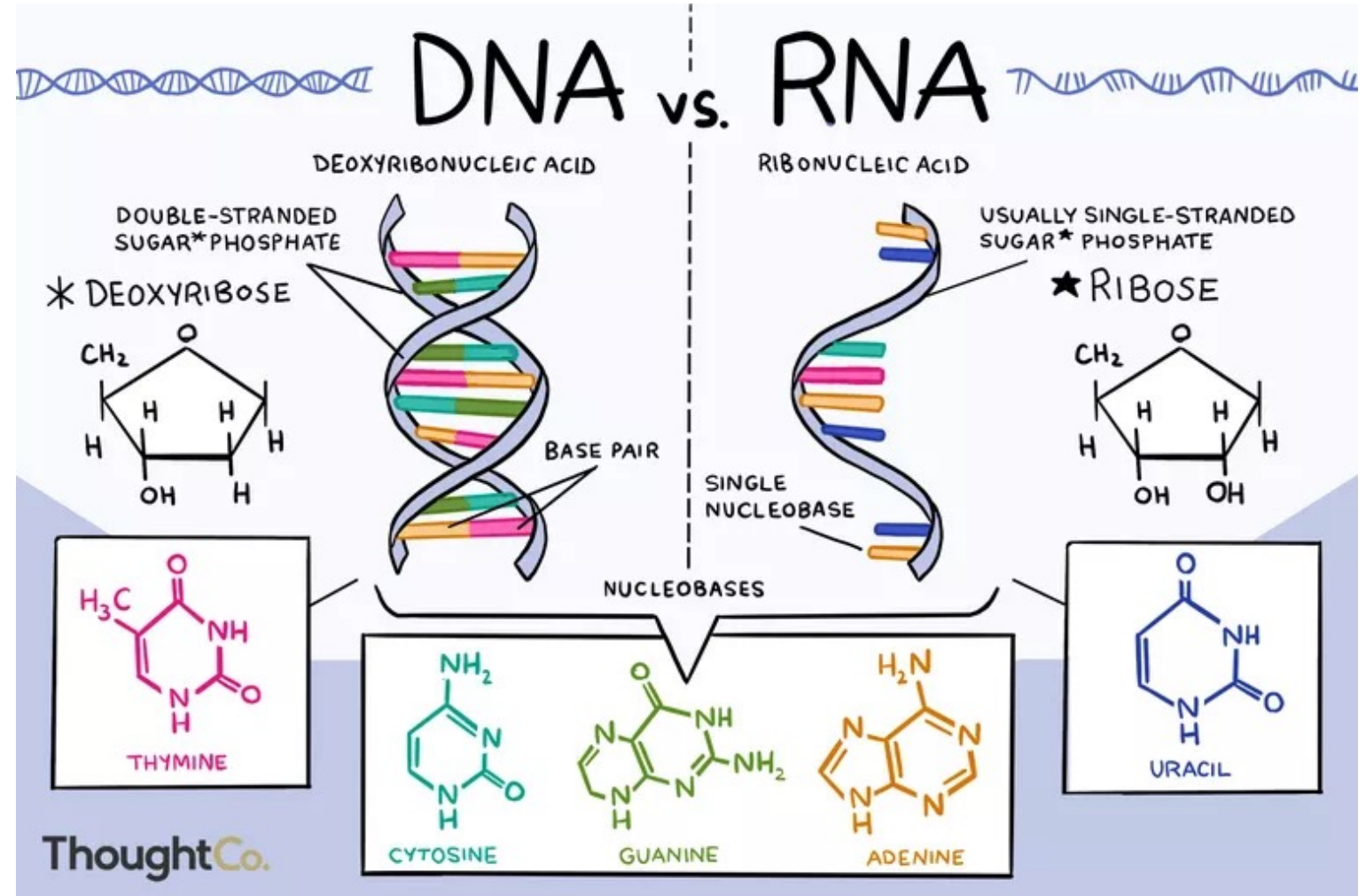**Nucleotides**

Phosphate

Sugar
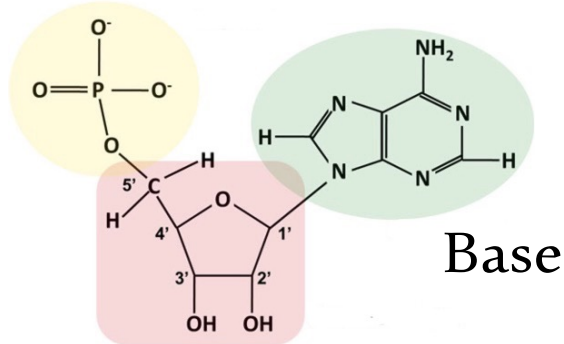
Base

# Domain Knowledge: DNA/RNA Structure
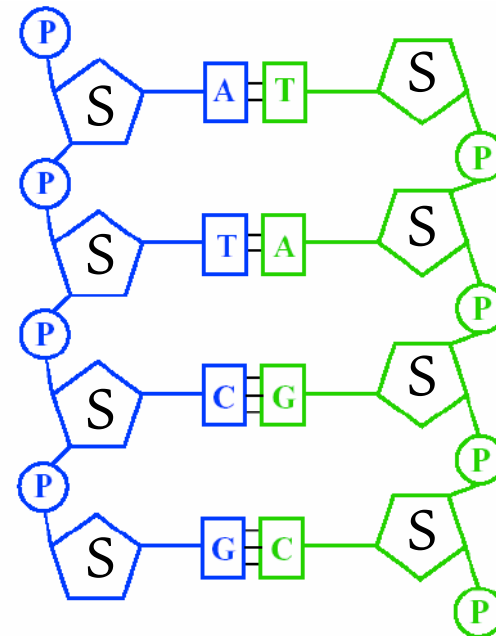


**Nucleotides**

**Phosphate-Sugar Backbone**

**Base Type**

**3D Structure**
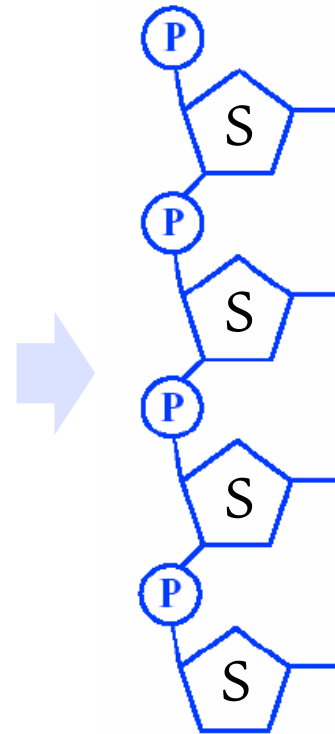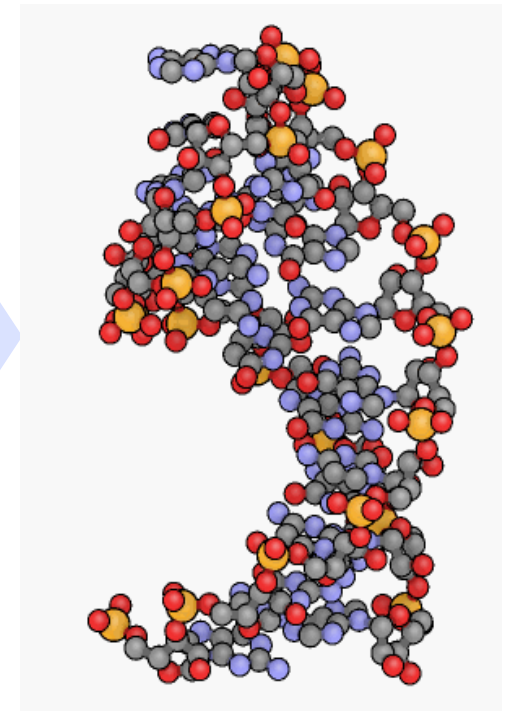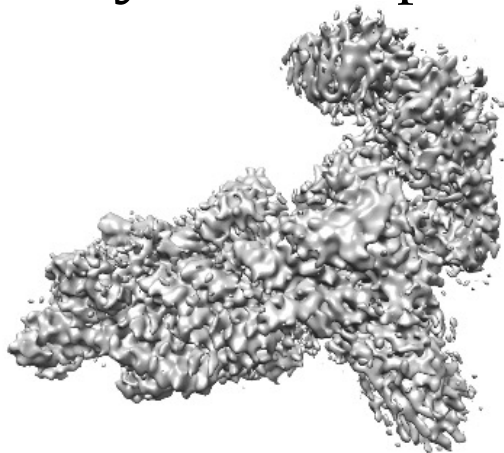
Phosphate

Base

Sugar

**Base Pairing**

Wang Xiao, et al. "Cryo-READ: De novo structure modeling for nucleic acids in cryo-EM maps using deep learning". Nature Methods. (2023)
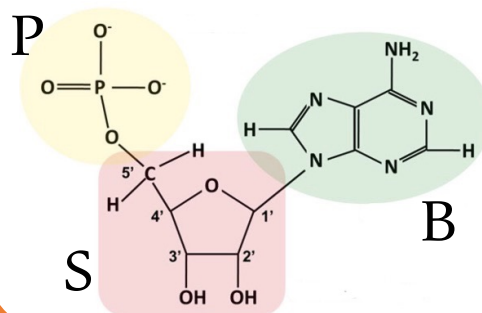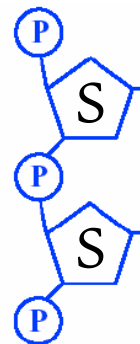
**Methods**

# Overview of CryoREAD



Cryo-EM Map

Nucleotides · Backbone · Base Type · 3D Structure

AI Segmentation → Backbone Tracing → Sequence Assignment → 3D Structure

Methods

Wang Xiao, et al. "Cryo-READ: De novo structure modeling for nucleic acids in cryo-EM maps using deep learning". Nature Methods. (2023)

# Step 1: Phos-Sugar-Base-Protein Segmentation



Nucleotides | Backbone | Base Type | 3D Structure

AI model

Wang Xiao, et al. "Cryo-READ: De novo structure modeling for nucleic acids in cryo-EM maps using deep learning". Nature Methods. (2023)

**Methods**

# Step 1: Phos-Sugar-Base-Protein Segmentation

Wang Xiao, et al. "Cryo-READ: De novo structure modeling for nucleic acids in cryo-EM maps using deep learning". Nature Methods. (2023)
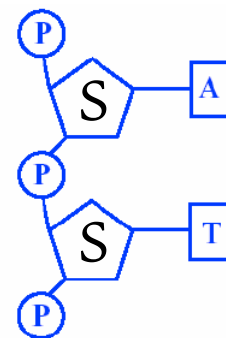
**Methods**

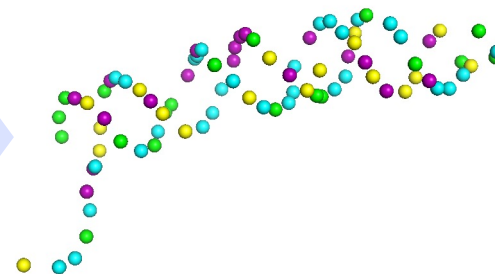# Step 2: Backbone Tracing



Nucleotides    Backbone    Base Type    3D Structure

P    B    S

Sugar Detection → Cluster → Sugar Node → Connect → Sugar Graph

Wang Xiao, et al. "Cryo-READ: De novo structure modeling for nucleic acids in cryo-EM maps using deep learning". Nature Methods. (2023)

**Methods**

# Step 2: Backbone Tracing



Nucleotides · Backbone · Base Type · 3D Structure

Sugar Graph

*Easy Job?*

Sugar Backbone

*How about This?*

Wang Xiao, et al. "Cryo-READ: De novo structure modeling for nucleic acids in cryo-EM maps using deep learning". Nature Methods. (2023)

**Methods**

# Step 2: Backbone Tracing



Nucleotides     Backbone     Base Type     3D Structure

## Vehicle Routing Problem



## Problem Reformulation

$$Cost = \sum_{i=1}^{N_p} \sum_{j=1}^{L(i)-1} w_{i,j} + P_{drop} \sum_{k=1}^{N} drop_k$$

Sugar Node: Locations

$N_p$ Chains: $N_p$ Vehicles

Drop Wrong Node: Drop Penalty

Wang Xiao, et al. "Cryo-READ: De novo structure modeling for nucleic acids in cryo-EM maps using deep learning". Nature Methods. (2023)

**Methods**

# : Backbone Tracing



Nucleotides — Backbone — Base Type — 3D Structure

Sugar Graph

**Vehicle Routing Problem**

Sugar Backbone

Wang Xiao, et al. "Cryo-READ: De novo structure modeling for nucleic acids in cryo-EM maps using deep learning". Nature Methods. (2023)

**Methods**

Nucleotides · Backbone · Base Type · 3D Structure

AI model

Base Type Detection

Wang Xiao, et al. "Cryo-READ: De novo structure modeling for nucleic acids in cryo-EM maps using deep learning". Nature Methods. (2023)

**Methods**

# Step 3: Sequence Assignment



Nucleotides | Backbone | Base Type | 3D Structure

Sugar Backbone **+** Base Type Detection **→** Initial Sequence Assignment

Wang Xiao, et al. "Cryo-READ: De novo structure modeling for nucleic acids in cryo-EM maps using deep learning". Nature Methods. (2023)

**Methods**

# Step 3: Sequence Assignment



Nucleotides   Backbone   Base Type   3D Structure

Initial Sequence Assignment

Limitations

Geometry Constraint

Distance?

Sequence Information

>chain A
UACCUGGUUGAUCC
UGCCAGUAGCAUAU
GCUUGUCUCAAAGA
UUAAGCCA

Wang Xiao, et al. "Cryo-READ: De novo structure modeling for nucleic acids in cryo-EM maps using deep learning". Nature Methods. (2023)

**Methods**

# Step 3: Sequence Assignment

Nucleotides

Backbone

Base Type

3D Structure

Initial Sequence Assignment

Scanned Fragments

Fragment Assignment

Sliding Window

DP

7.4    6.8    6.2

5.9    4.0    2.1

Input Sequence

....UACCUGGUUGAUCCUGCCAGUAGCA........

Wang Xiao, et al. "Cryo-READ: De novo structure modeling for nucleic acids in cryo-EM maps using deep learning". Nature Methods. (2023)

**Methods**

# Step 3: Sequence Assignment



Nucleotides | Backbone | Base Type | 3D Structure

**Multiple Assignments** 😭

$$G\text{-}TT....TC\text{-}T$$
$$G\text{-}TT....TT\text{-}T$$

**Conflict Assignments** 😭

**Constraint Programming**

① Traced Main Chain

12   19      16   20

②

T-T-A-C-A-C

T--....A-A-C

$$\max \sum_{i=1}^{N} x_i r_i$$

with        $x_i \in \{0,1\}$

s.t. $x_i + x_j + x_k + \cdots + x_l = 1$

**Methods**

# Step 3: Sequence Assignment



Nucleotides · Backbone · Base Type · 3D Structure

Initial Assignment → DP → Fragment-based Assignment → CP → Final Assignment

G-TT....TC-T
G-TT....TT-T

TG-....-A-A
TT-....-C-A

T-T....C-A-C
TT-....G-A-C

Wang Xiao, et al. "Cryo-READ: De novo structure modeling for nucleic acids in cryo-EM maps using deep learning". Nature Methods. (2023)

**Methods**

# Step 4: Atomic Structure Modeling



Nucleotides     Backbone     Base Type     3D Structure

Sequence Assignment     S-P-B Backbone     Atomic Structure

**Methods**

# RNA-protein Complex Example by CryoREAD

**Cryo-EM Map**

**Modeled Structure**

**Native Structure** (Ground Truth)



**Cryo READ** →

Automatically model **full** structure

Wang Xiao, et al. "Cryo-READ: De novo structure modeling for nucleic acids in cryo-EM maps using deep learning". Nature Methods. (2023)
Thoms, Matthias, et al. "Structural basis for translational shutdown and immune evasion by the Nsp1 protein of SARS-CoV-2." *Science* 369.6508 (2020): 1249-1255.

**Results**

# Method 2: DiffModeler: Protein Complex Structure Modeling

### Cryo-EM Map

### Protein Structure



**Structure Modeling**

**DiffModeler**

Wang Xiao, et al. "DiffModeler: Protein Complex Structure Modeling with Diffusion Model and AlphaFold in cryo-EM maps". In preparation. (2023)

**Overview**

# Background: Protein Structure

## Amino Acid



## Sequence



Met  Asp  Arg     Val     Gly  Ile  Lys     Val  Asp  Leu

N-terminus                                                              C-terminus

Phe     Ala Leu Gln     Ser Leu     Lys     Leu  Ala

## Structure



## Protein Complex Structure

**Background**

# Background: Template-based Structure Modeling



Sequence

SQETRKKCTEMKKKFKNCEVRCDESNHCVEVRCSDTKYTLC

Cryo-EM Map

Sequence Similarity Search

Protein BLAST
protein ▶ protein

Similar Structures

Fitting

**Challenges:** 1. *Where are the* templates(candidates)? 2. *Where to fit?*

Wang Xiao, et al. "DiffModeler: Protein Complex Structure Modeling with Diffusion Model and AlphaFold in cryo-EM maps". In preparation. (2023)

# Where are the templates(candidates): AlphaFold2

Sequence

SQETRKKCTEMKKKFKNCEVRCDESNHCVEVRCSDTKYTLC

**AlphaFold2**

Structure



Performance

Highly accurate structure by AlphaFold2 can serve as **reliable templates** for structure modeling in cryo-EM maps



TM Score

Predicted TM-Score

**Motivation**

# Where to fit?



Single-chain Structure

*Where?*

EM-Map

Corresponding Region

Wang Xiao, et al. "DiffModeler: Protein Complex Structure Modeling with Diffusion Model and AlphaFold in cryo-EM maps". In preparation. (2023)

**Motivation**

# Where to fit? Backbone tracing via diffusion model



Single-chain Structure

*Fitting*

Backbone Tracing

*Diffusion*

EM-Map

Wang Xiao, et al. "DiffModeler: Protein Complex Structure Modeling with Diffusion Model and AlphaFold in cryo-EM maps". In preparation. (2023)

**Motivation**

# Overview of DiffModeler



**Cryo-EM Map**

**Multi-Chain Sequence**

A  QMGYDRAITVFSPDGRLFQVEYAREA

1  TTTVGLVCKDGVVMATEKRATMGNFI

H  LLEKLKKLEEDYYKLRELYRRLEDEK

D  QMGYDRAITVFSPDGRLFQVEYAREA

**Diffusion Model**

**AlphaFold**

**Traced Backbone**

**Structure Pool**

**VESPER Fitting**

**Fitted Structure Pool**

**Assemble**

**Predicted Protein Complex**

**Method**

# DiffModeler Modeled Structure

**EM-Map and Native Structure**
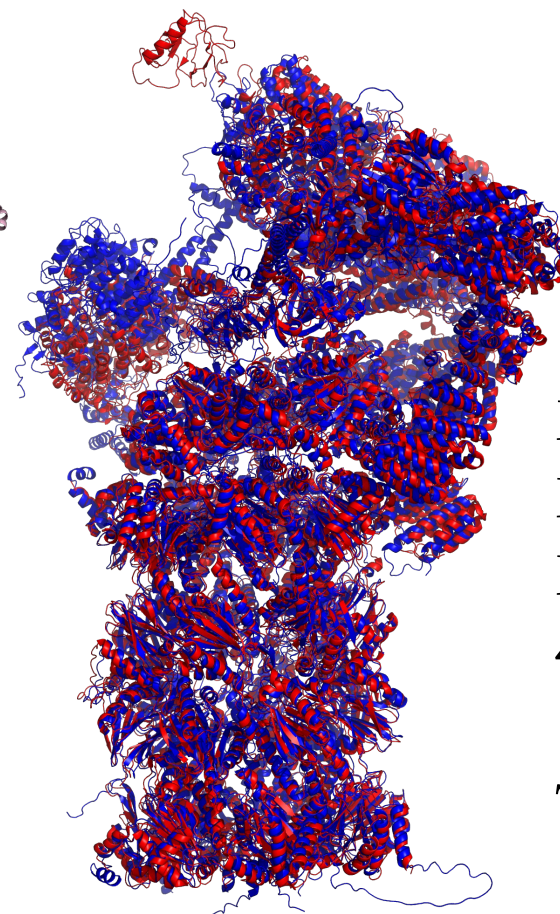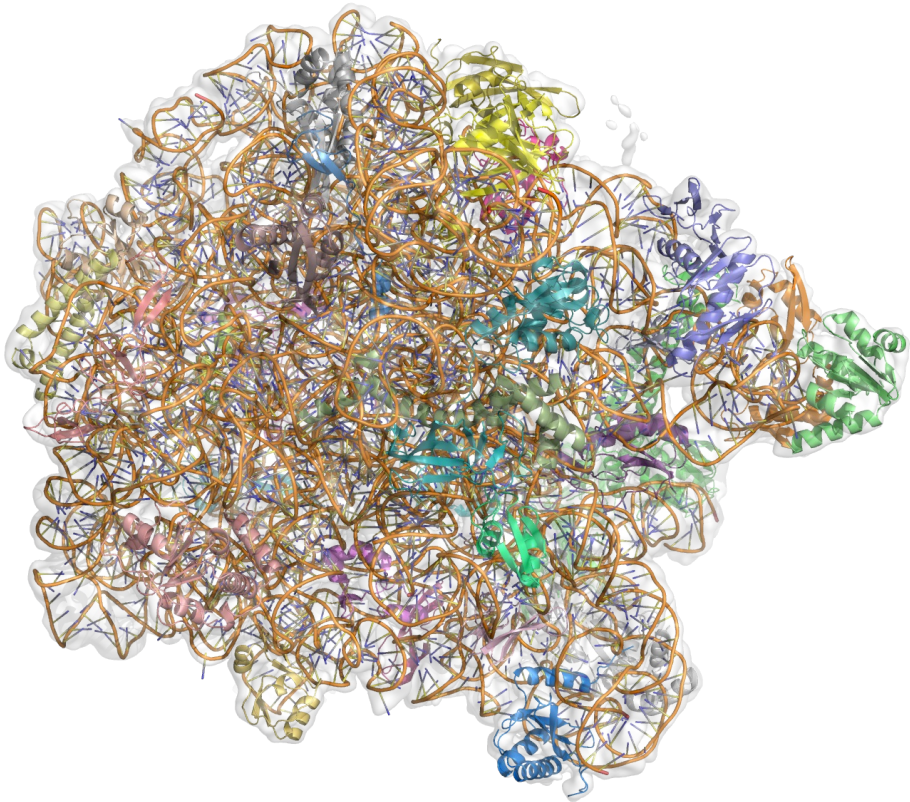
**Structure by DiffModeler**

**Structure Comparison**



**Ground Truth**

**Predicted Structure**
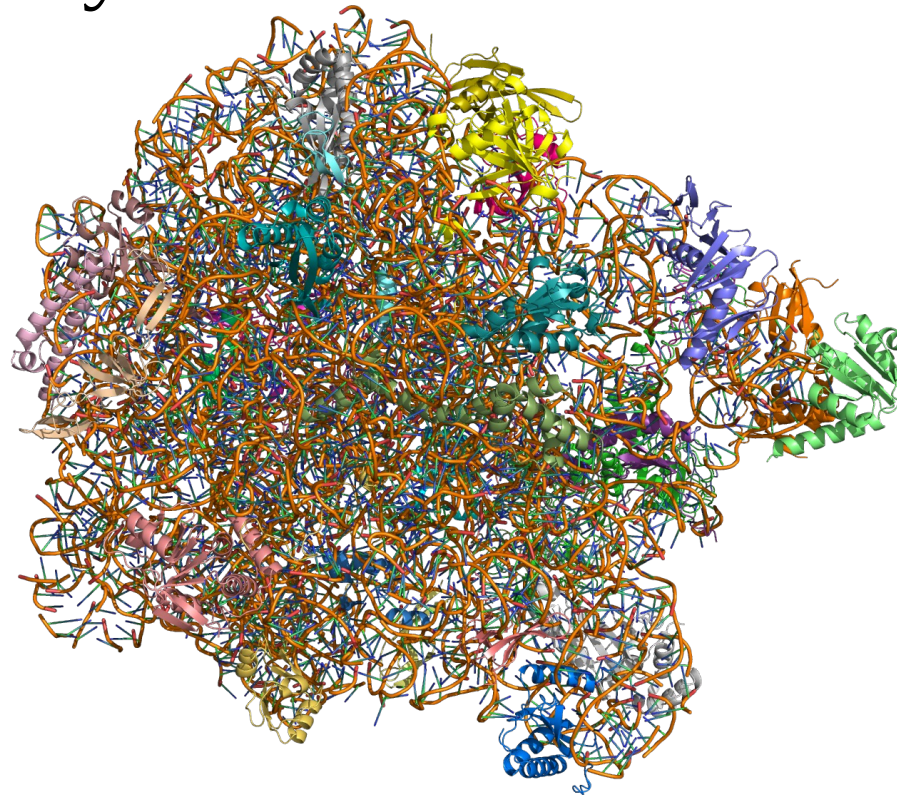
EMD-6693
PDB-ID: 5WVI
Resolution: 6.30 Å
47 chains
13,462 residues
TM-Score: **0.94**

Wang Xiao, et al. "DiffModeler: Protein Complex Structure Modeling with Diffusion Model and AlphaFold in cryo-EM maps". In preparation. (2023)

**Results**

# Protein-RNA complex by DiffModeler+CryoREAD

Cryo-EM Map and
Ground Truth Structure

Modeled Structure by
CryoREAD and DiffModeler



EMD-13017
PDB-ID: 7OPE
Resolution: 3.20 Å
3,818 residues
2,996 nucleotides
TM-Score(protein):
**0.92**
Backbone Recall
(RNA): **0.94**

**Results**

Wang Xiao, et al. "DiffModeler: Protein Complex Structure Modeling with Diffusion Model and AlphaFold in cryo-EM maps". In preparation. (2023)

# Acknowledgement