

Lecture 13:

AI Models and Biological Understanding

AI Revolution in Biology

Foundation Models

Scientific Discovery

Ho-min Park

<

Lecture Contents

Part 1: Foundation Models

Part 2: Biological Applications

Part 3: Design and Engineering

Part 1/3 - Foundation Models

Large-scale pretraining

Transfer learning

Emergent capabilities

Language Models in Biology

Biological sequences as text

DNA, RNA, Protein sequences → Text format

Tokenization strategies

K-mers, BPE, Character-level encoding

Pretraining objectives

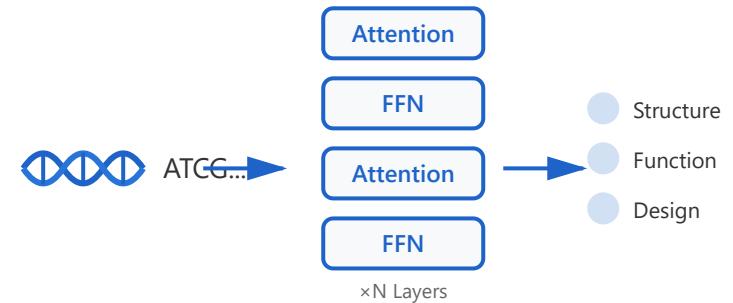
Masked LM, Next token prediction, Contrastive

Scale effects

Model size vs. performance trade-offs

Downstream tasks

Structure, Function, Design applications



1. Biological Sequences as Text

DNA Sequences

DNA consists of four nucleotide bases (A, T, G, C) that can be naturally represented as text strings, making them directly compatible with language model architectures.

Original: ATCGATCGTAGCTAGCTA

Tokenized: A T C G A T C G T A G C T A G C T A

Protein Sequences

Proteins use 20 amino acids represented by single-letter codes, creating a natural alphabet for language modeling similar to human language.

Original: MKTAYIAKQRQISFVKSH

Vocab size: 20 amino acids + special tokens

RNA Sequences

RNA sequences (A, U, G, C) can be treated similarly to DNA, with additional structural information from secondary structure annotations.

- ▶ Fixed vocabulary size makes tokenization straightforward
- ▶ Sequential nature enables transfer of NLP techniques
- ▶ Enables pre-training on massive unlabeled sequence databases

Sequence Representation

DNA: A T C G ...

Protein: M K T A ...

Text Format:

DNA: "ATCGATCGTAGCTA..."

Protein: "MKTAYIAKQRQI..."

RNA: "AUCGAUCGUAGCUA..."

2. Tokenization Strategies

Character-level Encoding

Each nucleotide or amino acid is treated as a single token. Simple and direct, but may miss important patterns spanning multiple positions.

Input: ATCGATCG
Tokens: [A] [T] [C] [G] [A] [T] [C] [G]

K-mer Tokenization

Sequences are split into overlapping or non-overlapping subsequences of length k. Captures local patterns and motifs effectively.

Input: ATCGATCG (k=3)
Tokens: [ATC] [TCG] [CGA] [GAT] [ATC] [TCG]

Byte Pair Encoding (BPE)

Data-driven approach that learns common subword units from the training corpus. Balances vocabulary size with sequence length.

Learns frequent patterns:
"ATG" → start codon
"TAA" → stop codon

- ▶ Choice affects model's ability to capture biological motifs
- ▶ K-mer size impacts computational efficiency and context
- ▶ BPE can discover biologically meaningful units

Tokenization Comparison

Original Sequence:

A T C G A T C G T A G C

Character-level (k=1):

A T C G A T ... Tokens: 12

K-mer (k=3):

ATC TCG CGA GAT ... Tokens: 10

K-mer (k=6):

ATCGAT TCGATC ... Tokens: 7

BPE (learned):

ATCG AT CGTA GC ... Tokens: 4

* Longer k-mers capture more context but increase vocab size

3. Pretraining Objectives

Masked Language Modeling (MLM)

Random tokens are masked, and the model learns to predict them using bidirectional context. Similar to BERT, enables learning rich representations.

Input: ATCG [MASK] TCGTA
Target: Predict 'A' using context
Model: ESM, ProtBERT

Next Token Prediction

Autoregressive training where the model predicts the next token given all previous tokens. Similar to GPT architecture, useful for generation tasks.

Input: ATCGATCG
Target: Predict next 'T'
Model: ProGen, ProtGPT2

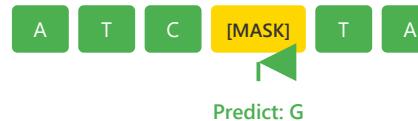
Contrastive Learning

Learns by contrasting positive pairs (e.g., sequence and structure) against negative pairs. Effective for multimodal alignment.

Positive: (Sequence, Structure)
Negative: (Sequence, Random Structure)
Model: ESM-IF, ProteinCLIP

Pretraining Objectives

Masked Language Modeling:



Next Token Prediction:



Contrastive Learning:



* All objectives learn from unlabeled sequence data at scale

- ▶ MLM: Best for understanding tasks (classification, prediction)
- ▶ Next token: Best for generation and design tasks
- ▶ Contrastive: Best for multimodal tasks and alignment

4. Scale Effects

Model Size Scaling

Larger models (more parameters) generally achieve better performance on downstream tasks, following similar scaling laws as in natural language models.

ESM-2: 8M → 150M → 650M → 3B → 15B params
Performance improves consistently with size

Data Scaling

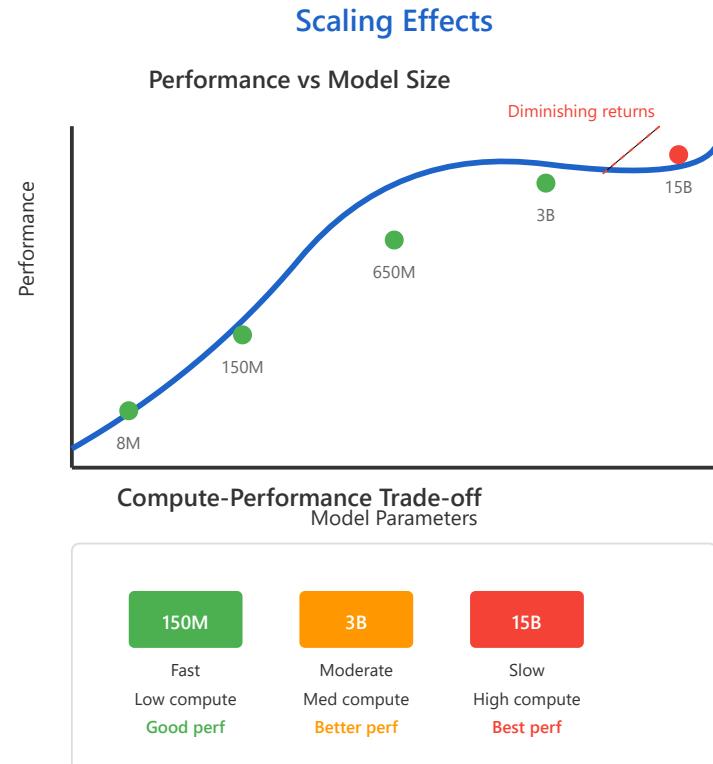
Training on larger sequence databases (UniProt, GenBank) provides richer representations. Models benefit from evolutionary diversity in training data.

ESM-2: Trained on 250M+ sequences
ProtT5: Trained on UniRef50 (45M sequences)

Compute-Performance Trade-offs

Larger models require more computational resources but provide diminishing returns. Need to balance accuracy gains with practical deployment constraints.

15B model: 100x compute of 150M model
Performance gain: ~15-20% on benchmarks



- ▶ Scaling laws similar to NLP models apply to biological sequences
- ▶ Emergent capabilities appear at certain scale thresholds
- ▶ Model selection depends on task complexity and resources

5. Downstream Tasks

Structure Prediction

Predicting 3D protein structures from sequences. Models learn structural constraints from sequence patterns. AlphaFold2 and ESMFold achieve near-experimental accuracy.

Input: Protein sequence
Output: 3D coordinates of all atoms
Applications: Drug design, protein engineering

Function Prediction

Predicting protein functions, subcellular localization, interactions, and enzymatic activity from sequence representations.

Tasks: GO term prediction, EC number
Active site identification
Protein-protein interaction prediction

Protein Design

Generating novel sequences with desired properties. Includes de novo design, optimization of existing proteins, and inverse folding (structure to sequence).

Input: Desired function/structure
Output: Novel protein sequence
Models: ProteinMPNN, ESM-IF, RFDiffusion

- ▶ Fine-tuning pretrained models dramatically improves performance
- ▶ Zero-shot capabilities emerge from large-scale pretraining

Downstream Applications

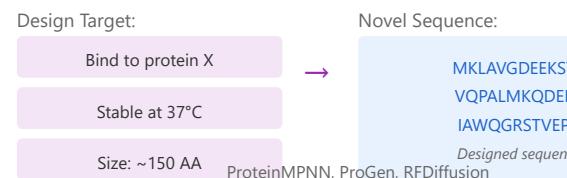
Structure Prediction:



Function Prediction:

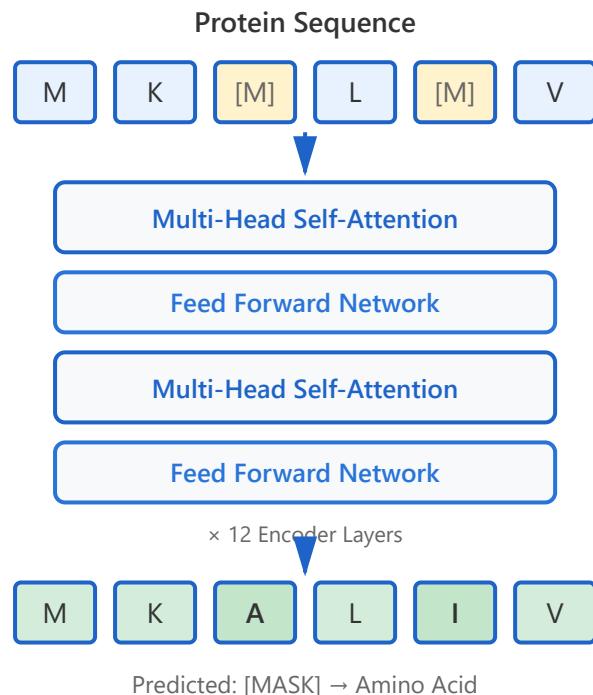


Protein Design:



- ▶ Multimodal models combine sequence, structure, and function

BERT for Proteins



ProtBERT architecture

12-layer bidirectional encoder

Masked language modeling

15% random masking strategy

Attention patterns

Learns residue interactions

Structural insights

Captures 3D contact maps

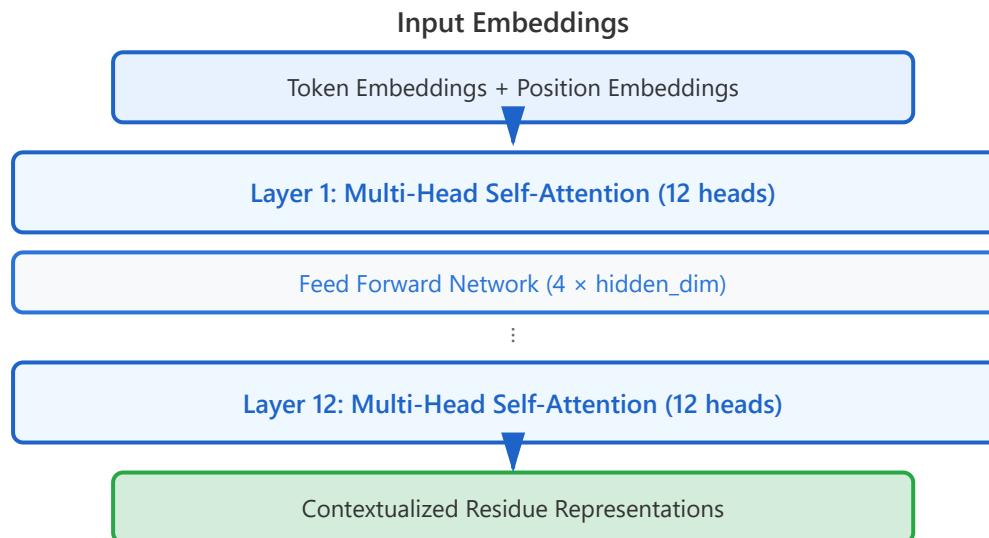
Function prediction

GO terms, EC numbers

ProtBERT Architecture

Architecture Overview

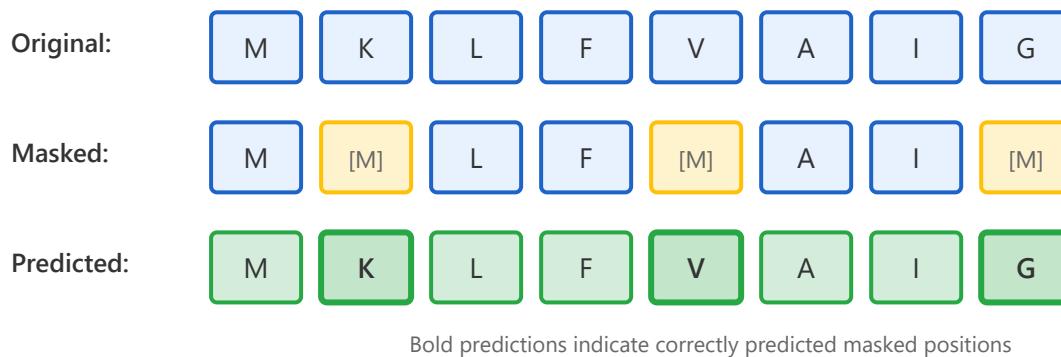
ProtBERT is based on the BERT architecture adapted for protein sequences. It employs a 12-layer bidirectional Transformer encoder that processes protein sequences to generate contextual embeddings for each amino acid residue.



Masked Language Modeling

Training Objective

ProtBERT is trained using Masked Language Modeling (MLM), where 15% of amino acids in the input sequence are randomly masked, and the model learns to predict the original amino acids based on bidirectional context.



Masking Strategy (15% of tokens)

- **80%**: Replace with [MASK] token
- **10%**: Replace with random amino acid

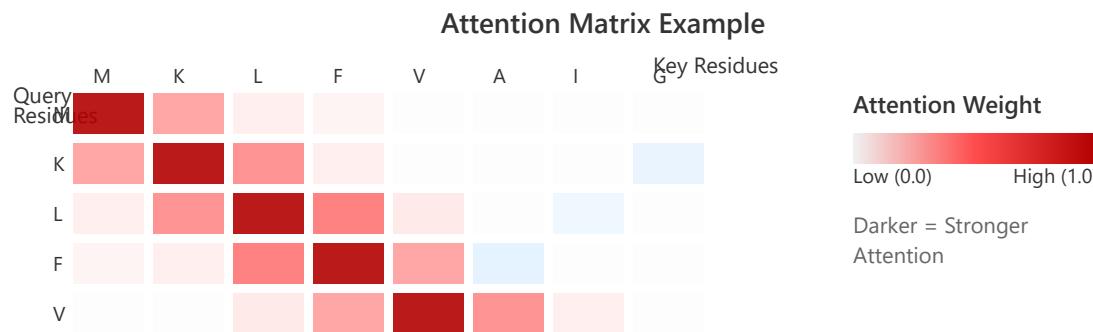
Training Data

- **UniRef100**: 217 million sequences
- **BFD**: 2.5 billion sequences

Attention Patterns

Learning Residue Interactions

The multi-head self-attention mechanism in ProtBERT learns to capture complex dependencies between amino acids. Different attention heads specialize in different types of interactions, from local sequential patterns to long-range contacts.



Types of Patterns Learned

- **Local patterns:** Adjacent residue correlations (α -helices, β -sheets)

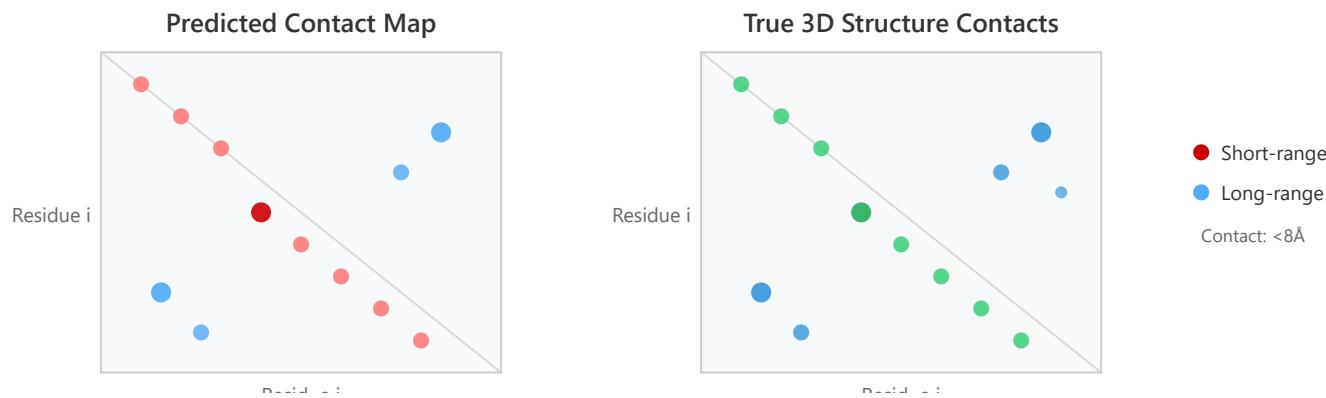
Head Specialization

- **Heads 1-4:** Focus on local sequence context

Structural Insights

Capturing 3D Contact Maps

ProtBERT learns implicit structural information from sequence data alone. The attention weights in deeper layers show strong correlation with actual 3D contacts in protein structures, enabling structure prediction tasks without explicit structural training.



Structural Features Learned

- **Secondary structure:** α -helices, β -sheets, loops with 85%+ accuracy

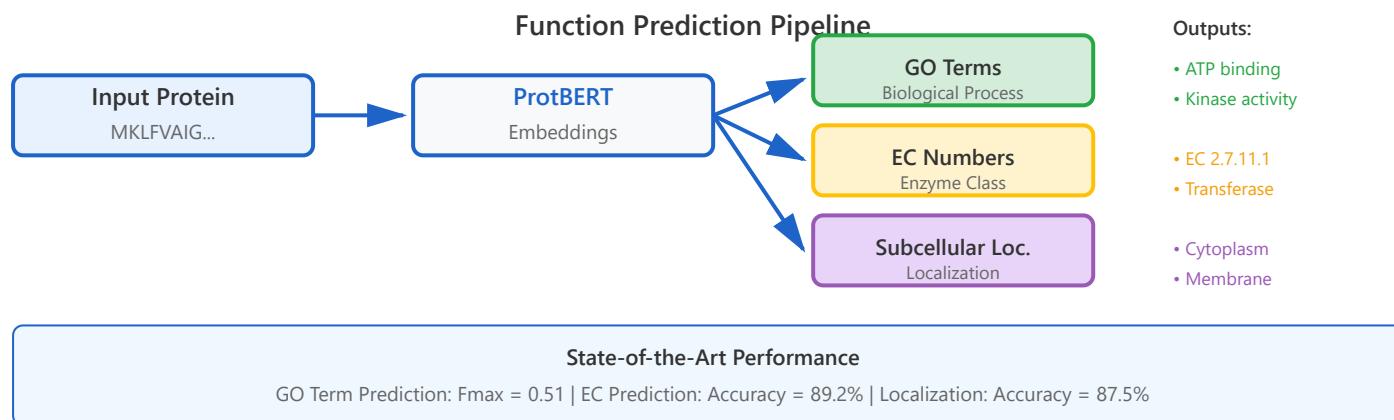
Applications

- **Structure prediction:** Input features for AlphaFold-like models
- **Protein design:** Guide mutations to maintain structure

Function Prediction

Predicting Protein Function

ProtBERT representations can be fine-tuned or used as features for various protein function prediction tasks, including Gene Ontology (GO) term annotation and Enzyme Commission (EC) number classification.



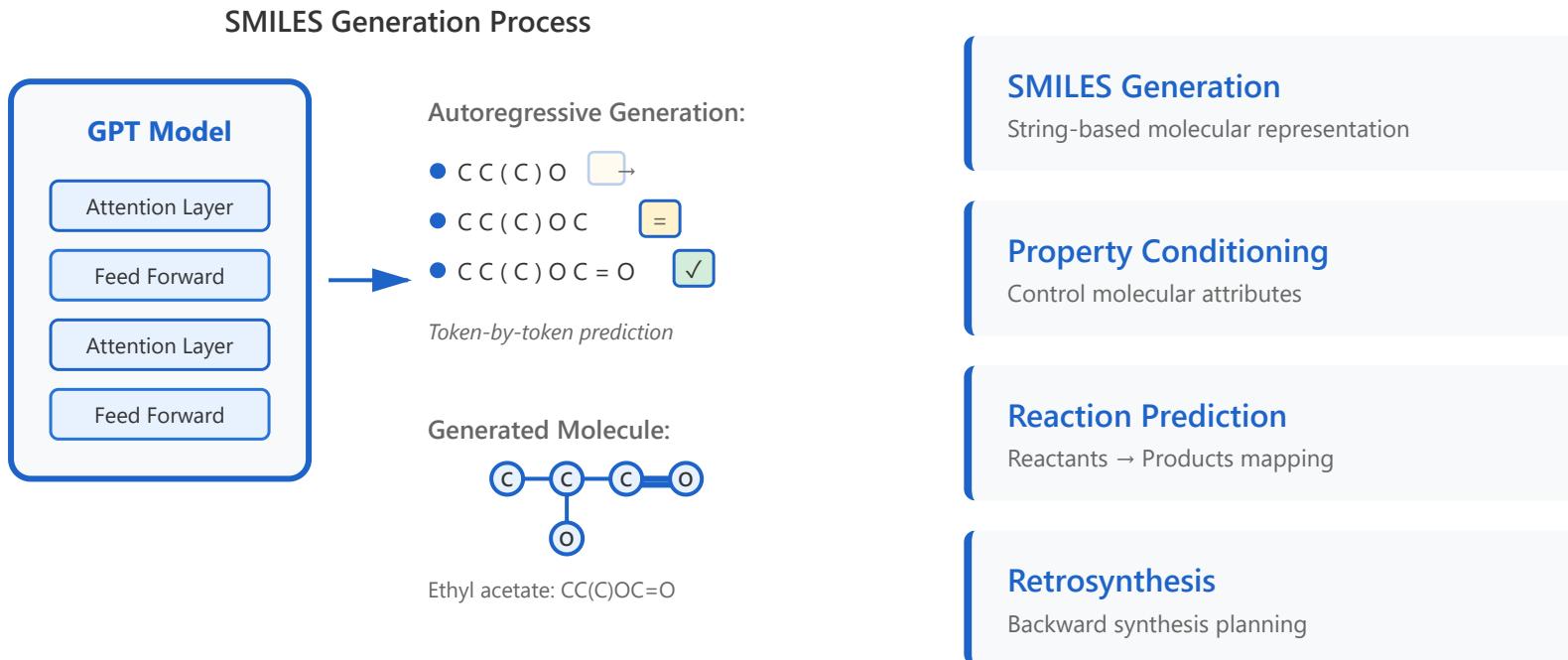
Gene Ontology (GO) Terms

GO annotations describe protein functions in three categories:

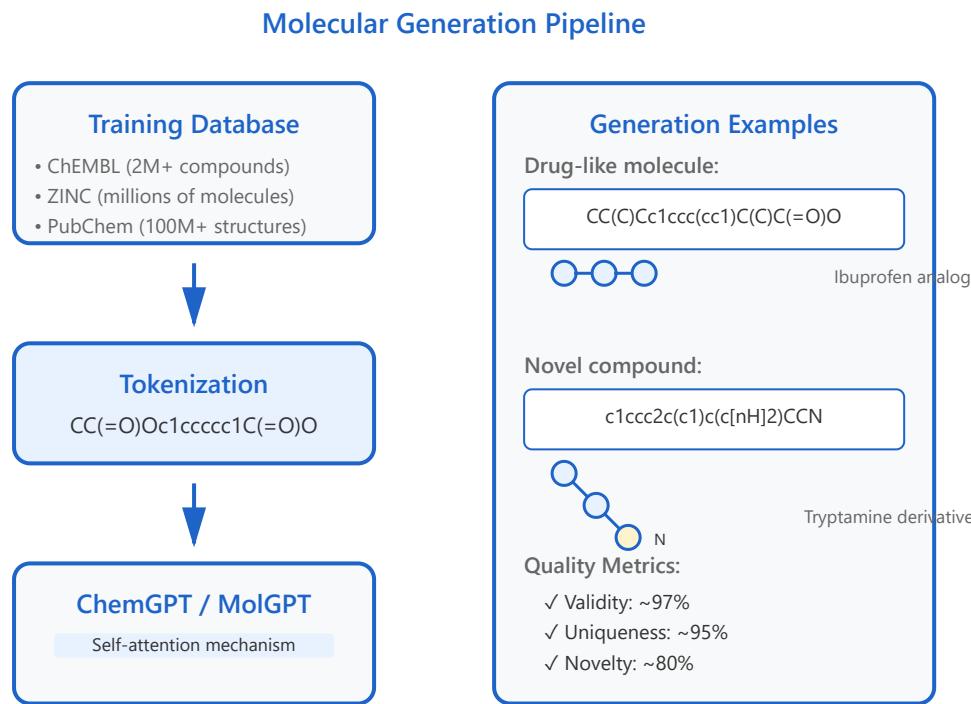
Enzyme Commission (EC) Numbers

EC numbers classify enzyme functions hierarchically:

GPT for Molecules



SMILES Generation with Chemical Language Models



SMILES Representation

SMILES (Simplified Molecular Input Line Entry System) converts molecular structures into sequential text strings, enabling language models to generate valid chemical structures.

Key Features

- ▶ Character-level or token-level encoding of molecular structure
- ▶ Autoregressive generation learns chemical syntax and rules
- ▶ Pre-training on massive molecular databases
- ▶ Fine-tuning for specific property targets

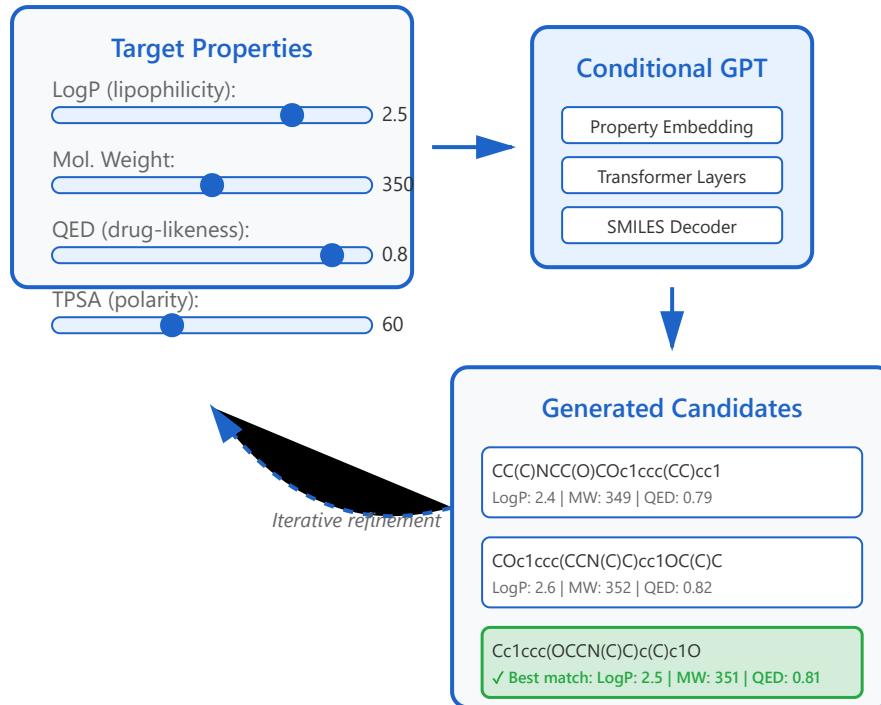
Applications

- ▶ De novo drug design and lead optimization
- ▶ Chemical space exploration
- ▶ Molecule optimization for desired properties

Models: ChemGPT, MolGPT, and SMILES-BERT leverage transformer architectures to understand chemical grammar and generate chemically valid molecules with high success rates.

Property-Conditioned Molecular Generation

Controlled Generation Framework



Controlled Generation

Property-conditioned models generate molecules that satisfy specific physicochemical or biological constraints by integrating target properties into the generation process.

Conditioning Approaches

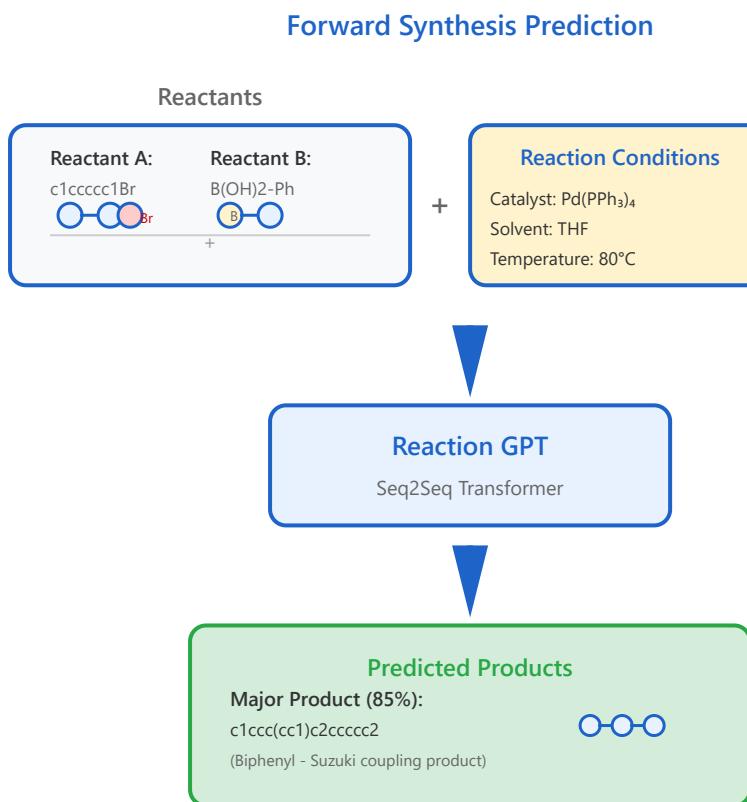
- ▶ **Prefix conditioning:** Property tokens prepended to SMILES
- ▶ **Latent conditioning:** Property embeddings in hidden layers
- ▶ **Reinforcement learning:** Reward-guided optimization
- ▶ **Multi-objective:** Balance multiple property constraints

Target Properties

- ▶ Physicochemical: LogP, molecular weight, TPSA
- ▶ Drug-likeness: QED, Lipinski's rule compliance
- ▶ Biological activity: Binding affinity, selectivity
- ▶ ADMET: Solubility, permeability, toxicity

Impact: Enables rational drug design by generating molecules with optimized pharmacokinetic properties, reducing experimental screening costs and accelerating lead discovery.

Chemical Reaction Prediction



Forward Reaction Prediction

GPT models learn to predict reaction outcomes by training on millions of reaction examples, mapping reactants and conditions to products using sequence-to-sequence architectures.

Model Architecture

- ▶ **Input:** Reactant SMILES + reaction conditions + reagents
- ▶ **Encoder:** Processes reactant molecular structure
- ▶ **Decoder:** Generates product SMILES sequentially
- ▶ **Attention:** Focuses on reactive sites and functional groups

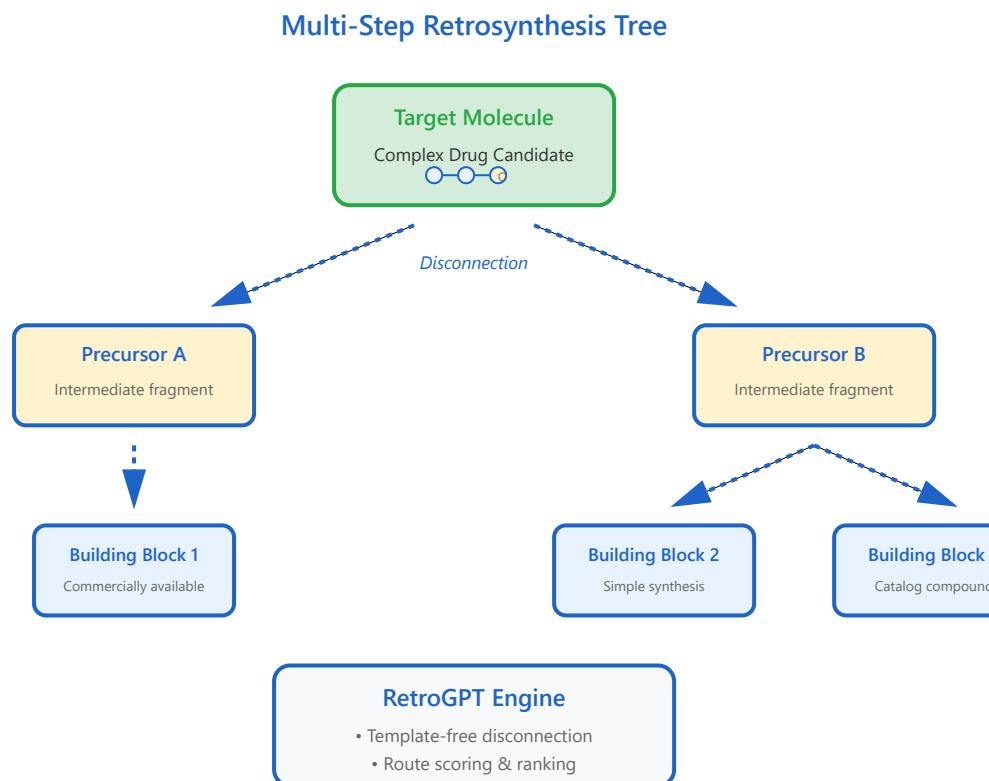
Key Capabilities

- ▶ Named reaction prediction (e.g., Suzuki, Grignard, Diels-Alder)
- ▶ Regioselectivity and stereochemistry prediction
- ▶ Side product and byproduct identification
- ▶ Yield estimation and reaction feasibility

Training Data: Models trained on USPTO (United States Patent and Trademark Office) dataset containing millions of experimentally validated

reactions, achieving >90% top-1 accuracy for common reaction types.

Retrosynthetic Planning with GPT Models



Retrosynthetic Analysis

Retrosynthesis works backwards from target molecules to identify synthetic routes using simpler, commercially available starting materials. GPT models automate this complex planning process.

Model Approaches

- ▶ **Template-free:** Direct SMILES transformation without predefined reaction rules
- ▶ **Template-based:** Apply learned reaction templates from databases
- ▶ **Hybrid:** Combine both approaches for robust predictions
- ▶ **Multi-step planning:** Build complete synthesis trees

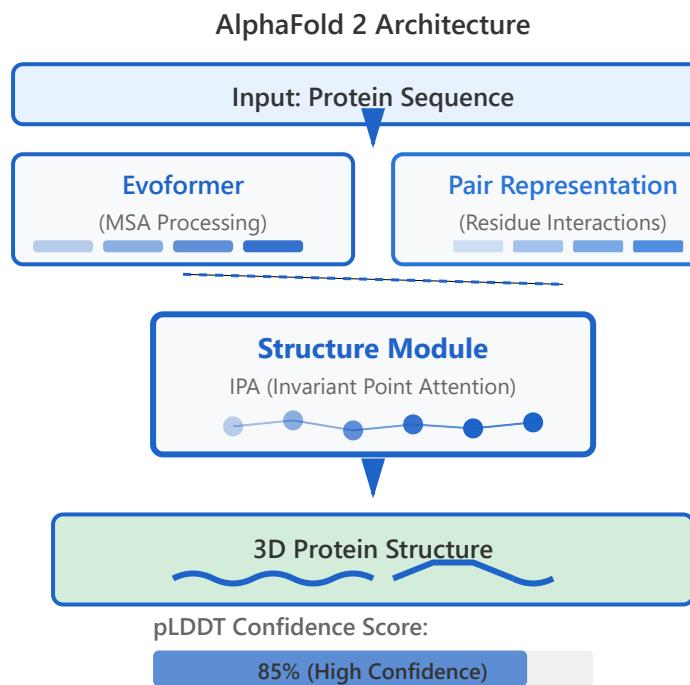
Key Features

- ▶ Automated disconnection site identification
- ▶ Route feasibility scoring and ranking
- ▶ Cost and availability optimization
- ▶ Stereoselective synthesis planning

Applications: Retrosynthetic GPT models like Molecular Transformer and RetroGPT are used in pharmaceutical companies to accelerate drug discovery, reducing synthesis planning time from

weeks to hours while suggesting novel synthetic routes.

AlphaFold Revolution



Architecture innovations

Evoformer + Structure module

MSA processing

Evolutionary information extraction

Structure module

IPA: SE(3)-equivariant attention

Confidence metrics

pLDDT per-residue scores

Database impact

200M+ structures predicted

1. Architecture Innovations

AlphaFold 2 introduced a revolutionary neural network architecture that combines the **Evoformer** module for processing evolutionary information with the **Structure**

Module for generating 3D coordinates. This two-stage pipeline represents a paradigm shift in protein structure prediction.

Evoformer Block

Processes multiple sequence alignments (MSA) and pairwise residue representations through 48 stacked blocks. Each block contains row/column attention mechanisms and transition layers that refine evolutionary patterns.

End-to-End Differentiable

Unlike traditional template-based methods, AlphaFold 2 is trained end-to-end, allowing gradient flow from 3D structure prediction back to sequence processing, enabling sophisticated feature learning.

Iterative Refinement

The structure module operates iteratively, refining the predicted structure through multiple cycles while maintaining geometric consistency through SE(3)-equivariant operations.

Two-Stage Architecture

Sequence + MSA

Evoformer (48 blocks)

MSA + Pair Processing

Structure Module (8 iterations)

IPA + Frame Updates

3D Structure

2. MSA Processing - Evolutionary Information

The **Multiple Sequence Alignment (MSA)** is crucial for AlphaFold's success. By analyzing thousands of related protein sequences from different species, the model extracts

evolutionary constraints that reveal which residues co-evolve, indicating structural proximity.

Evolutionary Co-variation

When two positions in a protein consistently mutate together across species, they are likely in close spatial proximity. AlphaFold learns these co-evolution patterns through MSA row and column attention.

MSA Representation

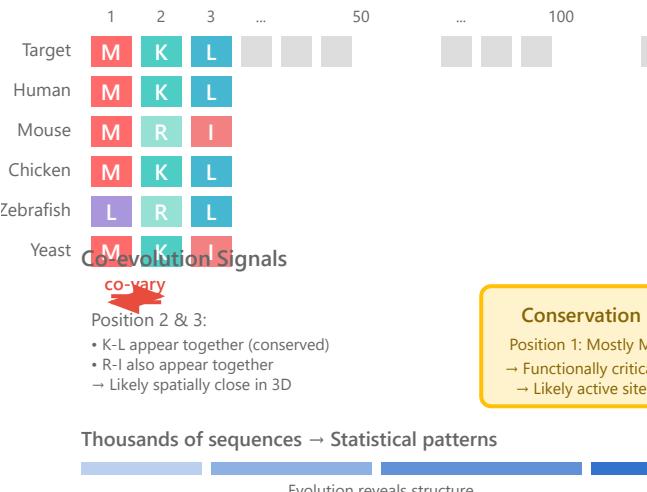
Each MSA row represents a homologous sequence. The Evoformer processes this matrix with specialized attention mechanisms that communicate both within sequences (row) and across positions (column).

Database Search

AlphaFold searches large databases (UniRef90, BFD, MGnify) to find homologous sequences, typically gathering thousands of related proteins to build a comprehensive evolutionary profile.

MSA: Evolutionary Patterns

Multiple Sequence Alignment



3. Structure Module - IPA & SE(3)-Equivariance

The Structure Module is AlphaFold's breakthrough component, featuring **Invariant Point Attention (IPA)**. This mechanism operates directly in 3D space while maintaining SE(3)-equivariance, meaning it respects rotations and translations of the protein structure.

Invariant Point Attention (IPA)

IPA computes attention in both the pair representation space and 3D coordinate space simultaneously. It measures geometric distances between points on local frames, making it rotation/translation invariant.

Local Reference Frames

Each residue has a local coordinate frame (backbone atoms N, C α , C). The structure module updates both frame orientations and translations iteratively, building the full 3D structure progressively.

Geometric Reasoning

Unlike previous methods that predict distance matrices, AlphaFold directly generates 3D coordinates. This enables natural modeling of chirality, angles, and other geometric constraints inherent to protein structures.

Structure Module: IPA Mechanism

Local Reference Frames



IPA Computation

$$\text{Attention} = f(\text{Pair Features}, \text{3D Geometry})$$

Pair Representation
 e_{ij} (learned features)

3D Point Distances
 $\|T_i(p) - T_j(q)\|^2$

SE(3)-Equivariance Property

Rotate/Translate Input →

The same transformation applies to output



4. Confidence Metrics - pLDDT Scores

AlphaFold provides **per-residue confidence scores (pLDDT)**

that indicate how reliable each predicted atom position is.

These scores are crucial for researchers to assess which parts of the structure are trustworthy and which regions might be disordered or incorrectly predicted.

pLDDT Definition

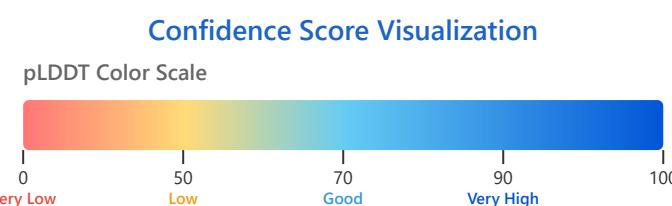
Predicted Local Distance Difference Test (pLDDT) scores range from 0-100, predicting the expected accuracy of C α atom positions. Scores >90 indicate very high confidence, 70-90 indicate good confidence, 50-70 indicate low confidence, and <50 indicate very low confidence.

Interpretation Guide

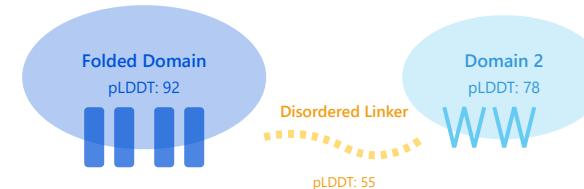
High pLDDT regions (blue) typically represent well-folded domains with strong evolutionary constraints. Low pLDDT regions (yellow/red) often correspond to disordered regions, linkers, or areas with insufficient evolutionary information.

PAE (Predicted Aligned Error)

AlphaFold also provides PAE matrices showing confidence in relative positions between residues. This is especially useful for multi-domain proteins to assess domain-domain orientations.



Example: Protein Structure colored by pLDDT



Interpretation Guidelines

- pLDDT > 90: Highly accurate backbone & side chains
- pLDDT 70-90: Accurate backbone, some side chain error
- pLDDT 50-70: Low confidence, possible disorder
- pLDDT < 50: Should not be interpreted

5. Database Impact - 200M+ Structures

AlphaFold has transformed structural biology by predicting **over 200 million protein structures**, covering nearly every known protein. The AlphaFold Protein Structure Database (AlphaFold DB) provides free access to these predictions, democratizing structural data for researchers worldwide.

Coverage Scale

Before AlphaFold: ~170,000 experimentally determined structures in PDB (50 years of work). After AlphaFold: 200+ million predicted

structures covering most of UniProt, representing a 1000x increase in structural knowledge.

Database Growth & Impact

Structural Coverage Timeline

Pre-AlphaFold (1970-2020)

PDB: ~170,000 structures

50 years of experimental work

170K

Post-AlphaFold (2021-2024)

AlphaFold DB: 200+ million structures

Nearly complete UniProt coverage

1000x increase in 3 years



Research Acceleration

Researchers can now instantly access predicted structures instead of waiting months/years for experimental determination. This has accelerated drug discovery, protein engineering, and fundamental biology research across all fields.

Organism Coverage

AlphaFold DB includes proteomes from model organisms (human, mouse, E. coli, yeast), plants, parasites, and environmental microbes. This enables comparative structural biology and evolutionary studies at unprecedented scale.

Research Impact Areas

Drug Discovery

Target identification
Binding site analysis

Protein Engineering

Rational design
Stability optimization

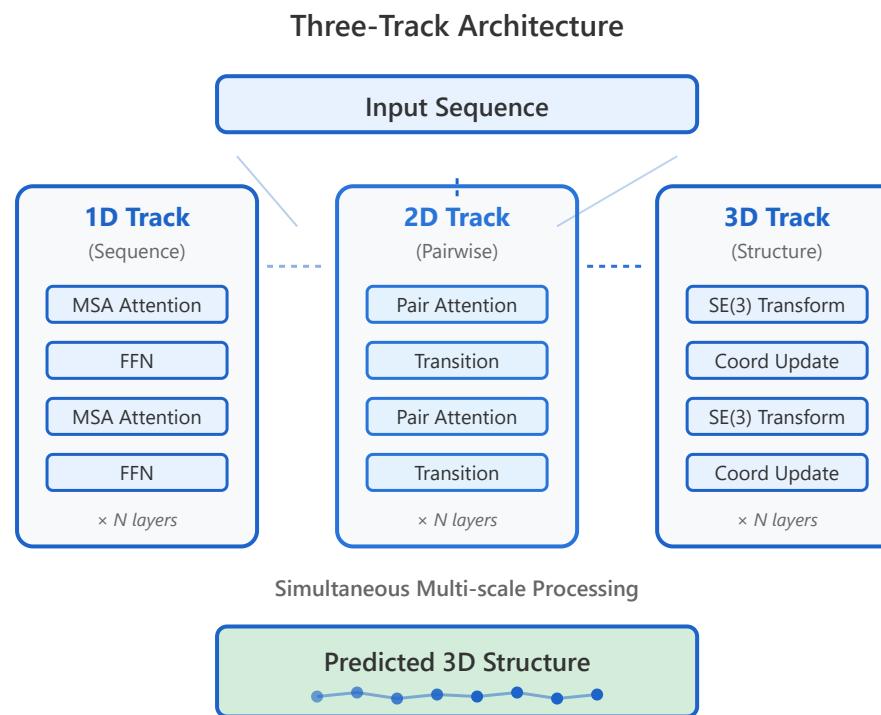
Disease Research

Mutation analysis
Pathway understanding

Evolution Studies

Comparative structures
Function prediction

RoseTTAFold: Accurate Protein Structure Prediction



Three-track architecture

1D, 2D, 3D parallel processing

End-to-end learning

Direct structure prediction

Complex prediction

Protein-protein interactions

Speed advantages

Faster than AlphaFold2

Applications

Structure, function, design

1

Three-Track Architecture

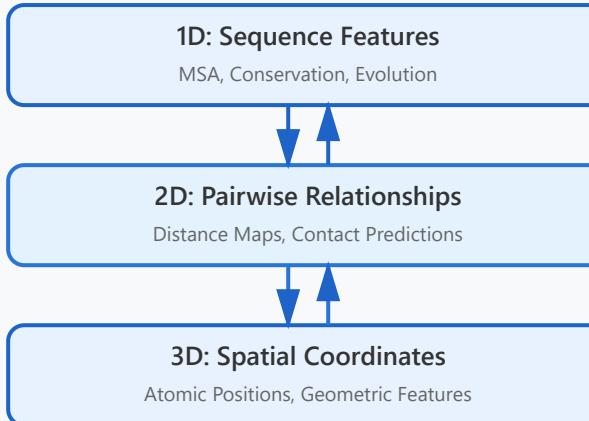
RoseTTAFold employs a unique three-track neural network architecture that simultaneously processes protein information at three different scales: 1D (sequence), 2D (pairwise distances), and 3D (coordinates). This parallel processing approach allows the model to capture multi-scale patterns and relationships that are crucial for accurate structure prediction.

- ▶ **1D Track (Sequence):** Processes multiple sequence alignments (MSA) to capture evolutionary information and identify conserved residues across homologous proteins
- ▶ **2D Track (Pairwise):** Models residue-residue relationships and distance constraints, capturing local and long-range interactions between amino acids
- ▶ **3D Track (Structure):** Directly operates on 3D coordinates using SE(3)-equivariant transformations, ensuring geometric consistency
- ▶ **Information Exchange:** The three tracks communicate bidirectionally at each layer, allowing features to flow between different representations

Key Innovation

Unlike traditional methods that process information sequentially, RoseTTAFold's parallel architecture enables simultaneous refinement across all three levels, leading to more coherent and accurate predictions.

Information Flow



Continuous Information Exchange

2

End-to-End Learning

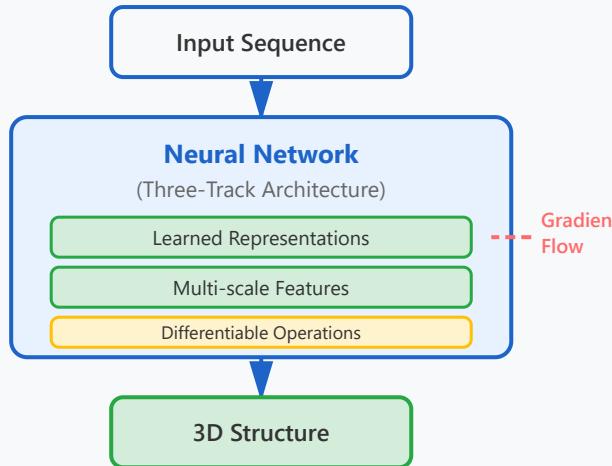
RoseTTAFold implements a fully differentiable end-to-end learning framework that directly maps from protein sequences to 3D structures without requiring intermediate steps or template-based modeling. This approach enables the network to learn complex structure-function relationships directly from data.

- ▶ **Direct Prediction:** Eliminates the need for fragment assembly or template-based modeling, which were standard in earlier methods
- ▶ **Gradient Flow:** Backpropagation flows through all three tracks simultaneously, allowing for holistic optimization
- ▶ **Learned Representations:** The model automatically learns relevant features at each level rather than relying on hand-crafted features
- ▶ **Structure Module:** Final layers convert learned representations directly into 3D atomic coordinates with associated confidence scores

Advantage

End-to-end learning allows the model to optimize for the final structural output rather than intermediate objectives, resulting in more accurate and physically realistic predictions.

End-to-End Pipeline



3

Complex Prediction

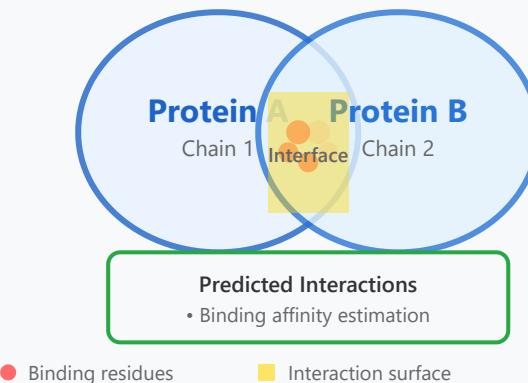
One of RoseTTAFold's most powerful capabilities is predicting protein-protein interaction complexes and multi-chain assemblies. By treating multiple chains simultaneously during inference, the model can capture inter-chain contacts, binding interfaces, and quaternary structure arrangements.

- ▶ **Multi-chain Support:** Processes multiple protein chains simultaneously, capturing inter-molecular interactions
- ▶ **Interface Prediction:** Accurately identifies binding sites and interaction surfaces between protein partners
- ▶ **Oligomer Assembly:** Can model homo- and hetero-oligomeric structures, including antibody-antigen complexes
- ▶ **Functional Insights:** Complex structures reveal mechanisms of protein function, regulation, and signaling pathways

Application Example

RoseTTAFold has successfully predicted structures of antibody-antigen complexes, enzyme-substrate interactions, and large multi-protein assemblies like the SARS-CoV-2 spike protein-ACE2 receptor complex.

Protein Complex Prediction



Multi-chain simultaneous prediction

4

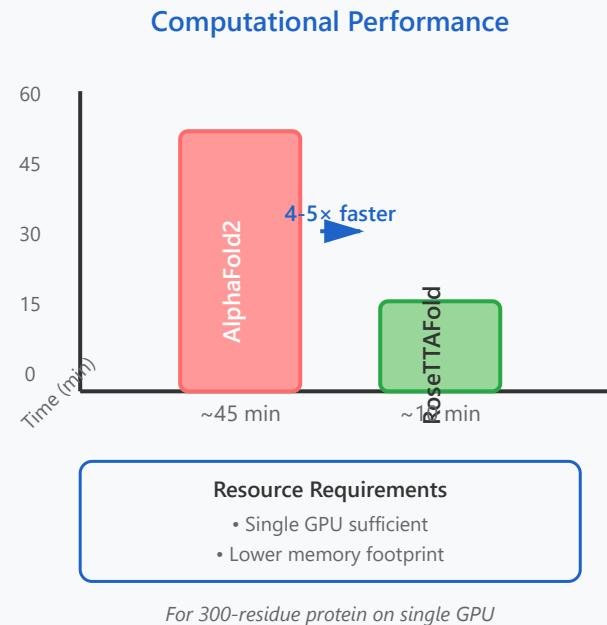
Speed Advantages

RoseTTAFold achieves significantly faster prediction times compared to AlphaFold2, making it practical for large-scale structural genomics projects and real-time applications. The speed improvement comes from architectural optimizations and efficient implementation without sacrificing accuracy.

- ▶ **Computational Efficiency:** Requires fewer computational resources and GPU memory compared to AlphaFold2
- ▶ **Faster Inference:** Typical predictions complete in minutes rather than hours for medium-sized proteins
- ▶ **Scalable Architecture:** Can process multiple proteins in parallel batches efficiently
- ▶ **Resource Accessibility:** Lower computational requirements make it accessible to more researchers without high-end computing clusters

Performance Comparison

For a 300-residue protein, RoseTTAFold typically completes prediction in 5-10 minutes on a single GPU, while maintaining accuracy comparable to AlphaFold2. This speed enables high-throughput screening of entire proteomes.

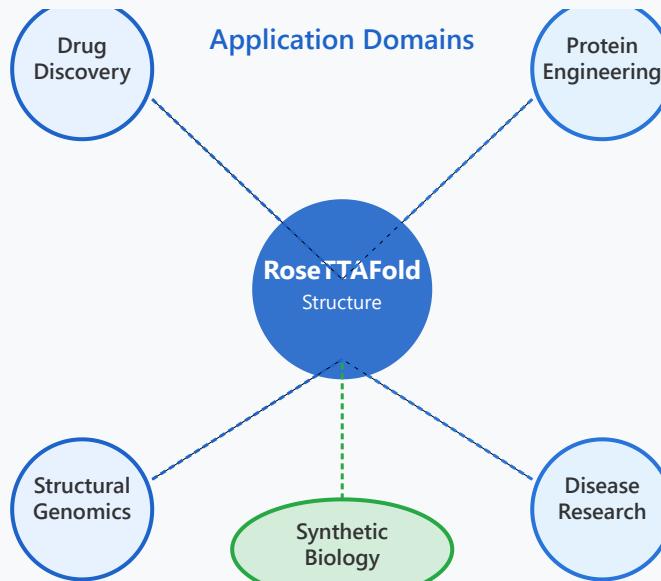


5

Applications

RoseTTAFold's capabilities extend across multiple domains in structural biology, drug discovery, and protein engineering. Its accuracy, speed, and ability to handle complex structures make it a versatile tool for both fundamental research and practical applications.

- ▶ **Drug Discovery:** Identifying binding pockets, predicting drug-target interactions, and virtual screening for therapeutic candidates
- ▶ **Protein Engineering:** Guiding rational design of proteins with enhanced stability, altered specificity, or novel functions
- ▶ **Structural Genomics:** Large-scale prediction of protein structures for entire genomes, filling gaps in structural databases
- ▶ **Disease Research:** Understanding structural basis of genetic diseases, identifying pathogenic variants, and designing therapeutic interventions
- ▶ **Synthetic Biology:** Designing novel protein folds, creating artificial enzymes, and engineering biosynthetic pathways



Real-world Impact

RoseTTAFold has been used to predict structures of orphan proteins, design COVID-19 therapeutic candidates, engineer novel enzymes for industrial applications, and accelerate

vaccine development by modeling antibody-antigen interactions.

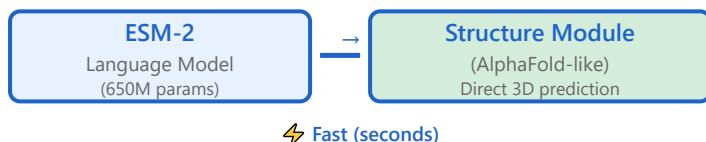
ESMFold

Language Model-Only Approach

Traditional (e.g., AlphaFold2):



ESMFold:



Key Innovation: No MSA Required

- Evolutionary info learned directly from 250M+ protein sequences
- 60× faster than AlphaFold2 (seconds vs minutes)
- Enables metagenomic-scale structure prediction

Language model only

ESM-2 pretrained transformer

No MSA required

Single sequence input

Speed benefits

60× faster inference

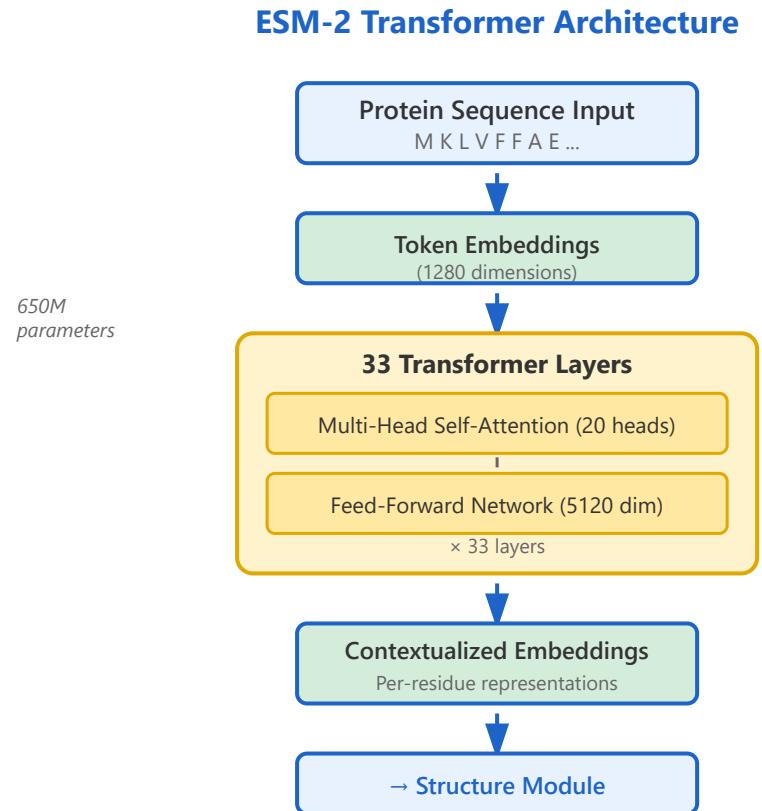
Metagenomic applications

Unknown protein discovery

Limitations

Lower accuracy on orphan proteins

1. Language Model Only: ESM-2 Architecture



ESM-2 Language Model

A transformer-based protein language model trained on 250 million protein sequences from UniRef. Uses masked language modeling to learn evolutionary patterns and structural constraints directly from sequence data.

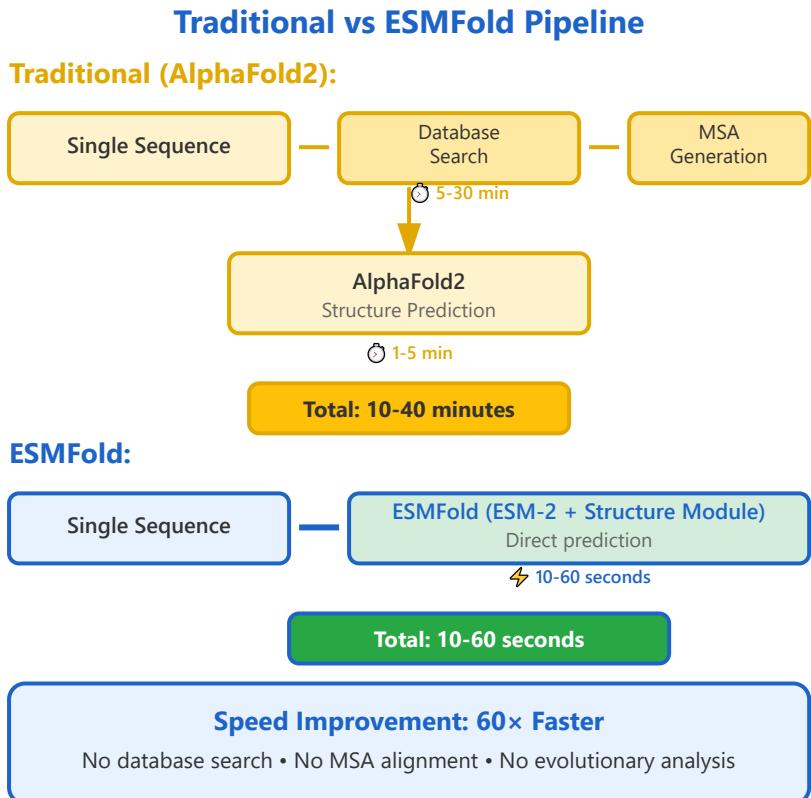
Model Architecture

- **33 transformer layers** with 20 attention heads each
- **1280-dimensional** residue embeddings
- **5120-dimensional** feed-forward layers
- **650 million** total parameters

Key Advantage

Captures evolutionary information implicitly through pretraining, eliminating the need for explicit MSA generation at inference time. This makes it drastically faster while maintaining competitive accuracy.

2. No MSA Required: Direct Sequence-to-Structure



What is MSA?

Multiple Sequence Alignment (MSA) aligns homologous protein sequences to identify conserved and variable regions, revealing evolutionary patterns critical for structure prediction.

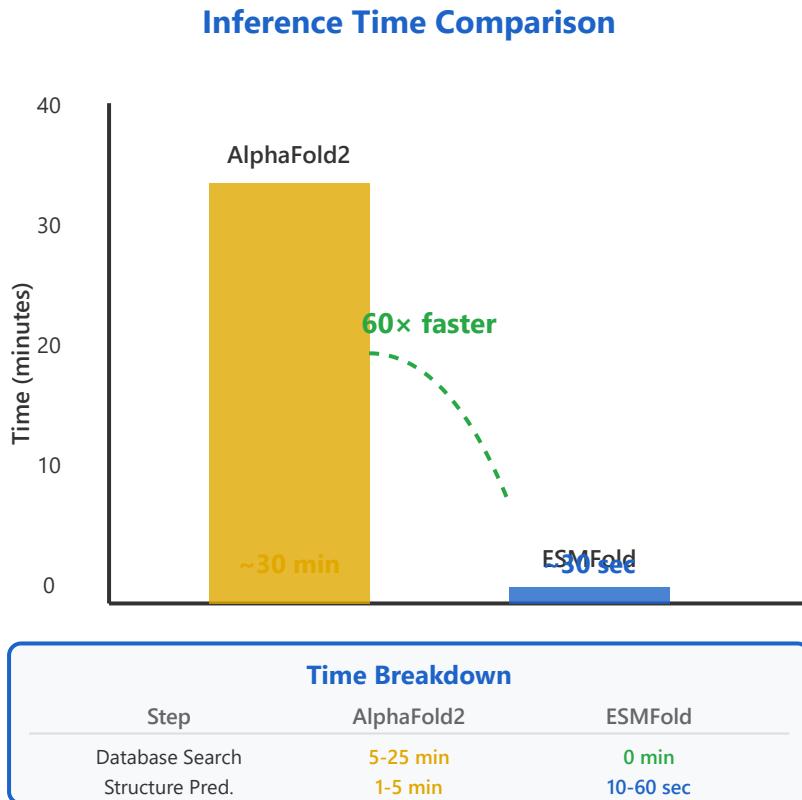
Why Skip MSA?

- **Database search** requires 5-30 minutes
- **Computational bottleneck** for large-scale predictions
- **Fails for orphan proteins** with no homologs
- **Not scalable** to metagenomic datasets

ESMFold's Solution

ESM-2 learns evolutionary patterns during pretraining on millions of sequences, embedding this knowledge directly into the model weights. At inference, only the single input sequence is needed.

3. Speed Benefits: Performance Comparison



Performance Metrics

60× faster than AlphaFold2

~30 seconds per protein (length 384)

No GPU required for database search

Linear scaling with sequence length

Throughput Advantage

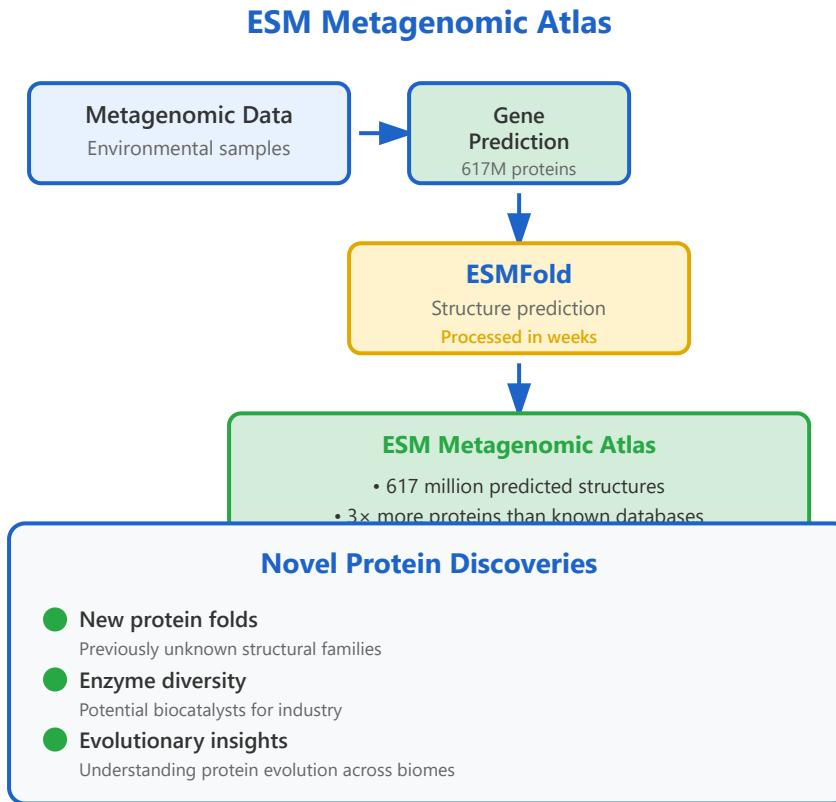
ESMFold can predict structures for thousands of proteins per day on a single GPU, enabling:

- Large-scale proteome analysis
- Real-time structure prediction in workflows
- High-throughput screening applications

⚡ Key Insight

The elimination of MSA generation removes the computational bottleneck, making ESMFold suitable for applications requiring rapid turnaround or processing millions of sequences.

4. Metagenomic Applications: Discovering Unknown Proteins



Why Metagenomics?

Metagenomic sequencing reveals millions of uncultured organisms from environmental samples (soil, ocean, human gut), containing proteins with no known homologs - perfect for ESMFold.

Scale Achievement

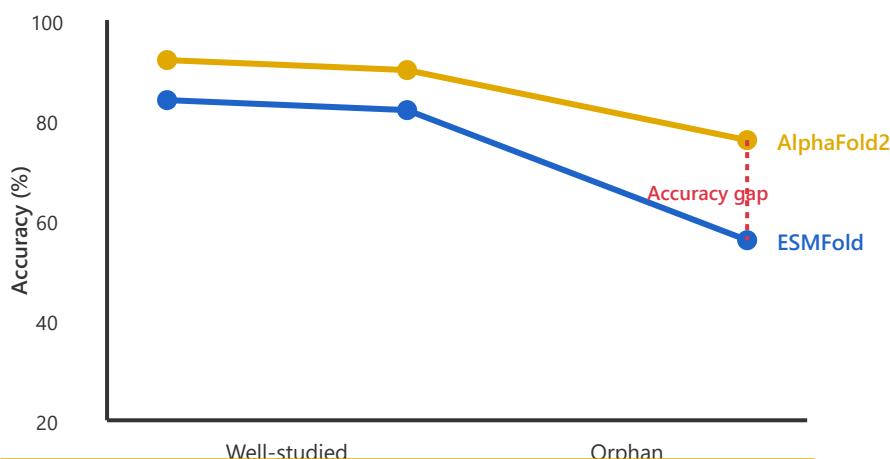
- **617 million structures** predicted
- **3x larger** than all known protein databases
- **Only possible with ESMFold** - AlphaFold2 would take decades
- **Public database** available for research

Impact Areas

Drug discovery, enzyme engineering, understanding microbial ecology, identifying novel antibiotic targets, and mapping the functional protein universe.

5. Limitations: Understanding the Trade-offs

Accuracy Comparison



Key Limitations

- Lower accuracy for orphan proteins
Proteins with few/no homologs show 5-10% lower accuracy
- Less reliable confidence scores
pLDDT scores less calibrated than AlphaFold2's
- Not ideal for protein complexes
Designed for single-chain predictions only

Orphan Proteins Challenge

Proteins without known homologs benefit most from MSA-based methods. ESMFold relies on patterns learned during pretraining, which may not cover rare protein families adequately.

When to Use Each Method

Use AlphaFold2 when:

- Maximum accuracy is critical
- Predicting protein complexes
- Working with orphan proteins

Use ESMFold when:

- Speed is essential
- Processing large datasets
- Working with metagenomic data

The Trade-off

ESMFold sacrifices 5-10% accuracy on difficult targets for a 60× speedup. For most applications, especially large-scale studies, this is an excellent trade-off.

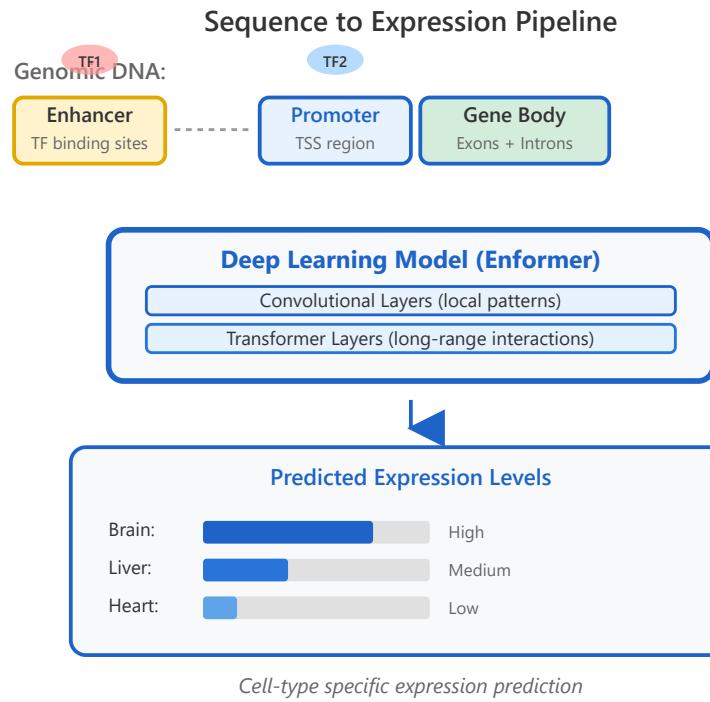
Part 2/3 - Biological AI

Predictive models

Interpretable AI

Biological insights

Gene Expression Prediction



Sequence to expression

DNA → RNA abundance mapping

Promoter models

TSS region activity prediction

Enhancer grammar

TF binding syntax learning

Cell type specificity

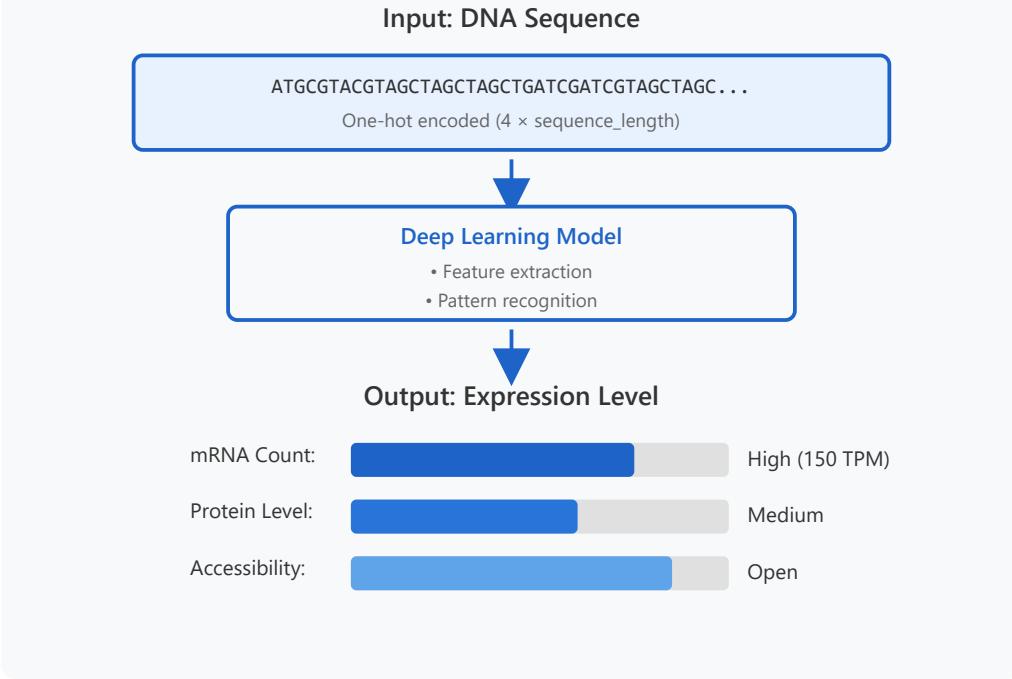
Context-dependent prediction

Enformer architecture

Transformer + CNN hybrid model

1. Sequence to Expression

Mapping DNA sequences to RNA abundance levels



Overview

Sequence-to-expression models predict gene activity directly from DNA sequence. These models learn the complex regulatory code that determines when and where genes are expressed.

Input Features

Models take raw DNA sequences as input, typically encoded as one-hot vectors representing the four nucleotides (A, T, G, C). The sequence context can span from hundreds to hundreds of thousands of base pairs.

Output Predictions

Models output quantitative predictions of gene expression levels, which can include mRNA abundance (e.g., TPM, FPKM), chromatin accessibility (e.g., ATAC-seq), histone modifications (e.g., ChIP-seq), or protein levels.

Key Applications & Challenges

- ▶ **Variant Effect Prediction:** Predict how genetic variants affect gene expression (eQTLs)
- ▶ **Therapeutic Design:** Design synthetic regulatory elements for gene therapy
- ▶ **Disease Mechanisms:** Understand regulatory disruptions in disease states
- ▶ **Challenge:** Long-range interactions can span megabases, requiring large context windows
- ▶ **Challenge:** Cell-type specific regulation requires integrated models with epigenetic features

2. Promoter Models

Predicting transcription start site (TSS) region activity

Promoter Architecture

Promoters are regulatory DNA regions located upstream of genes, typically spanning ~1kb around the transcription start site (TSS). They contain core elements (TATA box, Initiator, DPE) and binding sites for general transcription factors and RNA polymerase II.

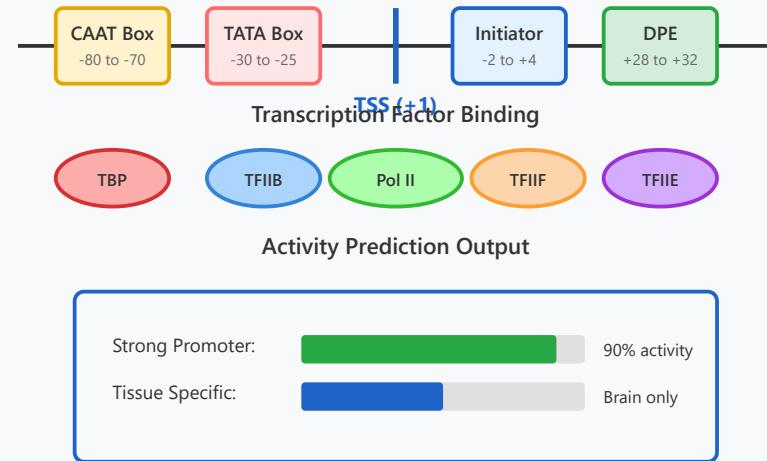
Model Approaches

Early models focused on position weight matrices (PWMs) for transcription factor binding sites. Modern deep learning approaches use CNNs to automatically learn motifs and their combinations, capturing complex syntax rules.

Prediction Tasks

Models predict TSS activity strength, directionality, tissue-specificity, and response to transcription factors. Advanced models can predict effects of promoter mutations and design synthetic promoters with desired properties.

Promoter Architecture



Key Insights & Applications

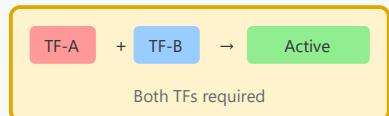
- Core Promoter Elements: TATA box, Initiator (lnr), and Downstream Promoter Element (DPE) determine basal transcription
- Proximal Elements: CAAT box and GC box enhance promoter activity
- Synthetic Biology: Design optimized promoters for gene expression systems
- Disease Variants: Predict how mutations in promoter regions affect gene expression
- Recent Models: ProCapNet, Xpresso, and ExPecto achieve high accuracy on human promoters

3. Enhancer Grammar

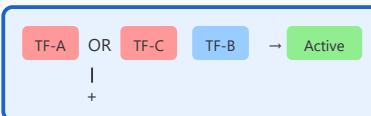
Learning the syntax of transcription factor binding

Enhancer Regulatory Logic

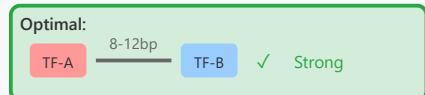
Simple: AND Logic



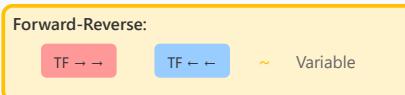
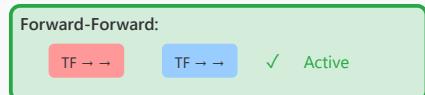
Complex: OR + AND Logic



Spacing Constraints Matter



Orientation Dependence



Deep Learning Discovers Grammar

Convolutional Neural Network

- Automatically learns motifs and their combinations
- Captures spacing, orientation, and order preferences
- Generalizes to predict activity of novel sequences

What is Enhancer Grammar?

Enhancer grammar refers to the rules governing how transcription factor binding sites combine to produce regulatory activity. Like linguistic grammar, it involves syntax (arrangement), semantics (meaning), and context-dependence.

Combinatorial Logic

Enhancers integrate signals from multiple transcription factors through boolean-like logic gates. Common patterns include AND gates (both TFs required), OR gates (either TF sufficient), and NOT gates (repressive interactions).

Spatial Constraints

The spacing and orientation between TF binding sites critically affects function. Optimal spacing allows protein-protein interactions, while poor spacing disrupts cooperative binding. Different TF pairs have characteristic preferred spacings.

Learning Grammar Rules

Deep learning models, particularly CNNs, automatically discover enhancer grammar from sequence and activity data. They learn motifs, their combinations, and higher-order syntax without explicit feature engineering.

Key Concepts & Applications

- ▶ **Motif Interactions:** TF binding sites work cooperatively or antagonistically based on their arrangement
- ▶ **Flexible Grammar:** Same TFs can produce different outputs depending on context and arrangement
- ▶ **Evolutionary Conservation:** Grammar rules are often conserved across species, indicating functional importance
- ▶ **Therapeutic Applications:** Design synthetic enhancers with predictable, cell-type specific activity
- ▶ **Key Models:** DeepSEA, Basset, and ChromBPNet excel at learning enhancer grammar

4. Cell Type Specificity

Context-dependent gene expression prediction

The Challenge

All cells in an organism share the same genome, yet express vastly different sets of genes. A neuron expresses different genes than a liver cell, despite having identical DNA sequences. This cell-type specificity arises from epigenetic regulation.

Epigenetic Context

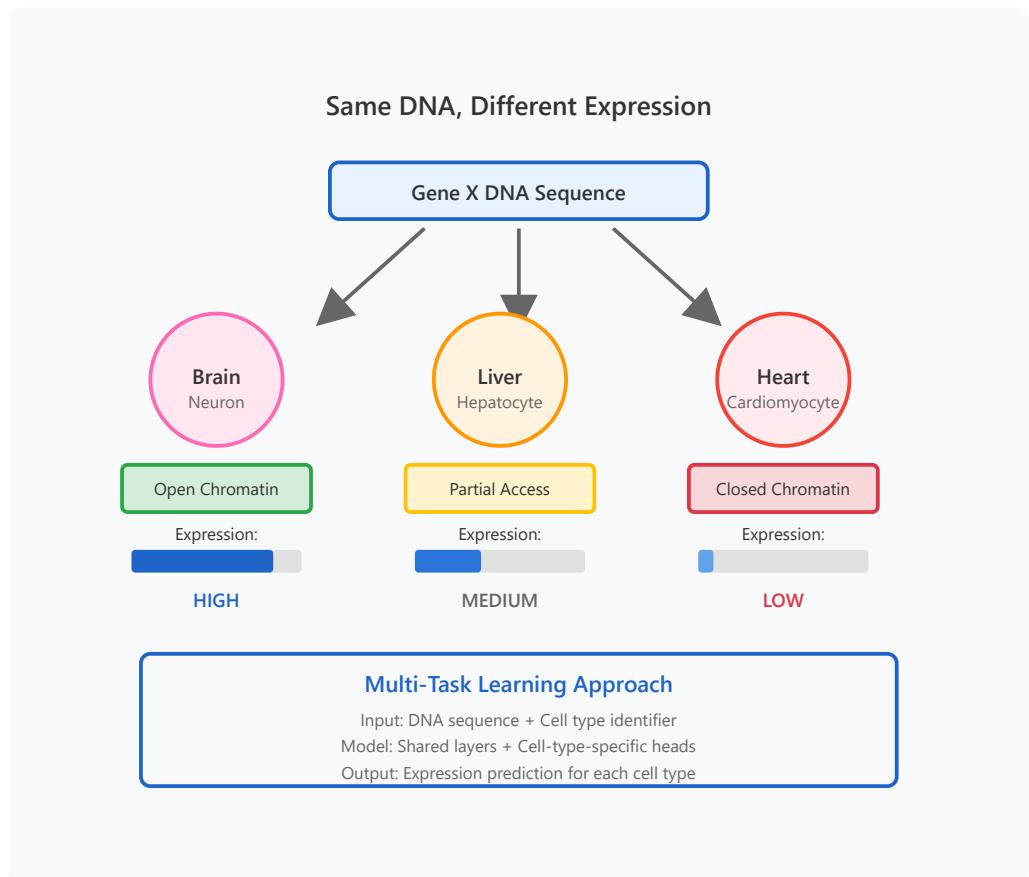
Cell type identity is encoded through chromatin accessibility, DNA methylation, and histone modifications. These epigenetic marks determine which regulatory elements are active in each cell type, creating a unique regulatory landscape.

Modeling Approaches

Modern models incorporate cell-type information through multi-task learning (predicting across cell types simultaneously), conditional models (cell type as input), or through learned cell-type embeddings that capture regulatory state.

Practical Impact

Understanding cell-type specificity enables prediction of how genetic variants affect different tissues, design of cell-type-specific gene therapies, and identification of regulatory elements driving cell fate decisions.



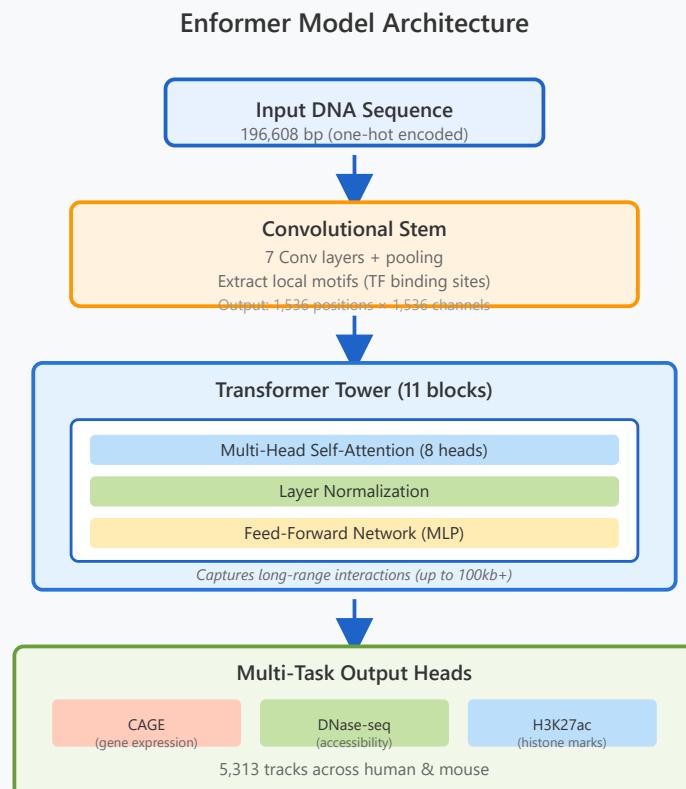
Key Mechanisms & Applications

- ▶ **Master Regulators:** Cell-type-specific transcription factors (e.g., MyoD in muscle, GATA1 in blood) drive expression programs
- ▶ **Chromatin Accessibility:** DNase-seq and ATAC-seq data reveal which regulatory elements are accessible in each cell type
- ▶ **Histone Marks:** H3K4me3 (promoters), H3K27ac (active enhancers), H3K27me3 (repression) mark regulatory states
- ▶ **Clinical Applications:** Predict tissue-specific effects of disease variants (e.g., heart vs brain)

- Notable Models: Basenji predicts cell-type-specific chromatin and expression across 200+ cell types

5. Enformer Architecture

Transformer + CNN hybrid for long-range regulatory prediction



Why Enformer?

Previous models (like Basenji) used only CNNs and were limited to ~40kb context windows. Enformer uses transformers to capture interactions across 200kb, dramatically improving predictions by modeling distal enhancers and TAD structures.

CNN Stem

The convolutional stem processes raw DNA sequence to extract local features like transcription factor binding motifs. This reduces the sequence length while enriching the representation with biologically meaningful patterns.

Transformer Tower

The transformer blocks use self-attention to model long-range interactions between regulatory elements. Unlike CNNs with limited receptive fields, attention can directly connect distant positions, capturing enhancer-promoter loops and chromatin interactions.

Multi-Task Learning

Enformer simultaneously predicts thousands of genomic tracks (CAGE, ChIP-seq, DNase-seq, etc.) across cell types. This multi-task approach helps the model learn generalizable regulatory principles and improves performance through shared representations.

Performance Gains

Enformer achieves state-of-the-art accuracy, explaining ~60% of variance in human expression data. It substantially outperforms previous models on variant effect prediction and can identify regulatory variants missed by GWAS.

Key Innovations in Enformer

1. Long Context Window

196,608 bp input
vs. 40kb in Basenji



Captures distal enhancers

2. Attention Mechanism

Self-attention layers
Direct long-range modeling



Better than convolution

3. Cross-Species Training

Human + Mouse data
Learns conserved principles
✓
Improved generalization

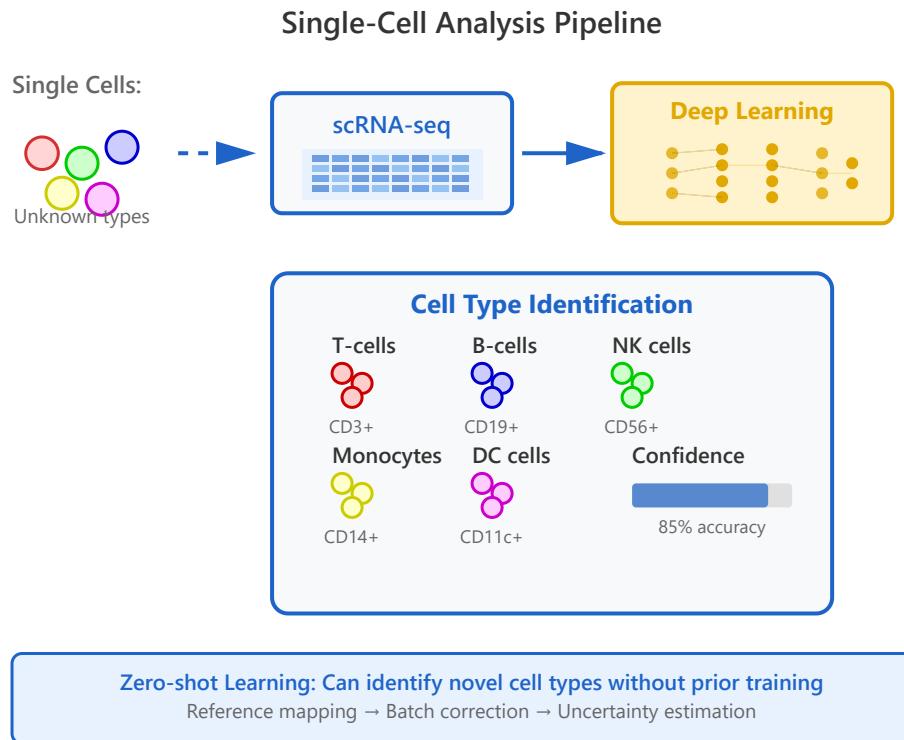
4. Variant Effect Prediction

In silico mutagenesis
Predict regulatory impact
★
Clinical applications

Technical Details & Impact

- ▶ **Parameters:** ~250M parameters, trained on TPUs for several weeks
- ▶ **Training Data:** Thousands of genomic assays from ENCODE, Roadmap Epigenomics, and GTEx
- ▶ **Architecture Benefit:** Attention mechanism provides computational efficiency and better gradient flow than deep CNNs
- ▶ **Interpretability:** Attention weights reveal which genomic regions interact to regulate expression
- ▶ **Applications:** Variant prioritization, synthetic biology, understanding disease mechanisms, drug target identification
- ▶ **Future Directions:** Single-cell predictions, 3D genome structure integration, protein sequence co-modeling

Cell Type Classification



Single-cell models

scBERT, Geneformer architectures

Reference mapping

Atlas-based annotation

Zero-shot learning

Novel cell type discovery

Batch correction

Remove technical variation

Uncertainty estimation

Confidence scoring

1

Single-cell Models

Overview

Single-cell foundation models are deep learning architectures specifically designed to understand and analyze gene expression patterns at the individual cell level. These models leverage transformer-based architectures, similar to those used in natural language processing, to learn meaningful representations of cellular states.

Key Architectures

scBERT (Single-cell BERT): Adapts the BERT architecture for single-cell RNA-seq data, treating genes as "words" and cells as "sentences". The model learns contextual relationships between genes through self-supervised pre-training on large-scale datasets.

Geneformer: A transformer model that processes genes ranked by expression level, enabling the model to capture gene regulatory networks and cellular hierarchies. It can predict cell fate and identify key regulatory genes.

Key Capabilities:

- Learn universal gene expression patterns across millions of cells
- Transfer learning to new datasets with minimal fine-tuning
- Identify cell type-specific gene signatures automatically
- Predict cell states and developmental trajectories

Single-cell Model Architecture

Input: Gene Expression



Transformer Layers

Multi-Head Self-Attention
Learning gene-gene relationships

Feed-Forward Network
Non-linear transformations

Layer Normalization
Stabilizing training

Output: Cell Embeddings



Why It Matters

Foundation models enable accurate cell type classification even with limited labeled data, dramatically reducing the need for manual annotation and improving consistency across different studies.

Overview

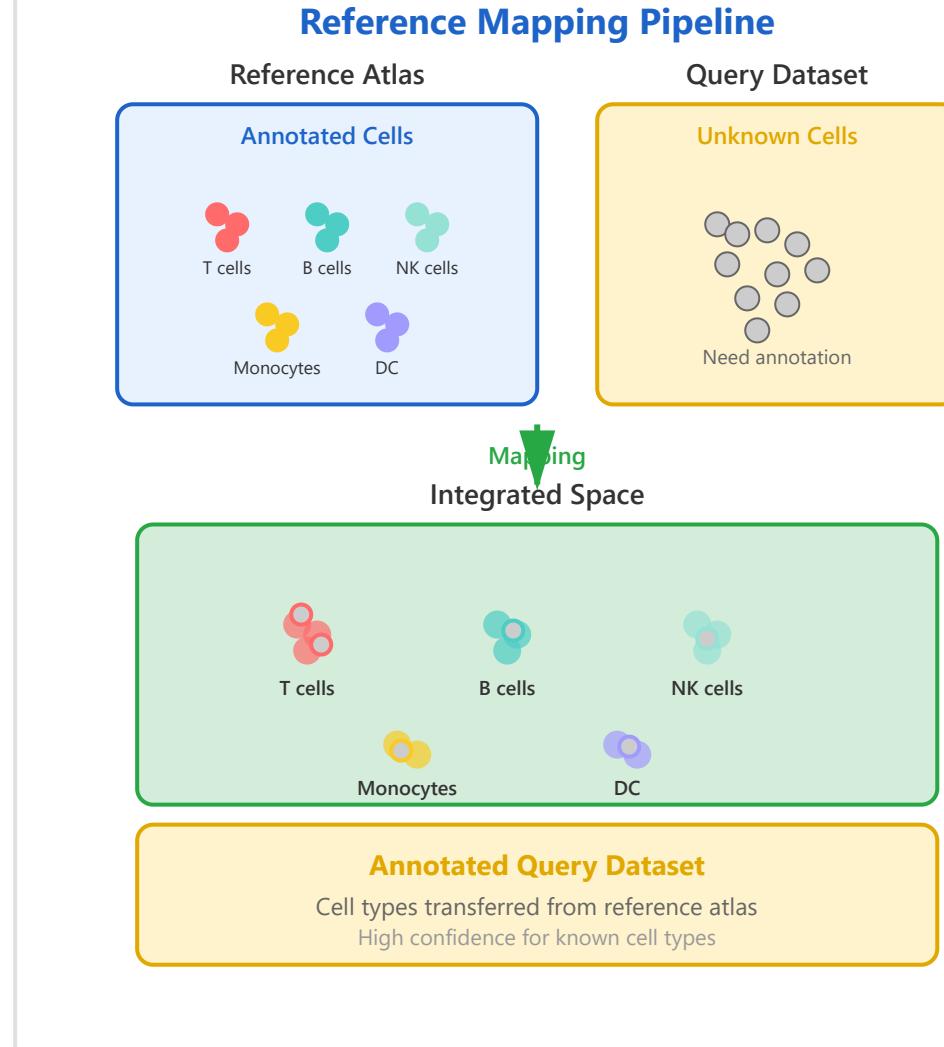
Reference mapping is a transfer learning approach that annotates new single-cell datasets by comparing them to well-characterized reference atlases. This method leverages comprehensive, manually curated cell type annotations from large-scale projects to automatically classify cells in new experiments.

How It Works

The process involves projecting query cells into the same embedding space as reference cells, then transferring labels based on similarity. Advanced methods like Seurat's reference mapping and Symphony use canonical correlation analysis (CCA) and harmony integration to align datasets while preserving biological variation.

Key Advantages:

- ▶ Leverages expert knowledge from reference atlases (e.g., Human Cell Atlas)
- ▶ Consistent annotations across different studies and laboratories
- ▶ Fast inference without requiring model training



- Works well for common, well-characterized cell types

Limitations:

- May struggle with novel or rare cell types not in reference
- Depends on quality and comprehensiveness of reference atlas
- Can be affected by batch effects between datasets

Real-World Application

Reference atlases like Tabula Sapiens contain millions of annotated human cells across 24 tissues, enabling researchers to classify cells in new disease studies rapidly and accurately.

3

Zero-shot Learning

Overview

Zero-shot learning enables AI models to identify and classify cell types that were never seen during training. This revolutionary capability is crucial for discovering novel cell populations, rare cell states, and disease-specific cell types that don't exist in healthy reference atlases.

Mechanism

Zero-shot models learn a semantic embedding space where cell types are represented based on their functional properties and gene expression characteristics rather than explicit labels. The model can recognize new cell types by understanding the relationships between genes and cellular functions, similar to how humans can identify an unfamiliar animal by recognizing its features.

Approaches

Semantic Embeddings: Models learn to associate gene expression patterns with cell type descriptions or functional annotations, enabling classification based on textual descriptions.

Zero-shot Learning Framework

Training Phase

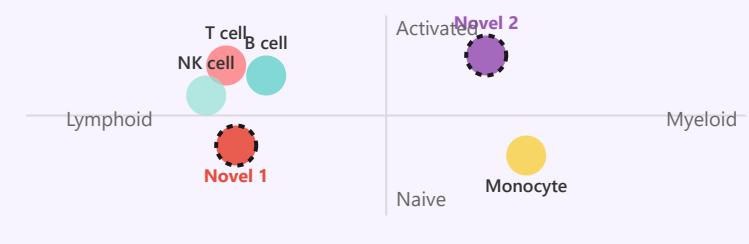
Known Cell Types



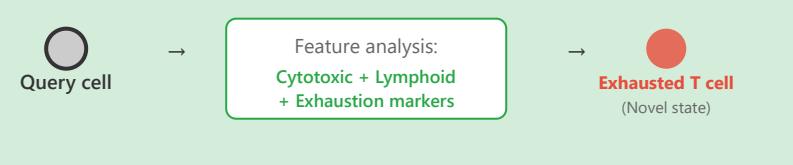
Learned Features

- Cytotoxicity markers
- Activation states
- Differentiation stages

Semantic Embedding Space



Zero-shot Inference



Compositional Learning: Breaking down cell types into fundamental properties (e.g., "cytotoxic" + "lymphocyte" = "cytotoxic T cell"), allowing recognition of novel combinations.

Key Capabilities:

- ▶ Identify disease-specific cell states not in healthy references
- ▶ Discover rare or transitional cell populations
- ▶ Classify cells in non-model organisms with limited annotations
- ▶ Adapt to emerging cell type nomenclature

Breakthrough Impact

Zero-shot learning was instrumental in identifying novel immune cell states in COVID-19 patients and discovering rare developmental intermediates in embryonic development studies.

Overview

Batch effects are systematic technical variations that arise from differences in experimental conditions, reagents, sequencing platforms, or processing times. These non-biological variations can obscure true biological signals and lead to incorrect cell type classifications if not properly addressed.

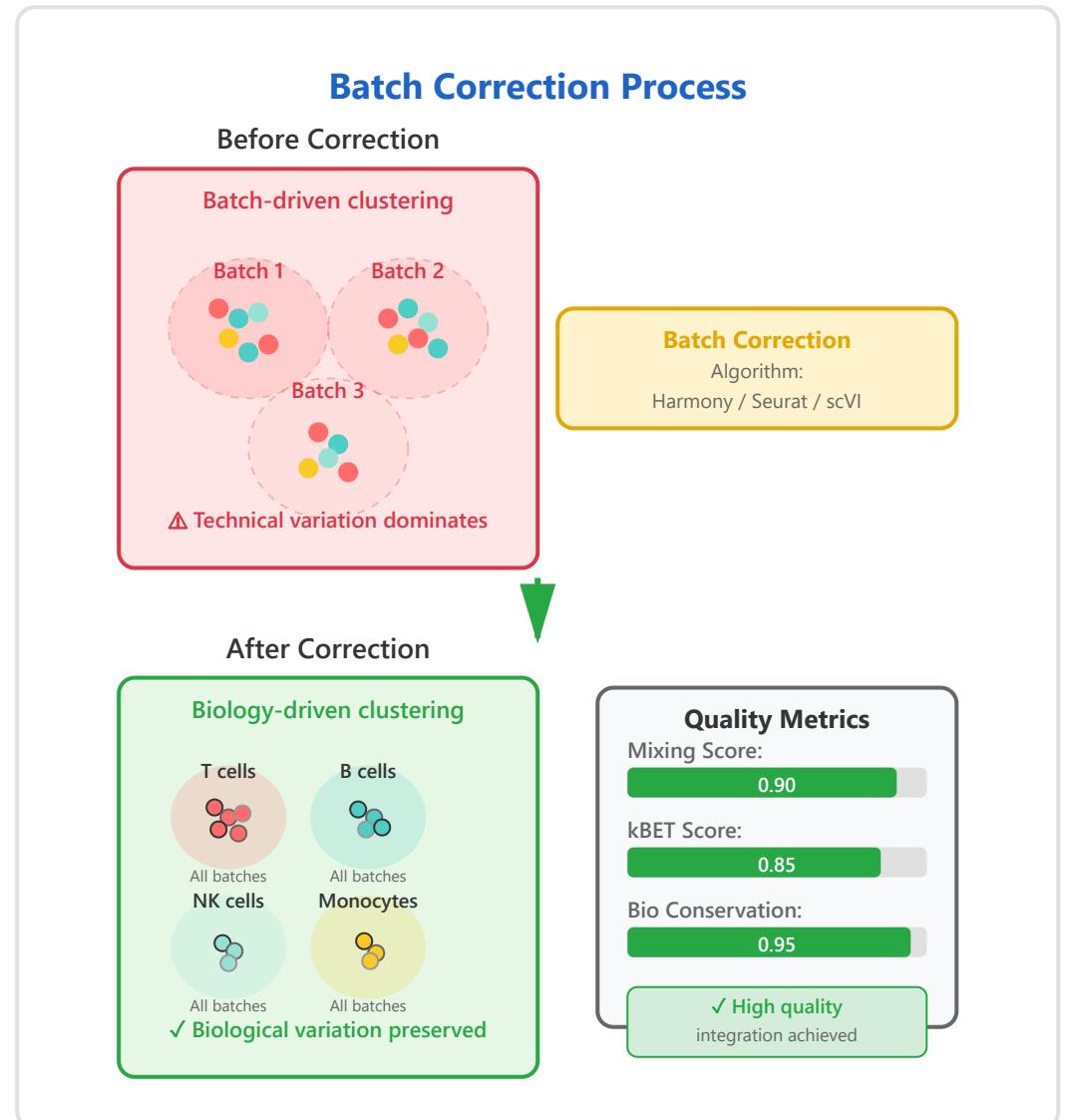
Sources of Batch Effects

Common sources include differences in cell capture efficiency, library preparation protocols, sequencing depth, ambient RNA contamination, and even the laboratory or technician performing the experiment. These effects can be so strong that cells cluster by experimental batch rather than biological cell type.

Correction Methods

Harmony: A fast integration method that iteratively corrects batch effects while preserving biological variation using soft k-means clustering in PCA space.

Seurat Integration: Uses canonical correlation analysis (CCA) to identify shared correlation structures across



batches, anchoring datasets together based on mutual nearest neighbors.

scVI (Single-cell Variational Inference): A deep learning approach using variational autoencoders to model both biological variation and batch effects simultaneously, learning a corrected latent representation.

Best Practices:

- ▶ Always visualize data before and after correction with UMAP/t-SNE
- ▶ Verify that biological variation is preserved, not removed
- ▶ Use multiple quality metrics (mixing metrics, kBET, LISI)
- ▶ Consider whether correction is necessary - some "batches" may have real biology

Critical Consideration

Over-correction can remove genuine biological differences. For example, disease-control comparisons should preserve disease-specific cell states while removing technical variation.

Overview

Uncertainty estimation quantifies the confidence of cell type predictions, distinguishing between cells that are confidently classified and those in ambiguous states. This is crucial for identifying transitional cells, doublets (two cells captured together), low-quality cells, and truly novel cell populations that require manual curation.

Types of Uncertainty

Aleatoric Uncertainty: Inherent noise in the data due to technical limitations like low gene capture efficiency or stochastic gene expression. This is irreducible uncertainty in the measurement itself.

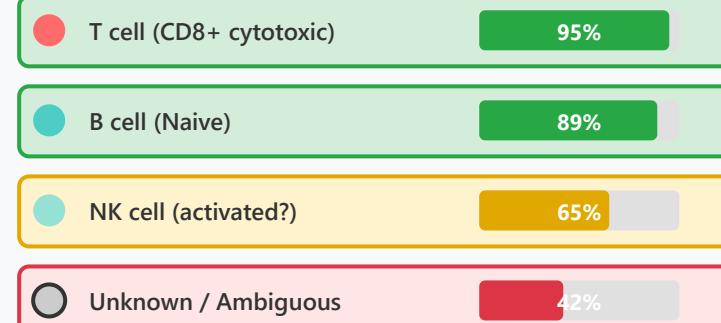
Epistemic Uncertainty: Uncertainty arising from model limitations or lack of training data. This can be reduced with more data or better models and indicates cells that are far from known training examples.

Estimation Methods

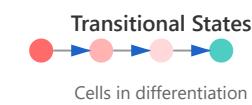
Probabilistic Classifiers: Models output probability distributions over cell types rather than hard labels, with entropy serving as an uncertainty measure.

Uncertainty Estimation Framework

Cell Classifications with Confidence



Sources of Uncertainty



Confidence-Based Decision Workflow

High Confidence
> 80%
✓ Auto-annotate

Med Confidence
50-80%
⚠ Flag for review

Transparent reporting builds trust in automated classification

Low Confidence
< 50%
💡 Expert review

Monte Carlo Dropout: Running multiple predictions with dropout enabled to create a distribution of predictions, estimating uncertainty from variance.

Ensemble Methods: Training multiple models and measuring disagreement between their predictions.

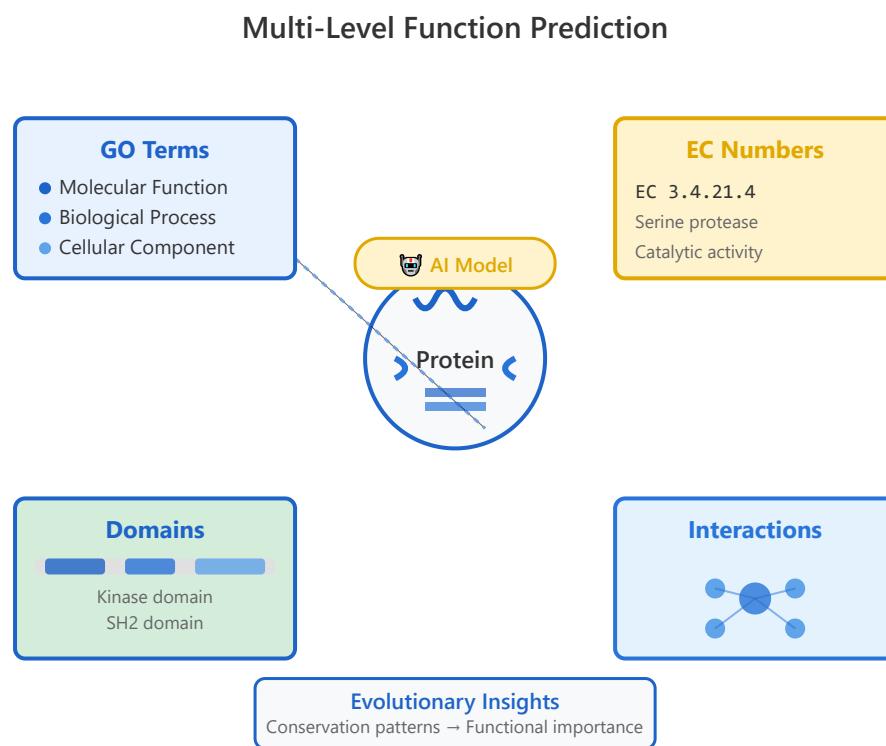
Practical Applications:

- ▶ Flag ambiguous cells for manual review by experts
- ▶ Identify potential novel cell states requiring further investigation
- ▶ Detect technical artifacts (doublets, damaged cells)
- ▶ Prioritize cells for validation experiments
- ▶ Provide honest assessment of classification reliability

Clinical Relevance

In clinical applications, uncertainty estimation is critical. High-confidence predictions can guide treatment decisions, while low-confidence cases can be flagged for additional testing or expert review.

Protein Function Prediction



GO term prediction

Molecular/biological/cellular

EC number classification

Enzyme commission numbers

Domain annotation

Functional regions identification

Interaction prediction

Protein-protein networks

Evolutionary insights

Conservation-function mapping

1

GO Term Prediction

What are GO Terms?

Gene Ontology (GO) terms provide a standardized vocabulary to describe protein functions across three main domains. They enable consistent annotation of protein characteristics across different organisms and databases.

Three GO Categories:

- ▶ **Molecular Function (MF):** Activities at the molecular level (e.g., catalytic activity, binding)
- ▶ **Biological Process (BP):** Larger processes accomplished by multiple molecular activities (e.g., signal transduction, metabolism)
- ▶ **Cellular Component (CC):** Location where a gene product is active (e.g., nucleus, mitochondrion)

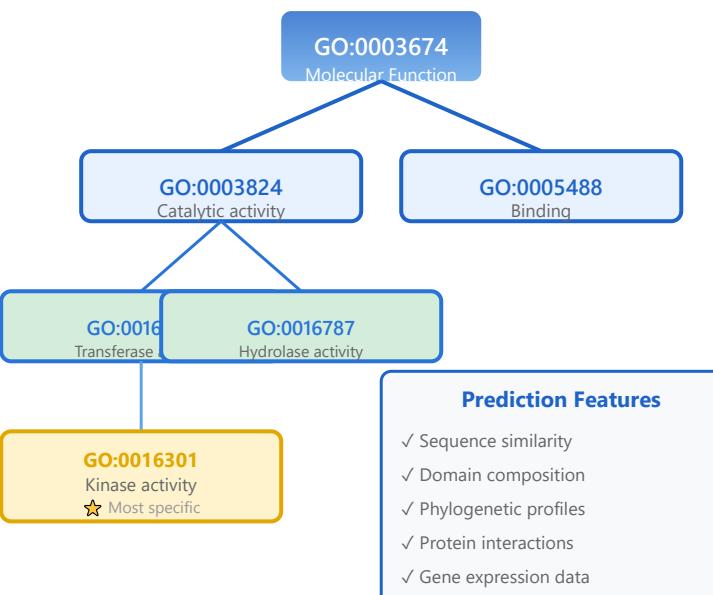
Example: Protein Kinase Annotation

GO:0004672 - protein kinase activity

GO:0006468 - protein phosphorylation

GO:0005737 - cytoplasm

GO Hierarchical Structure



Prediction Features

- ✓ Sequence similarity
- ✓ Domain composition
- ✓ Phylogenetic profiles
- ✓ Protein interactions
- ✓ Gene expression data

2 EC Number Classification

Enzyme Commission (EC) Numbers

EC numbers provide a hierarchical classification system for enzymes based on the chemical reactions they catalyze. Each EC number consists of four digits that progressively specify the enzyme's function.

EC Number Structure (EC a.b.c.d):

- ▶ **First digit (a):** Main enzyme class (1-7)
- ▶ **Second digit (b):** Subclass (substrate type)
- ▶ **Third digit (c):** Sub-subclass (specific substrate)
- ▶ **Fourth digit (d):** Serial number (specific enzyme)

Example: Trypsin Classification

EC 3.4.21.4 | | | L Trypsin (specific enzyme) | |
 L Serine endopeptidase | L Acting on
peptide bonds L Hydrolase

EC Classification System

Seven Main Enzyme Classes

EC 1 - Oxidoreductases

EC 2 - Transferases

EC 3 - Hydrolases

EC 4 - Lyases

EC 5 - Isomerases

EC 6 - Ligases

EC 7 - Translocases

Detailed Example: Protein Kinase

EC 2.7.11.1 - Non-specific serine/threonine protein kinase

Reaction catalyzed:

$\text{ATP} \rightarrow \text{ADP} + \text{phosphoprotein}$

Substrate specificity:

Serine or threonine residues on target proteins

Cofactor requirement: Mg^{2+} or Mn^{2+}

3 Domain Annotation

What are Protein Domains?

Protein domains are distinct structural and functional units within a protein sequence. They are evolutionarily conserved regions that can fold independently and often retain function even when separated from the rest of the protein.

Key Characteristics:

- ▶ **Modularity:** Can be mixed and matched in different proteins
- ▶ **Conservation:** Similar domains found across different species
- ▶ **Function:** Each domain typically has a specific function
- ▶ **Independence:** Can fold and function independently

Common Domain Databases

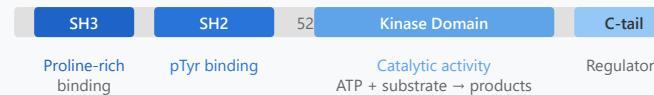
Pfam: Protein families database with HMM profiles

InterPro: Integrated resource of protein families

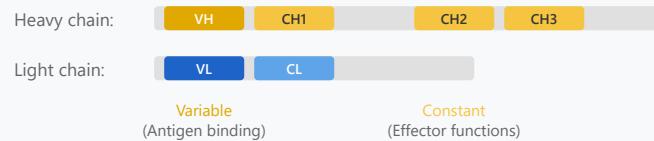
SMART: Simple Modular Architecture Research Tool

Multi-Domain Protein Architecture

Example 1: Src Kinase (c-Src)



Example 2: Immunoglobulin (Antibody)



Domain Prediction Methods

- 🔍 Sequence-based: HMM profiles, pattern matching
- 🌀 Structure-based: Fold recognition, structural alignment
- 🤖 Machine learning: Deep learning models (e.g., AlphaFold)
- 📊 Integrative: Combining multiple evidence sources
- ⌚ Context-aware: Domain co-occurrence patterns

4

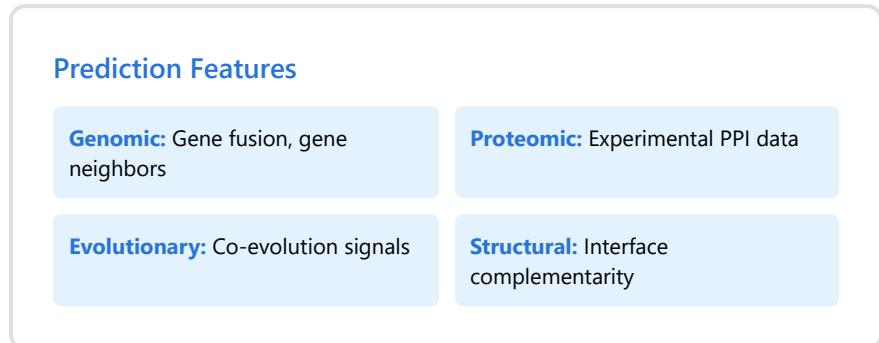
Protein Interaction Prediction

Protein-Protein Interactions (PPIs)

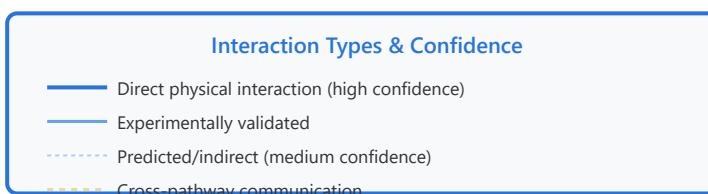
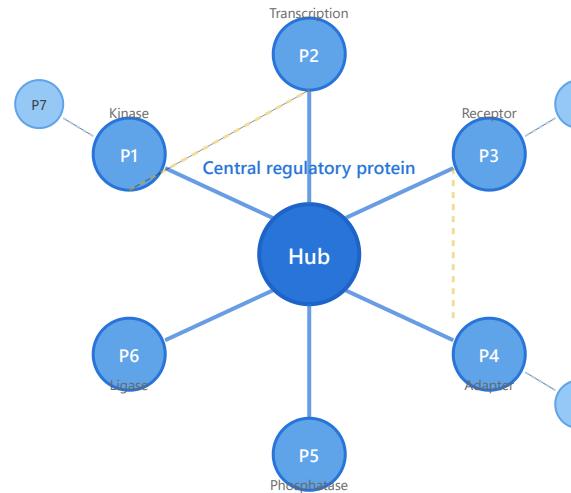
Proteins rarely act alone in cells. Understanding how proteins interact with each other is crucial for deciphering cellular mechanisms, signaling pathways, and disease processes. Computational prediction of PPIs helps identify potential interaction partners.

Types of Protein Interactions:

- ▶ **Stable complexes:** Long-lasting, often structural interactions
- ▶ **Transient interactions:** Brief contacts for signaling or catalysis
- ▶ **Direct binding:** Physical contact between protein surfaces
- ▶ **Indirect associations:** Mediated through other molecules



Protein Interaction Network



5

Evolutionary Conservation & Function

Conservation-Function Relationship

Evolutionary conservation analysis provides powerful insights into protein function. Highly conserved regions across species typically indicate functional importance, as mutations in these areas are often detrimental and eliminated by natural selection.

Key Principles:

- ▶ **Sequence conservation:** Similar amino acids across species
- ▶ **Structural conservation:** Preserved 3D structure despite sequence variation
- ▶ **Functional residues:** Highly conserved catalytic and binding sites
- ▶ **Co-evolution:** Correlated mutations reveal interaction partners

Conservation Score Interpretation

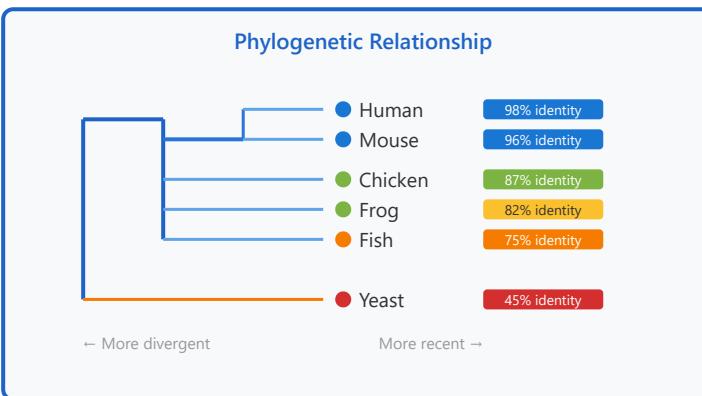
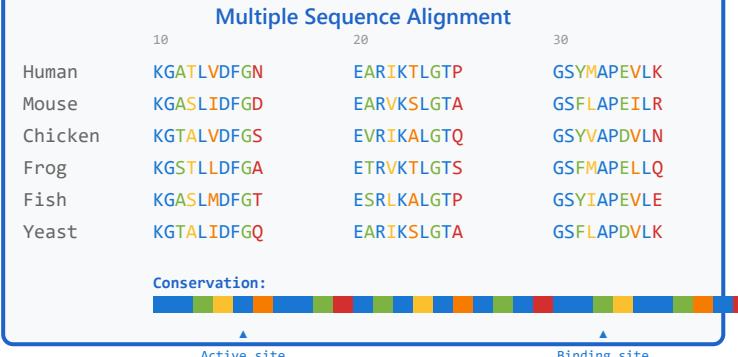


Low → High conservation

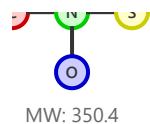
High conservation (blue): Functionally critical

Low conservation (red): Tolerates variation

Evolutionary Conservation Analysis

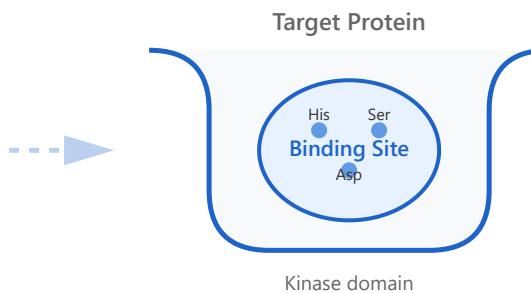


Drug-Target Affinity



MW: 350.4

Drug-Target Binding Prediction



AI Prediction

Binding Affinity (Kd): **2.3 nM**
IC50: **15.7 nM**
Selectivity Score: **0.92**

Advanced Predictions

● Allosteric sites ● Cryptic pockets ● Residence time ● Off-targets
Model Confidence:

84%

Binding prediction

Kd, Ki, IC50 values

Kinase selectivity

Off-target profiling

Allosteric sites

Non-competitive binding

Cryptic pockets

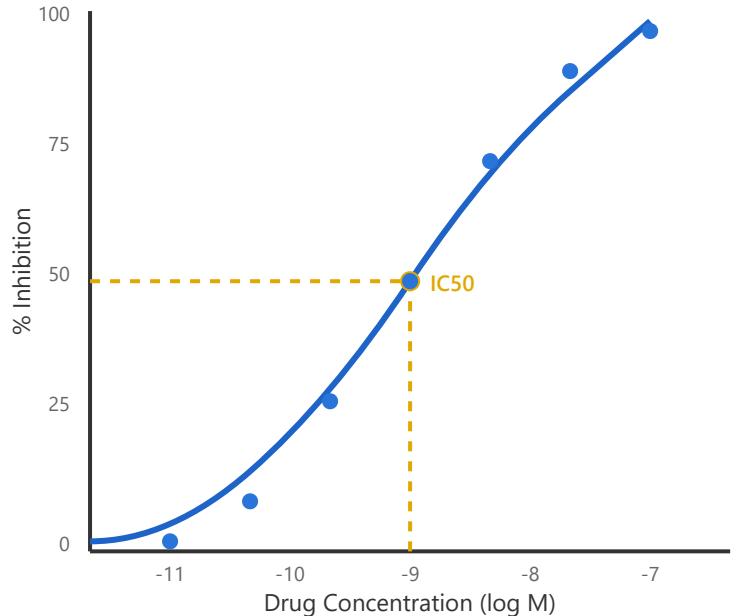
Hidden binding sites

Residence time

Drug-target kinetics

1. Binding Prediction (Kd, Ki, IC50)

Dose-Response Curve



Binding affinity quantifies the strength of interaction between a drug molecule and its target protein. Three key metrics are commonly used:

Key Parameters:

K_d (Dissociation Constant): The equilibrium constant for the dissociation of a drug-target complex. Lower K_d values indicate stronger binding.

$$K_d = [\text{Drug}] [\text{Target}] / [\text{Drug-Target Complex}]$$

K_i (Inhibition Constant): Measures the affinity of an inhibitor for its target enzyme. Represents the concentration needed to produce half-maximum inhibition.

IC₅₀: The concentration of drug required to inhibit 50% of the target activity. Most commonly used in drug screening.

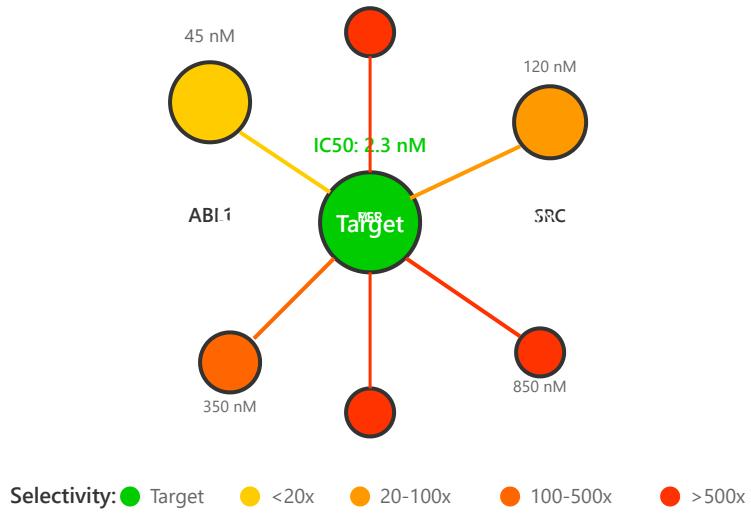
Typical ranges:

- Strong binders: K_d/K_i < 1 nM
- Moderate binders: 1-100 nM
- Weak binders: > 100 nM

AI Prediction Advantages: Machine learning models can predict binding affinity from molecular structure, reducing the need for extensive experimental screening and accelerating drug discovery.

2. Kinase Selectivity & Off-Target Profiling

Kinase Selectivity Heatmap



Kinase selectivity is critical for developing safe and effective kinase inhibitors. The human kinome contains over 500 protein kinases with similar ATP-binding sites, making selectivity a major challenge.

Why Selectivity Matters:

Poor selectivity can lead to:

- Off-target toxicity and adverse effects
- Reduced therapeutic window
- Unpredictable drug interactions
- Clinical trial failures

Selectivity Metrics:

$$\text{Selectivity Score} = \frac{\text{IC50 (off-target)}}{\text{IC50 (target)}}$$

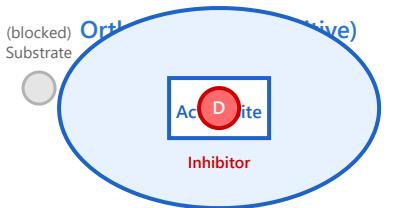
Selectivity Index: A score typically ranging from 0 to 1, where values closer to 1 indicate higher selectivity for the intended target.

Example: A drug with $\text{IC50} = 2.3 \text{ nM}$ for target kinase and $\text{IC50} = 120 \text{ nM}$ for SRC kinase has a selectivity ratio of ~ 52 -fold, indicating good selectivity.

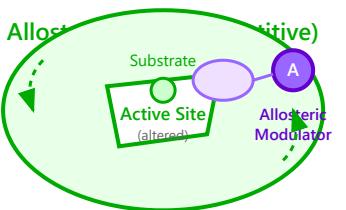
AI-Powered Profiling: Machine learning models can predict binding to hundreds of kinases simultaneously, identifying potential off-targets early in development and enabling structure-based optimization for improved selectivity.

3. Allosteric Sites & Non-Competitive Binding

Allosteric vs Orthosteric Binding



VS



Advantages: Higher selectivity • Unique binding sites • Modulates rather than blocks

Allosteric sites are binding pockets located away from the active site that regulate protein function through conformational changes. Unlike orthosteric inhibitors that compete with natural substrates, allosteric modulators offer unique therapeutic advantages.

Key Characteristics:

- **Non-competitive binding:** Does not compete with natural substrate/ligand
- **Conformational regulation:** Induces structural changes that affect activity
- **Can be activators or inhibitors:** Positive or negative allosteric modulators (PAMs/NAMs)

Therapeutic Advantages:

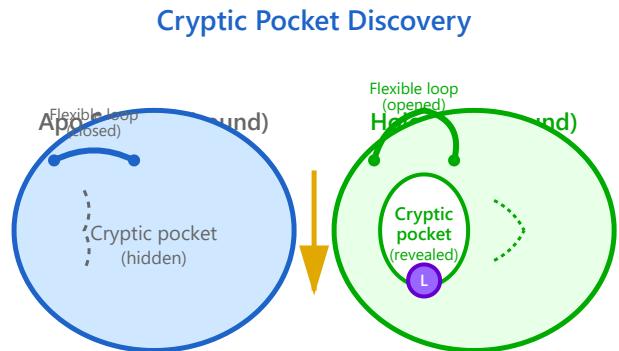
1. Enhanced Selectivity: Allosteric sites are often less conserved than active sites across protein families, enabling development of highly selective drugs.

2. Tunable Efficacy: Can fine-tune protein function rather than completely blocking it, potentially reducing side effects.

3. Overcome Resistance: Mutations in active sites that confer drug resistance may not affect allosteric sites.

Clinical Examples: MEK inhibitors (trametinib), BCR-ABL allosteric inhibitors (asciminib), and EGFR allosteric

4. Cryptic Pockets & Hidden Binding Sites



Cryptic pockets are binding sites that are not visible in the native protein structure but can be revealed through conformational changes induced by ligand binding or protein dynamics. These hidden pockets represent untapped opportunities for drug discovery.

Characteristics of Cryptic Pockets:

- **Transient nature:** Form through protein breathing motions and conformational fluctuations
- **Induced fit:** Ligand binding stabilizes the open pocket conformation
- **Often druggable:** Can provide binding sites in "undruggable" targets

Methods for Cryptic Pocket Discovery

MD Simulations

- Sample conformations
- Identify transient pockets
- Predict druggability

AI/ML Prediction

- Deep learning models
- Structural analysis
- Binding site prediction

Discovery Approaches:

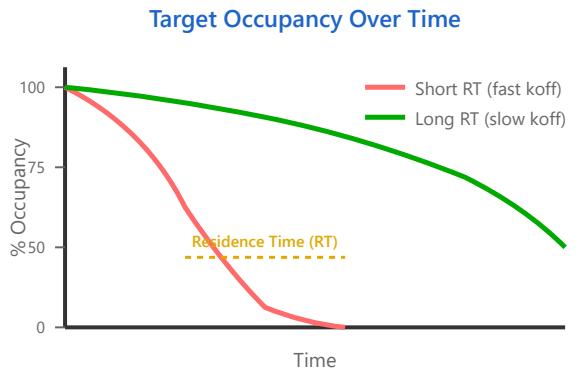
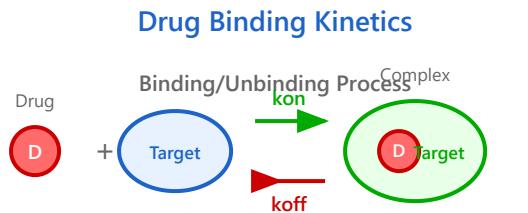
Molecular Dynamics (MD) Simulations: Computational simulations that sample protein conformational space to identify transient pockets that appear during protein motion.

AI-Based Prediction: Machine learning models trained on protein structures and dynamics can predict the location and druggability of cryptic pockets without extensive MD simulations.

Success Story: Cryptic pockets have been successfully targeted in proteins previously considered "undruggable," including K-Ras (AMG 510/sotorasib) and BCR-ABL (asciminib).

Impact on Drug Discovery: Cryptic pocket identification expands the druggable genome by 30-40%, opening new therapeutic opportunities for difficult targets such as transcription factors and scaffold proteins.

5. Residence Time & Drug-Target Kinetics



Residence time (RT) is the average duration a drug molecule remains bound to its target. It has emerged as a critical parameter in drug design, often correlating better with *in vivo* efficacy than binding affinity alone.

Key Concepts:

$$\text{Residence Time (RT)} = 1 / k_{off}$$

Where **koff** is the dissociation rate constant. A smaller k_{off} means slower unbinding and longer residence time.

$$K_d = k_{off} / k_{on}$$

Two drugs can have identical K_d but very different residence times depending on their k_{on} and k_{off} values.

Clinical Importance:

Extended Efficacy: Long residence time can provide sustained target engagement even when plasma drug concentrations drop, allowing for less frequent dosing.

Selectivity Enhancement: Longer residence time on target vs. off-targets can improve the therapeutic window.

Example: Alectinib (ALK inhibitor) has a residence time of ~200 minutes compared to crizotinib's ~30 minutes, contributing to superior efficacy and ability to overcome resistance.

Prediction Challenges:

Residence time is influenced by multiple factors including binding pathway, conformational changes, and rebinding events, making it computationally challenging to predict. Advanced AI models now incorporate molecular dynamics and transition state modeling to estimate RT.

Mutation Effects

Mutation Impact Analysis

Wild Type:

M K L V F F A R G I L S D N Q K Y Position 234

R234W
↓

Mutant:

M K L V F F A W G I L S D N Q K Y

Predicted Effects

Structural Impact

- $\Delta\Delta G$: +3.2 kcal/mol
- Destabilizing

Functional Impact

- Activity: 12% WT
- Loss of function

Pathogenicity:

Clinical Interpretation

90% (Likely Pathogenic)

Conservation Score: 0.98 (Highly Conserved)

Pathogenicity prediction

Disease association scoring

Stability changes

$\Delta\Delta G$ calculation

Function impact

Activity & binding changes

Evolutionary constraints

Conservation analysis

Clinical interpretation

Variant classification

1. Pathogenicity Prediction

Pathogenicity Assessment Framework

Input Features:

Conservation
Score: 0.98

Amino Acid Change
R → W (charge loss)

Structural Context
Active site

Prediction Algorithms:

PolyPhen-2

Score: 0.956
Probably Damaging

SIFT

Score: 0.002
Deleterious

Meta-Predictor Integration

REVEL, CADD, MetaSVM

PATHOGENICITY PREDICTION

Classification:

LIKELY PATHOGENIC

High confidence

90%

Overview

Pathogenicity prediction assesses whether a genetic variant is likely to cause disease. This computational approach integrates multiple lines of evidence to estimate the probability that a mutation contributes to pathology.

Key Prediction Tools

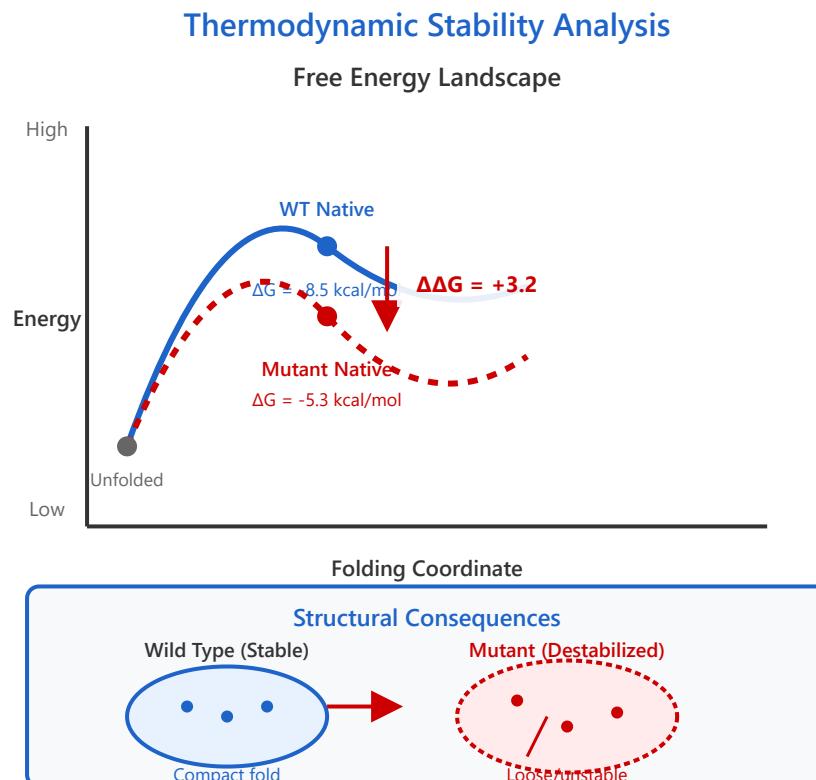
- **PolyPhen-2:** Uses sequence conservation and structural information to predict impact on protein function
- **SIFT:** Predicts whether an amino acid substitution affects protein function based on sequence homology
- **CADD:** Integrates diverse annotations into a single deleteriousness score
- **REVEL:** Ensemble method combining 13 individual tools for improved accuracy

Clinical Application: These predictions help prioritize variants for experimental validation and guide clinical decision-making when interpreting genetic test results.

Interpretation Guidelines

- Scores > 0.8: Likely pathogenic
- Scores 0.4-0.8: Uncertain significance
- Scores < 0.4: Likely benign

2. Protein Stability Changes ($\Delta\Delta G$)



Overview

The change in Gibbs free energy ($\Delta\Delta G$) quantifies how a mutation affects protein stability. Positive $\Delta\Delta G$ values indicate destabilization, while negative values suggest stabilization of the protein structure.

Calculation Methods

- **FoldX:** Empirical force field-based approach using high-resolution structures
- **Rosetta:** Energy function combining physics-based and knowledge-based terms
- **DynaMut:** Considers protein dynamics and flexibility changes
- **I-Mutant:** Machine learning approach trained on experimental data

Formula: $\Delta\Delta G = \Delta G_{\text{mutant}} - \Delta G_{\text{wild-type}}$
A value of +3.2 kcal/mol indicates significant destabilization, often leading to protein misfolding or degradation.

Biological Consequences

- **$\Delta\Delta G > +2 \text{ kcal/mol}$:** Significant destabilization, likely protein degradation
- **$\Delta\Delta G +0.5 \text{ to } +2$:** Moderate instability, temperature-sensitive phenotypes
- **$\Delta\Delta G -0.5 \text{ to } +0.5$:** Minimal impact on stability

3. Functional Impact Assessment

Molecular Function Analysis

Enzymatic Activity

Wild Type

Full Activity

100%

Mutant (R234W)

Loss

12%

88% reduction in catalytic efficiency

Substrate Binding Affinity

K_M

High
Low

2.5 μM

7.4x increase

18.6 μM

Molecular Mechanism

Arg234 → Trp: Loss of positive charge disrupts substrate coordination and catalytic geometry

Overview

Functional impact assessment evaluates how a mutation affects the molecular activities of a protein, including catalytic activity, binding affinity, and interaction with other biomolecules.

Key Functional Parameters

- **Catalytic Efficiency (k_{cat}/K_M)**: Overall measure of enzyme performance
- **Binding Affinity (K_D or K_M)**: Strength of protein-ligand interactions
- **V_{max}** : Maximum reaction velocity, reflects enzyme concentration
- **Protein-Protein Interactions**: Changes in binding to partner proteins

Example R234W Impact:

- Activity reduced to 12% of wild-type
- K_M increased 7.4-fold (weaker substrate binding)
- Charge loss at position 234 disrupts active site geometry

Experimental Approaches

- **Enzyme kinetics:** Measure k_{cat} , K_M , V_{max}
- **Binding assays:** SPR, ITC, fluorescence polarization
- **Cell-based assays:** Functional readouts in cellular context
- **Structural studies:** X-ray crystallography, cryo-EM

4. Evolutionary Constraints & Conservation

Evolutionary Conservation Analysis

Multiple Sequence Alignment

Species	Sequence (Position 230-240)									
Human	L	V	F	F	A	R	G	I	L	S D
Mouse	L	V	F	F	A	R	G	I	L	S D
Rat	L	V	F	F	A	R	G	I	L	T D
Dog	L	V	F	F	A	R	G	V	L	S D
Cow	L	V	F	F	A	R	G	I	L	S E
Chicken	L	V	Y	F	A	R	G	I	L	S D
Zebrafish	M	V	F	F	A	K	G	I	L	S N

Conservation Metrics

Shannon Entropy:	<div style="width: 20px; height: 10px; background-color: blue;"></div>	0.06 (Low variability)
Conservation Score:	<div style="width: 90%; height: 10px; background-color: red;"></div>	0.98 (Highly conserved)

Evolutionary Time Scale



Arg234 conserved across 600 million years

Overview

Evolutionary conservation analysis examines how well a protein position is preserved across species. Highly conserved positions are typically functionally important, and mutations at these sites are more likely to be deleterious.

Conservation Metrics

- Shannon Entropy:** Measures amino acid variability at each position (0 = identical, higher = variable)
- Conservation Score:** Quantifies evolutionary constraint (0-1 scale)
- GERP Score:** Identifies positions under strong selective pressure
- PhyloP:** Measures evolutionary conservation based on multiple alignments

Position 234 Analysis:

- Conserved as Arginine in all mammals
- Only conservative substitution (Lysine) in distantly related species
- Conservation score: 0.98/1.00
- Indicates critical functional role

Interpretation

- Score > 0.9:** Extremely conserved, mutations likely deleterious
- Score 0.7-0.9:** Well conserved, mutations often harmful
- Score 0.4-0.7:** Moderately conserved, variable tolerance

- **Score < 0.4:** Poorly conserved, mutations often tolerated

5. Clinical Interpretation & Variant Classification

ACMG/AMP Classification Framework

Lines of Evidence

Pathogenic Evidence (Supporting)

- PS3 (Strong): Well-established functional studies show damaging effect
- PM1 (Moderate): Located in critical functional domain
- PM2 (Moderate): Absent from controls in large databases
- PP3 (Supporting): Multiple computational predictions support damage
- PP2 (Supporting): Missense in gene with low tolerance to variation

Benign Evidence

Classification Decision Tree

Evidence
PS3 + 2xPM + 2xPP

ACMG Rules
Apply criteria

Likely
Pathogenic

**CLINICAL CLASSIFICATION
LIKELY PATHOGENIC (Class 4)**

Overview

Clinical interpretation follows standardized guidelines from the American College of Medical Genetics (ACMG) and Association for Molecular Pathology (AMP) to classify variants into five categories based on their likelihood of causing disease.

ACMG Classification Tiers

- **Pathogenic (Class 5):** Sufficient evidence of disease causation
- **Likely Pathogenic (Class 4):** Strong but not conclusive evidence
- **Uncertain Significance (Class 3):** Insufficient or conflicting evidence
- **Likely Benign (Class 2):** Strong evidence against pathogenicity
- **Benign (Class 1):** Established as non-pathogenic

Evidence Strength Levels:

- **Very Strong (PVS):** Null variants, proven functional effects
- **Strong (PS):** Well-established in vitro/in vivo

studies

- **Moderate (PM):** Computational predictions, location
- **Supporting (PP):** Conservation, multiple algorithms

Clinical Action

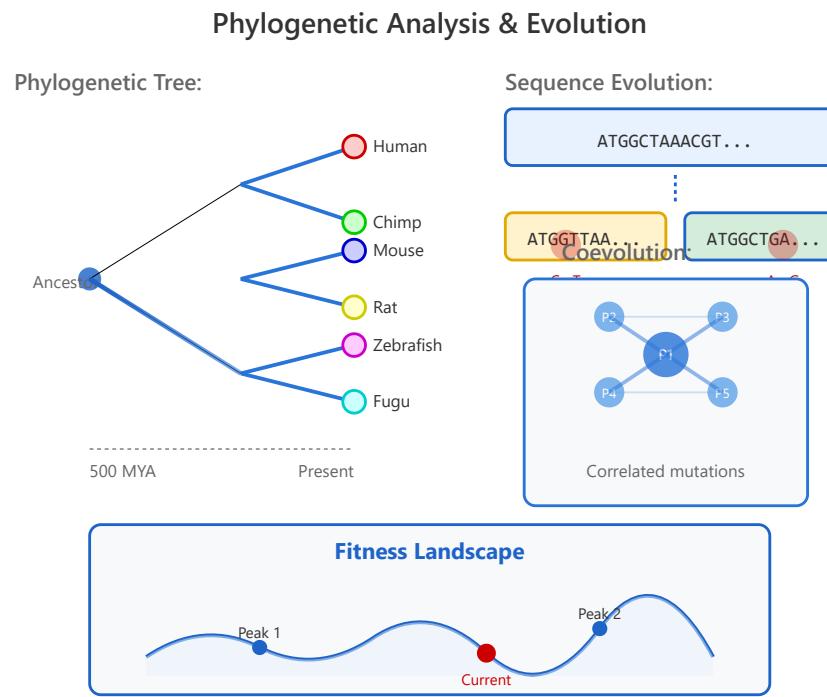
- **Genetic counseling:** Discuss implications with patients
- **Cascade testing:** Test family members for variant
- **Management:** Enhanced surveillance or preventive measures
- **Research:** Further functional validation if needed

Databases & Resources

- **ClinVar:** Repository of variant interpretations
- **gnomAD:** Population frequency data
- **OMIM:** Gene-disease relationships
- **HGMD:** Human gene mutation database

Evolution Modeling

Computational Approaches to Understanding Biological Evolution



Sequence Evolution

Substitution models & rates

Phylogenetic Inference

Tree reconstruction methods

Ancestral Reconstruction

Ancient sequence prediction

Coevolution

Correlated mutations analysis

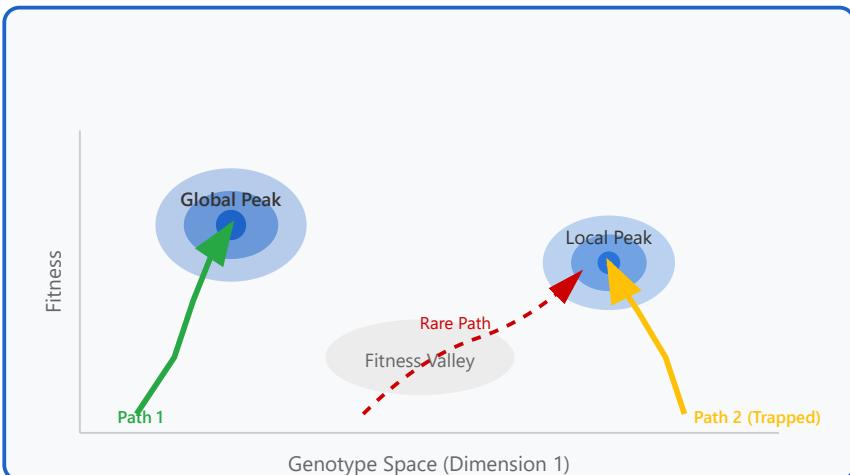
Fitness Landscapes

Adaptive evolution mapping

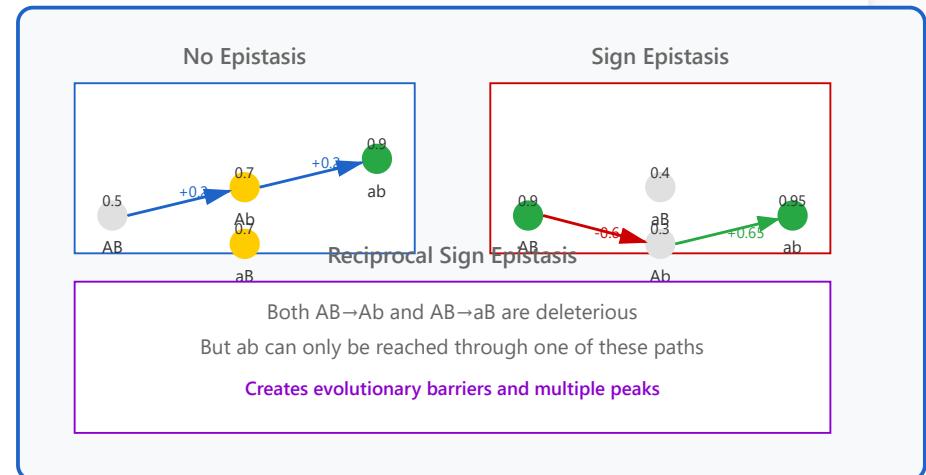
Fitness landscapes represent the relationship between genotypes or phenotypes and their reproductive success. These multidimensional spaces help us understand evolutionary trajectories, constraints on adaptation, and the accessibility of beneficial mutations through the complex topology of sequence space.

Fitness Landscape Topology and Evolutionary Paths

3D Fitness Landscape (2D Projection)

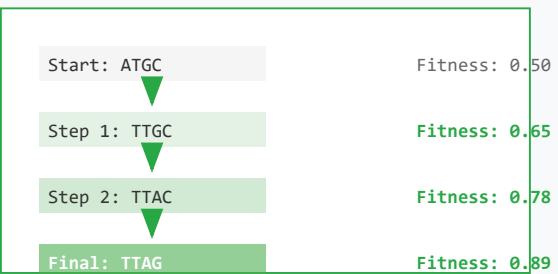


Epistatic Interactions



Adaptive Walks and Evolutionary Dynamics

Greedy Adaptive Walk



Population Dynamics



Key Factors

- Population size (N)
Large N → selection dominates
- Mutation rate (μ)
High μ → more exploration
- Selection strength (s)
Strong s → fast adaptation
- Recombination (r)
 $r > 0$ → breaks linkage

Key Concepts:

- **Fitness Peaks:** Genotypes with maximum fitness in their local neighborhood; populations tend to evolve toward peaks
- **Fitness Valleys:** Low-fitness regions between peaks that can trap populations and prevent access to higher peaks
- **Ruggedness:** The complexity of the landscape; rugged landscapes have many local optima
- **Epistasis:** Interactions between mutations where the effect of one depends on genetic background
- **Sign Epistasis:** When a mutation is beneficial in one background but deleterious in another
- **Evolutionary Accessibility:** Which genotypes can be reached through single mutational steps
- **Adaptive Walks:** Trajectories through sequence space following fitness gradients

Wright-Fisher Model for Fixation Probability:

$$P(\text{fixation}) \approx (1 - e^{(-2Ns)}) / (1 - e^{(-4Ns)})$$

where N = population size, s = selection coefficient

Example: Antibiotic Resistance Evolution

Studies of beta-lactamase evolution revealed a rugged fitness landscape where the path to high-level antibiotic resistance requires crossing fitness valleys. Some highly resistant variants can only be accessed through specific mutational paths involving 5+ mutations, explaining why certain resistance mechanisms emerge more frequently than others in clinical settings. Sign epistasis between mutations means that intermediate genotypes have lower fitness, creating evolutionary constraints.

Protein Engineering

Navigate sequence space to design proteins with desired properties

Drug Resistance

Predict evolutionary trajectories of pathogens under drug pressure

Crop Improvement

Guide breeding strategies to optimize multiple traits simultaneously

 **Synthetic Biology**

Design genetic circuits with predictable
evolutionary stability

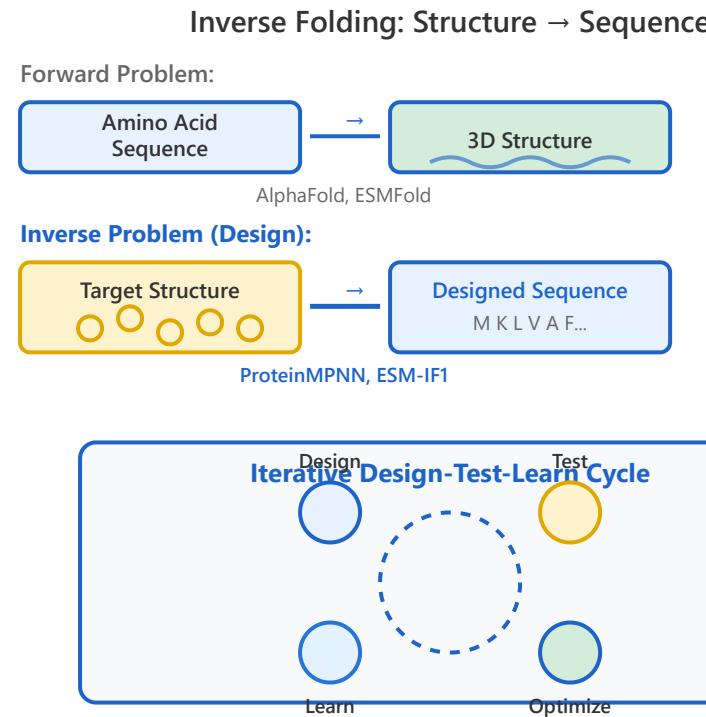
Part 3/3 - Applications

Design problems

Engineering solutions

Therapeutic development

Protein Design



Inverse folding

Structure → sequence prediction

Scaffold design

De novo backbone generation

Interface design

Protein-protein interactions

De novo binders

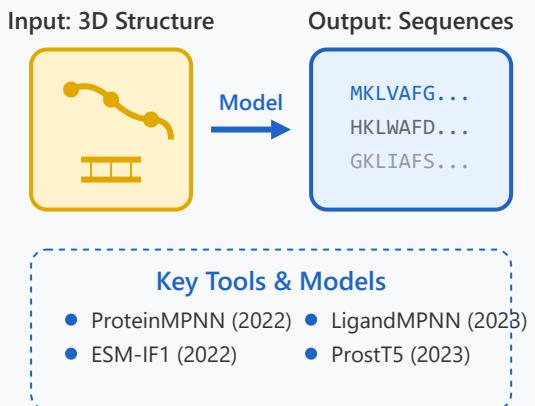
Target-specific protein design

Stability optimization

Thermostability enhancement

Detailed Method Descriptions

Inverse Folding: Structure → Sequence Prediction



Inverse folding is the fundamental problem in computational protein design where we predict amino acid sequences that will fold into a desired three-dimensional structure. Unlike forward folding (structure prediction from sequence), inverse folding solves the reverse problem: given a target structure, what sequences could produce it?

Modern deep learning models like ProteinMPNN and ESM-IF1 have revolutionized this field by learning from vast databases of protein structures. These models encode the geometric constraints of the backbone structure and predict sequences that satisfy both structural and biochemical requirements.

Key Features:

- ▶ Graph neural networks encode backbone geometry and residue environments
- ▶ Generates multiple diverse sequences for the same structure
- ▶ Considers rotamer preferences and side-chain packing
- ▶ Can condition on specific residue constraints or motifs
- ▶ High success rates (60-90%) in experimental validation

Applications:

- ▶ Redesigning protein scaffolds for improved stability
- ▶ Engineering enzyme active sites

- ▶ Creating novel protein folds
- ▶ Protein humanization for therapeutics

Example:

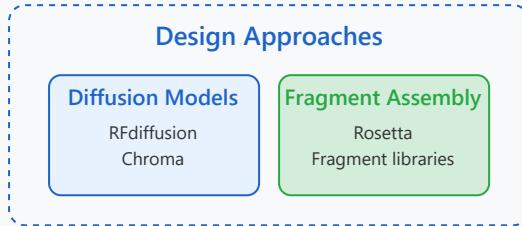
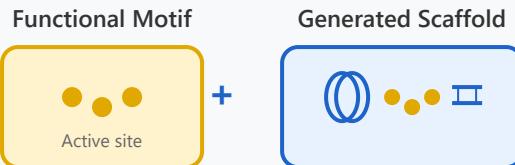
Given a TIM barrel structure, ProteinMPNN can generate 100+ diverse sequences that all fold into the same barrel topology, enabling exploration of sequence space while maintaining function.

2

Scaffold Design: De Novo Backbone Generation

Scaffold design involves creating entirely new protein backbones that can support specific functional motifs or binding sites. Rather than modifying existing proteins, this approach generates novel three-dimensional architectures from scratch that position key residues in precise geometric arrangements.

Recent breakthroughs using diffusion models (like RFdiffusion) have dramatically improved scaffold design. These models learn to generate protein backbones by reversing a noise-adding process, allowing them to create diverse, designable structures that incorporate functional constraints while maintaining protein-like geometry and secondary structure composition.



Key Features:

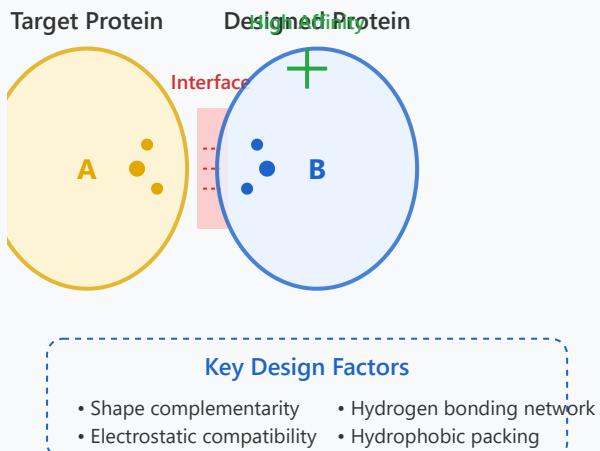
- Generates completely novel protein folds not found in nature
- Can incorporate motif constraints (binding sites, active sites)
- Controls secondary structure composition (helices, sheets, loops)
- Symmetry constraints for multi-subunit assemblies
- Iterative refinement through structure prediction validation

Applications:

- Creating novel enzyme scaffolds with tailored active sites
- Designing protein binders with specific epitope recognition
- Building protein cages and nanomaterials
- Engineering vaccines with optimized antigen presentation

Example:

RFdiffusion was used to design novel protein binders against SARS-CoV-2 spike protein by generating scaffolds that position key binding residues to match the ACE2 receptor interface, resulting in nanomolar-affinity binders.



Interface design focuses on engineering the interaction surfaces between proteins to create or enhance protein-protein interactions. This involves optimizing the complementarity between two protein surfaces through shape, electrostatics, and chemical interactions to achieve high-affinity binding.

Computational approaches model the interface region at atomic detail, considering backbone conformational changes, side-chain rotamers, and the balance between binding affinity and specificity. Modern machine learning methods can predict interface residues and suggest mutations that improve binding while maintaining specificity.

Key Features:

- ▶ Shape complementarity scoring using surface geometry analysis
- ▶ Electrostatic potential matching across interfaces
- ▶ Hydrogen bond network optimization
- ▶ Buried surface area maximization
- ▶ Hot spot residue identification and enhancement
- ▶ Specificity design to avoid off-target interactions

Applications:

- ▶ Antibody-antigen interface optimization for therapeutics
- ▶ Protein complex stabilization in structural biology
- ▶ Designing protein inhibitors for disease targets
- ▶ Engineering synthetic signaling pathways
- ▶ Creating self-assembling protein materials

Example:

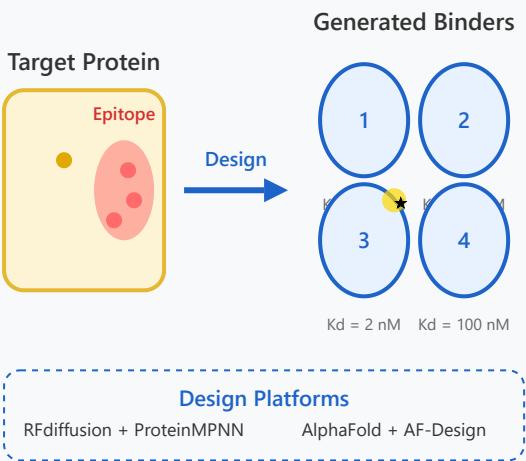
Redesigning the interface between an antibody and a viral protein to increase binding affinity from micromolar to picomolar range through computational optimization of 5-10 key interface residues.

4

De Novo Binders: Target-Specific Protein Design

De novo binder design creates entirely new proteins from scratch that bind to specific target molecules with high affinity and specificity. Unlike antibodies or natural binding proteins, these are computationally designed proteins optimized for a particular binding task.

This approach combines scaffold generation with interface design, using diffusion models to create backbones that present binding residues in optimal geometric arrangements. The process typically involves generating thousands of



candidates, screening them computationally, and experimentally validating the top designs, achieving success rates of 10-50% for nanomolar binders.

Key Features:

- ▶ Target-conditioned generation of binding proteins
- ▶ Multiple scaffold topologies (helical bundles, repeat proteins, mini-proteins)
- ▶ Epitope-specific targeting on protein surfaces
- ▶ Integration with AlphaFold for structure prediction validation
- ▶ Rapid design-test cycles (weeks instead of months)
- ▶ Generates highly specific binders with minimal off-target binding

Applications:

- ▶ Therapeutic protein development (antibody alternatives)
- ▶ Diagnostic tools and biosensors
- ▶ Research tools for target validation
- ▶ Blocking viral entry (e.g., COVID-19 inhibitors)
- ▶ Creating new signaling molecules in synthetic biology

Example:

Researchers at University of Washington used RFdiffusion to design mini-protein binders against the SARS-CoV-2 receptor binding domain, achieving picomolar affinity binders that neutralize the virus *in vitro* within 3 weeks of design.

5

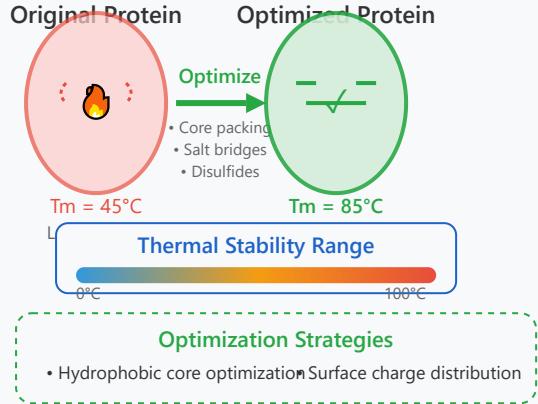
Stability Optimization: Thermostability Enhancement

Stability optimization focuses on enhancing protein thermostability and resistance to denaturation through rational design and computational prediction. This involves identifying and modifying residues that contribute to protein unfolding, improving hydrophobic core packing, introducing stabilizing interactions, and removing destabilizing elements.

Machine learning models can now predict the effects of mutations on stability ($\Delta\Delta G$) with high accuracy, enabling systematic exploration of stability-enhancing variants. This is particularly valuable for therapeutic proteins, industrial enzymes, and any application requiring proteins to function under harsh conditions.

Key Features:

- ▶ $\Delta\Delta G$ prediction for mutation effects on folding stability
- ▶ Core residue packing optimization using rotamer libraries
- ▶ Introduction of disulfide bonds for covalent stabilization
- ▶ Salt bridge network design for electrostatic stabilization
- ▶ Loop rigidification and proline substitutions



- ▶ Removal of thermolabile residues (asparagine, glutamine deamidation sites)
- ▶ Consensus design from homologous sequences

Applications:

- ▶ Industrial enzyme optimization for high-temperature processes
- ▶ Therapeutic protein formulation stability
- ▶ Vaccine antigen stabilization
- ▶ Biosensor proteins for harsh environments
- ▶ Extending shelf-life of protein-based products

Example:

Engineering a thermostable variant of T4 lysozyme by introducing 5 key mutations identified through computational stability prediction, increasing melting temperature from 42°C to 72°C while maintaining full enzymatic activity.

Antibody Design

Comprehensive Guide to Therapeutic Antibody Engineering

- CDR Optimization

- Humanization

- Affinity Maturation

- Specificity Engineering

- Developability

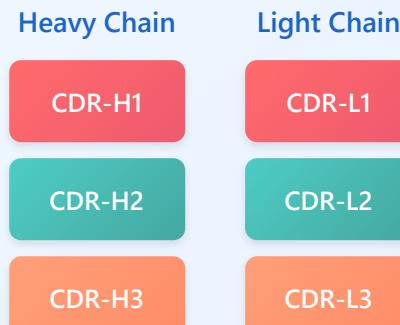
1

CDR Optimization

Complementarity-Determining Regions (CDRs) are the hypervariable loops within the antibody variable domains that directly interact with the antigen. CDR optimization focuses on modifying these critical regions to enhance binding properties, reduce immunogenicity, and improve overall antibody performance. The three CDR loops in both heavy and light chains (CDR-H1, CDR-H2, CDR-H3, CDR-L1, CDR-L2, CDR-L3) form the antigen-binding site.

Key Approaches:

- ▶ Rational design based on structural analysis of antibody-antigen complexes
- ▶ Computational modeling to predict favorable mutations
- ▶ Focused mutagenesis libraries targeting hotspot residues
- ▶ Structure-guided optimization using crystallography or cryo-EM data
- ▶ Conservation of framework integrity while modifying CDR loops
- ▶ Balancing affinity improvements with stability maintenance



Six CDR loops form the antigen-binding site

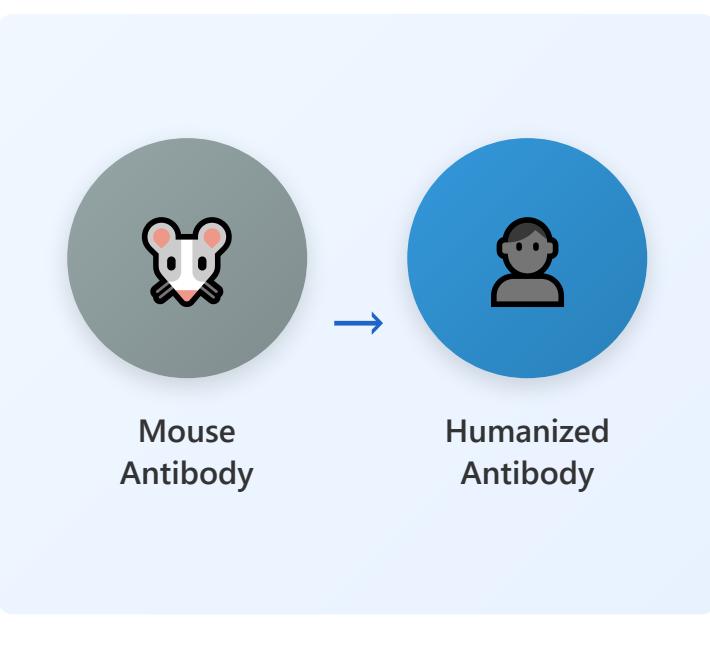
2

Humanization

Humanization is the process of converting non-human antibodies (typically from mice or other rodents) into forms that closely resemble human antibodies. This critical step reduces immunogenicity in patients, minimizing the risk of anti-drug antibodies (ADAs) that can neutralize therapeutic efficacy or cause adverse reactions. The goal is to maintain the antigen-binding properties while replacing most of the antibody sequence with human-derived sequences.

Key Methods:

- ▶ CDR grafting: Transfer CDRs from mouse antibody to human framework
- ▶ Framework region selection from human germline sequences
- ▶ Back-mutation of key residues to restore binding affinity
- ▶ Vernier zone optimization to support CDR conformation
- ▶ Computational tools for identifying critical residues
- ▶ Validation of reduced immunogenicity through T-cell epitope prediction



Conversion from murine to human-compatible format

3 Affinity Maturation

Affinity maturation is the process of improving antibody binding strength to its target antigen. This mimics the natural immune process where B cells undergo somatic hypermutation to produce higher-affinity antibodies. In therapeutic development, artificial affinity maturation enhances antibody potency, potentially reducing required doses and improving efficacy. The goal is to achieve low nanomolar to picomolar binding affinities (Kd values).

Key Strategies:

- ▶ Random mutagenesis followed by high-throughput screening
- ▶ Site-directed mutagenesis at hotspot positions
- ▶ Display technologies (phage, yeast, ribosome display)
- ▶ Error-prone PCR to generate diversity libraries
- ▶ Computational design using molecular dynamics simulations
- ▶ Iterative rounds of mutation and selection for incremental improvements



Progressive improvement in binding affinity (lower Kd = stronger binding)

4

Specificity Engineering

Specificity engineering ensures that antibodies recognize and bind exclusively to their intended target while avoiding off-target interactions. High specificity is crucial for therapeutic safety, preventing adverse effects from cross-reactivity with similar proteins or unintended tissues. This involves careful selection and engineering to discriminate between highly similar molecules, including closely related family members or post-translational variants.

Key Considerations:

- ▶ Differential screening against related proteins and homologs
- ▶ Epitope mapping to identify unique binding sites
- ▶ Cross-reactivity testing across species and protein variants
- ▶ Negative selection to eliminate polyreactive clones
- ▶ Engineering selectivity for specific conformational states

✓
Target
Antigen

X
Off-
Target
1

X
Off-
Target
2

Selective binding to intended target only

- ▶ Computational analysis of potential off-target binding

5 Developability

Developability refers to the pharmaceutical and biophysical properties that enable an antibody to be successfully manufactured, formulated, and administered as a drug product. Even highly potent and specific antibodies can fail in development due to poor biophysical properties.

Developability assessment evaluates stability, solubility, aggregation propensity, viscosity, and manufacturing efficiency to ensure the molecule can progress through clinical development.

Critical Parameters:

- ▶ Thermal stability (T_m , Tagg) for storage and handling
- ▶ High-concentration formulation capability (>100 mg/mL)



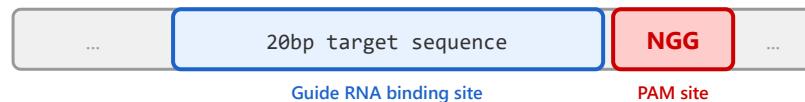
Key biophysical properties for successful development

- ▶ Low aggregation tendency and colloidal stability
- ▶ Acceptable viscosity for subcutaneous administration
- ▶ Chemical stability against degradation pathways
- ▶ Expression yield in manufacturing cell lines (CHO, HEK293)
- ▶ Post-translational modification profiles
- ▶ Absence of self-association or polyreactivity

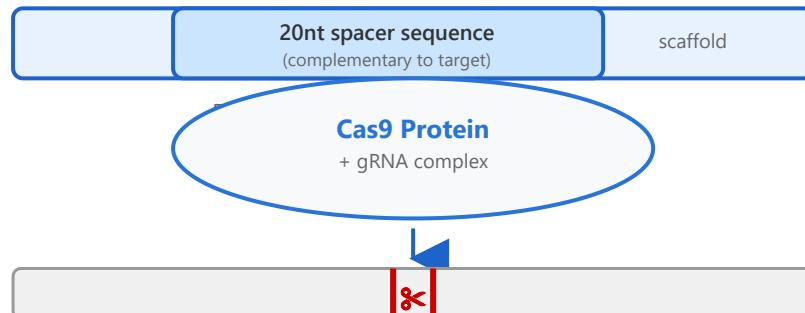
CRISPR Optimization

CRISPR-Cas9 Guide RNA Design

Target DNA Sequence:



Guide RNA (gRNA):



Guide RNA design

20nt spacer + scaffold optimization

Off-target prediction

Minimize unintended cuts

Efficiency scoring

On-target activity models

Prime editing

Precise base substitutions

Base editing

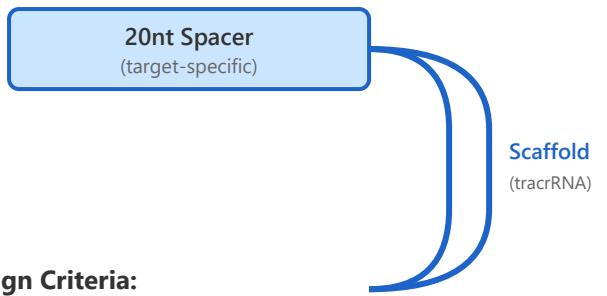
C→T, A→G conversions

1

Guide RNA Design

Optimizing the 20nt spacer and scaffold structure for maximum specificity

gRNA Structure & Optimization



Design Criteria:

- ✓ **GC Content:** 40-60% (optimal stability)
Too low: weak binding | Too high: off-targets
- ✓ **Avoid Poly-T:** ≥4 consecutive T's (U6 terminator)
- ✓ **Start with G:** Enhances U6 promoter transcription
- ✓ **Secondary Structure:** Minimize hairpins in spacer

Guide RNA (gRNA) design is the foundation of CRISPR efficiency. The gRNA consists of a 20-nucleotide spacer sequence that is complementary to the target DNA, fused to a scaffold sequence that binds Cas9 protein. Proper design ensures high on-target activity while minimizing off-target effects.

The spacer sequence must be carefully selected based on multiple biochemical and structural criteria. GC content affects binding stability—too low results in weak DNA binding, while too high can lead to increased off-target activity. The position within the target gene is also critical, with exons near the N-terminus being preferred for gene knockout studies.

AI-Enhanced Design

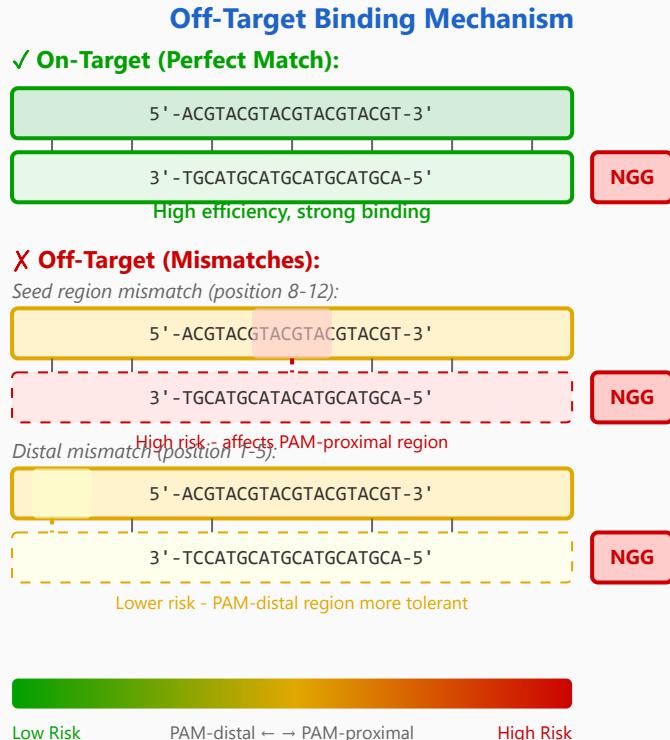
Modern tools use machine learning algorithms trained on thousands of validated gRNAs to predict efficiency scores. These models consider sequence context, chromatin accessibility, and epigenetic marks to recommend optimal designs.

Key Optimization Strategies:

- ▶ Use computational tools (CRISPOR, Benchling, IDT) for initial screening
- ▶ Select gRNAs in constitutively open chromatin regions
- ▶ Avoid SNPs and repetitive sequences in spacer region
- ▶ Test 3-5 gRNAs per target for optimal results

2 Off-Target Prediction

Identifying and minimizing unintended genomic cuts



Off-target effects occur when Cas9 binds and cuts at genomic sites similar to the intended target. These unintended edits can cause mutations in critical genes, leading to cellular dysfunction or confounding experimental results. The risk depends on the number, position, and nature of mismatches between the gRNA and off-target site.

Not all mismatches are equal. The PAM-proximal region (seed sequence, positions 8-12 from PAM) is most critical for specificity. Mismatches here significantly reduce binding, while mismatches in the PAM-distal region (positions 1-7) are more tolerated. Modern prediction algorithms use these principles along with genomic context to calculate off-target risk scores.

Prediction Tools

Leading tools include Cas-OFFinder (comprehensive genome-wide search), GUIDE-seq (experimental validation), and deep learning models like DeepCRISPR that predict cutting likelihood at potential off-target sites. Many tools now integrate chromatin accessibility data for more accurate predictions.

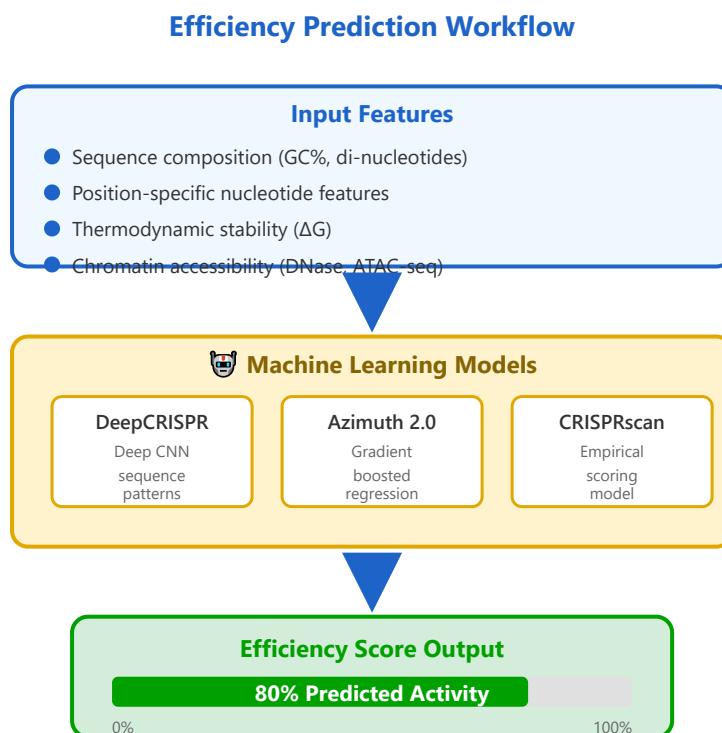
Minimizing Off-Targets:

- ▶ Select gRNAs with no 0-1 mismatch off-targets in coding regions

- ▶ Use high-fidelity Cas9 variants (eSpCas9, HiFi Cas9)
- ▶ Reduce Cas9 exposure time through RNP delivery
- ▶ Validate edited clones with whole-genome sequencing

3 Efficiency Scoring

Predicting on-target cutting activity with machine learning models



Not all correctly designed gRNAs are equally effective. Cutting efficiency can vary from <5% to >90% depending on subtle sequence features and genomic context. Efficiency scoring uses machine learning models trained on thousands of empirically tested gRNAs to predict activity levels before experimental validation.

Modern scoring algorithms combine multiple layers of information including sequence composition, position-specific nucleotide preferences, thermodynamic parameters, and epigenetic features. Deep learning models can capture complex, non-linear relationships that traditional rule-based methods miss, significantly improving prediction accuracy.

Model Performance

State-of-the-art models achieve Spearman correlations of 0.65-0.75 with experimental data. While not perfect, these predictions help prioritize gRNA candidates and reduce experimental

screening burden. Ensemble methods combining multiple algorithms often provide the most robust predictions.

Key Efficiency Factors:

- ▶ Position 16-20 nucleotide identity strongly influences activity
- ▶ Local chromatin state affects Cas9 accessibility
- ▶ DNA repair pathway availability impacts editing outcomes
- ▶ Cell type-specific models improve prediction accuracy

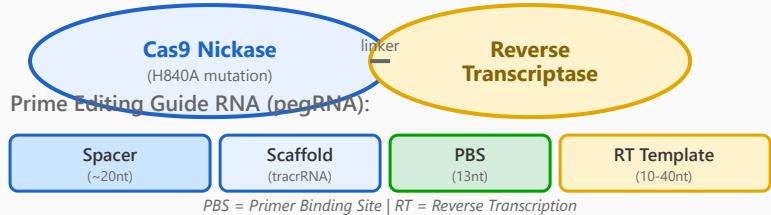
4

Prime Editing

Precise genomic modifications without double-strand breaks

Prime Editing Mechanism

Prime Editor (PE) Components:



Prime Editing Steps:

- 1 pegRNA guides PE to target site
- 2 Cas9(H840A) nicks non-target strand **Nick**
- 3 PBS binds, RT synthesizes new DNA (3' flap)
New DNA + edit
- 4 Flap integration & DNA repair → edited product
✓ Edited DNA sequence (permanent)

Prime editing represents a major advancement in genome editing precision. Unlike traditional CRISPR that creates double-strand breaks, prime editors use a Cas9 nickase fused to reverse transcriptase. The system directly writes new genetic information into the genome using an extended guide RNA (pegRNA) that carries a template for the desired edit.

The pegRNA contains not only the targeting spacer but also a primer binding site (PBS) and a reverse transcription template (RTT) encoding the desired edit. After nicking the DNA, the PBS anneals to the exposed strand, allowing the reverse transcriptase to synthesize new DNA containing the programmed changes. This mechanism enables all 12 types of point mutations plus small insertions and deletions without requiring double-strand breaks or donor DNA templates.



Advantages Over Standard CRISPR

Prime editing achieves insertions and deletions up to 80bp, all 12 base-to-base conversions, and reduced off-target activity. The single-strand nick is less genotoxic than DSBs, minimizing unwanted indels and large deletions. Efficiency ranges from 0-60% depending on edit type and genomic context.

Optimization Strategies:

- ▶ PBS length 10-17nt; RTT length optimized for each edit (10-40nt)
- ▶ Second-strand nick (PE3) increases efficiency 2-6 fold
- ▶ Position edit 1-10bp from nick site for best results
- ▶ Use enhanced PE variants (ePE, PEmax) for higher activity

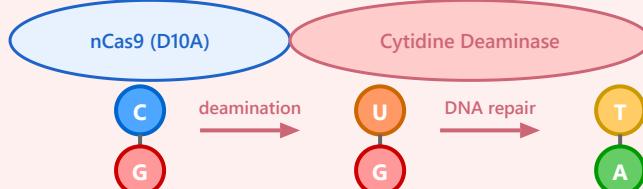
5 Base Editing

Direct chemical conversion of DNA bases without cutting

Base Editor Types & Mechanisms

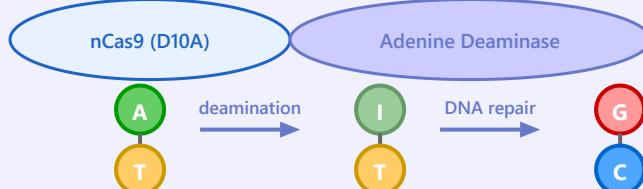
Cytosine Base Editor (CBE)

C → T (or G → A on complementary strand)



Adenine Base Editor (ABE)

A → G (or T → C on complementary strand)



U = Uracil (intermediate) | I = Inosine (intermediate) | Editing window: typically positions 4-8 from

Base editors enable precise single-nucleotide changes without creating double-strand breaks or requiring donor DNA templates. These molecular machines fuse a catalytically impaired Cas9 (nickase) to a deaminase enzyme that chemically converts one base to another. This approach achieves high-efficiency point mutations with minimal indel formation.

Cytosine Base Editors (CBEs) convert C•G to T•A base pairs through cytidine deamination, creating a uracil intermediate that is processed by cellular DNA repair machinery. Adenine Base Editors (ABEs) perform the reverse transition, converting A•T to G•C through adenosine deamination to inosine, which is read as guanine by polymerases. Together, these editors enable four of the 12 possible base transitions, representing approximately 50% of known pathogenic point mutations.



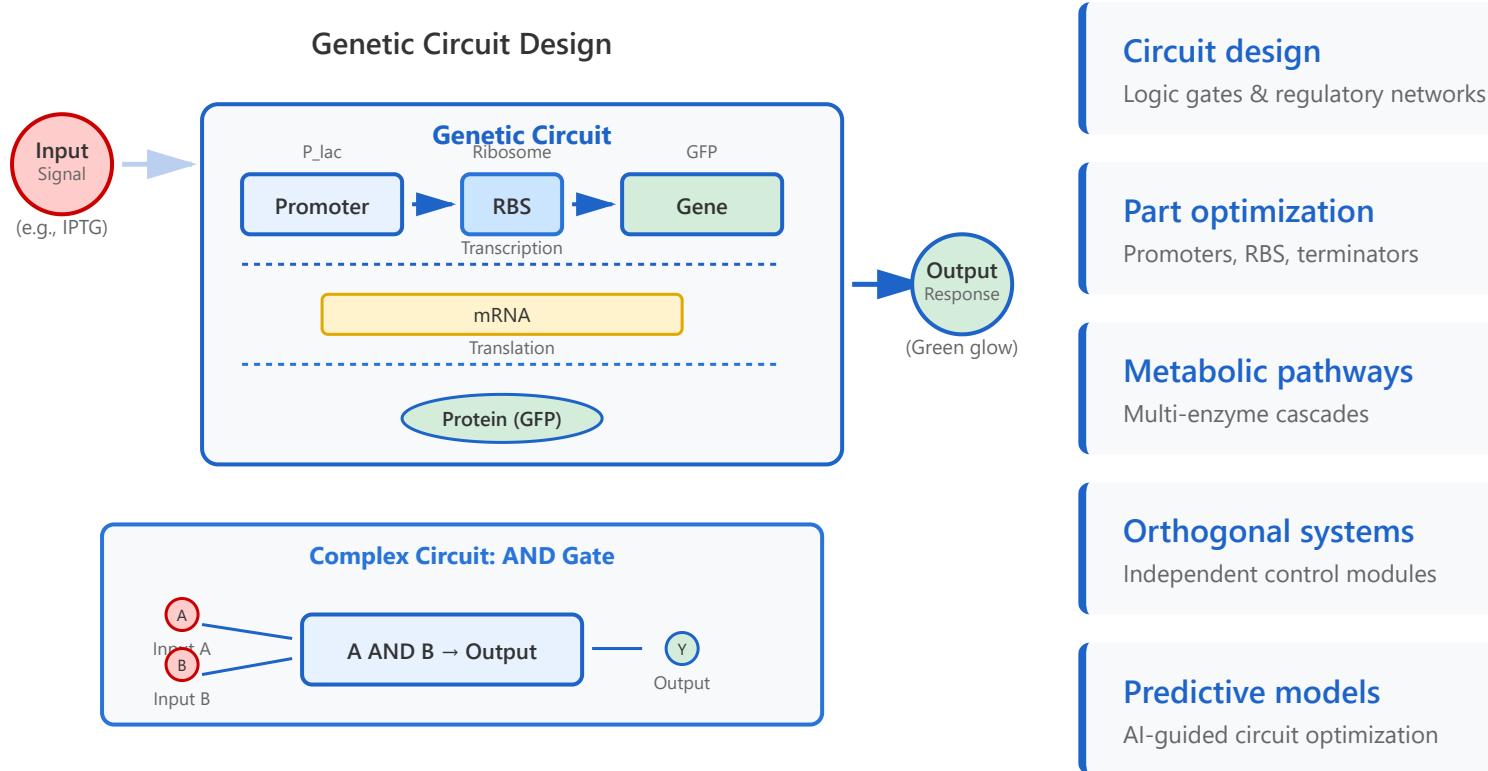
Clinical Applications

Base editors are being developed to correct disease-causing mutations including those in sickle cell disease (HBB E6V), progeria (LMNA G608G), and hereditary hemochromatosis (HFE C282Y). Their precision and low indel rates make them ideal candidates for therapeutic applications where unintended mutations must be minimized.

Design Considerations:

- ▶ Editing window typically spans positions 4-8 (CBE) or 4-7 (ABE) from PAM
- ▶ Check for bystander bases in editing window (may cause unwanted edits)
- ▶ Use narrow-window variants (BE4max, ABE8e) to minimize bystanders
- ▶ Consider RNA off-targets for cytidine deaminases (CBE > ABE)

Synthetic Biology

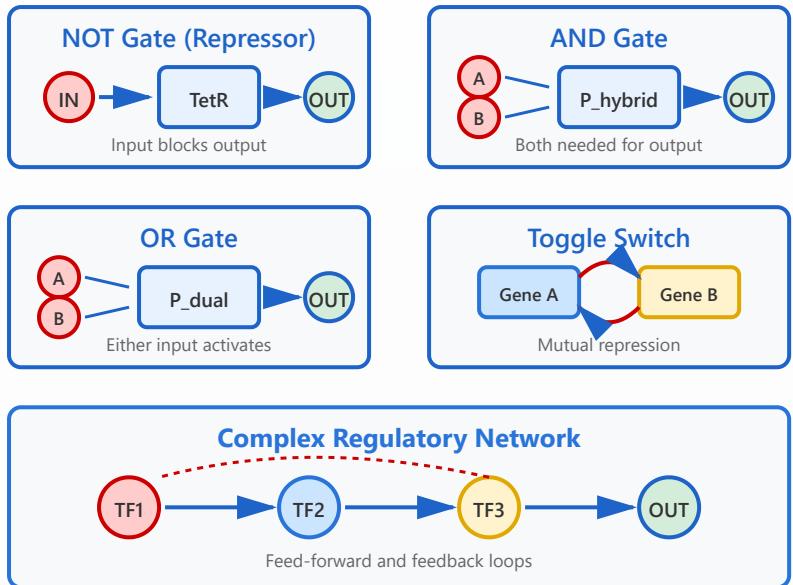


1

Circuit Design

Logic gates & regulatory networks

Biological Logic Gates



Circuit design in synthetic biology involves creating genetic constructs that process inputs and generate outputs, similar to electronic circuits. By combining promoters, repressors, and activators, researchers engineer biological logic gates that perform computational functions within living cells.

These circuits enable cells to make decisions based on environmental conditions, such as sensing chemicals, responding to light, or detecting disease markers. The complexity ranges from simple NOT gates using transcriptional repressors to sophisticated multi-layer networks with feedback and feed-forward loops.

Key Design Principles

- **Modularity:** Genetic parts can be combined like building blocks to create complex behaviors
- **Orthogonality:** Circuit components should not interfere with host cell processes
- **Tunability:** Expression levels can be adjusted through promoter and RBS strength
- **Reliability:** Circuits must function consistently across different conditions

Real-World Application

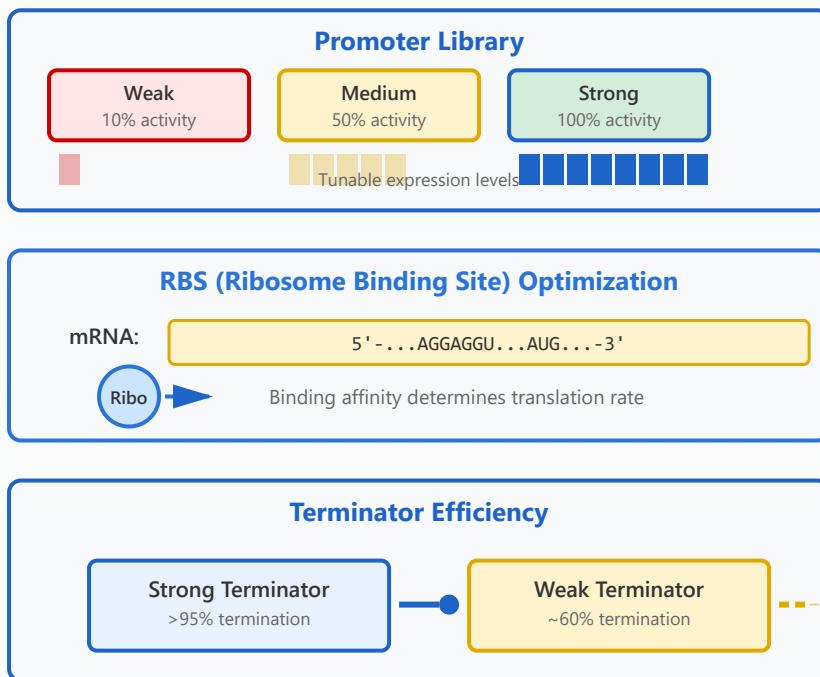
The toggle switch (genetic bistable system) was one of the first synthetic circuits, creating a memory element in *E. coli*. This circuit uses two genes that mutually repress each other, allowing cells to "remember" which state they're in. Such switches are now used in biosensors, therapeutic delivery systems, and cellular computation.

2

Part Optimization

Promoters, RBS, terminators

Genetic Part Standardization



Part optimization focuses on standardizing and characterizing genetic components to ensure predictable circuit behavior. The Registry of Standard Biological Parts (BioBricks) catalogs thousands of characterized DNA sequences including promoters, ribosome binding sites (RBS), coding sequences, and terminators.

Promoters control transcription initiation and can be constitutive (always on) or inducible (regulated by signals). RBS sequences determine translation efficiency by affecting ribosome binding to mRNA. Terminators stop transcription and prevent read-through. Optimizing these parts allows fine-tuned control over gene expression levels.

Optimization Strategies

- **Sequence libraries:** Collections of variants with characterized strengths enable precise tuning
- **Computational prediction:** Tools like RBS Calculator predict translation rates from sequence
- **Context dependency:** Part behavior can vary based on surrounding sequences and host organism

- **Standardization:** Common assembly standards (BioBrick, Golden Gate, Gibson) facilitate part sharing

Engineering Example

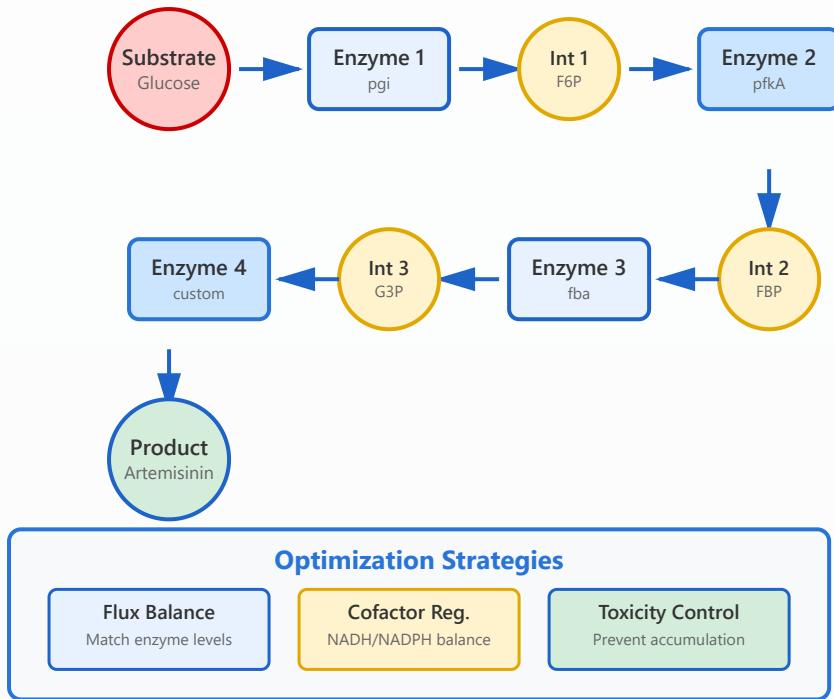
The Anderson Promoter Collection provides a library of constitutive promoters with characterized relative strengths from 1% to 100% of a reference promoter. Researchers can select appropriate promoters to balance pathway enzymes, preventing metabolic bottlenecks and optimizing production of valuable compounds like artemisinin or biofuels.

3

Metabolic Pathways

Multi-enzyme cascades

Engineered Metabolic Pathway



Metabolic pathway engineering involves reconstructing or modifying multi-step biochemical reactions to produce valuable compounds. By introducing heterologous enzymes from different organisms and optimizing their expression, researchers can create entirely new biosynthetic routes or enhance existing ones.

Success requires careful balancing of enzyme levels to prevent metabolic bottlenecks and toxic intermediate accumulation. Computational models help predict optimal enzyme ratios, while directed evolution and protein engineering improve enzyme efficiency and specificity. Integration with host metabolism must be carefully managed to maintain cell viability.

Engineering Challenges

- **Metabolic burden:** Overexpression of heterologous enzymes can reduce growth rate
- **Cofactor availability:** Pathways requiring NADPH or ATP need balanced regeneration
- **Intermediate toxicity:** Some pathway intermediates may inhibit cell growth
- **Enzyme compatibility:** Optimal pH and temperature may differ between enzymes

Breakthrough Achievement

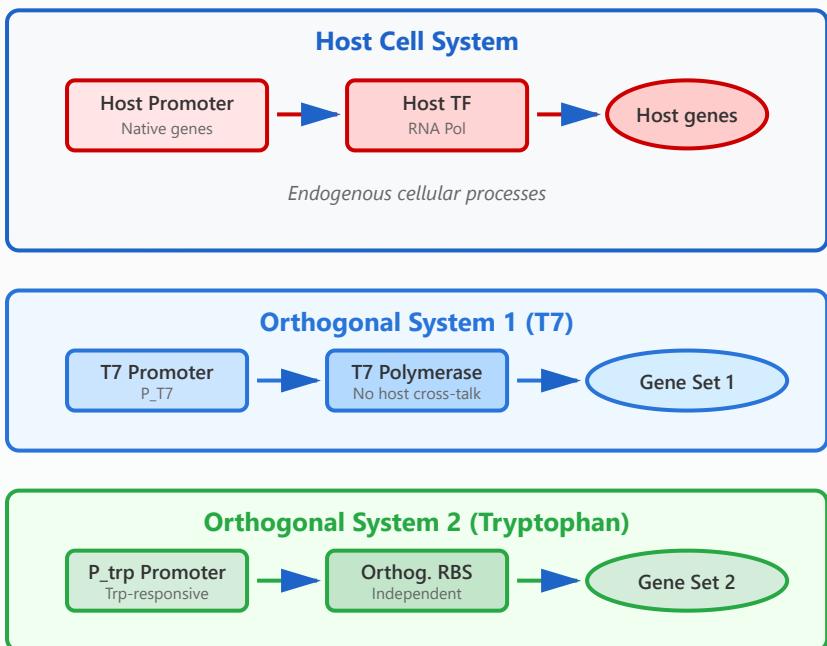
Artemisinin, an antimalarial drug, was traditionally extracted from sweet wormwood plants at low yields. Scientists engineered yeast with a 10-enzyme pathway from three different organisms to produce artemisinic acid, a precursor that's chemically converted to artemisinin. This approach

dramatically increased accessibility and reduced cost of this life-saving medicine.

4 Orthogonal Systems

Independent control modules

Orthogonal Control Systems



Orthogonal systems are genetic circuits designed to function independently of the host cell's native regulatory machinery. This independence prevents unintended interactions between synthetic circuits and endogenous cellular processes, enabling more predictable and robust engineered behaviors.

Key examples include T7 RNA polymerase systems, which use a viral polymerase that recognizes only T7 promoters, and orthogonal ribosomes with modified rRNA that translate only specially designed mRNAs. These systems allow researchers to create multiple independent layers of gene regulation within a single cell.

Design Principles

- **Minimal cross-talk:** Orthogonal parts don't interact with host machinery or other circuits
- **Predictable behavior:** Isolated from cellular context, making design more reliable
- **Scalability:** Multiple orthogonal systems can operate simultaneously in one cell

- **Resource partitioning:** Separate resource pools prevent competition with host

Advanced Application

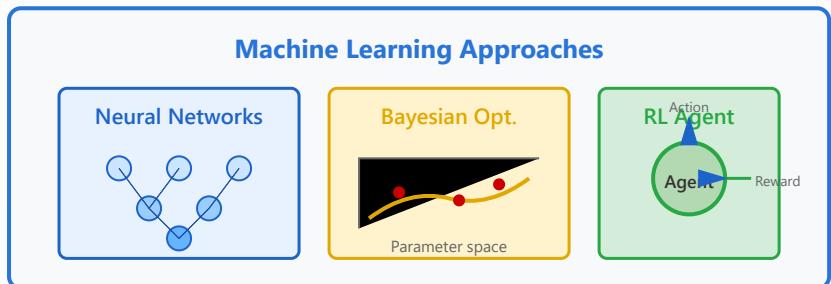
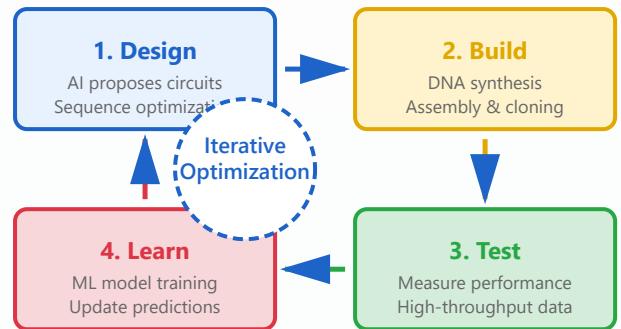
Orthogonal ribosome systems enable parallel translation channels in *E. coli*. Scientists evolved ribosomes that recognize unique Shine-Dalgarno sequences not used by natural ribosomes. This allows simultaneous production of different proteins with independent control, useful for complex metabolic engineering and preventing burden from overexpression on native protein synthesis.

5

Predictive Models

AI-guided circuit optimization

AI-Driven Design Cycle



Predictive modeling combines machine learning with synthetic biology to accelerate the design-build-test cycle. AI algorithms learn from experimental data to predict circuit behavior, optimize DNA sequences, and suggest design improvements, dramatically reducing the time and cost of circuit development.

Deep learning models can predict promoter strength, RBS efficiency, and protein expression from sequence alone. Bayesian optimization helps navigate vast design spaces to find optimal parameter combinations. Reinforcement learning agents discover novel circuit architectures by exploring and learning from simulation or experimental results.

Applications & Benefits

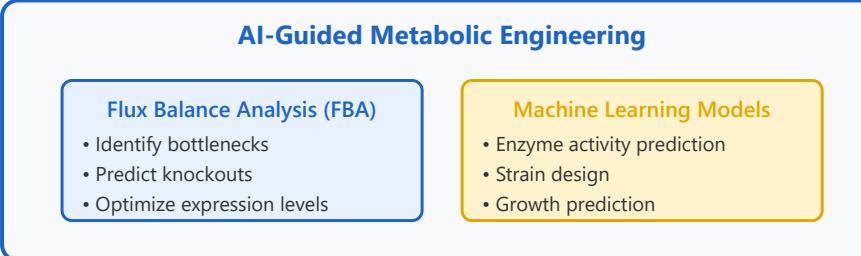
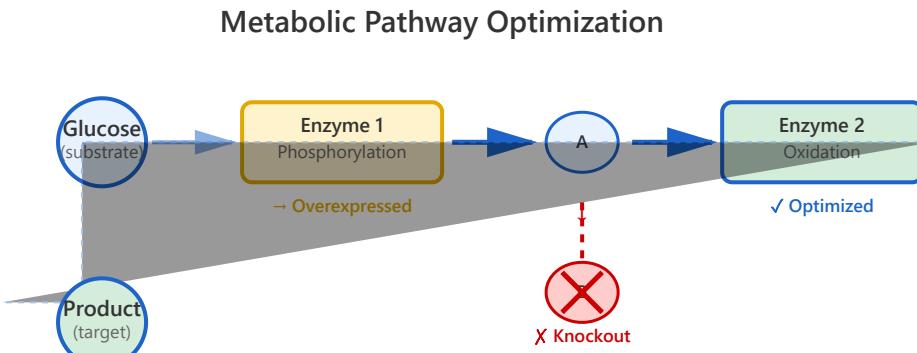
- **Sequence optimization:** Predict functional DNA sequences without extensive testing
- **Circuit modeling:** Simulate dynamic behavior and identify design flaws early
- **Parameter tuning:** Find optimal expression levels and regulatory strengths
- **Knowledge extraction:** Discover design rules from successful circuits

Cutting-Edge Research

Recent work uses transformer-based language models pretrained on millions of DNA sequences to predict gene expression with unprecedented accuracy. These models can design entirely new promoters with specific expression profiles, and even propose novel

genetic circuits that human designers hadn't considered. This AI-guided approach reduced design iterations from dozens to just a few attempts.

Metabolic Engineering



Flux optimization

Balance metabolic flow

Enzyme engineering

Improve catalytic efficiency

Pathway design

Novel biosynthetic routes

Strain optimization

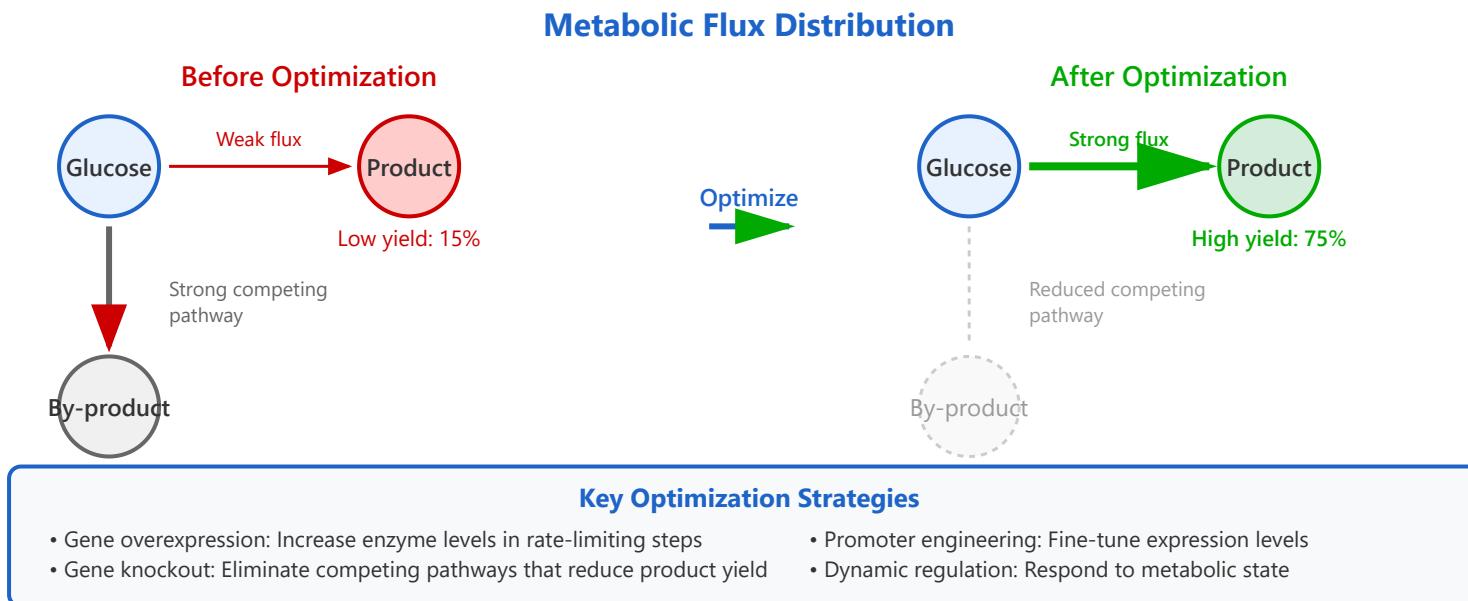
Host organism engineering

Scale-up prediction

Lab → production modeling

1. Flux Optimization

Flux optimization involves redistributing metabolic flux through cellular pathways to maximize the production of desired compounds. This strategy focuses on balancing the flow of metabolites to prevent bottlenecks and accumulation of toxic intermediates while maximizing product yield.



Real-World Example: 1,3-Propanediol Production

DuPont engineered E. coli to produce 1,3-propanediol (a polymer precursor) from glucose. By redirecting flux from glycerol to 1,3-propanediol through overexpression of dhaB and dhaT genes while knocking out competing pathways, they achieved industrial-scale production with yields exceeding 130 g/L, making it commercially viable for producing Sorona® polymer.

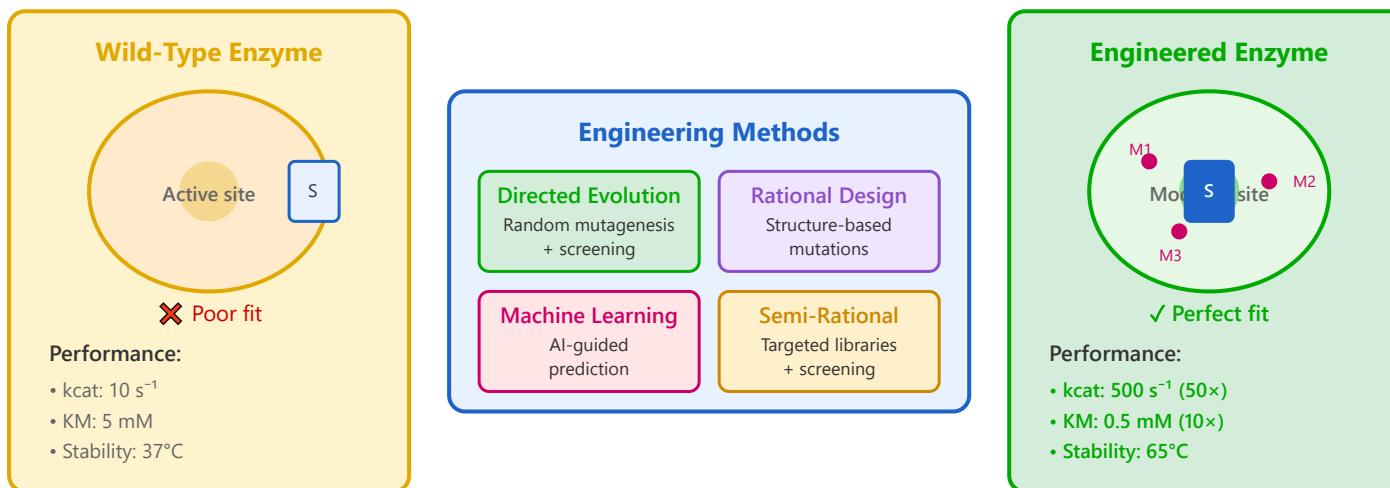
- Computational tools like Flux Balance Analysis (FBA) predict optimal gene modifications
- ¹³C metabolic flux analysis experimentally validates flux distributions

- ▶ Balancing growth rate and production rate is critical for industrial applications

2. Enzyme Engineering

Enzyme engineering focuses on improving the catalytic properties of enzymes through directed evolution, rational design, or computational methods. This approach enhances reaction rates, substrate specificity, stability, and overall pathway efficiency.

Enzyme Engineering Approaches



Real-World Example: Artemisinic Acid Production

Researchers at UC Berkeley and Amyris engineered cytochrome P450 enzymes in yeast to produce artemisinic acid, a precursor to the anti-malaria drug artemisinin. Through directed evolution and rational design, they improved enzyme activity by over 200-fold, enabling economical semi-synthetic production of this life-saving medication that was previously only available from plant extraction.

- ▶ AlphaFold and RoseTTAFold enable structure-guided enzyme design
- ▶ Deep learning models predict beneficial mutations with >80% accuracy

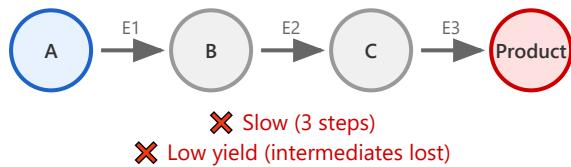
- ▶ Combination approaches (rational + directed evolution) often yield best results
- ▶ Enzyme thermostability is crucial for industrial bioprocesses

3. Pathway Design

Pathway design involves constructing novel biosynthetic routes by introducing heterologous genes or creating entirely synthetic pathways. This strategy enables production of compounds not naturally made by the host organism or improves efficiency through shortened pathways.

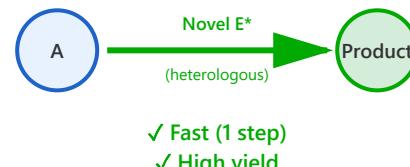
Pathway Design Strategies

Natural Pathway (Long)



✗ Slow (3 steps)
✗ Low yield (intermediates lost)

Engineered Pathway (Optimized)



Modular Pathway Construction

Module 1: Uptake



Transporter



Kinase

Module 2: Conversion



Enzyme A



Enzyme B

Module 3: Export



Efflux pump



Regulator

Parts Library

- ❑ Promoters (strong/weak)
- ❑ RBS variants
- ❑ Terminators
- ❑ Regulatory elements

Real-World Example: Taxol (Paclitaxel) Precursor Production

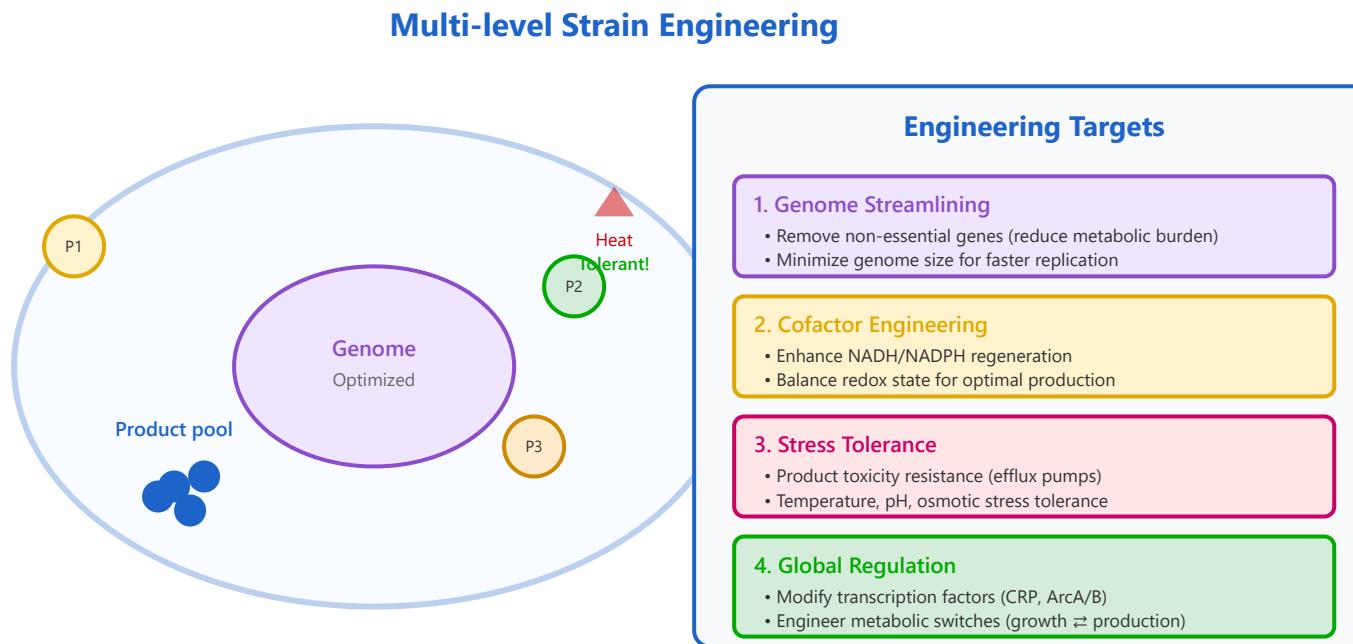
Researchers assembled a 13-gene heterologous pathway in yeast to produce taxadiene, a key precursor to the anti-cancer drug Taxol. This pathway combined genes from Pacific yew trees with engineered yeast enzymes, creating a biosynthetic route that would be impossible in any natural organism. The modular approach allowed rapid optimization of each pathway segment independently.

- ▶ Retrosynthetic analysis identifies optimal enzymatic routes to target molecules

- ▶ Codon optimization ensures proper expression of heterologous genes
- ▶ Enzyme scaffolding and spatial organization reduce intermediate diffusion
- ▶ Standardized genetic parts (BioBricks, MoClo) accelerate pathway construction

4. Strain Optimization

Strain optimization encompasses the systematic improvement of host organisms to enhance production capabilities. This includes modifying cellular machinery, improving tolerance to stress, and engineering global regulatory networks for optimal biosynthesis.



Real-World Example: Tolerance-Enhanced *E. coli* for Biofuel Production

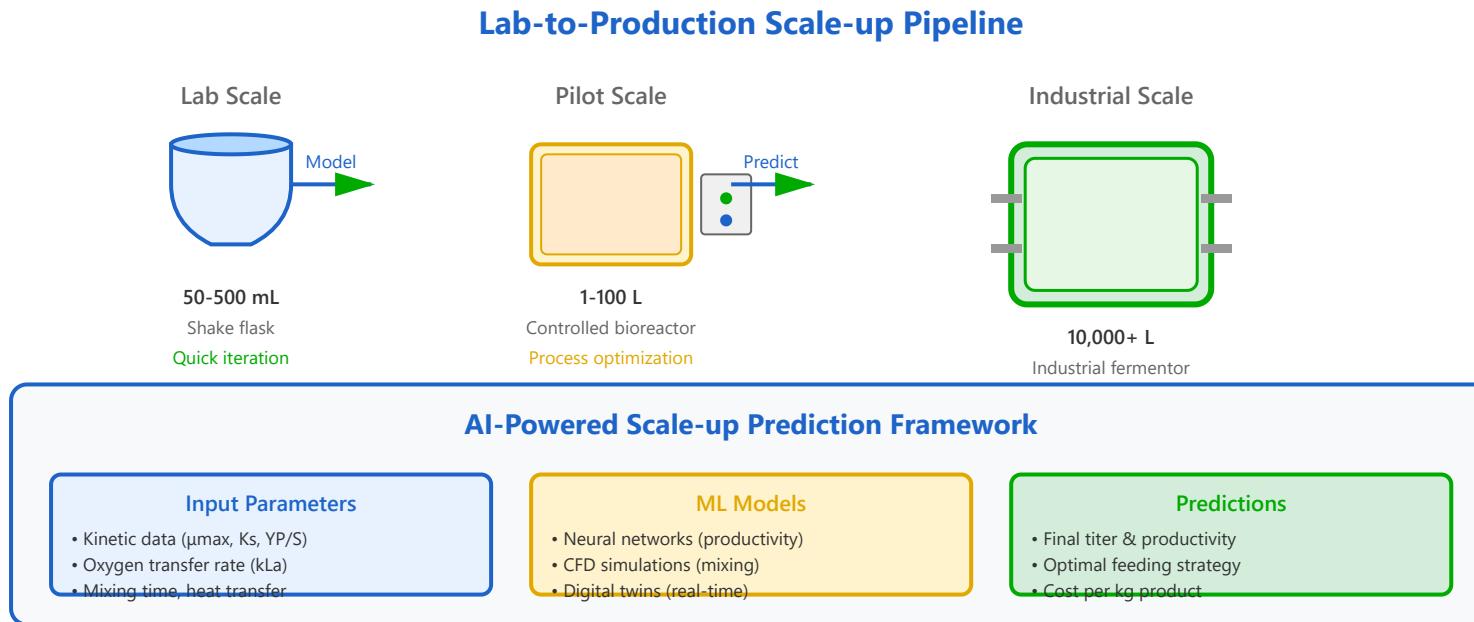
Scientists at MIT engineered *E. coli* for advanced biofuel production by incorporating multiple modifications: deletion of 15 non-essential gene clusters to reduce metabolic burden, overexpression of stress response genes (*groESL*, *dnaKJ*), and modification of membrane composition to tolerate high concentrations of toxic alcohols and fatty acids. The optimized strain showed 3-fold improvement in titer and could operate continuously in fed-batch fermentation for over 200 hours.

- Adaptive laboratory evolution (ALE) complements rational engineering for complex phenotypes

- ▶ CRISPR interference (CRISPRi) enables fine-tuning of gene expression without knockouts
- ▶ Systems biology approaches identify hidden bottlenecks in cellular metabolism
- ▶ Chassis organisms (*E. coli*, *S. cerevisiae*, *B. subtilis*) each offer unique advantages

5. Scale-up Prediction

Scale-up prediction uses computational models and machine learning to forecast how laboratory-scale bioprocesses will perform at industrial scale. This strategy reduces costly trial-and-error in pilot plants and enables rapid translation from bench to commercial production.



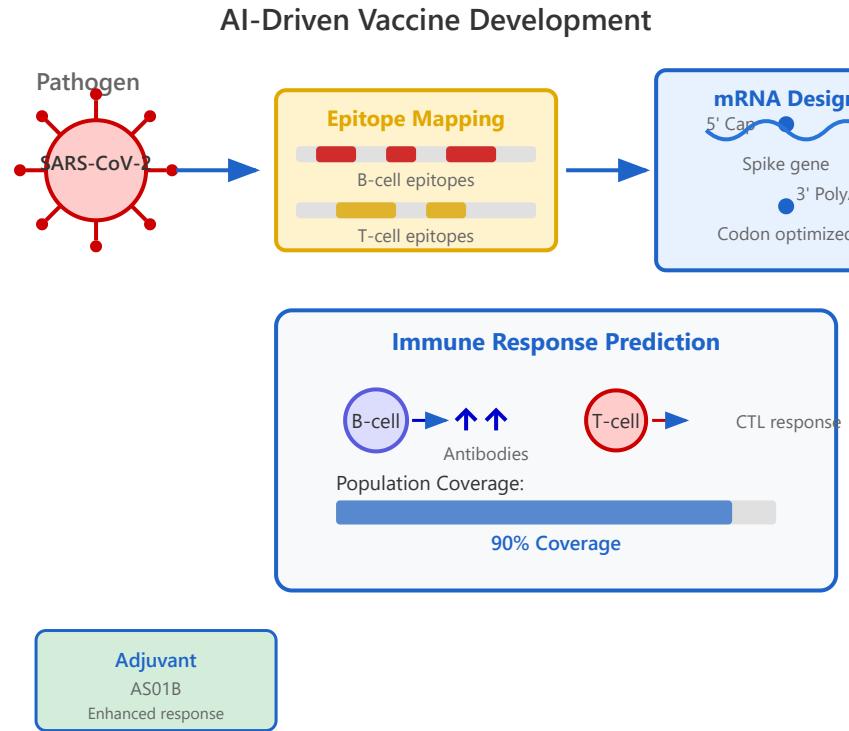
Real-World Example: Ginkgo Bioworks Scale-up Platform

Ginkgo Bioworks developed a machine learning platform that predicts fermentation performance at 10,000L+ scale from 96-well plate data. By training models on thousands of fermentation runs, their system predicts titer within 15% accuracy and recommends optimal media composition and feeding strategies. This reduced their time from strain design to commercial production from 18 months to under 6 months for multiple products including fragrance molecules and pharmaceutical intermediates.

- ▶ Computational Fluid Dynamics (CFD) models predict mixing and mass transfer at scale

- ▶ Machine learning on historical data outperforms mechanistic models for complex systems
- ▶ Digital twins enable real-time process optimization during production
- ▶ Scale-down simulators validate predictions before committing to large-scale runs
- ▶ Techno-economic analysis (TEA) integrated with models guides cost-effective scale-up

Vaccine Design



Epitope prediction

B-cell & T-cell epitopes

Immunogenicity

Immune response modeling

Coverage optimization

Population HLA diversity

Adjuvant selection

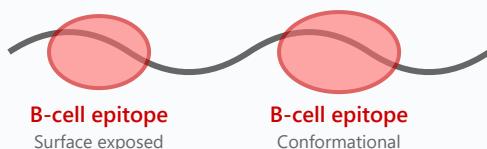
Enhance immune response

mRNA design

Codon optimization & stability

1. Epitope Prediction

Pathogen Antigen



Processed Peptides

FLKDCVMYV	AQFAPSASAFFGMS
CD8+ T-cell epitope (MHC Class I)	CD4+ T-cell epitope (MHC Class II)

AI-Based Prediction Tools

BepiPred 3.0 B-cell epitopes	NetMHCPan MHC binding
IEDB Analysis Immunogenicity	AlphaFold 3D structure

Overview

Epitope prediction identifies specific regions on pathogen proteins that are recognized by the immune system. These epitopes serve as the primary targets for vaccine-induced immunity.

B-Cell Epitopes

B-cell epitopes are recognized by antibodies and can be:

- **Linear epitopes:** Continuous amino acid sequences (5-15 residues)
- **Conformational epitopes:** Discontinuous sequences brought together by protein folding
- **Surface accessibility:** Must be exposed on the pathogen surface

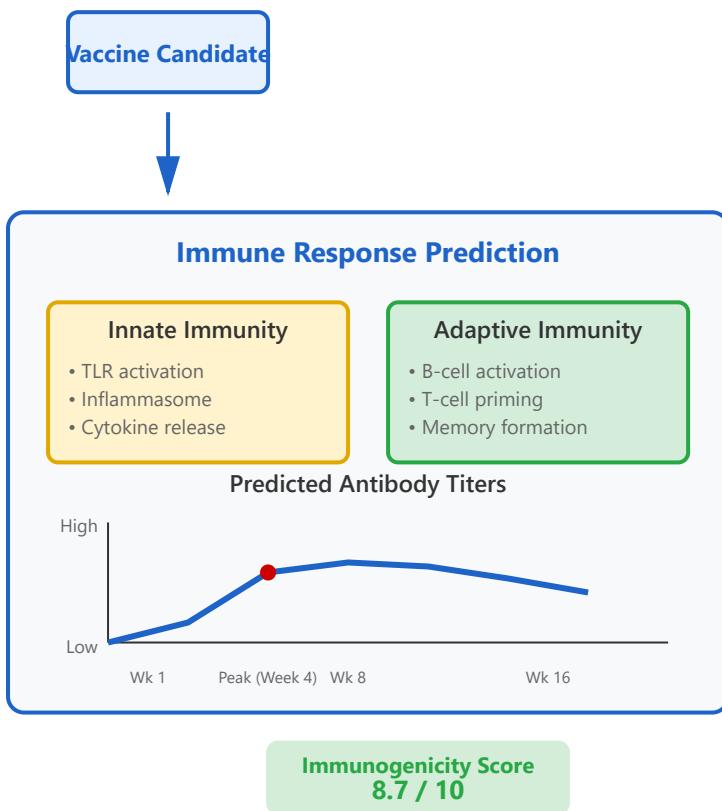
T-Cell Epitopes

T-cell epitopes are short peptide sequences presented by MHC molecules:

- **MHC Class I (8-11 amino acids):** Activates CD8+ cytotoxic T-cells
- **MHC Class II (13-25 amino acids):** Activates CD4+ helper T-cells
- **Processing requirements:** Must be cleaved and loaded properly

Key Tools: Modern epitope prediction uses machine learning algorithms like BepiPred 3.0, NetMHCPan 4.1, and IEDB tools, achieving >85% accuracy in identifying immunogenic epitopes.

2. Immunogenicity Prediction



Overview

Immunogenicity prediction assesses how strongly a vaccine candidate will stimulate the immune system. This involves modeling both innate and adaptive immune responses.

Key Factors

- **Antigen dose:** Optimal concentration for immune activation without tolerance
- **Route of administration:** Intramuscular, subcutaneous, or intradermal delivery
- **Adjuvant effects:** Enhancement of immune recognition and response
- **Epitope density:** Number and spacing of immunogenic sites

Response Modeling

AI models predict multiple immune parameters:

- **Antibody titers:** Concentration and kinetics over time
- **T-cell response:** CD4+ and CD8+ activation levels
- **Cytokine profiles:** Type and magnitude of inflammatory response
- **Memory formation:** Long-term protection durability

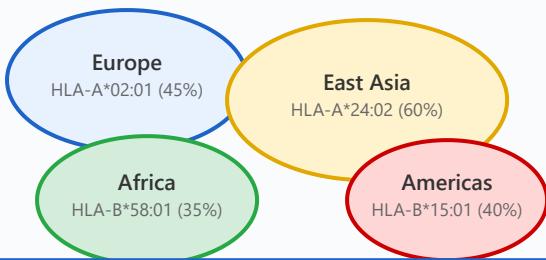
Clinical Validation: Immunogenicity predictions are validated against clinical trial data, with modern algorithms achieving >75% accuracy in predicting successful vaccine candidates.

Safety Considerations

Models also evaluate potential adverse reactions including autoimmunity risk, excessive inflammation, and allergic responses to ensure a favorable benefit-risk profile.

3. Population Coverage Optimization

Global HLA Diversity



Multi-Epitope Vaccine Design

Selected Epitopes:

FLKDCVMYV	HLA-A*02:01	DPFLGVYY	HLA-B*15:01
RYPANSIVR	HLA-A*24:02	YQAGSTPCN	HLA-A*11:01
KQIYKTPPIK	HLA-B*58:01		

Population Coverage Analysis

Global Coverage:

90% Coverage

Europe: 94% | Asia: 88% | Africa: 87% | Americas: 92%

Based on 12 major HLA alleles covering 95% of world population

Overview

Population coverage optimization ensures that a vaccine will be effective across diverse human populations, accounting for genetic variation in immune response genes (HLA alleles).

HLA Diversity Challenge

Human Leukocyte Antigen (HLA) genes are the most polymorphic in the human genome:

- Thousands of alleles:** Over 28,000 HLA alleles identified globally
- Geographic variation:** Different populations have distinct HLA frequency distributions
- Peptide binding specificity:** Each HLA allele binds different peptide sequences

Optimization Strategy

- Multi-epitope approach:** Include 8-15 epitopes targeting multiple HLA alleles
- Frequency weighting:** Prioritize alleles that are common across populations
- Conserved regions:** Select epitopes from pathogen regions with low mutation rates
- Redundancy:** Multiple epitopes per HLA allele for robust coverage

Coverage Targets: Modern vaccines aim for ≥90% population coverage globally. The COVID-19 mRNA vaccines achieve ~95% coverage by targeting highly conserved spike protein epitopes.

Computational Tools

Tools like IEDB Population Coverage, OptiVax, and Vaxign use algorithms to select optimal epitope combinations that maximize coverage while minimizing the number of epitopes needed.

4. Adjuvant Selection

Adjuvant Classification

Aluminum Salts

- Alum ($\text{Al(OH}_3\text{)}$)
 - Aluminum phosphate
- Most widely used

TLR Agonists

- CpG oligonucleotides
 - Monophosphoryl lipid A
- Strong innate activation

Oil Emulsions

- MF59 (squalene)
 - AS03
- Enhanced uptake

Liposome-based

- AS01 (liposome + MPL)
 - Virosomes
- Targeted delivery

Adjuvant Mechanisms

Depot Formation

Sustained antigen release at injection site

Immune Cell Activation

Recruitment and activation of APCs

Controlled Inflammation

Cytokine production and immune signaling

Selection Criteria

- ✓ Antigen compatibility and stability
- ✓ Desired immune response type (Th1/Th2 balance)
- ✓ Safety profile and regulatory approval status

Overview

Adjuvants are substances added to vaccines to enhance and direct the immune response. Proper adjuvant selection is critical for vaccine efficacy and safety.

Functions of Adjuvants

- **Immunopotentiation:** Increase magnitude of immune response
- **Dose-sparing:** Achieve protection with less antigen
- **Response shaping:** Direct toward Th1 or Th2 response
- **Duration enhancement:** Prolong immune memory

Major Adjuvant Classes

Aluminum-based adjuvants: Most commonly used, promote Th2 responses and antibody production. Safe track record with over 70 years of use.

TLR agonists: Activate pattern recognition receptors, inducing strong innate immunity and Th1 responses. Examples include CpG-ODN (TLR9) and MPL (TLR4).

Emulsion-based: Oil-in-water emulsions like MF59 and AS03 enhance antigen uptake and presentation, particularly effective for influenza vaccines.

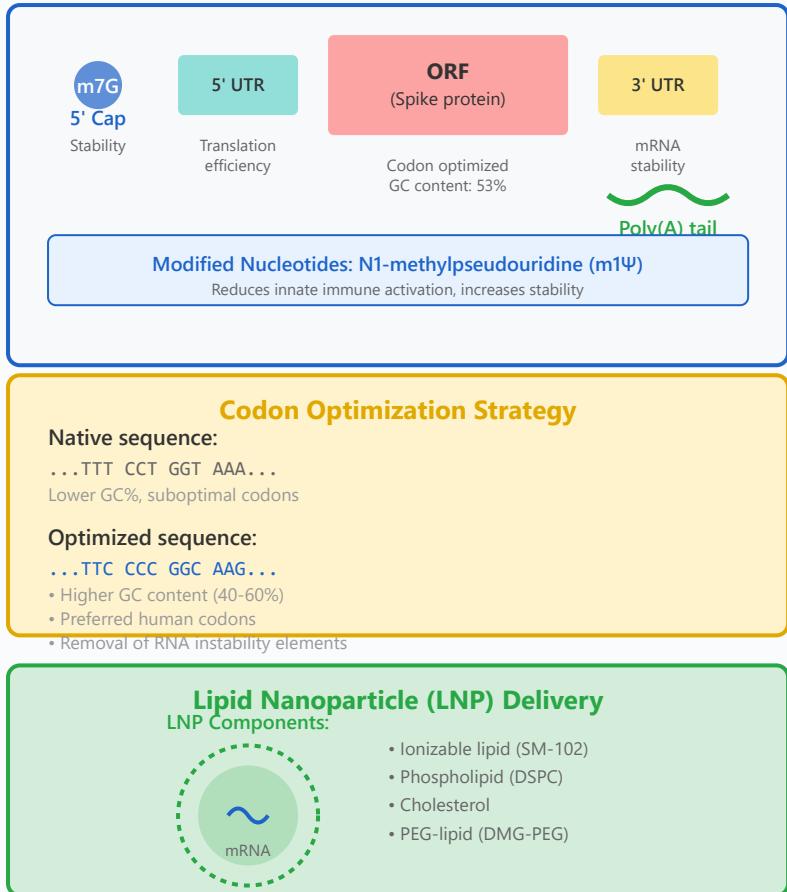
Case Study: The AS01B adjuvant used in the Shingrix vaccine combines liposomes with MPL and QS-21 saponin, achieving >90% efficacy in elderly populations by strongly activating both innate and adaptive immunity.

Selection Strategy

Adjuvant choice depends on target pathogen, patient population (age, immune status), desired response type, and regulatory considerations. AI models can predict optimal adjuvant-antigen combinations based on immunological data.

5. mRNA Vaccine Design

mRNA Vaccine Architecture



Overview

mRNA vaccines represent a revolutionary platform that instructs cells to produce antigens directly. Successful design requires optimization of multiple molecular features for stability, translation efficiency, and immunogenicity.

Key Structural Elements

5' Cap structure: Modified guanosine cap (m7G or cap1) protects against degradation and enables ribosome binding for translation initiation.

Untranslated Regions (UTRs): The 5' UTR contains regulatory elements for translation efficiency, while the 3' UTR provides stability signals and poly(A) binding sites.

Open Reading Frame (ORF): Encodes the target antigen with extensive codon optimization to maximize expression while maintaining protein structure.

Codon Optimization

- **GC content balance:** Target 50-60% for optimal stability and translation
- **Codon usage:** Replace rare codons with frequently used human codons
- **Secondary structure:** Minimize hairpins and self-complementary regions
- **Immune evasion:** Remove CpG dinucleotides and uridine-rich motifs

Modified Nucleotides: COVID-19 mRNA vaccines use N1-methylpseudouridine instead of uridine. This modification reduces innate immune detection (TLR activation), increases translation efficiency by 10-fold, and improves mRNA stability.

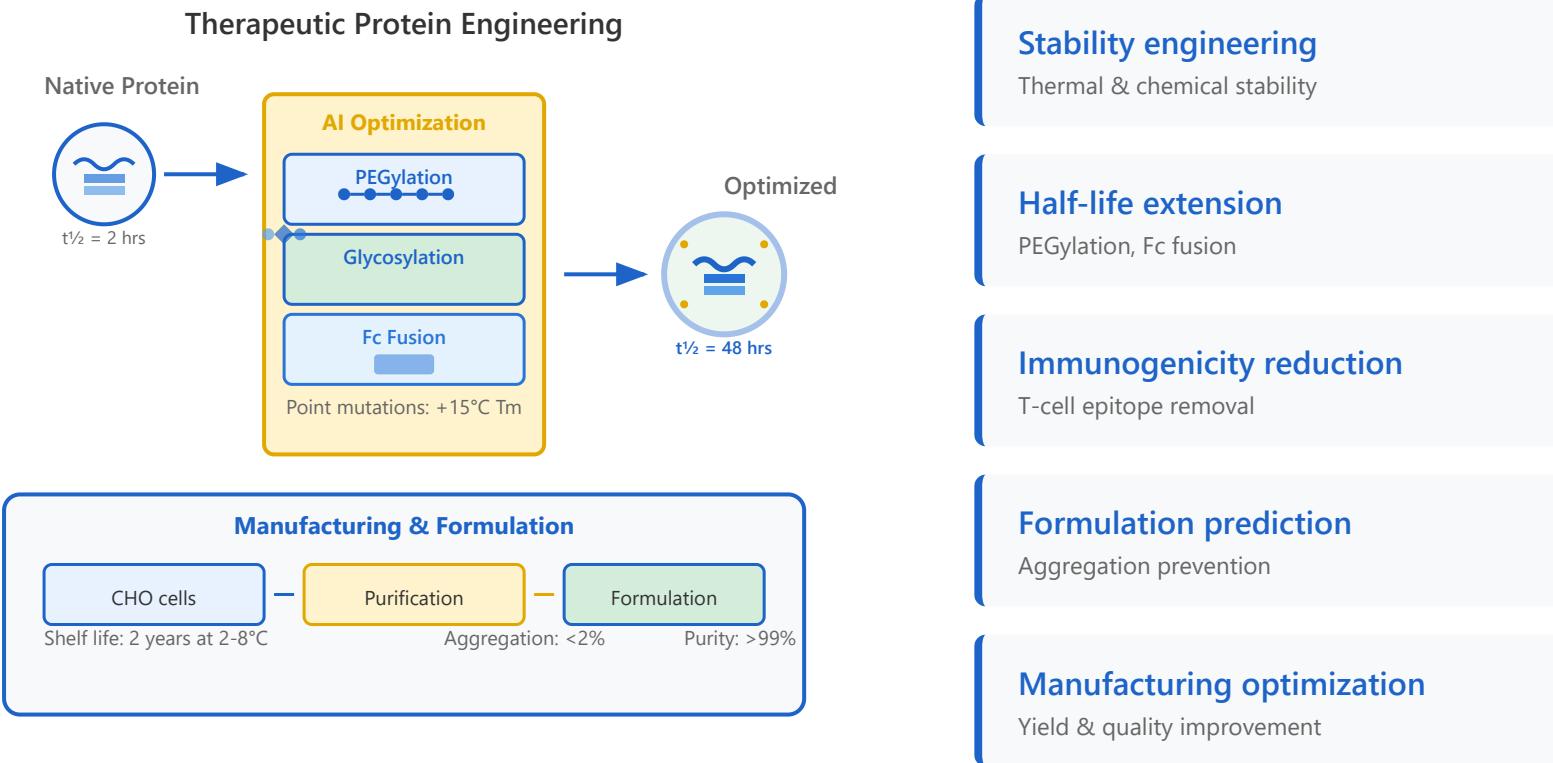
Delivery System

Lipid nanoparticles (LNPs) encapsulate mRNA for protection and cellular delivery. The ionizable lipid component enables endosomal escape, releasing mRNA into the cytoplasm where translation occurs.

Quality Control

Critical parameters include mRNA integrity (>80%), encapsulation efficiency (>90%), particle size (80-100 nm), and endotoxin levels. AI models predict optimal sequences and formulations before synthesis.

Therapeutic Proteins



1. Stability Engineering

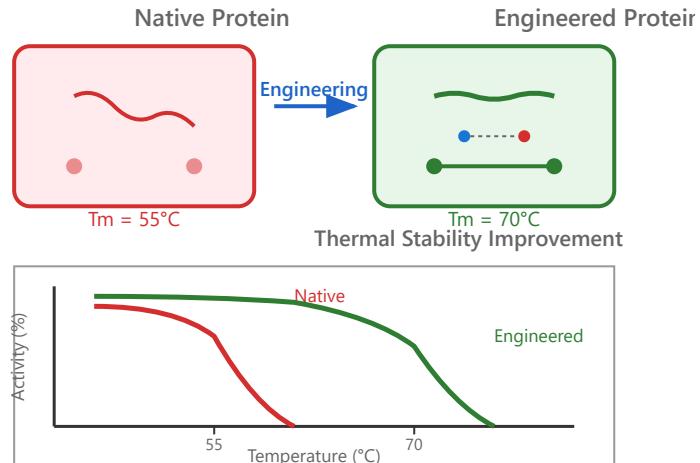
Thermal & Chemical Stability Enhancement

Objective: Improve protein resistance to temperature, pH changes, and chemical degradation to extend shelf life and maintain therapeutic efficacy.

Key Strategies:

- **Disulfide bond engineering** - Introduce strategic cysteine pairs to stabilize protein structure
- **Salt bridge optimization** - Enhance electrostatic interactions between charged residues
- **Hydrophobic core packing** - Improve interior residue arrangement to prevent unfolding
- **Surface charge modification** - Reduce aggregation-prone patches

Stability Engineering Approaches



Clinical Example: Enzyme Replacement Therapy

Recombinant human α -glucosidase (Pompe disease treatment) was engineered with strategic mutations increasing thermal stability from 55°C to 70°C , enabling room temperature storage and reducing cold-chain logistics requirements by 40%.

2. Half-life Extension

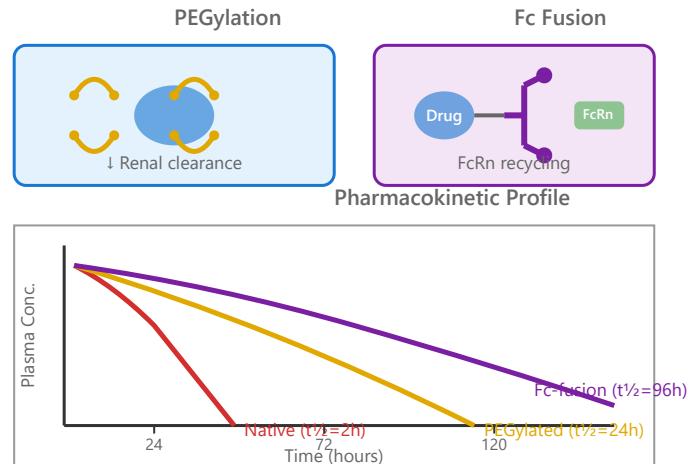
Prolonging Therapeutic Circulation Time

Objective: Extend protein residence time in circulation to reduce dosing frequency and improve patient compliance.

Major Approaches:

- **PEGylation** - Attachment of polyethylene glycol (PEG) chains increases hydrodynamic radius
- **Fc fusion** - Linking to IgG Fc domain enables FcRn-mediated recycling
- **Albumin fusion** - Leverage albumin's long half-life (19 days)
- **Polysialylation** - Attach polysialic acid chains for stealth effect

Half-life Extension Strategies



Clinical Examples:

Pegfilgrastim (Neulasta®): PEGylated G-CSF extends half-life from 3.5 hours to 42 hours, reducing injections from daily to once per chemotherapy cycle.

Etanercept (Enbrel®): TNF receptor-Fc fusion achieves 102-hour half-life, enabling twice-weekly dosing for rheumatoid arthritis treatment.

3. Immunogenicity Reduction

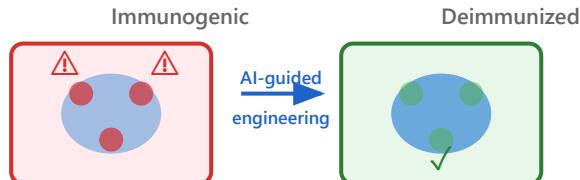
Minimizing Unwanted Immune Response

Objective: Reduce the risk of anti-drug antibodies (ADAs) that can neutralize therapeutic effect or cause adverse reactions.

Key Strategies:

- **T-cell epitope removal** - Identify and eliminate MHC-II binding sequences using computational tools
- **Humanization** - Replace non-human sequences with human framework regions
- **Deimmunization** - Strategic mutations to disrupt epitope binding without affecting function
- **Glycosylation engineering** - Shield immunogenic surfaces with glycan

Immunogenicity Reduction Process



T-cell Epitope Elimination



Clinical Success Story:

Factor VIII for Hemophilia A: Computational deimmunization reduced T-cell epitopes from 26 to 4, decreasing immunogenicity by 73% in preclinical models. Clinical trials showed a 5-fold reduction in inhibitor antibody formation compared to first-generation products, significantly improving treatment outcomes.

4. Formulation Prediction

Aggregation Prevention & Stability Optimization

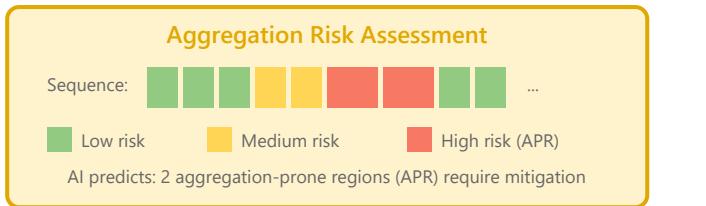
Objective: Design optimal formulation conditions to prevent protein aggregation, maintain stability during storage, and ensure consistent drug product quality.

AI-Driven Approaches:

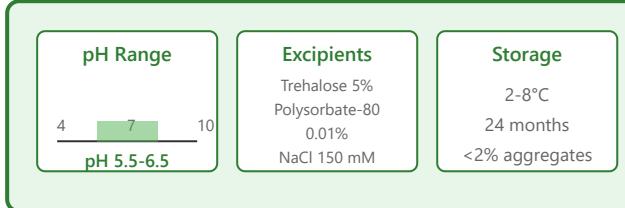
- **Aggregation propensity prediction** - Identify aggregation-prone regions using ML models
- **pH optimization** - Predict optimal pH range for maximum stability

- **Excipient screening** - AI-guided selection of stabilizing agents (sugars, salts, surfactants)
- **Concentration optimization** - Balance high concentration needs with aggregation risk

Formulation Optimization Workflow



Optimized Formulation Conditions



Real-World Application:

Monoclonal Antibody (150 mg/mL): Machine learning models predicted optimal formulation with histidine buffer (pH 6.0), 8% sucrose, and 0.02% polysorbate-80. This formulation achieved 36-month stability at 5°C with less than 1% high molecular weight species, eliminating the need for extensive empirical screening and saving 18 months of development time.

5. Manufacturing Optimization

Yield & Quality Improvement through Process Engineering

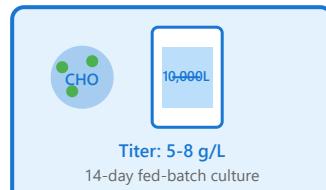
Objective: Maximize protein production efficiency, ensure consistent quality, and reduce manufacturing costs through bioprocess optimization.

Optimization Targets:

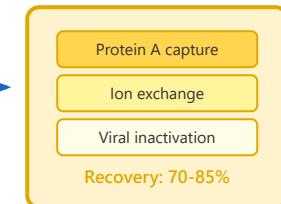
- **Cell line engineering** - CHO, HEK293, E. coli strain optimization for higher titers
- **Expression system** - Codon optimization, signal peptide design, secretion enhancement
- **Culture conditions** - Media composition, temperature, pH, dissolved oxygen
- **Purification strategy** - Chromatography sequences, yield optimization
- **Process analytics** - Real-time monitoring and AI-driven process control

Manufacturing Process Flow

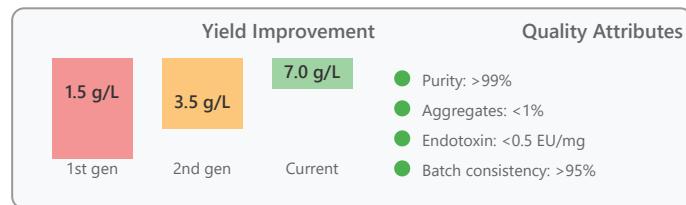
Upstream Processing



Downstream



Key Performance Indicators

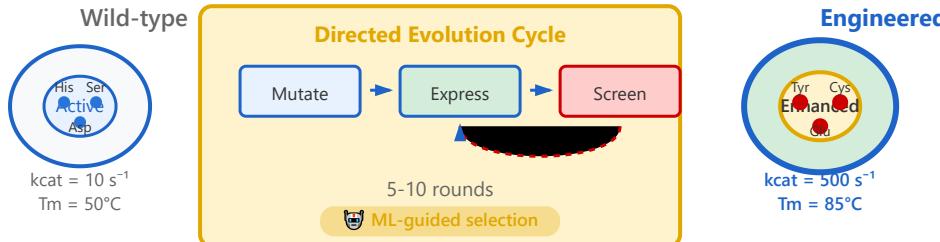


Industry Example:

Adalimumab Biosimilar Production: Process optimization through AI-guided media design and cell line engineering increased volumetric productivity from 2.5 g/L to 7.2 g/L in CHO cells. This 3-fold improvement reduced cost of goods by 60% while maintaining product quality matching the reference product across all critical quality attributes. The optimized process enabled commercial-scale production in smaller bioreactors, significantly reducing capital expenditure requirements.

Enzyme Engineering

Directed Evolution & Rational Design



Catalytic Reaction



Industrial Applications

Biofuel
Cellulase

200% activity

Pharma
Transaminase

99% ee

Detergent
Protease

pH 10, 60°C

Food
Amylase

High temp

Activity improvement

k_{cat}/K_m optimization

Substrate specificity

Promiscuity engineering

Thermostability

High temperature operation

Solvent tolerance

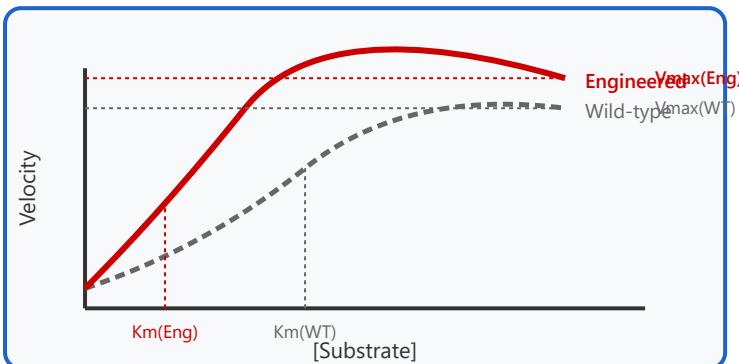
Organic solvent resistance

Directed evolution

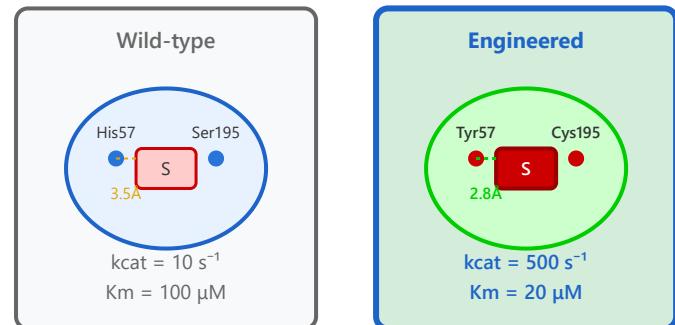
Iterative improvement cycles

1. Activity Improvement (kcat/Km Optimization)

Catalytic Efficiency Enhancement



Active Site Optimization



Objective

Enhance the catalytic efficiency (kcat/Km) of enzymes to increase reaction rates and substrate binding affinity. This is crucial for industrial processes requiring high throughput.

Key Strategies

- Transition state stabilization:** Modify active site residues to better stabilize the transition state
- Substrate binding optimization:** Engineer binding pocket geometry for improved substrate fit
- Product release enhancement:** Reduce product inhibition by facilitating product dissociation
- Catalytic triad engineering:** Optimize spatial arrangement and pKa of catalytic residues

Case Study: Subtilisin Protease

Wild-type: $k_{cat}/K_m = 1.0 \times 10^5 \text{ M}^{-1}\text{s}^{-1}$

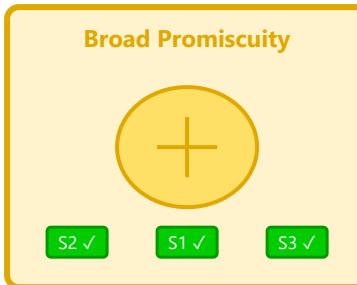
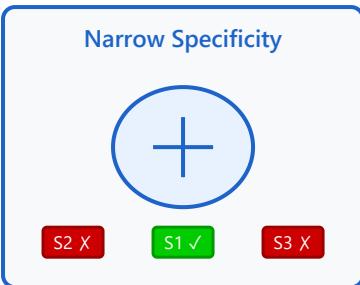
Engineered (N62D/G166D): $k_{cat}/K_m = 5.2 \times 10^6 \text{ M}^{-1}\text{s}^{-1}$

Improvement: 52-fold increase in catalytic efficiency through rational design of the oxyanion hole and substrate binding pocket.

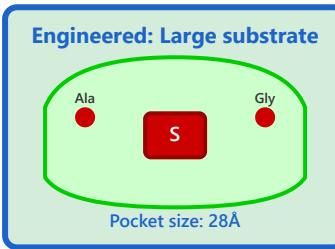
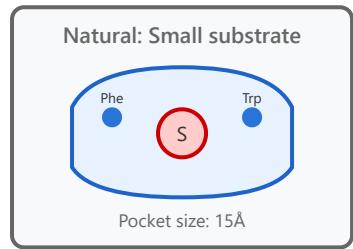
Impact: Activity improvements of 10-1000 fold are achievable through directed evolution combined with computational design, enabling reduced enzyme loading and lower production costs.

2. Substrate Specificity & Promiscuity Engineering

Substrate Recognition Engineering



Binding Pocket Reshaping



Engineering Applications

Stereospecificity
R S
ee > 99%

Regioselectivity
C-2 C-4
C-2:C-4 = 95:5

New Activity
A→B → C→D
Novel pathway

Objective

Modify substrate binding specificity to either narrow selectivity for a single substrate or broaden promiscuity to accept multiple substrates. This enables enzymes to process non-natural substrates or improve stereoselectivity.

Engineering Approaches

- Binding pocket reshaping:** Alter size and geometry through mutations (e.g., Phe→Ala for pocket enlargement)
- Electrostatic tuning:** Change charge distribution to favor specific substrate classes
- Hydrophobic interactions:** Engineer aromatic residues for π-stacking with substrates
- Gatekeeper residue modification:** Control substrate entry and selectivity

Case Study: P450 BM3 Hydroxylase

Wild-type: Hydroxylates C12-C16 fatty acids

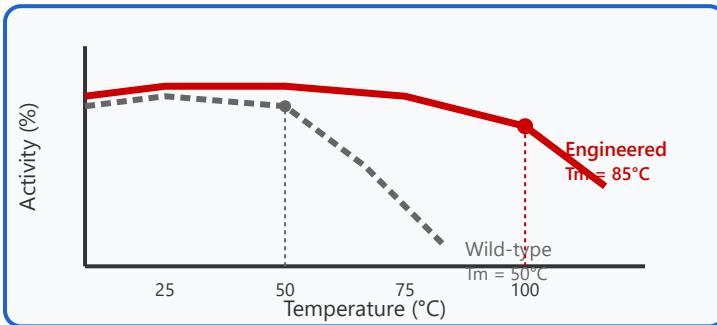
Engineered (9 mutations): Accepts propane and alkanes (C2-C8)

Achievement: Complete substrate scope inversion - from long-chain to short-chain hydrocarbons, enabling production of valuable chemical intermediates.

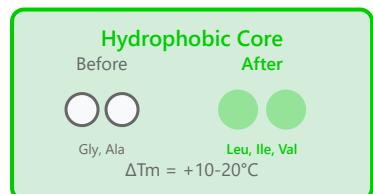
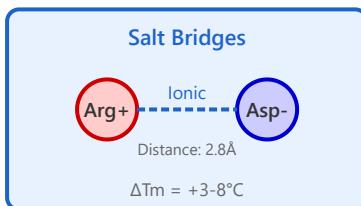
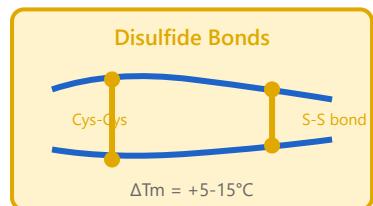
Applications: Pharmaceutical synthesis (>99% ee), biocatalytic cascades, plastic degradation, and production of non-natural amino acids and chemicals.

3. Thermostability Engineering

Thermal Stability Enhancement



Molecular Stabilization Mechanisms



Objective

Increase enzyme thermal stability (T_m) to enable operation at elevated temperatures, which improves reaction rates, reduces contamination risks, and extends enzyme shelf life in industrial processes.

Stabilization Strategies

- Disulfide bonds:** Introduce Cys-Cys bridges to constrain structure (5-15°C increase)
- Salt bridges:** Engineer ionic interactions between charged residues (3-8°C increase)
- Hydrophobic core packing:** Replace small residues (Gly, Ala) with bulky hydrophobic ones (Leu, Ile, Val) for tighter packing (10-20°C increase)
- Proline substitution:** Reduce loop flexibility by inserting proline (5-12°C increase)
- N/C-terminal modifications:** Add stabilizing residues or tags
- Glycosylation:** Attach sugar moieties for protection

Case Study: *Bacillus α-Amylase*

Wild-type T_m : 55°C (half-life: 15 min at 90°C)

Engineered (15 mutations): $T_m = 95^\circ\text{C}$ (half-life: 7 hours at 90°C)

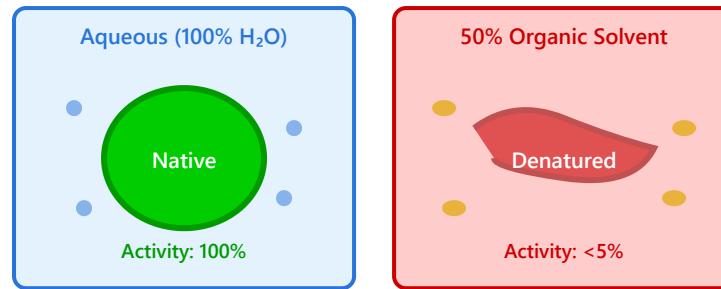
Key mutations: Introduction of 3 disulfide bonds, 5 salt bridges, and 7 core packing improvements. Now used in high-temperature starch processing.

Industrial Impact: Thermostable enzymes enable higher process temperatures (70-100°C), reducing viscosity, increasing mass

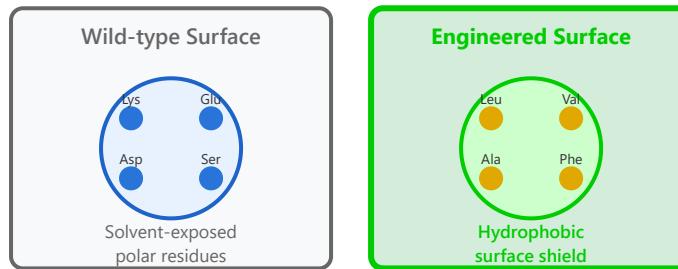
transfer, and preventing microbial contamination without sterilization.

4. Solvent Tolerance Engineering

Organic Solvent Resistance



Surface Engineering for Solvent Tolerance



Solvent Compatibility Enhancement

Methanol	WT: 2%	DMSO	WT: <1%
	Eng: 65%		Eng: 40%
Acetonitrile	WT: 5%	Acetonitrile	WT: 5%
	Eng: 70%		Eng: 35%
Toluene	WT: 0%	Toluene	WT: 0%
	Eng: 35%		Eng: 35%

Objective

Engineer enzymes to maintain activity and stability in organic solvents, enabling reactions with hydrophobic substrates and products that are poorly soluble in water.

Engineering Strategies

- Surface hydrophobicity:** Replace charged surface residues (Lys, Glu, Asp) with hydrophobic ones (Leu, Val, Ala, Phe)
- Core stabilization:** Strengthen hydrophobic core to resist solvent penetration
- Removal of water-binding sites:** Eliminate surface pockets that trap destabilizing water molecules
- Increased rigidity:** Reduce conformational flexibility through proline substitutions and disulfide bonds
- Active site protection:** Shield catalytic residues from solvent deactivation

Case Study: *Candida antarctica* Lipase B (CALB)

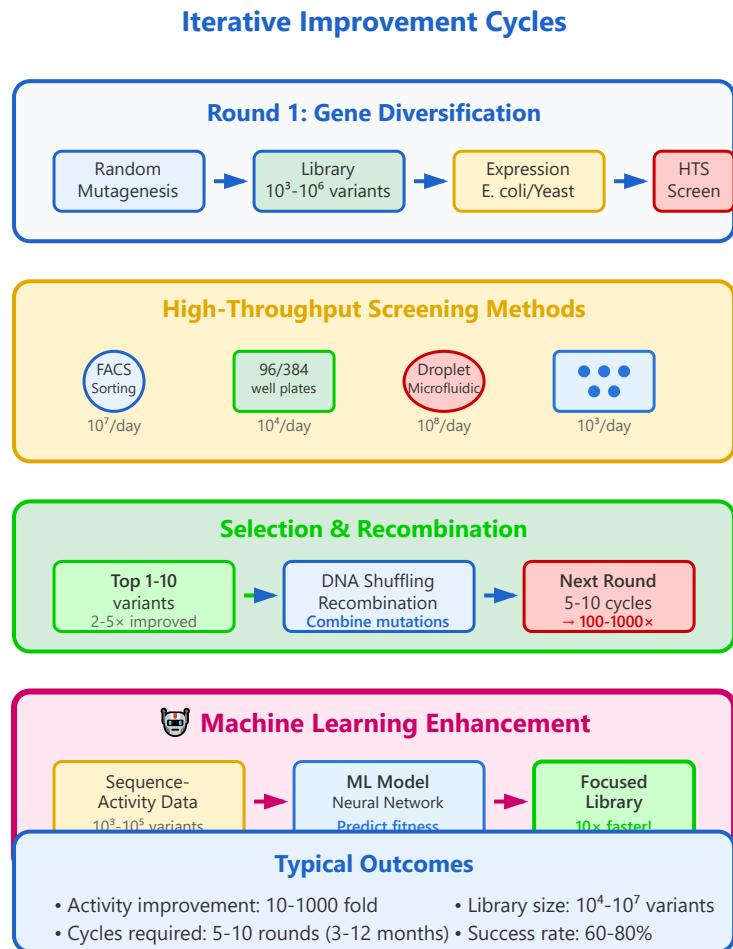
Wild-type: 30% activity in 30% methanol

Engineered (K26L/D223A/E226V): 85% activity in 70% methanol

Application: Used in biodiesel production and synthesis of pharmaceutical esters in high solvent concentrations, dramatically improving product yields.

Industrial Benefits: Enables biphasic reactions, increases substrate/product solubility, reduces water activity for reversing hydrolysis, and facilitates downstream product recovery.

5. Directed Evolution Strategy



Objective

Use iterative rounds of random mutagenesis, recombination, and selection to evolve enzymes with desired properties without requiring detailed structural knowledge. This approach mimics natural evolution but accelerated 1000-fold.

Key Components

- **Mutagenesis methods:** Error-prone PCR (0.1-1% mutation rate), DNA shuffling, saturation mutagenesis
- **Library construction:** Generate 10³-10⁷ variants with diverse mutations
- **High-throughput screening:** FACS (10⁷/day), microfluidics (10⁸/day), or plate-based assays (10⁴/day)
- **Selection criteria:** Activity, stability, specificity, or multiple properties simultaneously
- **Recombination:** DNA shuffling to combine beneficial mutations from different variants
- **ML-guidance:** Machine learning models predict promising variants, reducing screening by 10-fold

Nobel Prize Example: Frances Arnold's P450 Evolution

Goal: Evolve P450 for propane hydroxylation (non-natural activity)

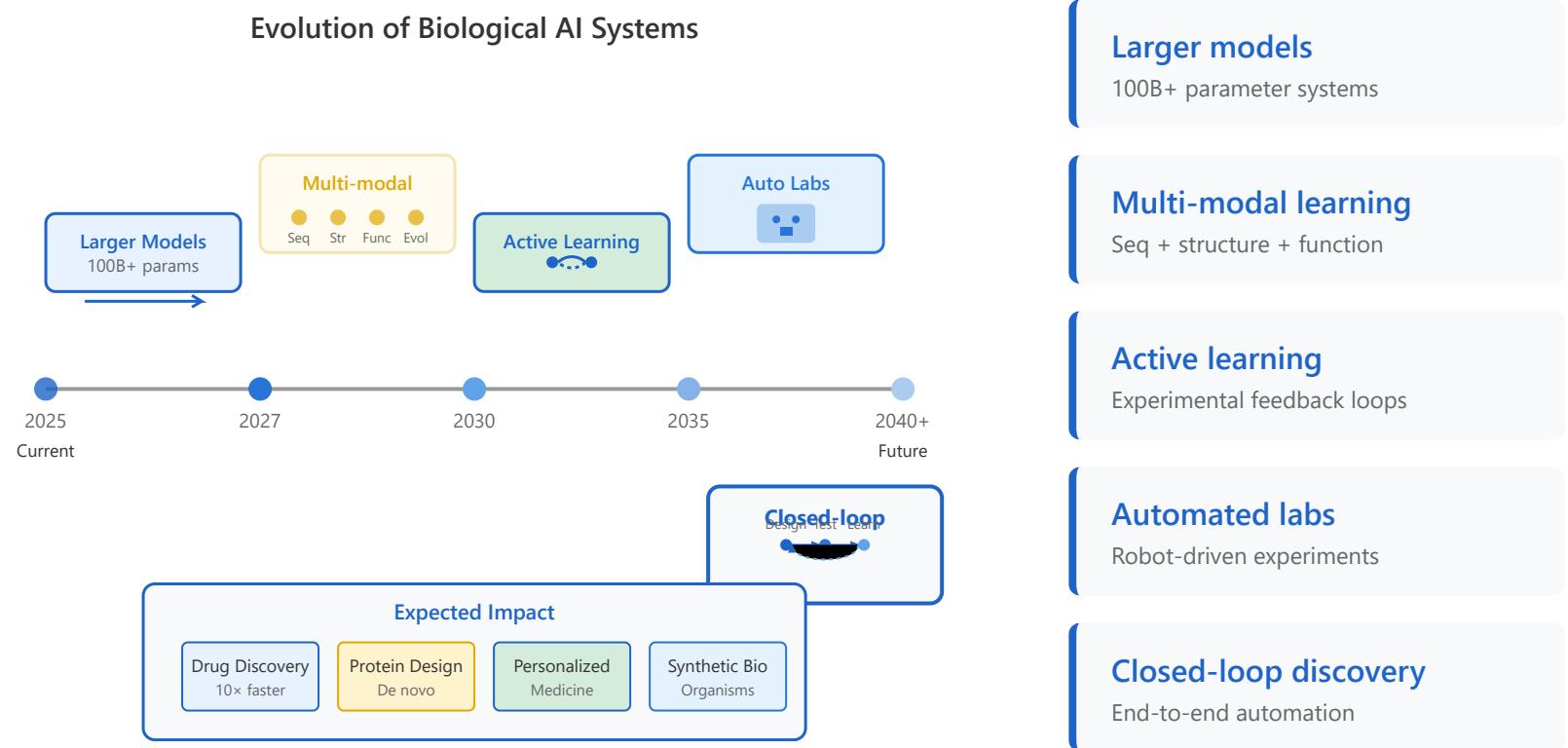
Starting: 0% activity on propane

After 5 rounds: 300,000 turnovers with 98% selectivity

Impact: Enabled sustainable biocatalytic production of valuable chemicals from cheap alkanes, awarded 2018 Nobel Prize in Chemistry.

Advantages: No structural knowledge required, can optimize multiple properties simultaneously, discovers unexpected beneficial mutations, and integrates easily with computational design and ML prediction.

Future Perspectives



1

Larger Models: Scaling to 100B+ Parameters

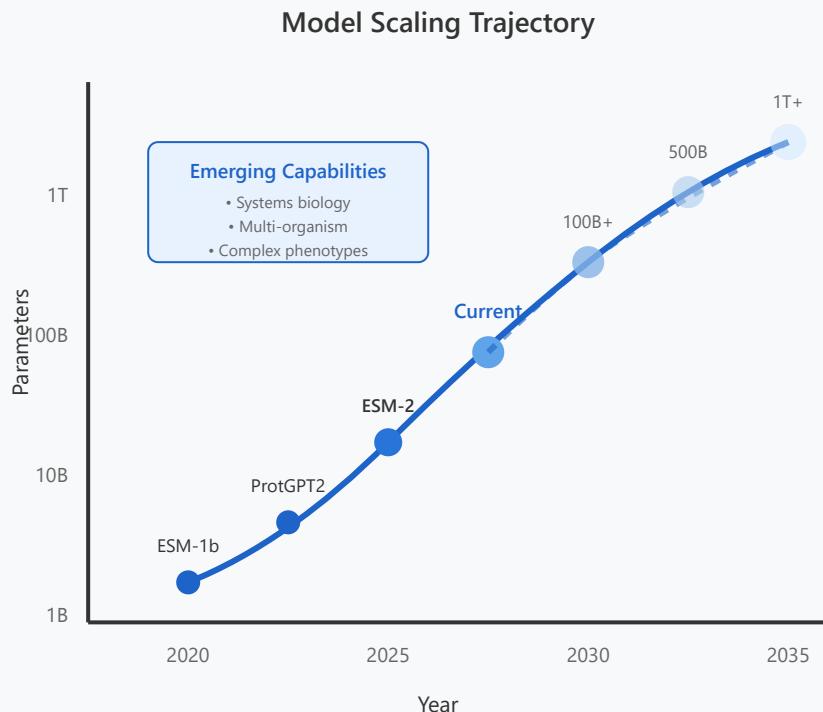
The next generation of biological AI models will scale beyond current architectures, reaching 100 billion parameters or more. This scaling enables models to capture increasingly complex biological patterns, from molecular interactions to systems-level behavior.

Key Capabilities

- ▶ Understanding protein-protein interaction networks at proteome scale
- ▶ Predicting complex phenotypes from genomic sequences
- ▶ Modeling cellular pathways and metabolic networks
- ▶ Cross-species transfer learning for rare organisms
- ▶ Integration of evolutionary information across phylogenies

Expected Impact

Larger models will enable prediction of complex biological phenomena that current models cannot address, such as multi-gene disease mechanisms, organism-level responses to perturbations, and emergent properties in synthetic biological systems.



2 Multi-modal Learning: Integration Across Data Types

Overview

Future AI systems will seamlessly integrate multiple biological data modalities including sequences, 3D structures, functional annotations,

and evolutionary information. This holistic approach mirrors how biologists naturally reason about biological systems.

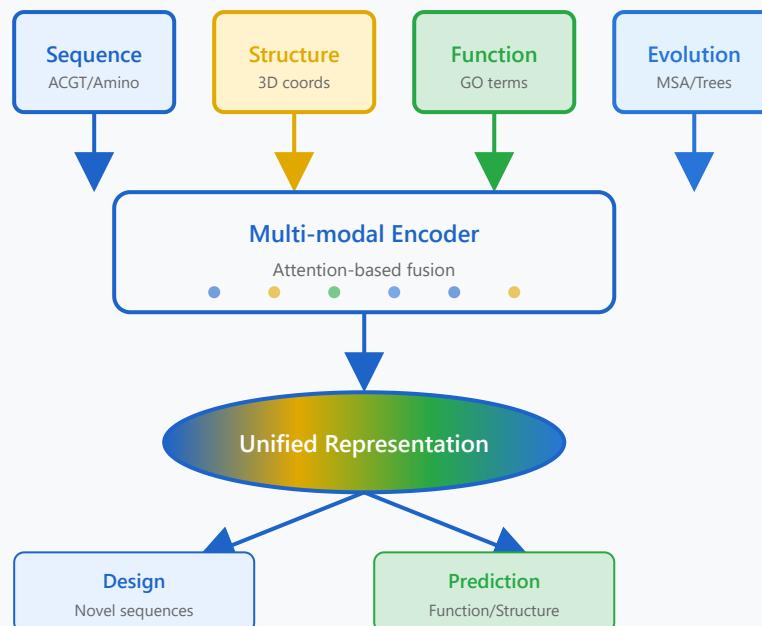
Data Modalities

- ▶ **Sequence:** Genomic and protein sequences with contextual information
- ▶ **Structure:** 3D conformations, dynamics, and structural ensembles
- ▶ **Function:** Biochemical activities, cellular localization, interactions
- ▶ **Evolution:** Phylogenetic relationships and conservation patterns
- ▶ **Expression:** Temporal and spatial gene expression profiles

Expected Impact

Multi-modal models will provide comprehensive understanding of biological entities, enabling accurate prediction of functional effects from sequence alone and facilitating the design of proteins with specified structures and functions.

Multi-modal Integration Architecture



3 Active Learning: Experimental Feedback Loops

Overview

Active learning strategies enable AI systems to identify the most informative experiments to conduct next, dramatically improving data efficiency. The model learns from experimental results and iteratively refines its predictions through strategic sampling.

Core Components

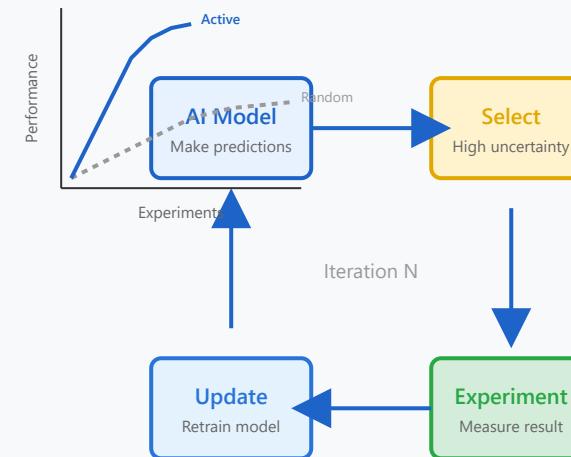
- ▶ **Uncertainty estimation:** Quantifying confidence in model predictions
- ▶ **Acquisition functions:** Selecting maximally informative experiments
- ▶ **Batch optimization:** Planning multiple parallel experiments
- ▶ **Transfer learning:** Leveraging knowledge across related tasks
- ▶ **Cost-aware selection:** Balancing information gain with experimental cost

Expected Impact

Active learning can reduce the number of required experiments by 10-100x, accelerating discovery cycles and enabling exploration of vast sequence spaces that would be prohibitively expensive with traditional approaches.

Active Learning Cycle

Efficiency Comparison



4

Automated Laboratories: Robot-Driven Experiments

Overview

Automated laboratories combine robotic systems, microfluidics, and AI control to execute thousands of experiments in parallel with minimal human intervention. These systems can operate continuously, generating high-quality data at unprecedented scales.

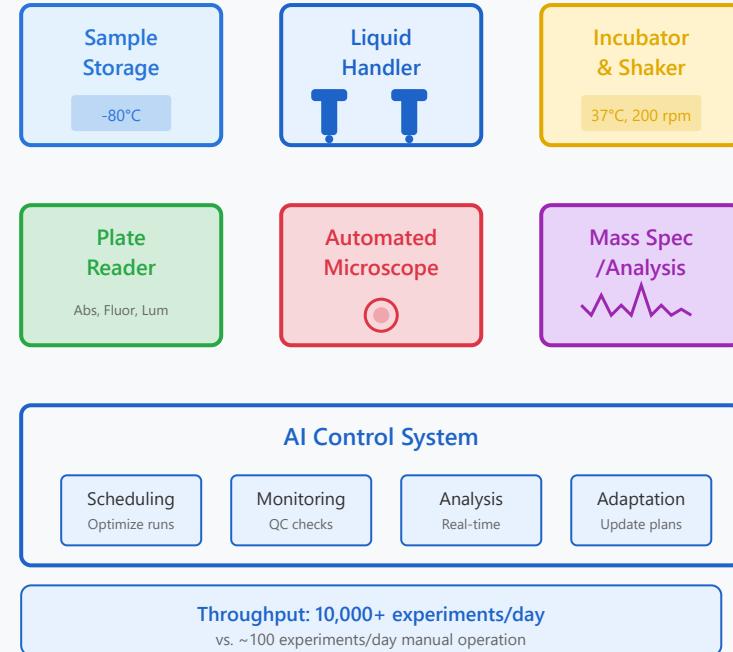
Key Technologies

- ▶ **Liquid handling robots:** Precise pipetting and sample preparation
- ▶ **High-throughput screening:** Parallel assays in microplate formats
- ▶ **Microfluidics:** Miniaturized reactions and cell manipulation
- ▶ **Real-time monitoring:** Automated imaging and spectroscopy
- ▶ **Adaptive protocols:** AI-driven experimental adjustments

Expected Impact

Automated labs will enable 24/7 experimentation with throughput 1000× higher than manual approaches, while maintaining reproducibility and reducing costs. This will democratize access to advanced experimental capabilities.

Automated Laboratory Architecture



5 Closed-loop Discovery: End-to-End Automation

Overview

Closed-loop discovery represents the ultimate integration of AI and automated experimentation, where the entire scientific discovery process operates autonomously. The system designs experiments, executes them, analyzes results, updates models, and iterates without human intervention.

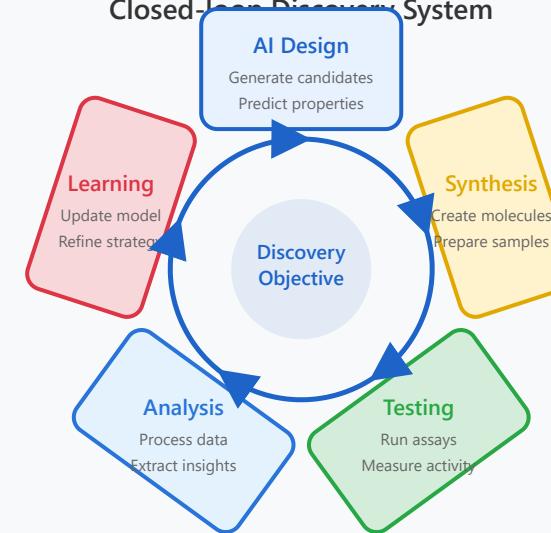
System Components

- ▶ **Autonomous design:** AI generates novel hypotheses and experiments
- ▶ **Robotic execution:** Automated labs perform experiments
- ▶ **Real-time analysis:** Instant processing of experimental data
- ▶ **Model updating:** Continuous learning from new results
- ▶ **Goal optimization:** Multi-objective optimization toward targets

Expected Impact

Closed-loop systems will compress discovery timelines from years to weeks, enable exploration of combinatorially vast design spaces, and accelerate the pace of innovation in drug discovery, materials science, and synthetic biology by 100× or more.

Closed-Loop Discovery System

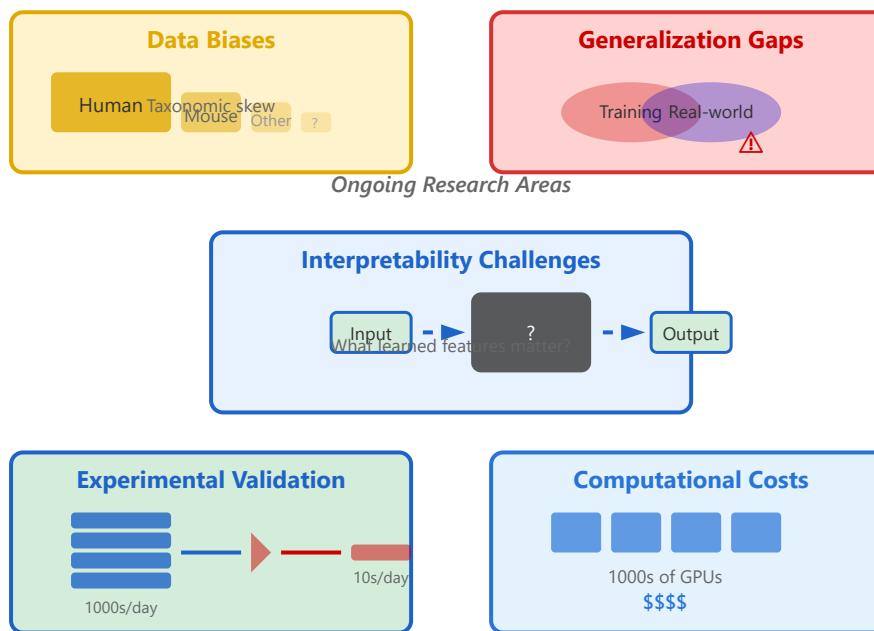


Key Performance Indicators

Cycle Time: 1-7 days Throughput: 100s/cycle Improvement: 10-100×

Limitations

Current Challenges in Biological AI



Data biases

Taxonomic & functional skew

Generalization gaps

Out-of-distribution failures

Interpretability challenges

Black box models

Experimental validation

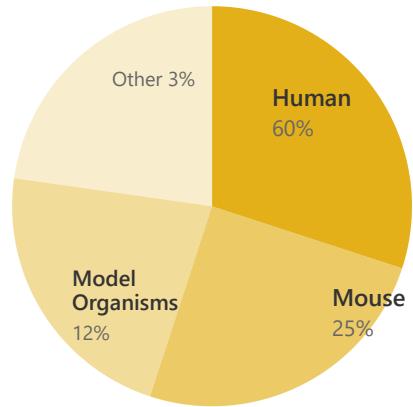
Lab throughput bottleneck

Computational costs

Training & inference expense

1. Data Biases

Taxonomic Distribution in Biological Databases



Underrepresented:

- Non-model organisms
- Rare diseases
- Non-coding regions

The Problem

Biological AI models are trained on highly imbalanced datasets that overrepresent certain species, tissues, and biological processes while underrepresenting others. This creates systematic biases in model predictions.

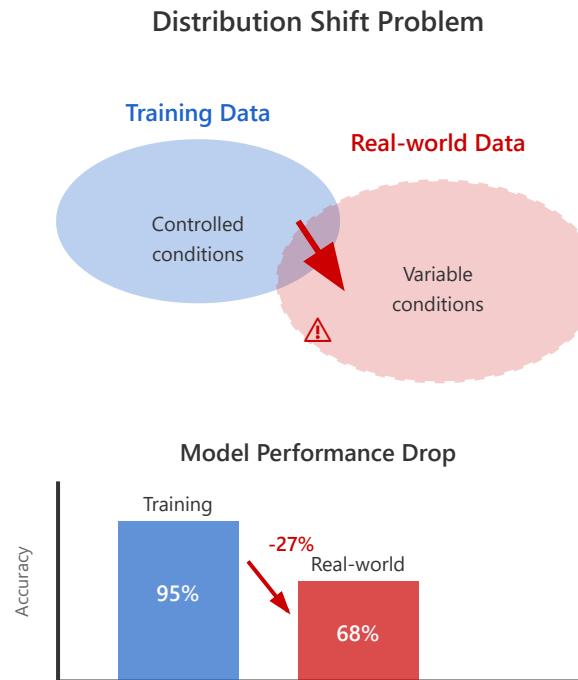
Key Issues

- **Taxonomic skew:** Human and mouse data dominate, while other organisms are severely underrepresented
- **Functional bias:** Well-studied pathways and genes receive more attention than novel or rare functions
- **Tissue bias:** Easily accessible tissues are overrepresented in training data
- **Disease bias:** Common diseases studied more than rare diseases

Example Impact

A protein function prediction model trained primarily on human proteins may fail to accurately predict functions in non-model organisms like extremophiles or plant species, limiting its utility for agricultural or environmental applications.

2. Generalization Gaps



The Problem

Models trained on carefully curated datasets often fail when confronted with real-world data that differs from training conditions. This out-of-distribution (OOD) problem leads to unreliable predictions in practical applications.

Key Issues

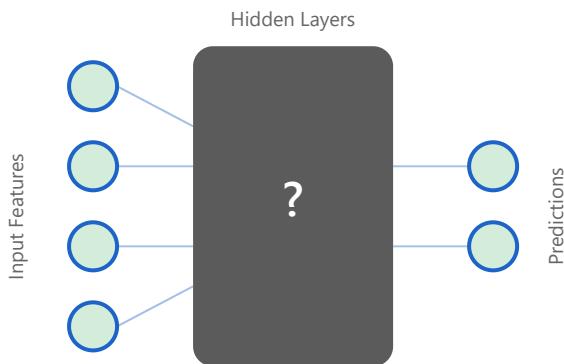
- **Batch effects:** Different experimental protocols and platforms create systematic variations
- **Population diversity:** Models trained on one population may not generalize to others
- **Environmental factors:** Lab conditions differ from natural environments
- **Temporal shifts:** Biological systems evolve and change over time

Example Impact

A drug response prediction model trained on cell lines from European populations may show significantly reduced accuracy when applied to patients of African or Asian ancestry due to genetic and environmental differences.

3. Interpretability Challenges

The Black Box Problem



Key Questions:

- ❓ Which features does the model rely on?
- ❓ Why did it make this specific prediction?
- ❓ Are the learned patterns biologically meaningful?
- ❓ Can we trust it for critical decisions?

The Problem

Deep learning models, while powerful, function as "black boxes" where the relationship between inputs and outputs is opaque. This lack of transparency poses serious challenges for biological applications requiring mechanistic understanding.

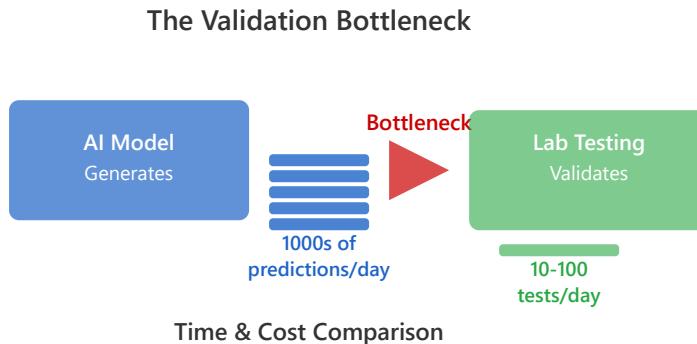
Key Issues

- **Feature attribution:** Difficulty identifying which input features drive predictions
- **Biological plausibility:** Model patterns may not reflect known biological mechanisms
- **Trust and adoption:** Clinicians and researchers hesitant to use unexplainable models
- **Debugging:** Hard to identify and fix systematic errors in model reasoning

Example Impact

A deep learning model predicts a protein will bind to a specific drug target with high confidence, but cannot explain which structural features matter most. Researchers cannot determine if the prediction is based on relevant biochemistry or spurious correlations.

4. Experimental Validation



Consequences:

- Only fraction of predictions can be tested
- Delayed feedback for model improvement
- Risk of deploying unvalidated predictions
- Selection bias in which predictions to test

The Problem

AI models can generate predictions orders of magnitude faster than they can be experimentally validated. This creates a severe bottleneck that limits the practical utility of computational predictions and slows the research cycle.

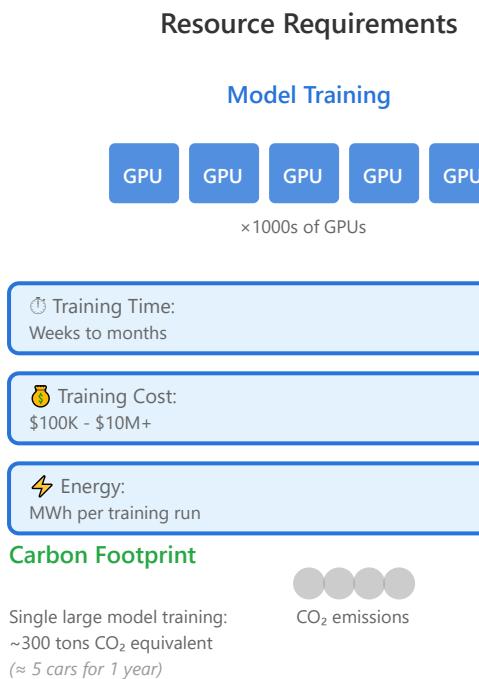
Key Issues

- **Throughput mismatch:** Models predict thousands of hypotheses daily while labs test tens
- **Resource constraints:** Limited lab space, equipment, and trained personnel
- **Time delays:** Experiments take days to months while predictions are instant
- **Cost barriers:** Each validation experiment costs significantly more than predictions

Example Impact

A protein engineering model suggests 10,000 potentially beneficial mutations. The lab can only test 50 mutations per month due to time and cost constraints. It would take over 16 years to validate all predictions, by which time the model may be outdated.

5. Computational Costs



The Problem

Training and deploying state-of-the-art biological AI models requires massive computational resources, creating barriers to entry and raising concerns about sustainability and accessibility of these technologies.

Key Issues

- **Hardware requirements:** Need for thousands of specialized GPUs or TPUs
- **Financial barriers:** Training costs can exceed millions of dollars per model
- **Energy consumption:** Significant electricity use with environmental impact
- **Accessibility:** Only well-funded institutions can develop large models
- **Inference costs:** Running predictions at scale remains expensive

Example Impact

Training a large protein language model like ESM-2 (650M parameters) requires thousands of GPU-hours and costs approximately \$200,000. This puts such models out of reach for most academic labs and smaller biotech companies, concentrating power in well-funded institutions.

Addressing the Limitations: Current Approaches

Mitigation Strategies

Data Biases

Solutions:

- Active data collection from underrepresented species
- Transfer learning techniques
- Data augmentation strategies
- Meta-learning approaches

Generalization Gaps

Solutions:

- Domain adaptation methods
- Robust training with diverse data
- Uncertainty quantification
- Ensemble approaches
- Causal inference methods

Interpretability

Solutions:

- Attention visualization
- Integrated gradients
- Concept activation vectors
- Mechanistic interpretability
- Post-hoc explanation tools

Experimental Validation

Solutions:

- Active learning for efficient experiment selection
- High-throughput screening
- Automated labs (lab-on-chip)
- Simulation-based validation

Computational Costs

Solutions:

- Model compression techniques
- Knowledge distillation
- Efficient architectures
- Cloud-based democratization
- Green AI initiatives

Key Takeaways

While biological AI faces significant limitations, active research is addressing these challenges through methodological innovations, improved experimental design, and technological advances. Success requires interdisciplinary collaboration between computational scientists, experimentalists, and domain experts.

The field is moving toward more robust, interpretable, and accessible AI systems that can reliably contribute to biological discovery and clinical applications.

Hands-on: AlphaFold Usage



- Structure prediction
- Confidence interpretation
- Complex modeling
- Mutation analysis
- Drug discovery applications

Hands-on: Bio Transformers



- Model loading
- Sequence encoding
- Fine-tuning
- Embedding extraction
- Downstream tasks

Thank You!

Scientific breakthroughs

Drug discoveries

Future potential

Career opportunities

Questions? Contact: homin.park@ghent.ac.kr

Thank You!

Scientific breakthroughs

Drug discoveries

Future potential

Career opportunities

Questions? Contact: