

Multimodal deep learning for protein engineering

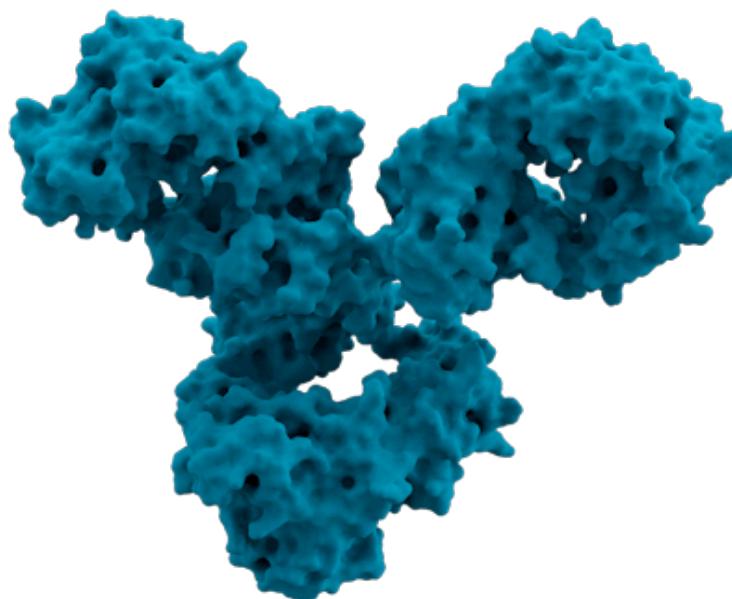
Kevin Kaichuang Yang
Microsoft Research New England
 @KevinKaichuang

Proteins are biology's actuators

- Human cells contain 1-3 billion proteins each
- Structure, metabolism, and signaling

Proteins are biology's actuators

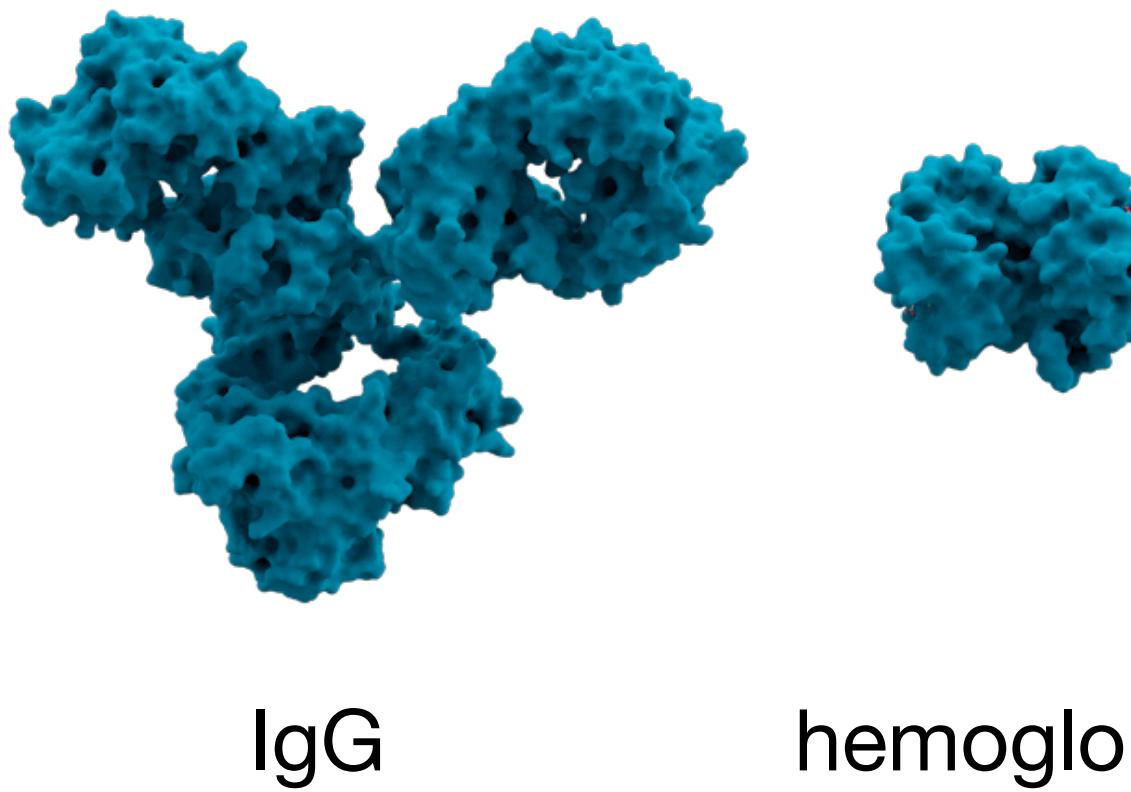
- Human cells contain 1-3 billion proteins each
- Structure, metabolism, and signaling



IgG

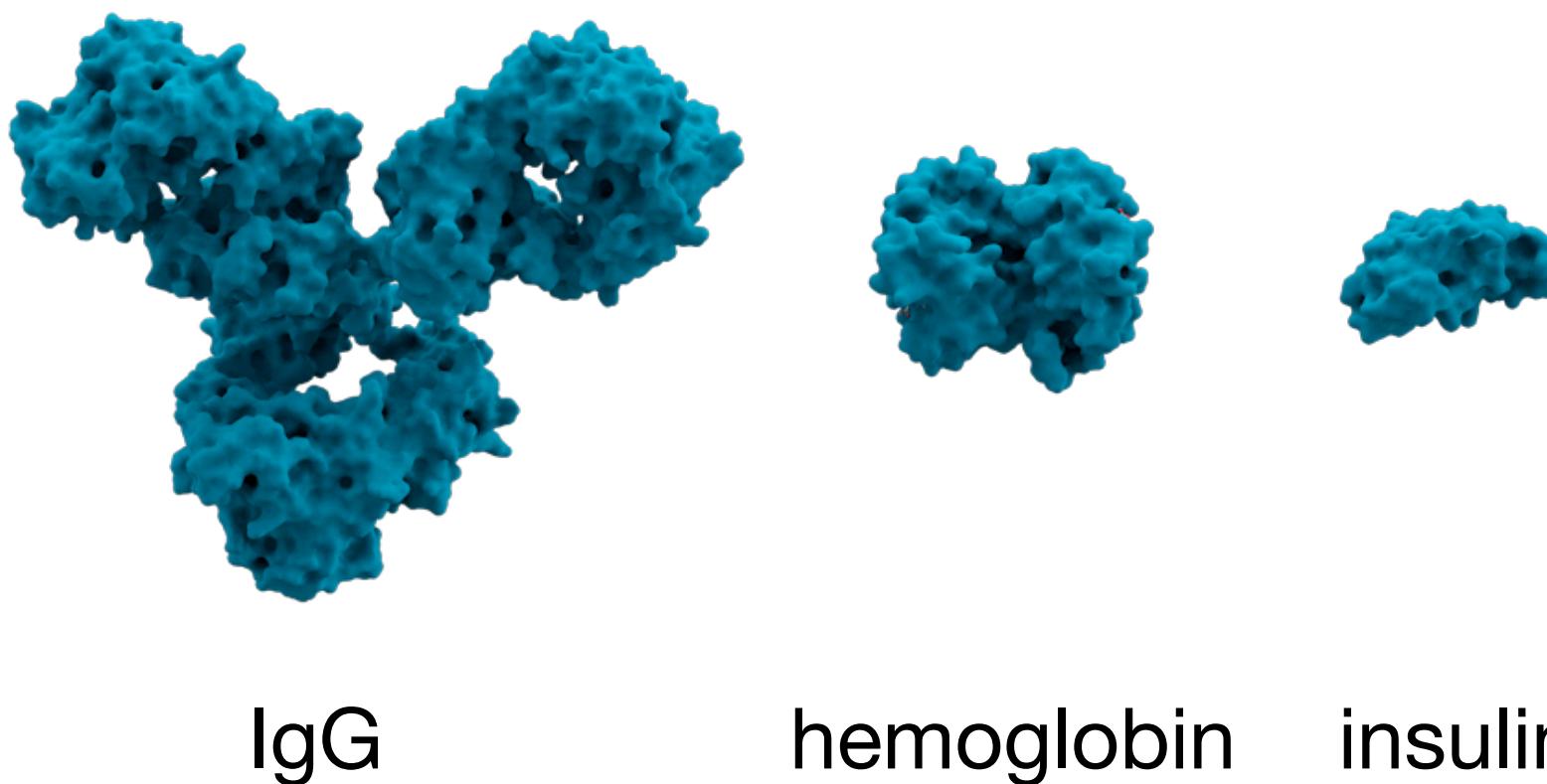
Proteins are biology's actuators

- Human cells contain 1-3 billion proteins each
- Structure, metabolism, and signaling



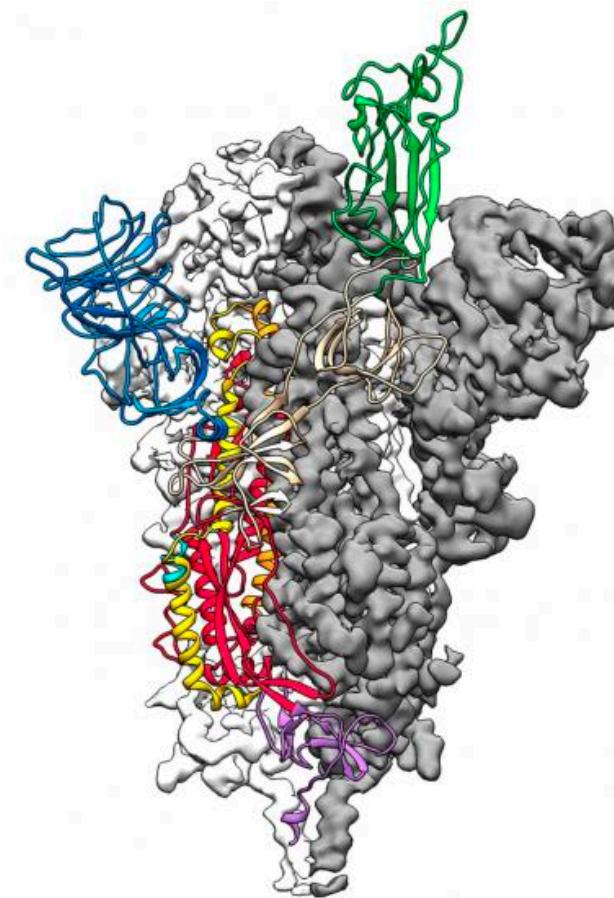
Proteins are biology's actuators

- Human cells contain 1-3 billion proteins each
- Structure, metabolism, and signaling

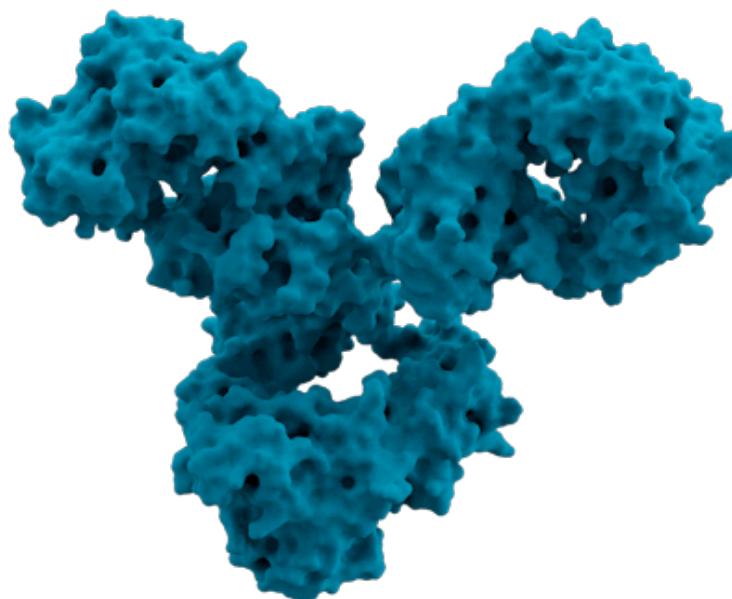


Proteins are biology's actuators

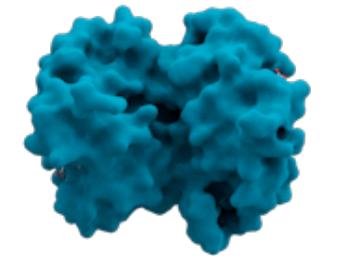
- Human cells contain 1-3 billion proteins each
- Structure, metabolism, and signaling



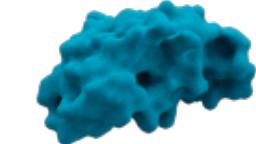
coronavirus spike protein



IgG



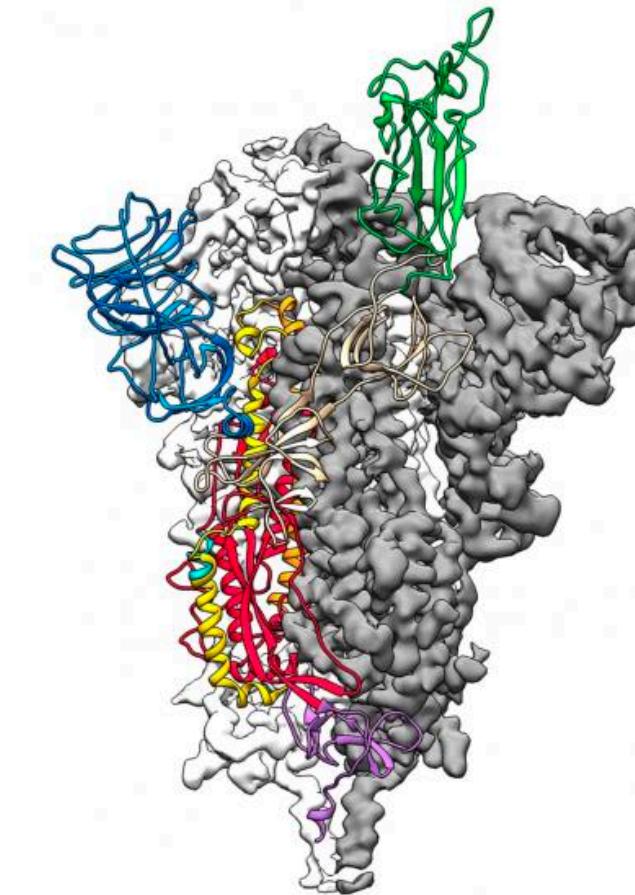
hemoglobin



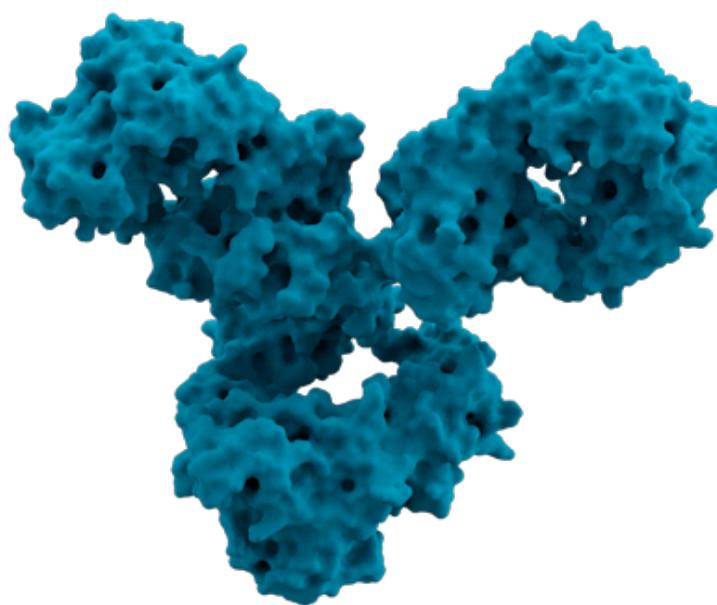
insulin

Proteins are biology's actuators

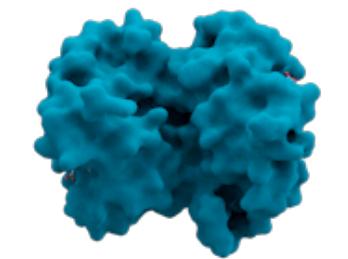
- Human cells contain 1-3 billion proteins each
- Structure, metabolism, and signaling



coronavirus spike protein



IgG



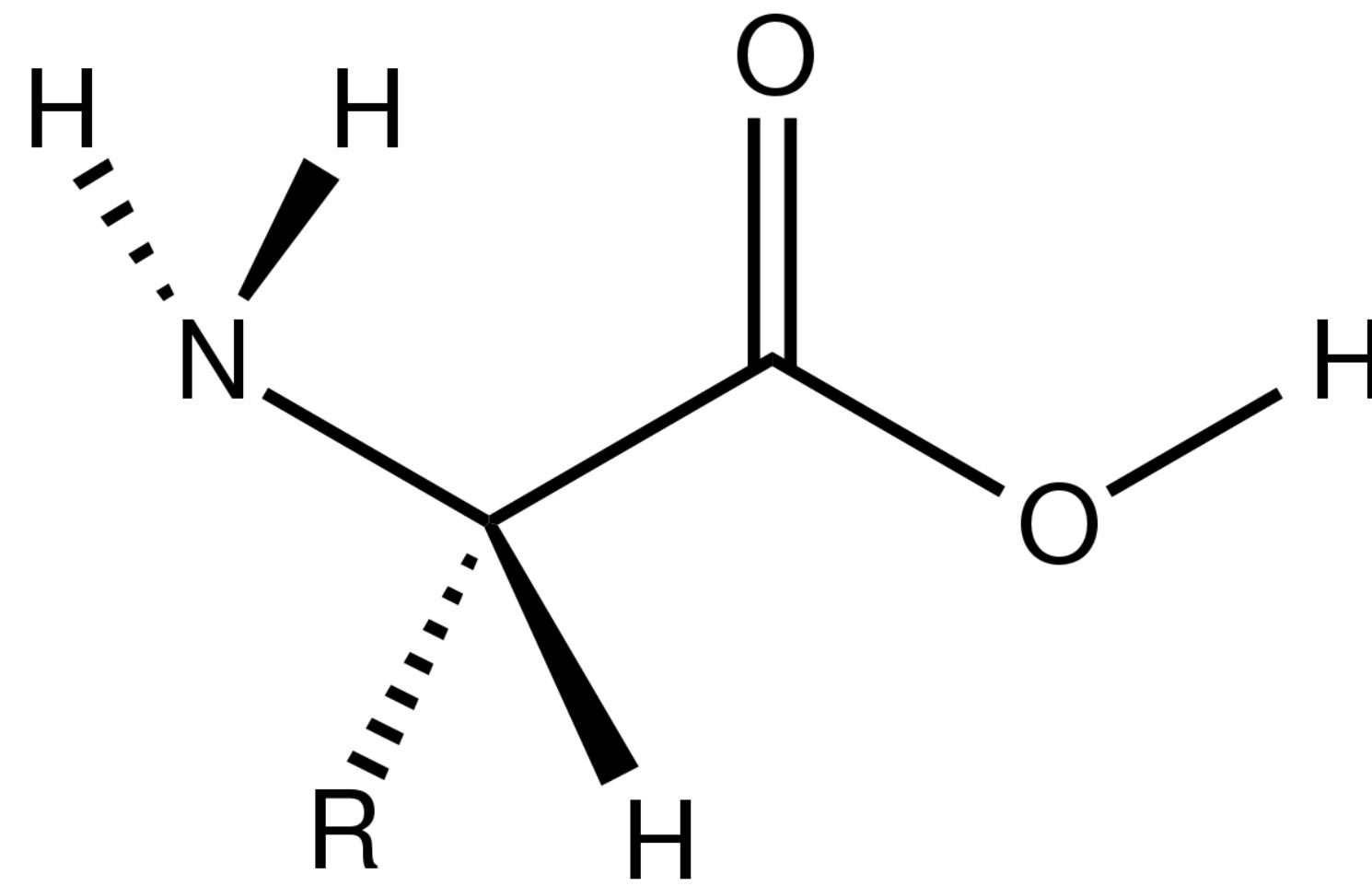
hemoglobin insulin

2

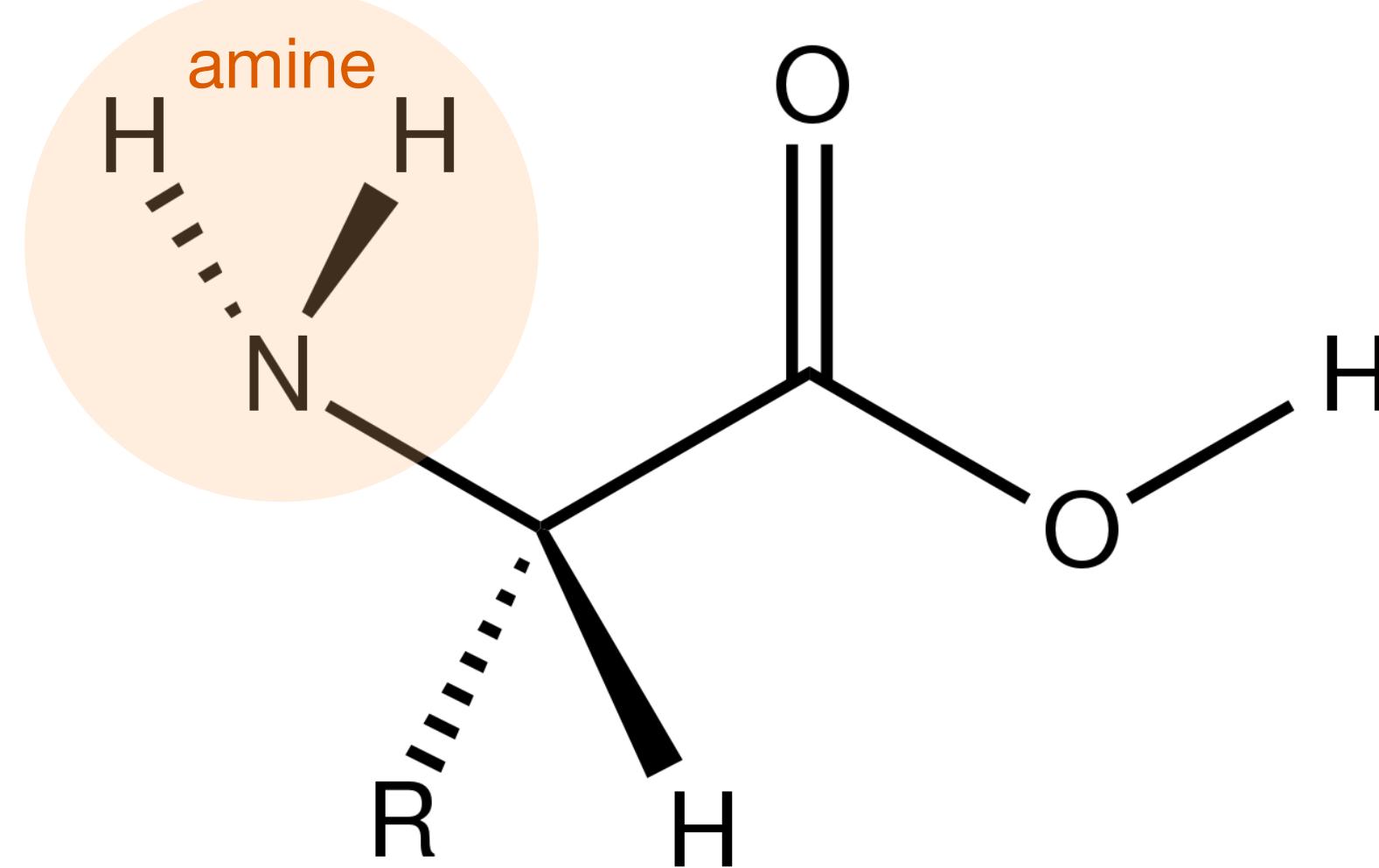


luciferase

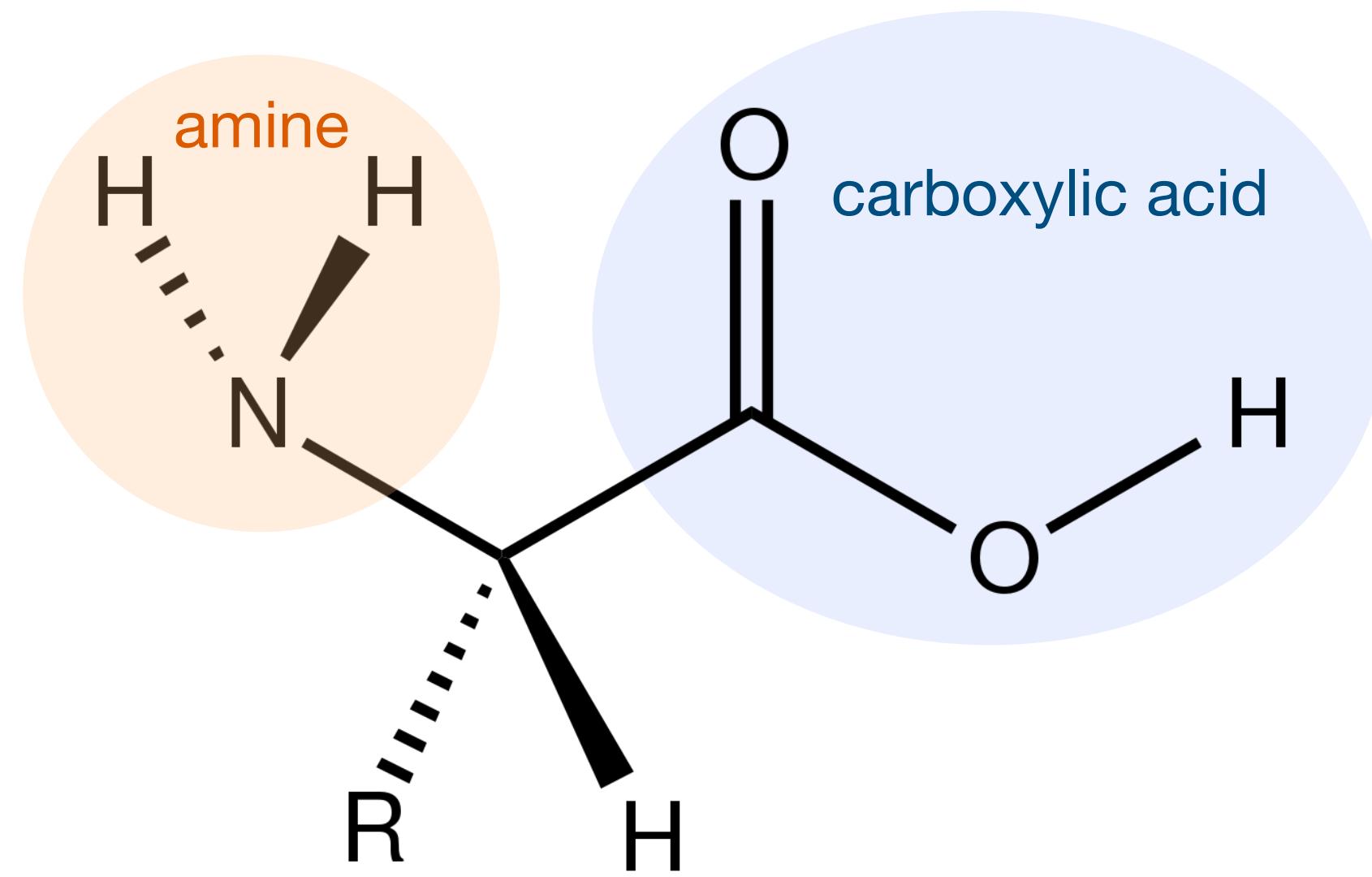
Diversity arises from 20 building blocks



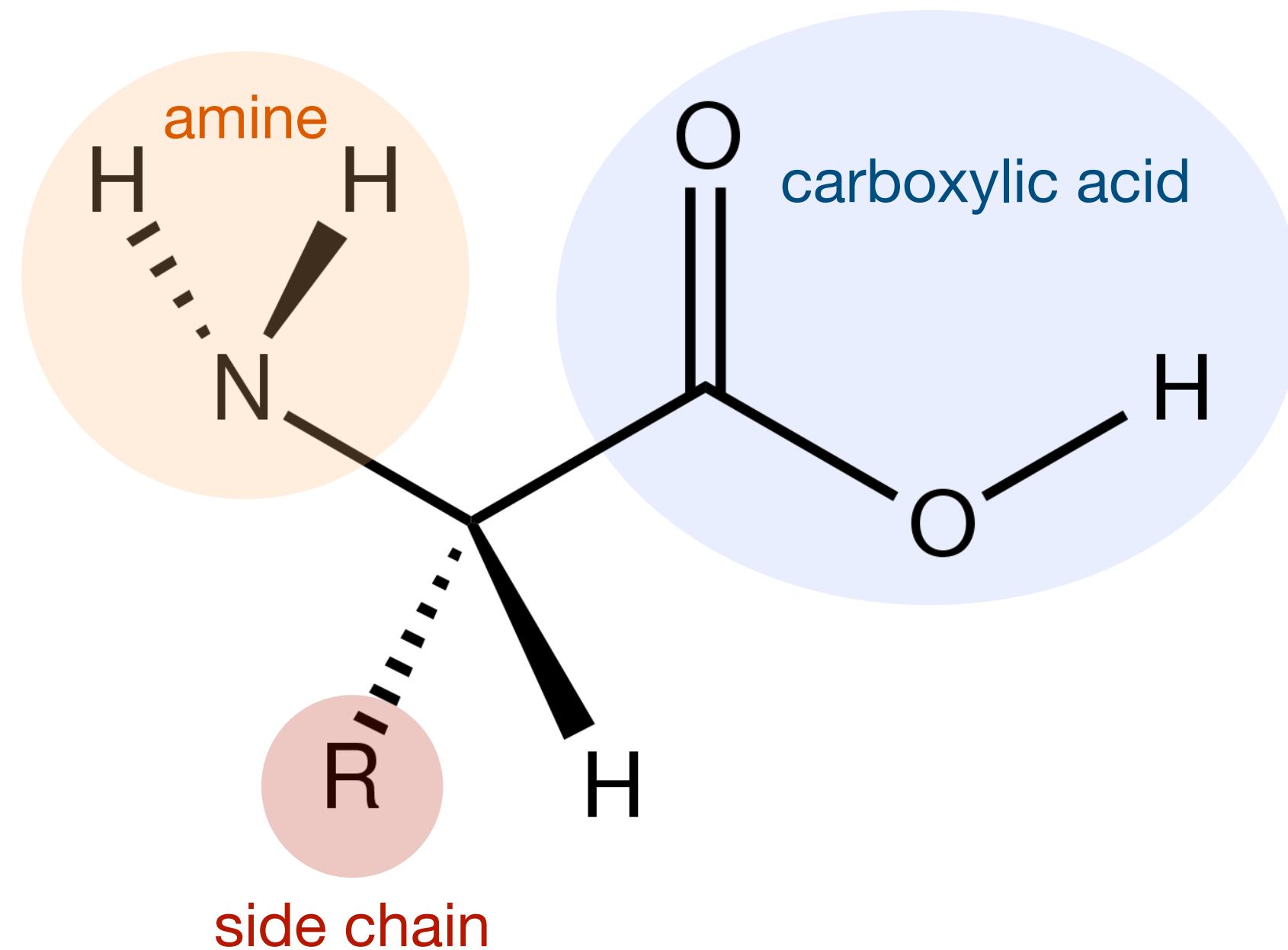
Diversity arises from 20 building blocks



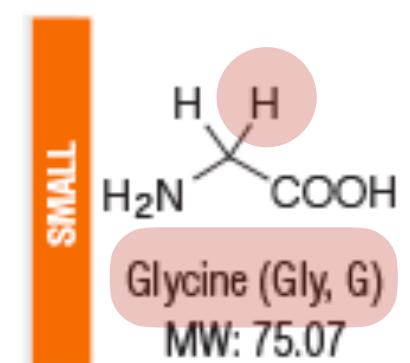
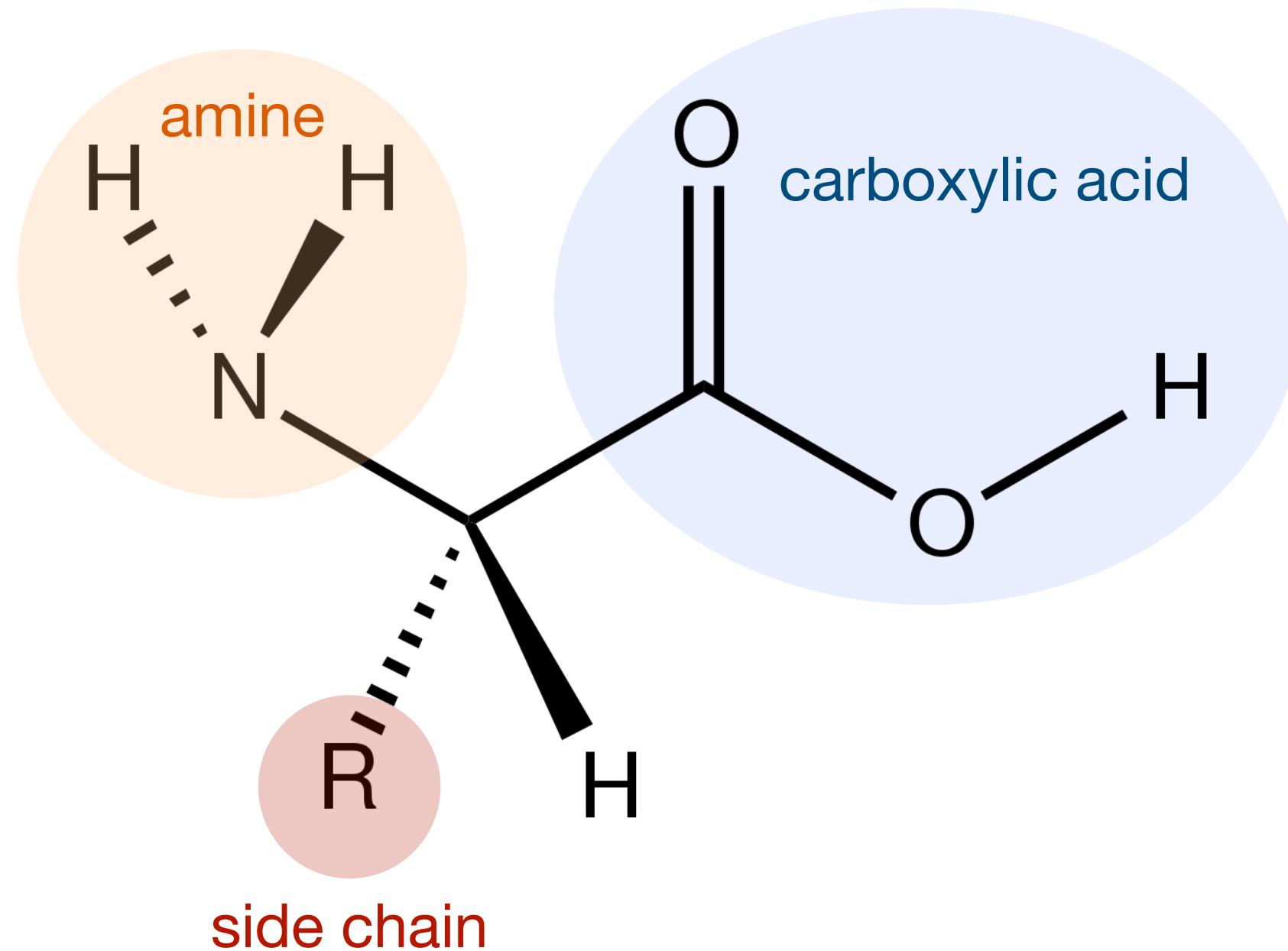
Diversity arises from 20 building blocks



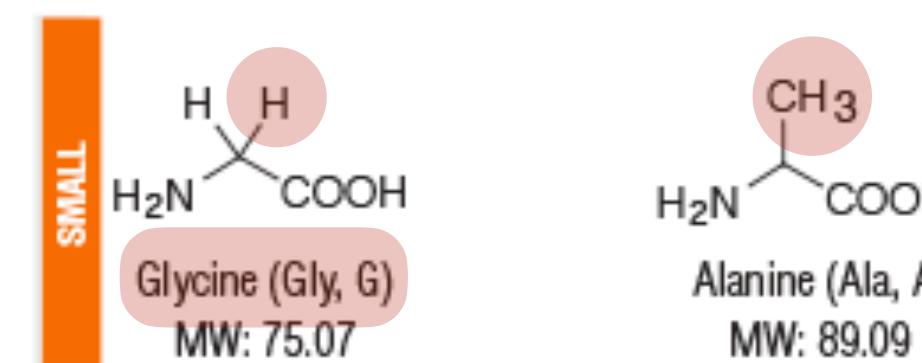
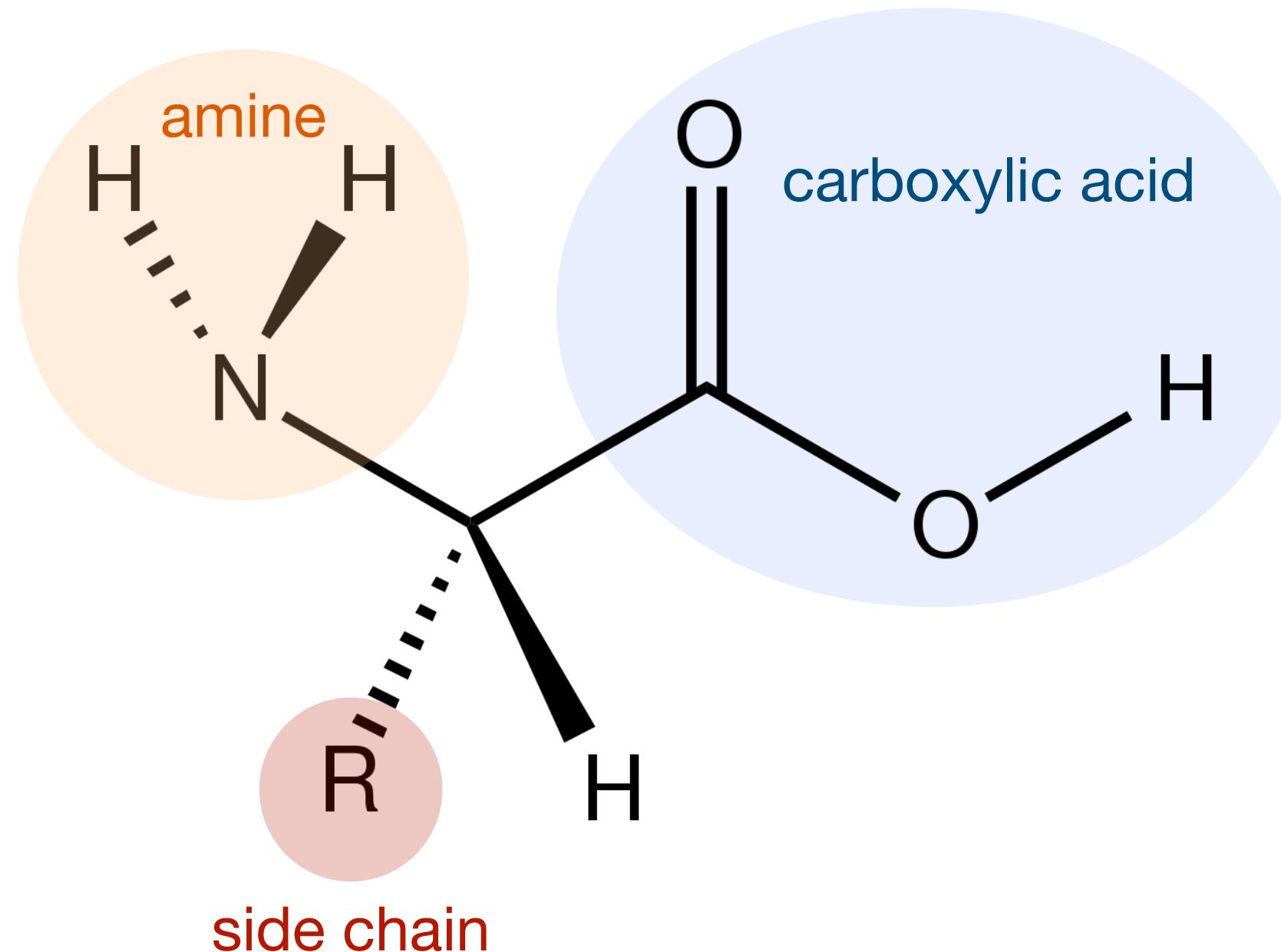
Diversity arises from 20 building blocks



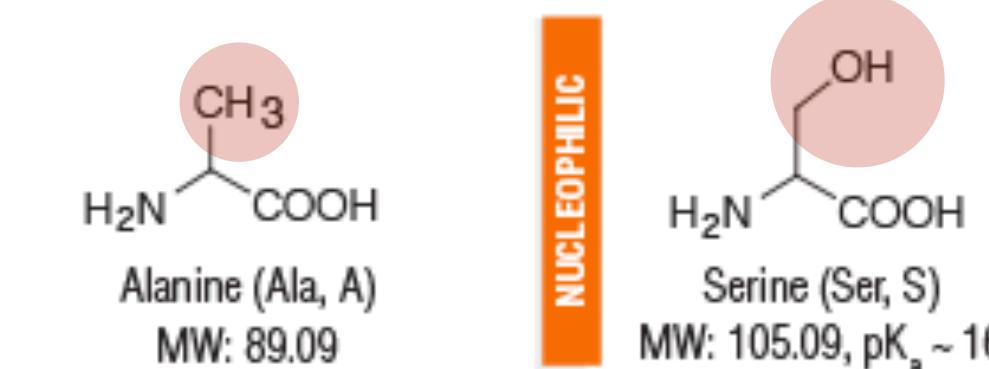
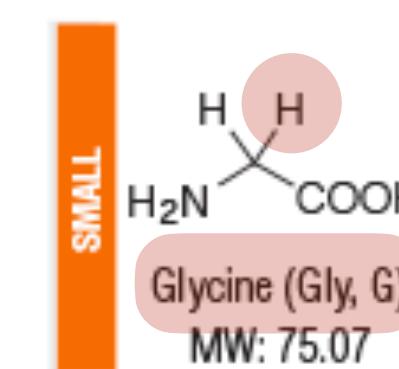
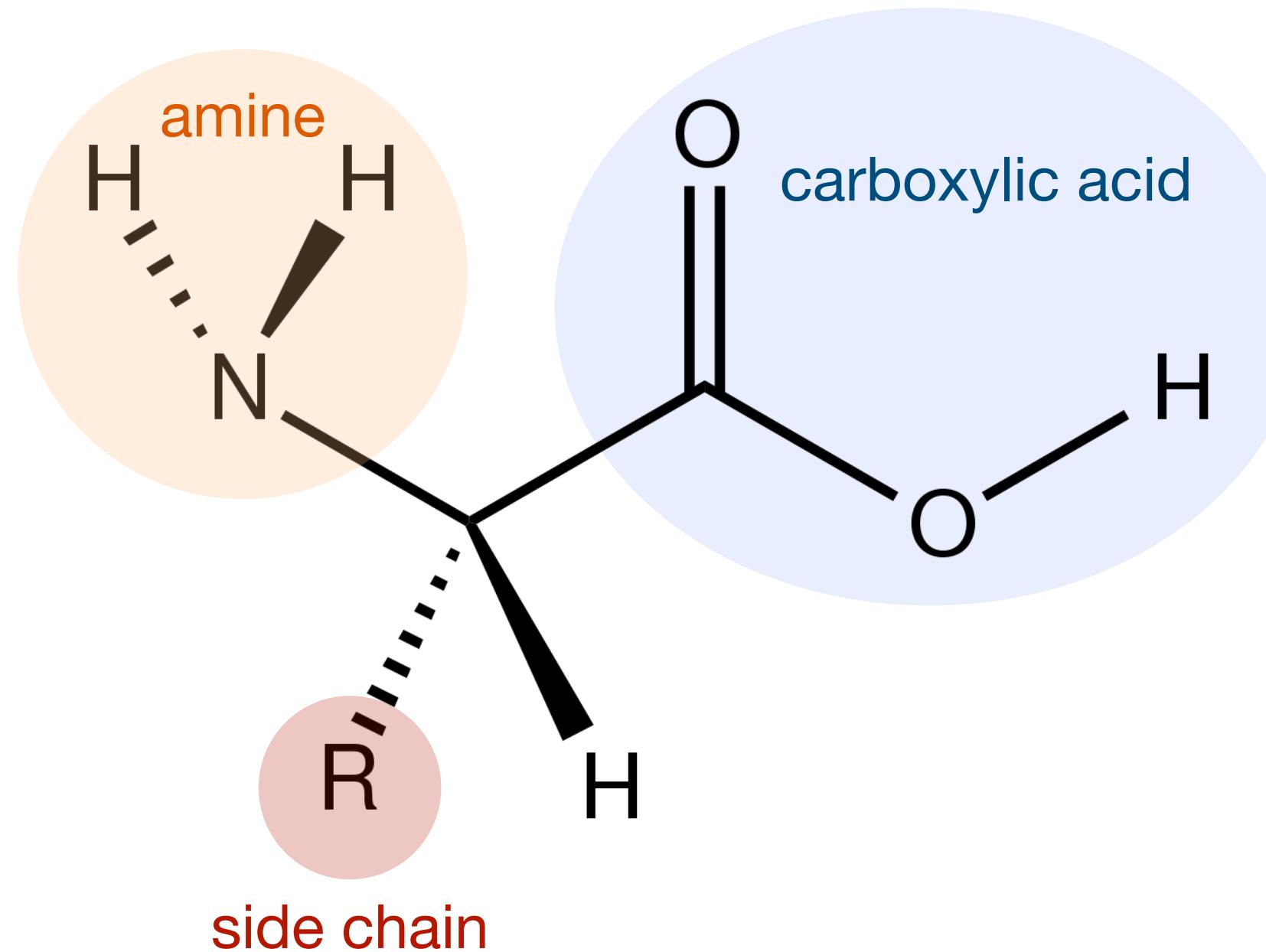
Diversity arises from 20 building blocks



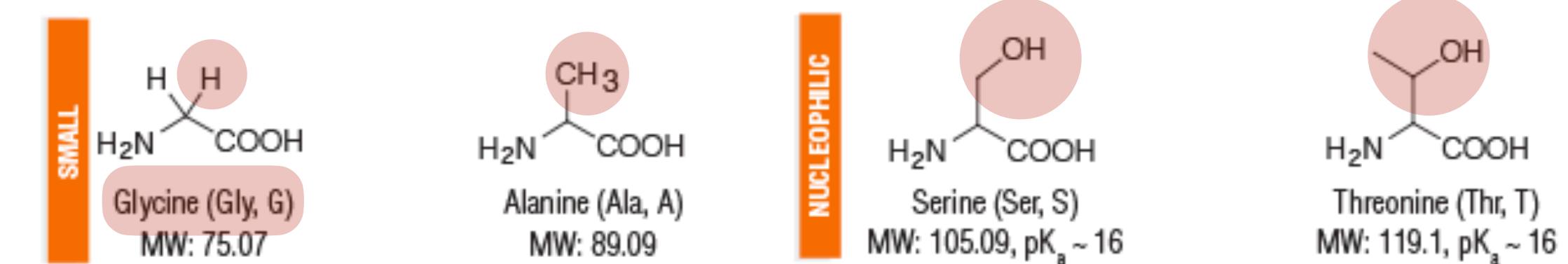
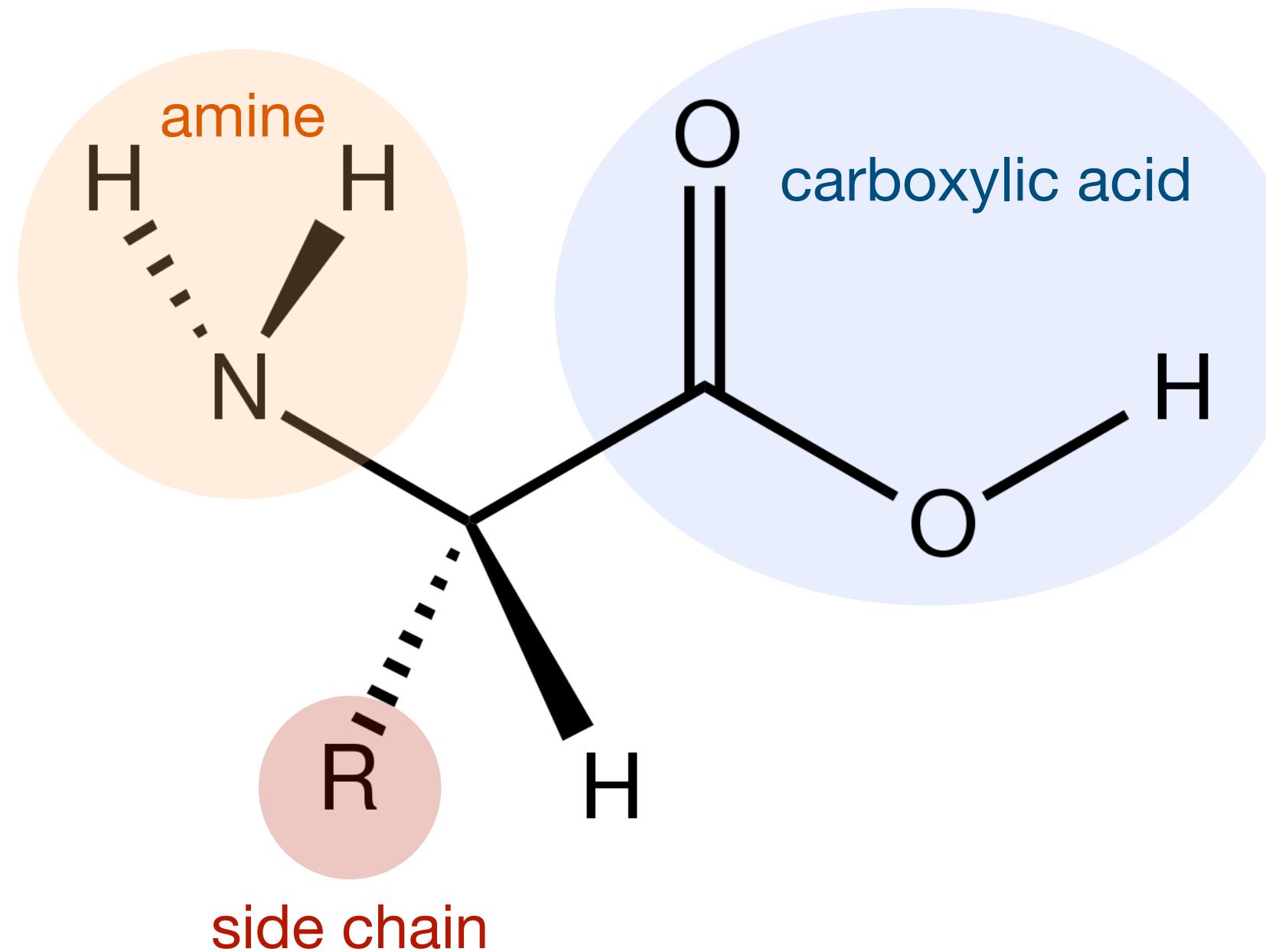
Diversity arises from 20 building blocks



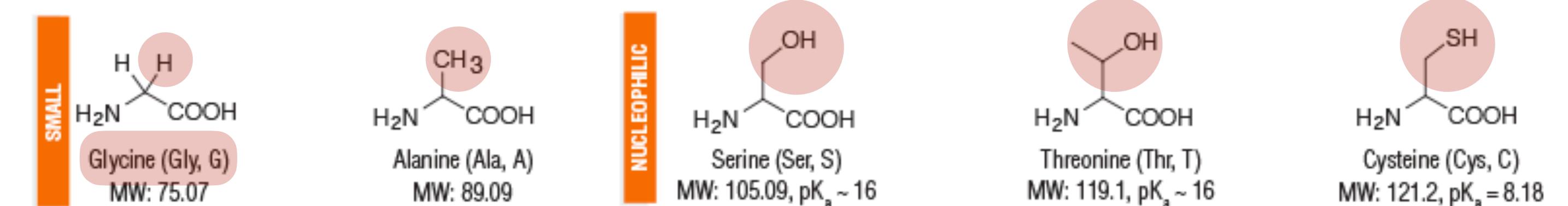
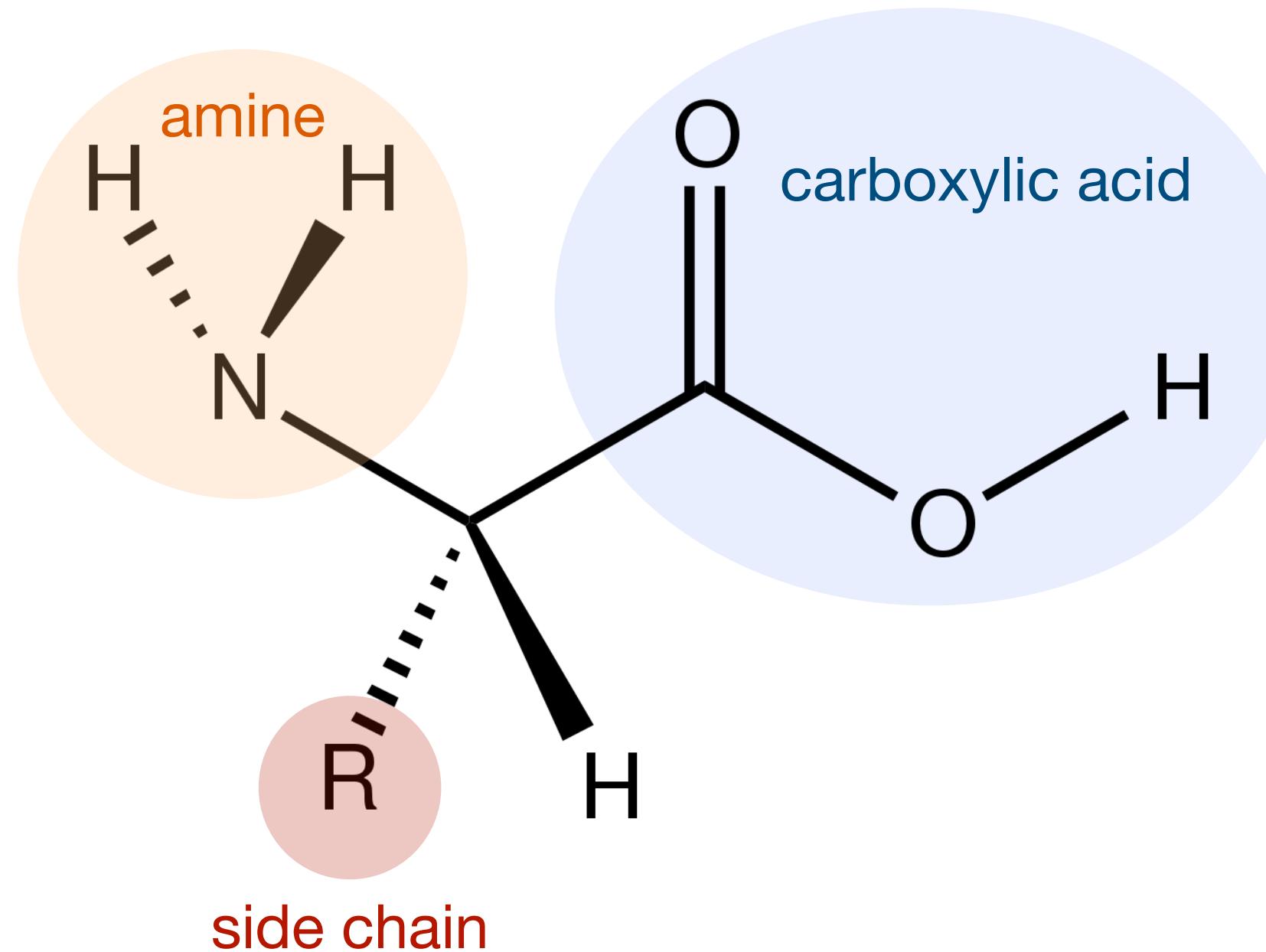
Diversity arises from 20 building blocks



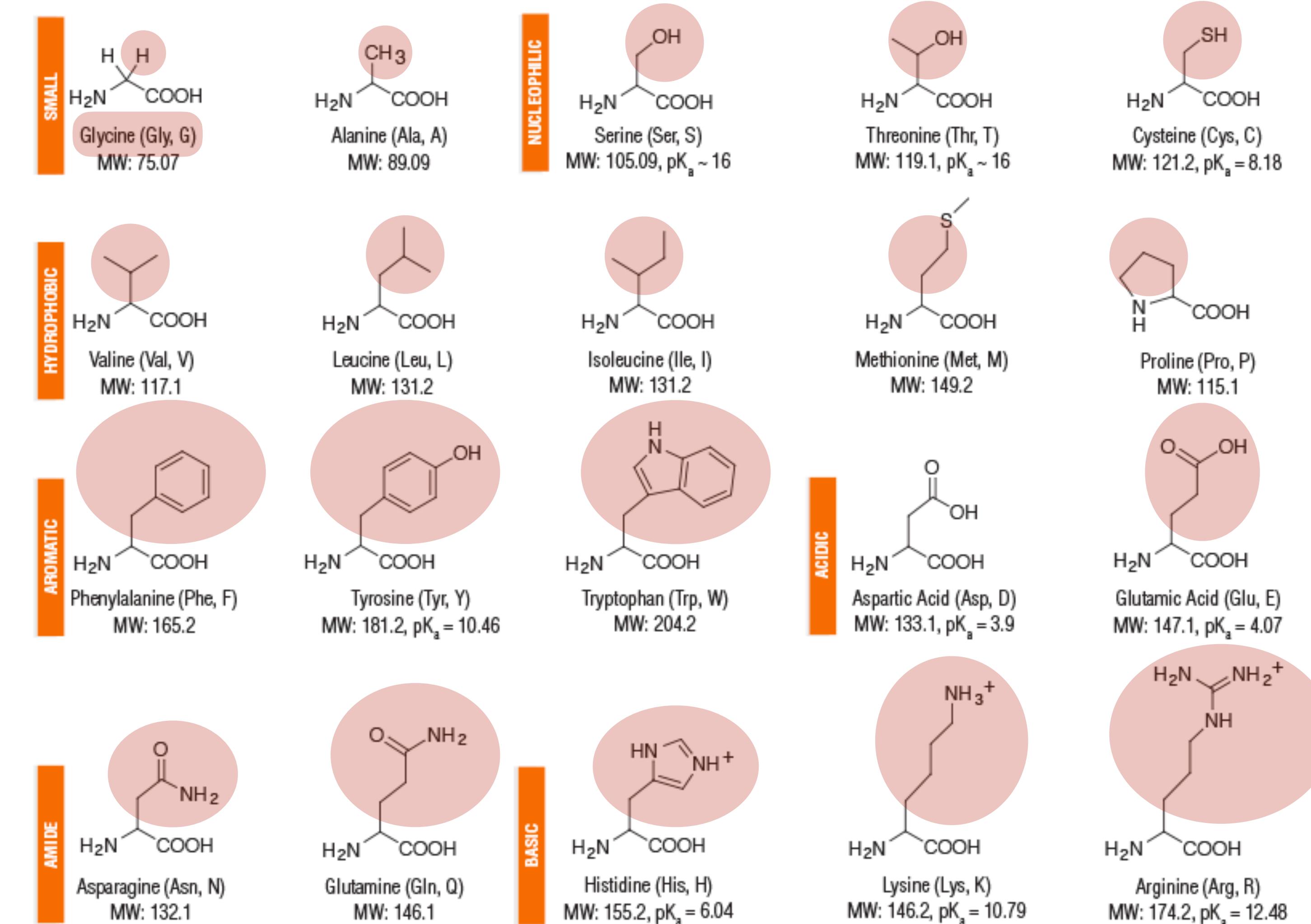
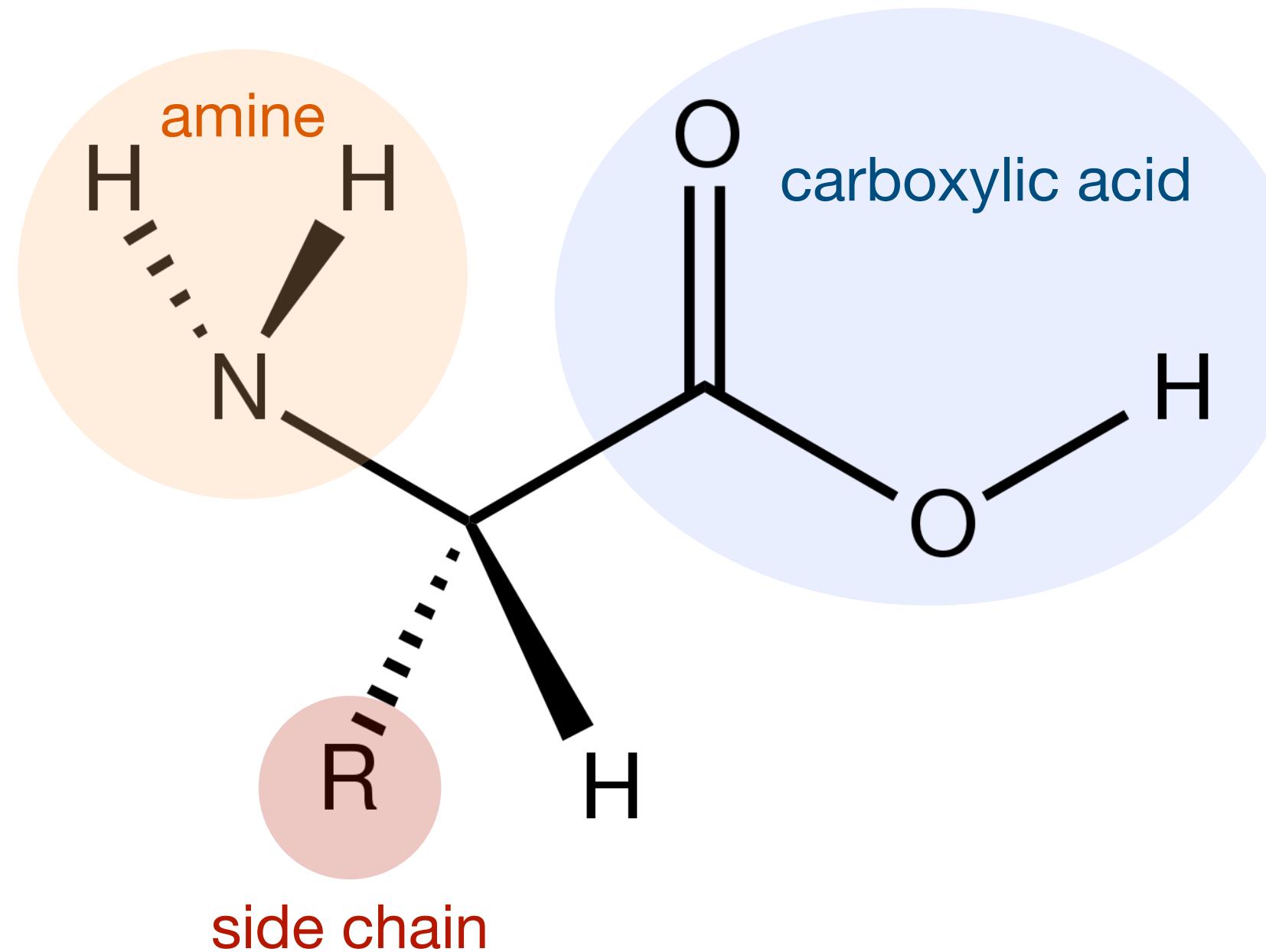
Diversity arises from 20 building blocks



Diversity arises from 20 building blocks

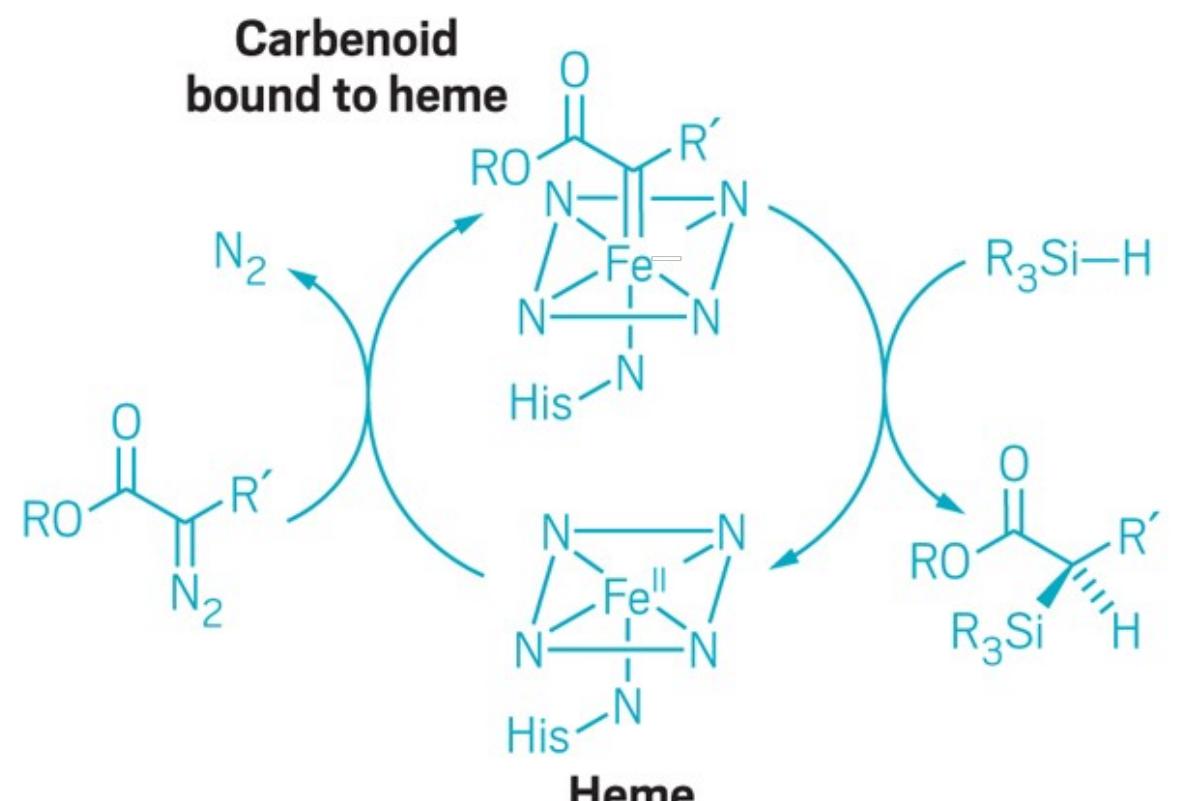


Diversity arises from 20 building blocks



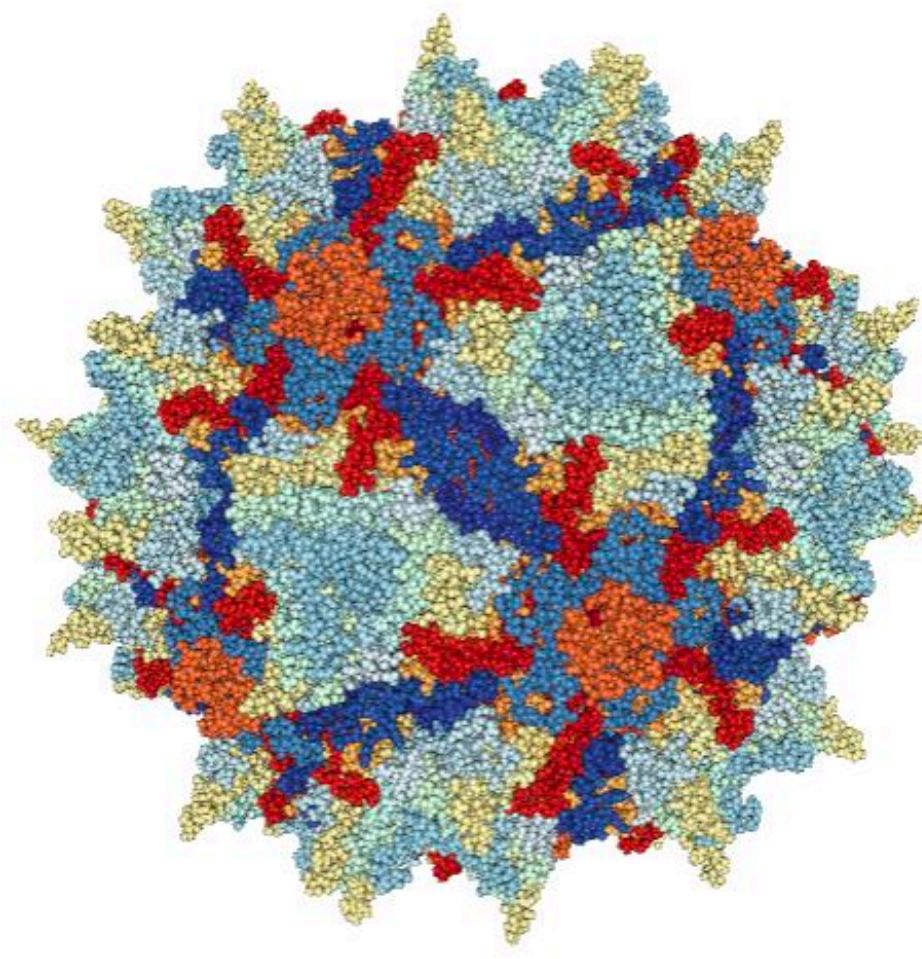
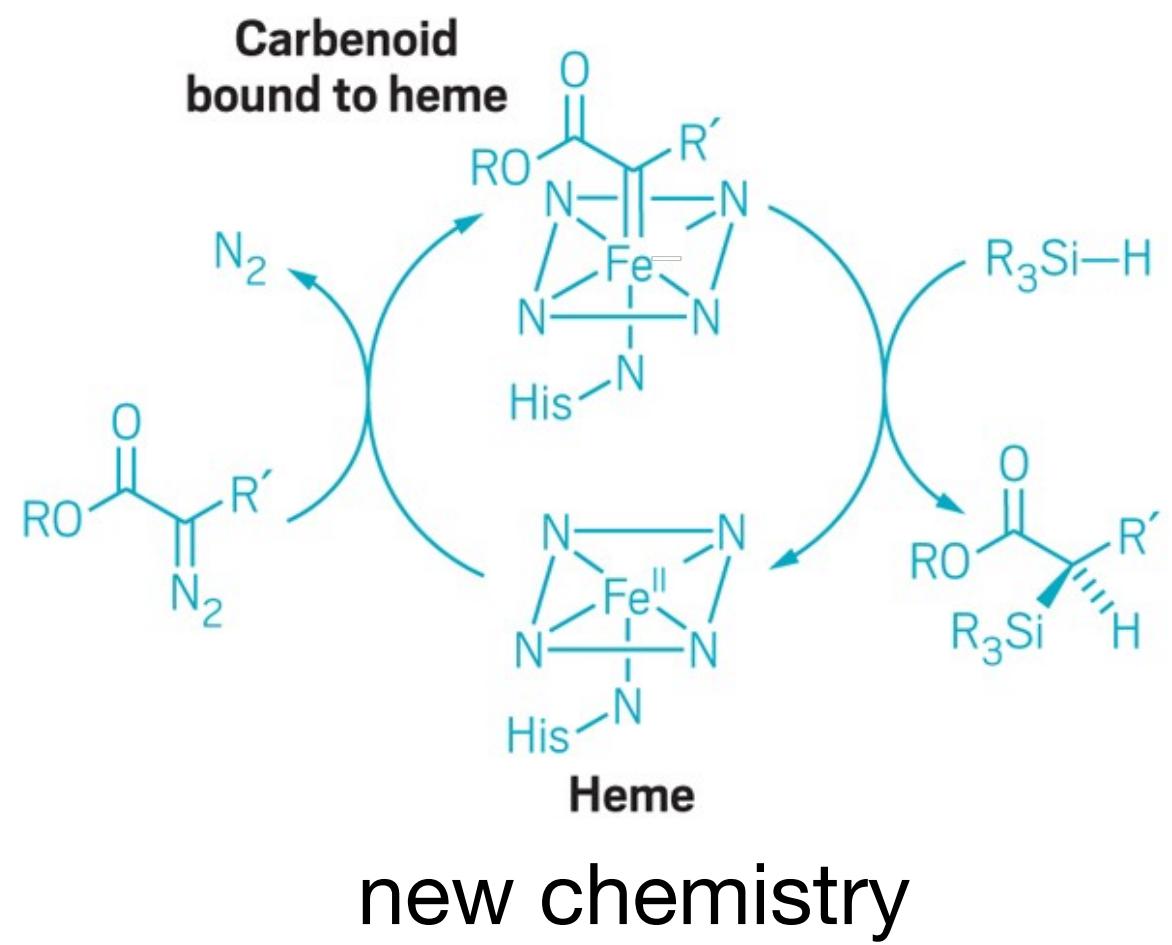
Why design proteins?

Why design proteins?

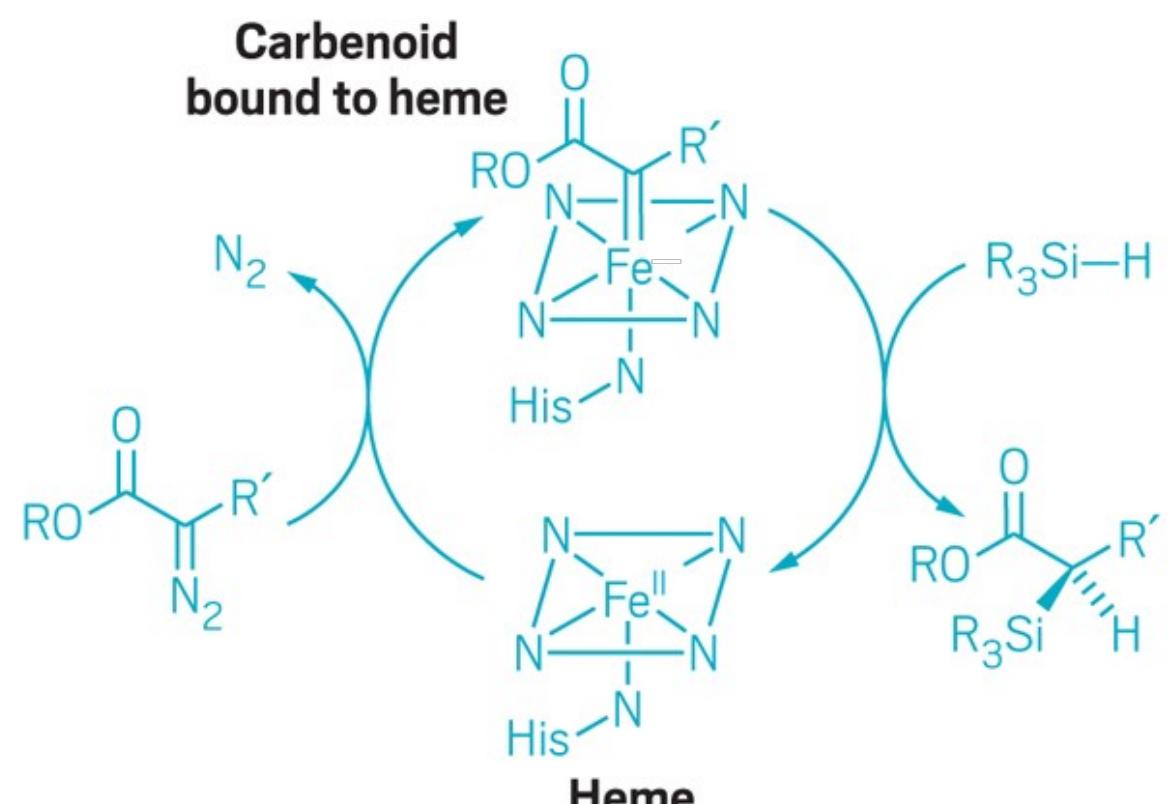


new chemistry

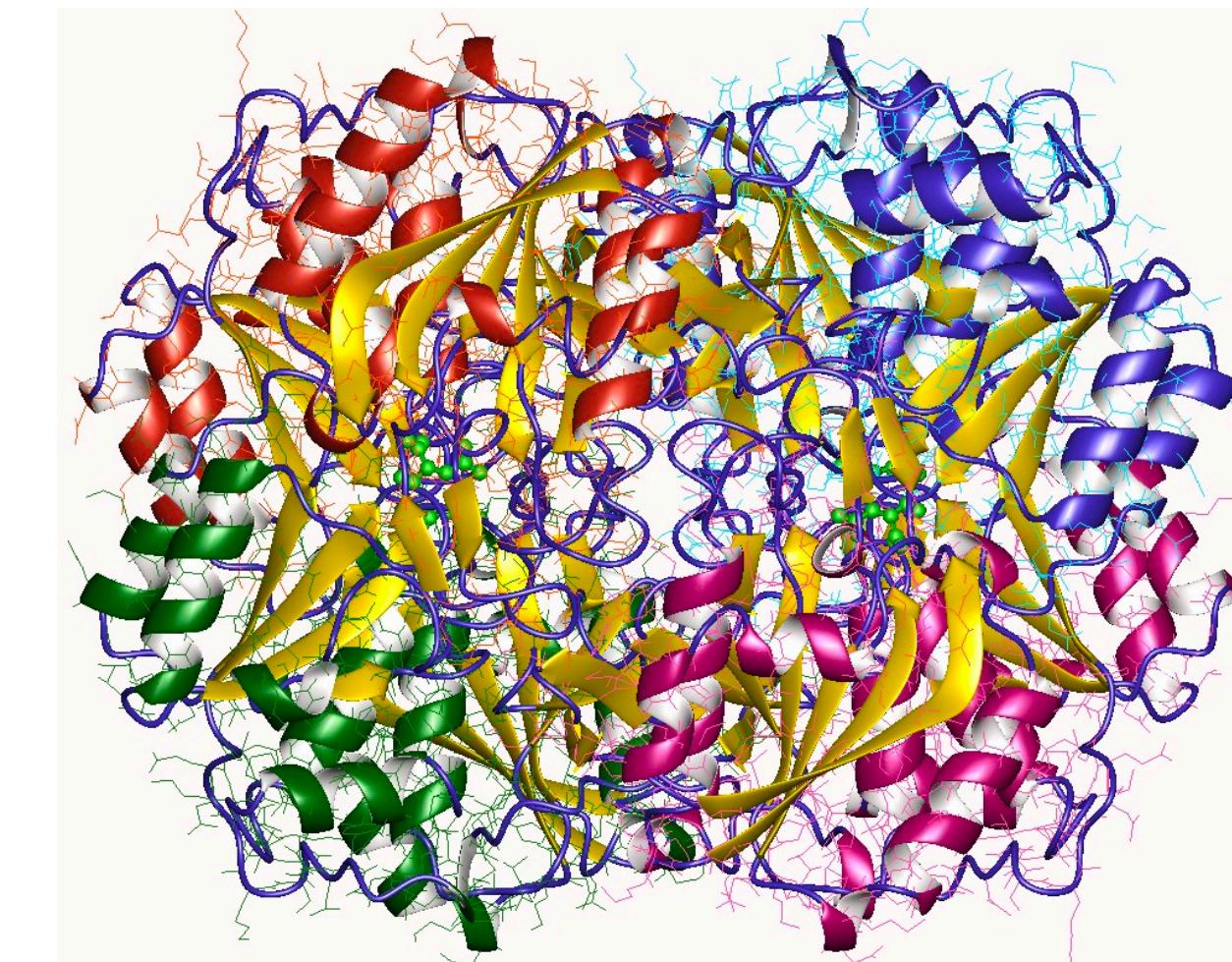
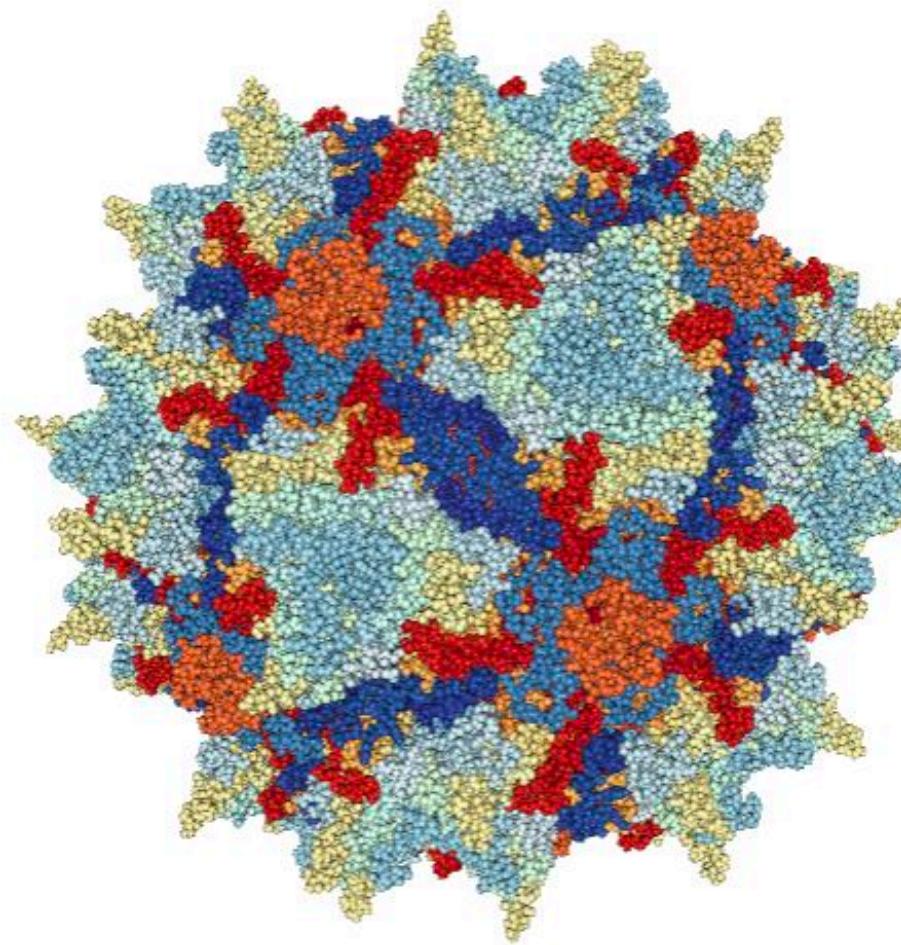
Why design proteins?



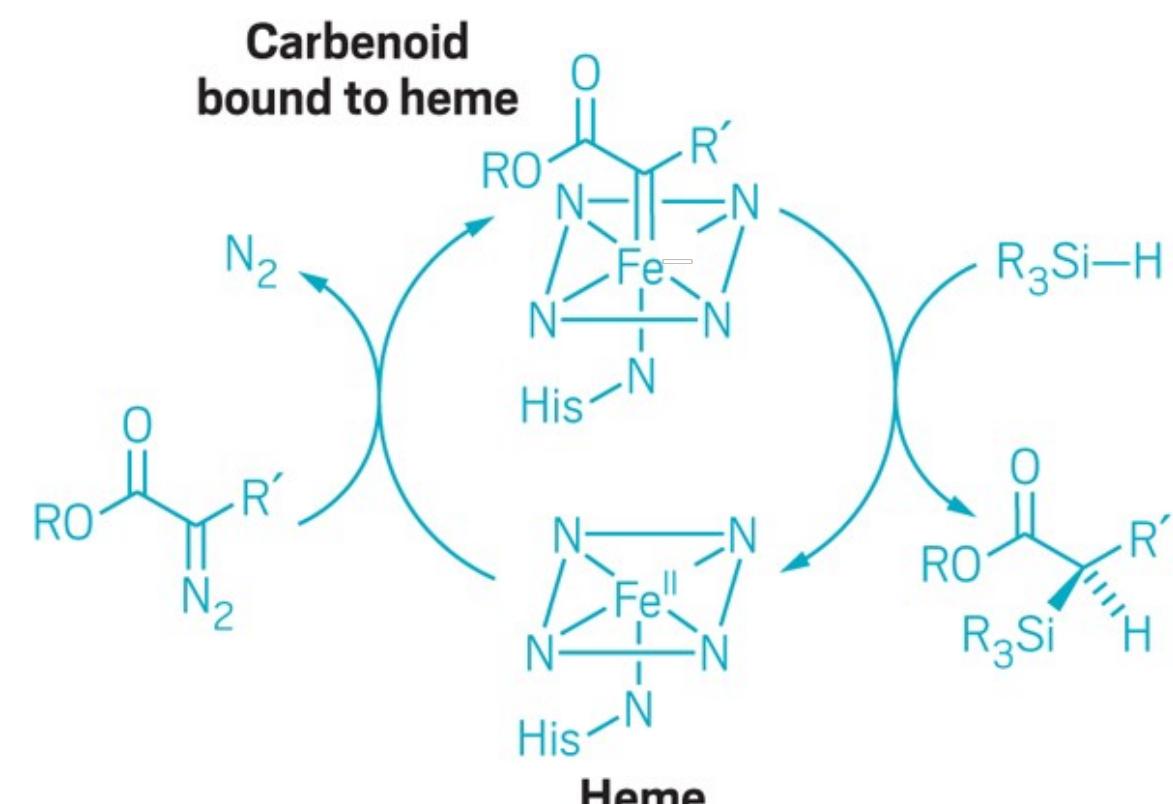
Why design proteins?



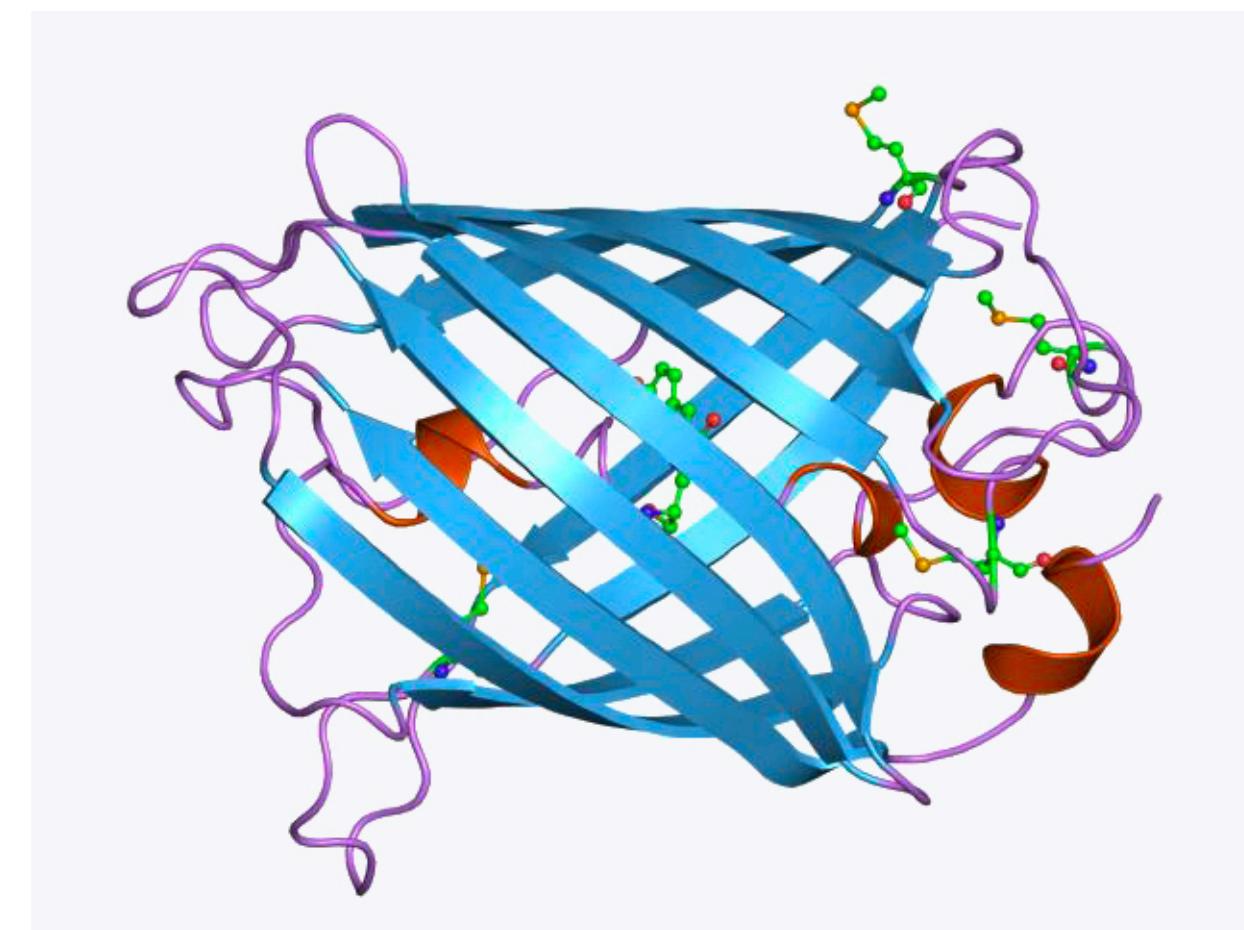
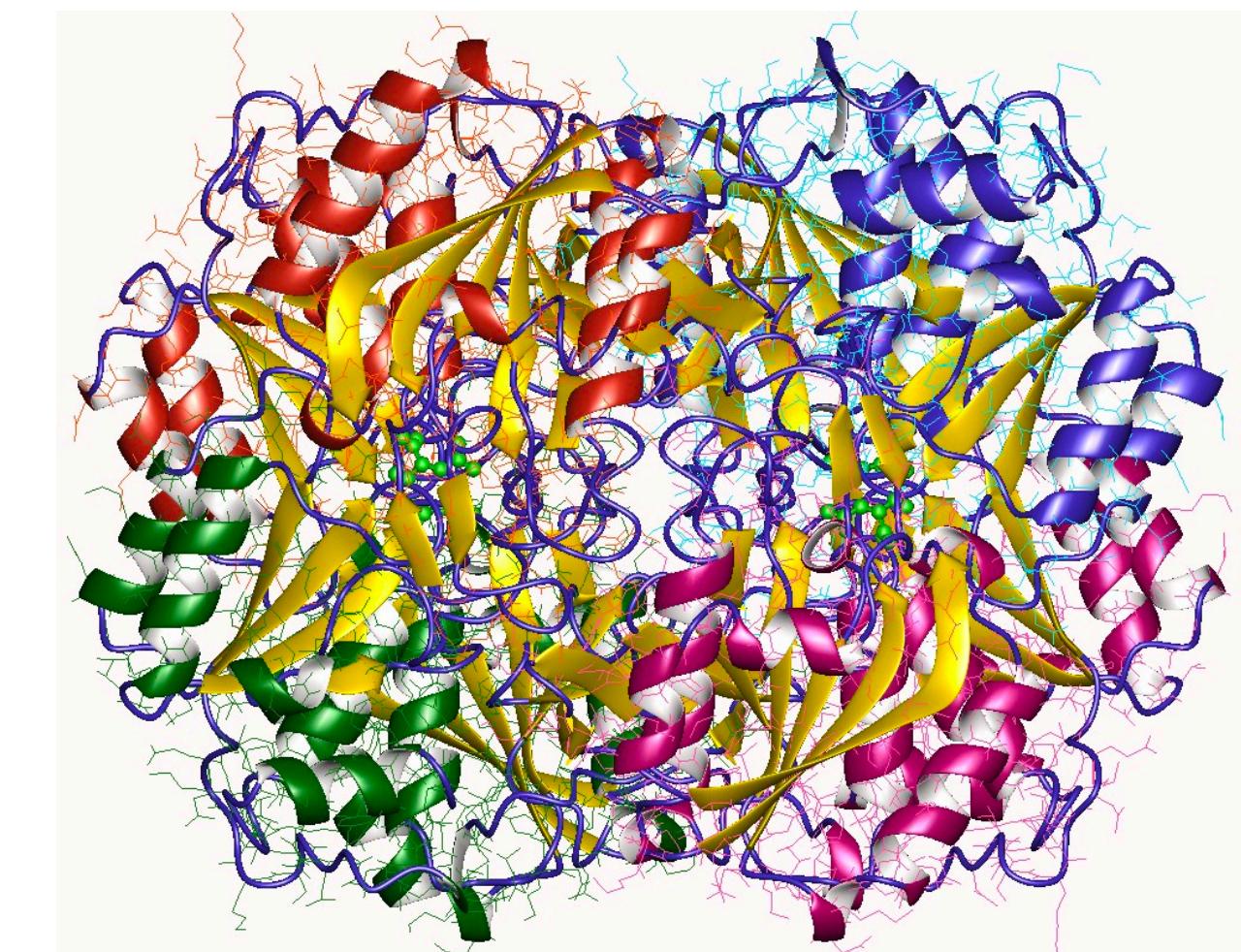
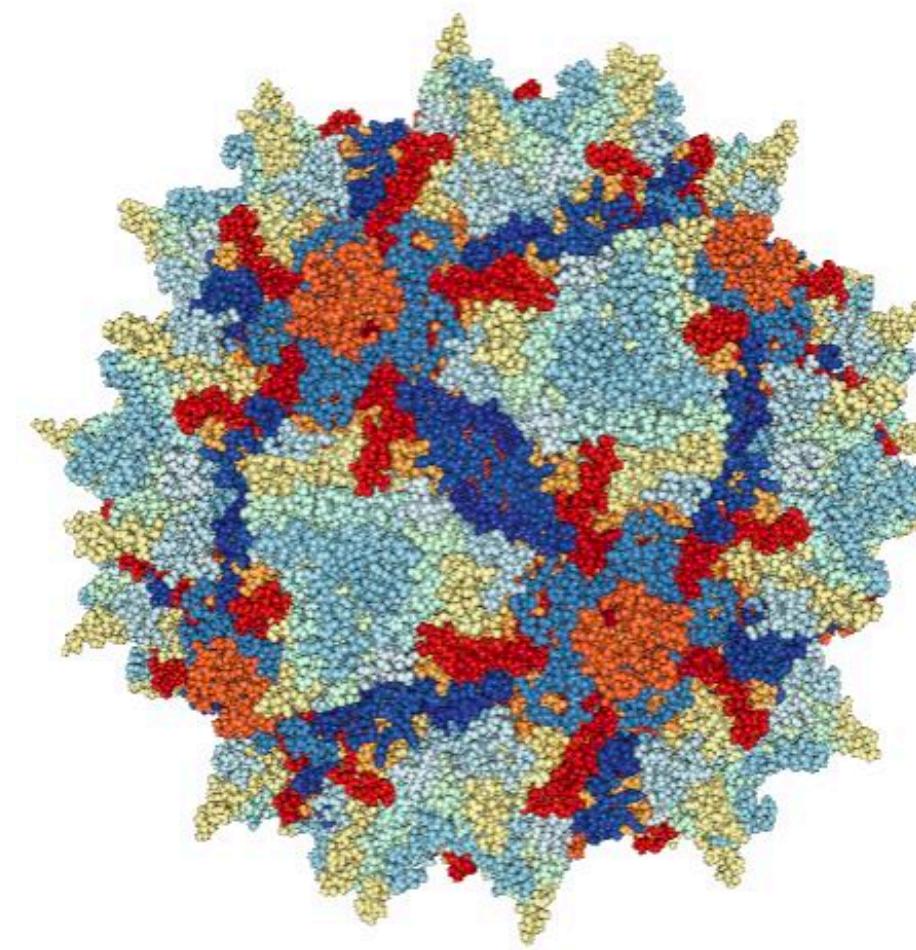
new chemistry



Why design proteins?



new chemistry

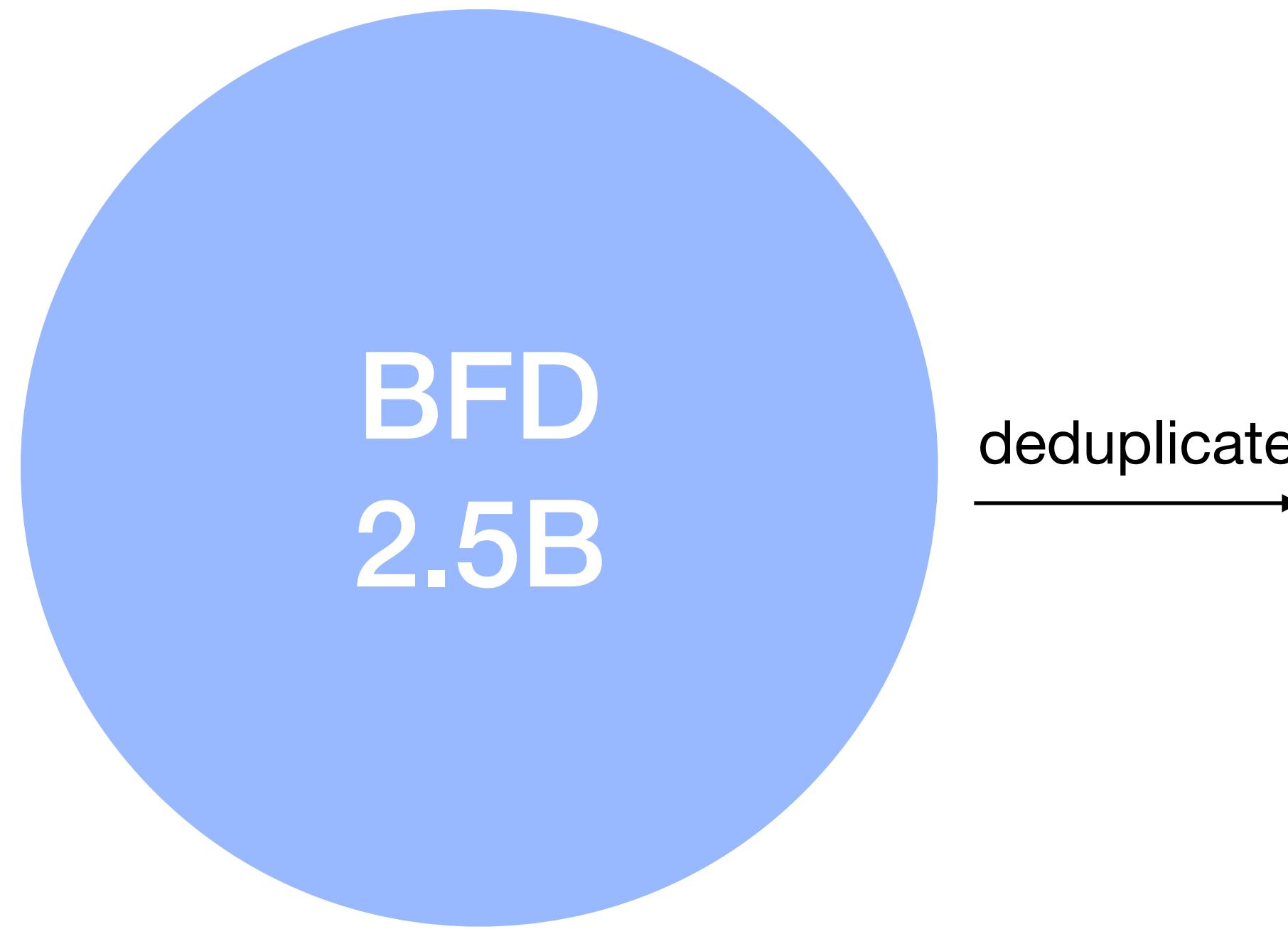


We have access to large protein databases

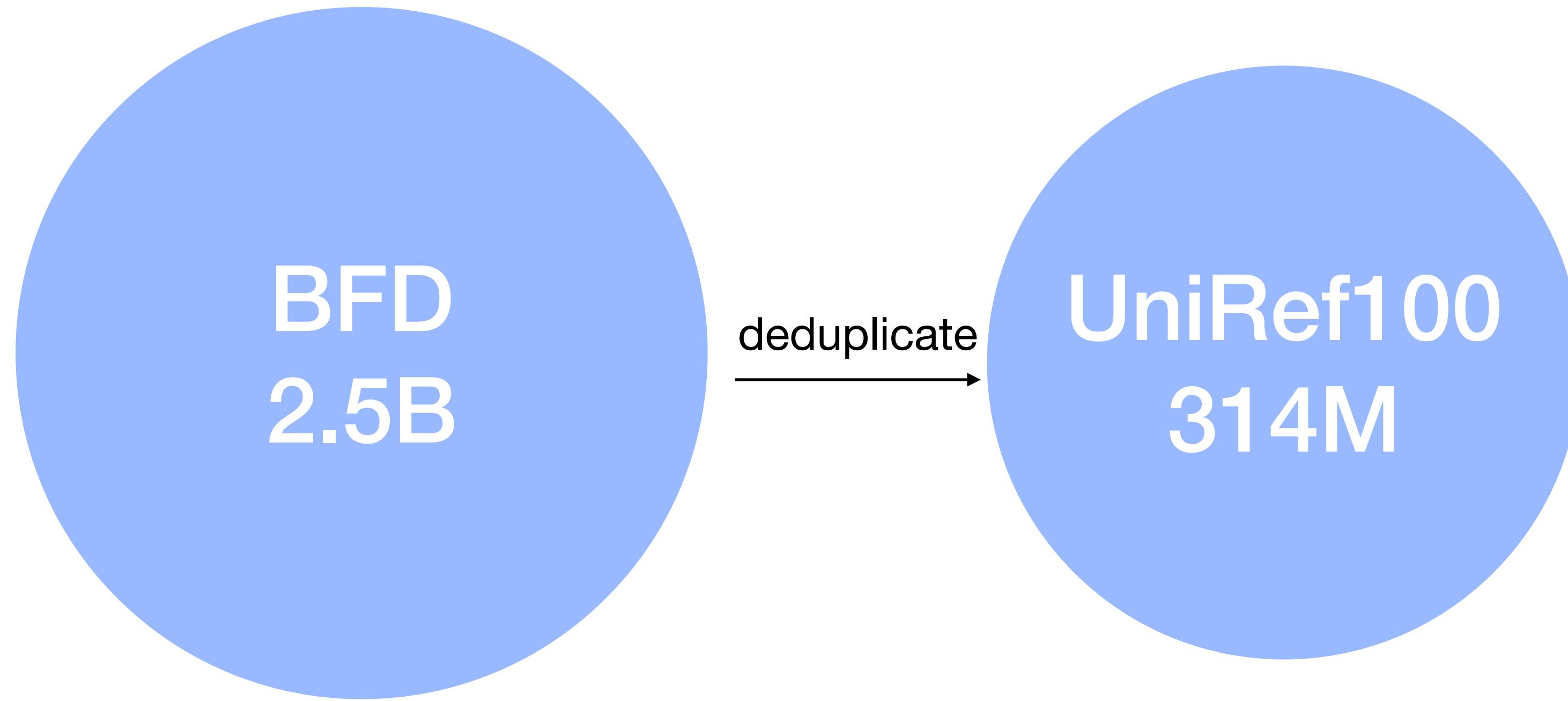
We have access to large protein databases



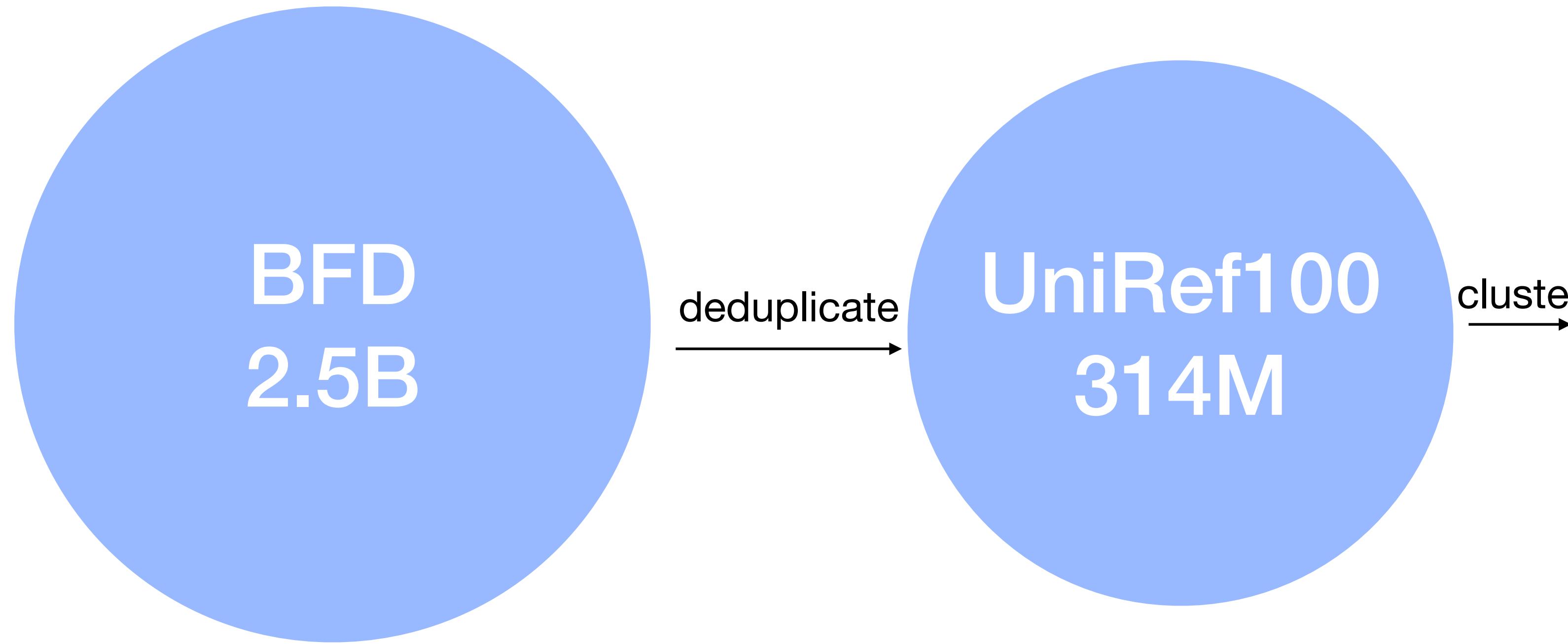
We have access to large protein databases



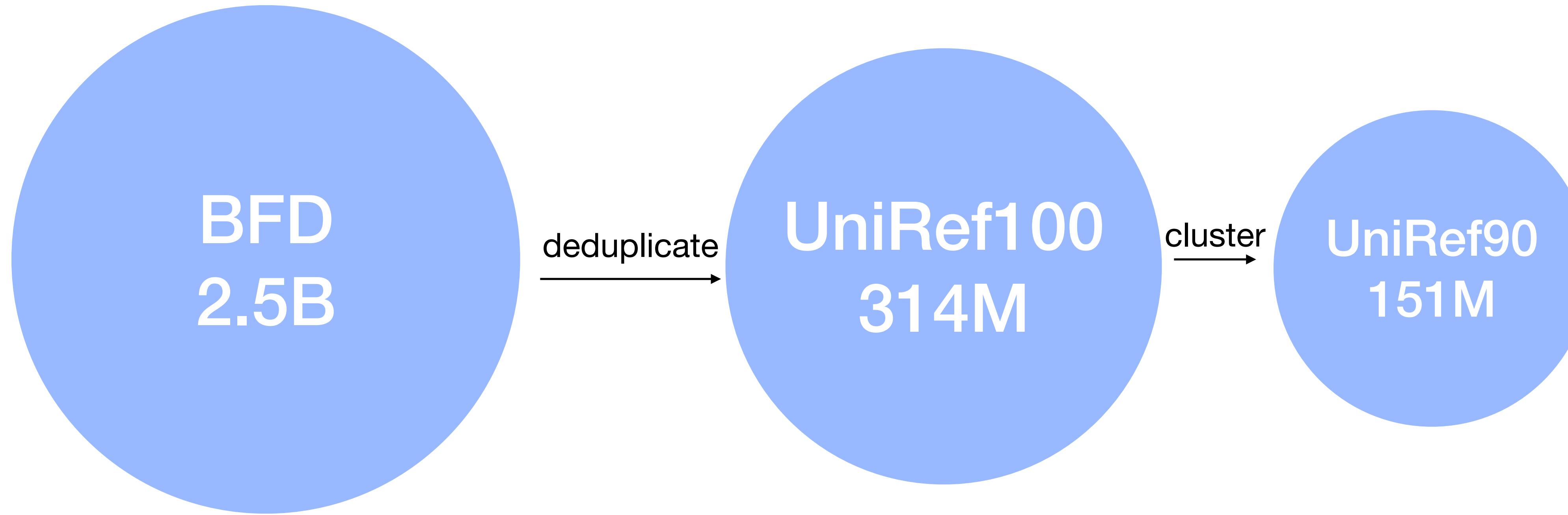
We have access to large protein databases



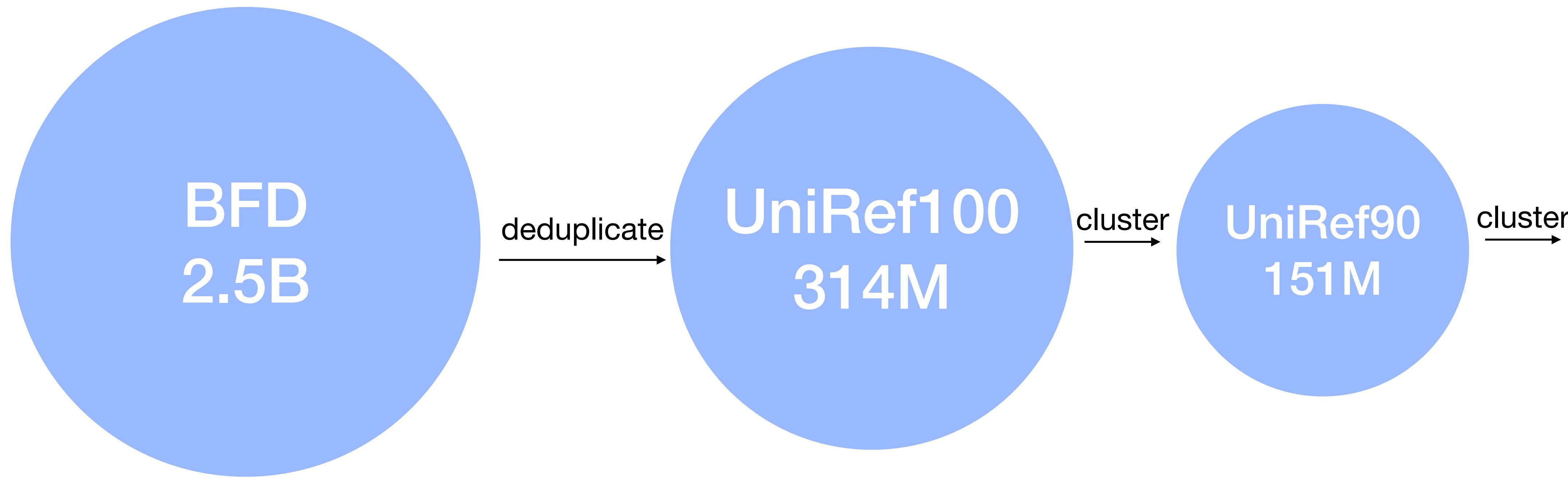
We have access to large protein databases



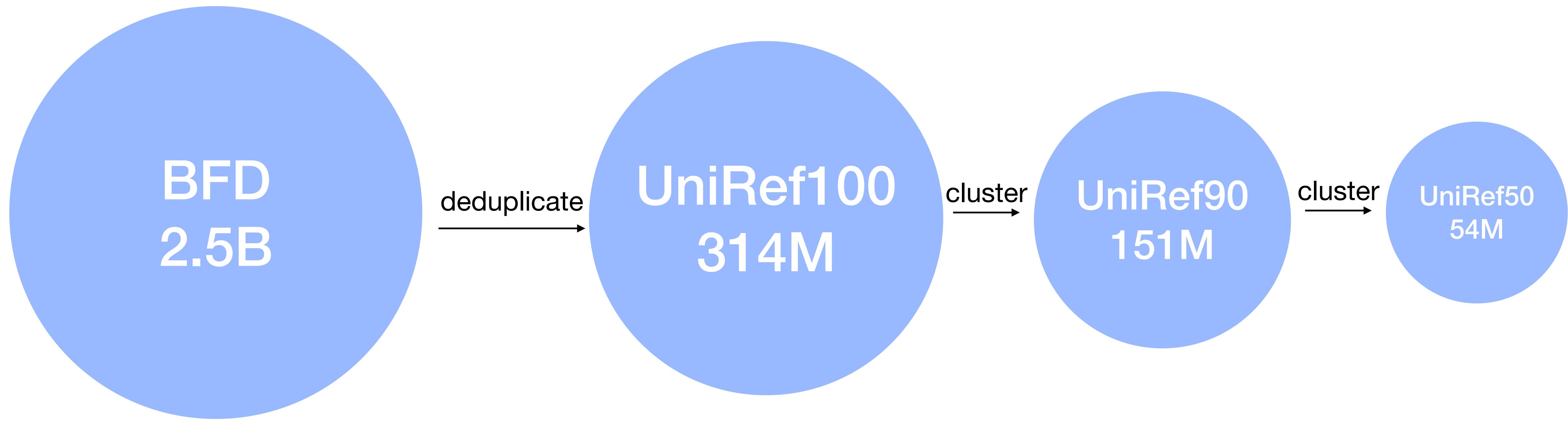
We have access to large protein databases



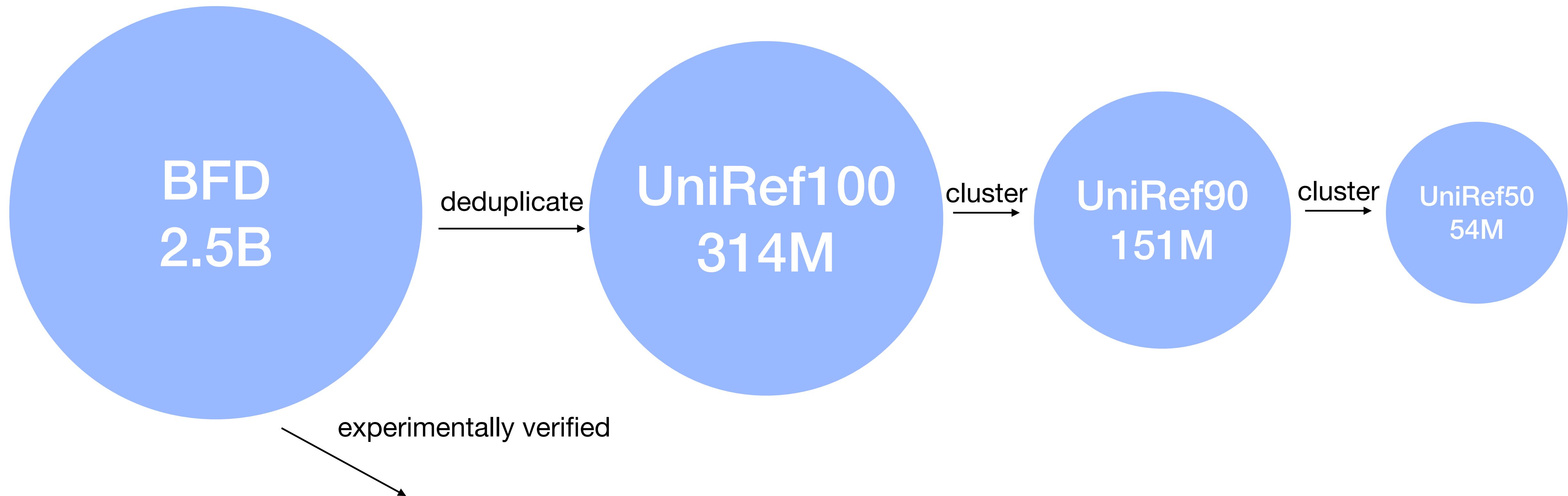
We have access to large protein databases



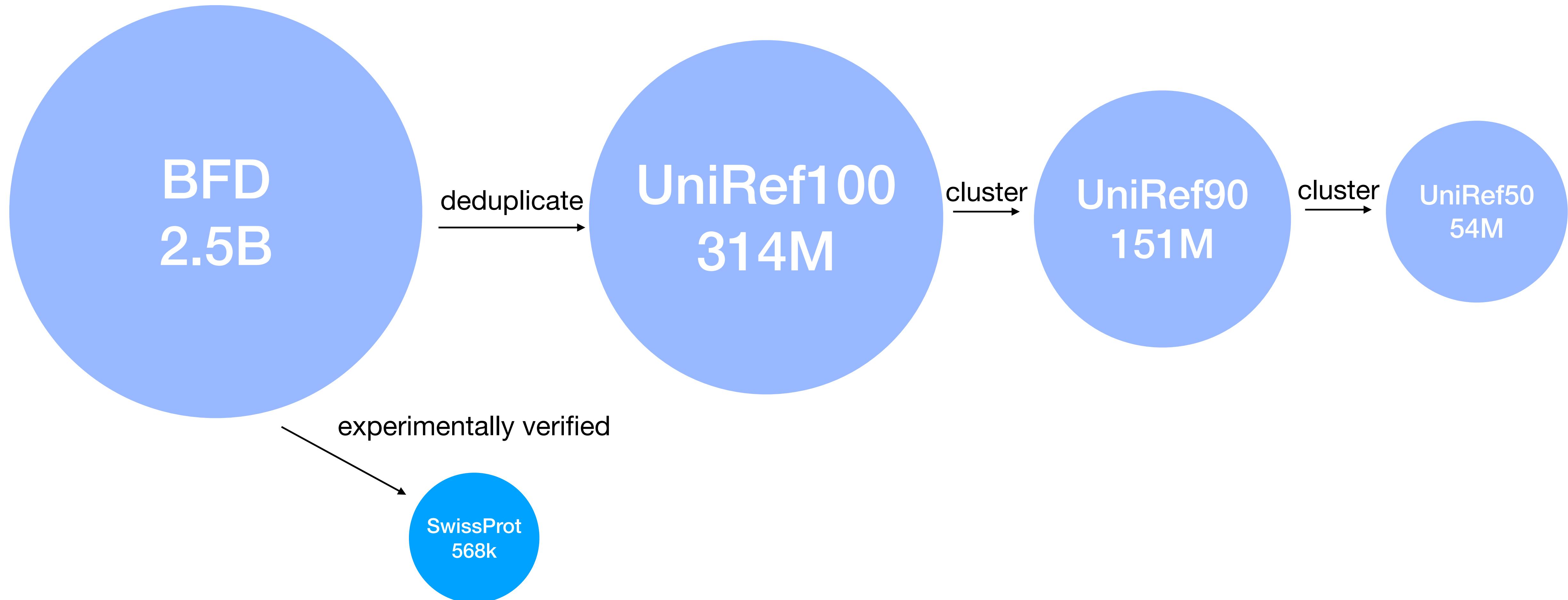
We have access to large protein databases



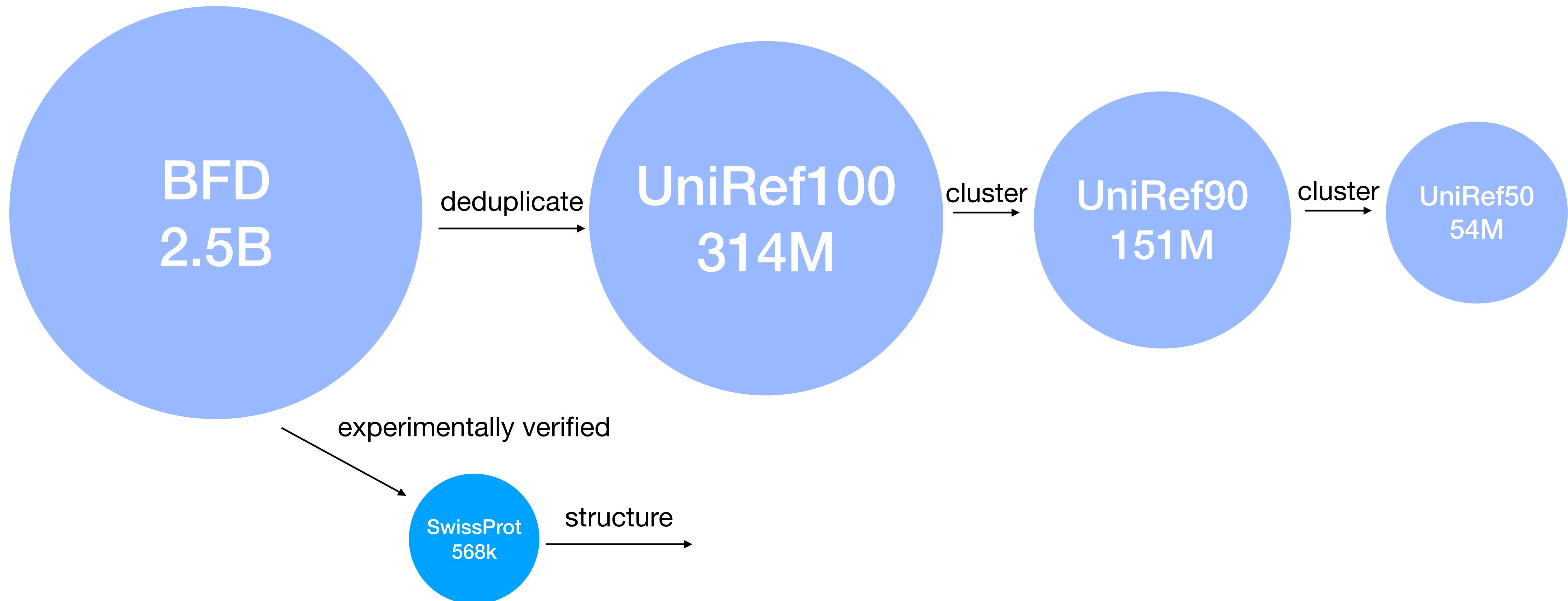
We have access to large protein databases



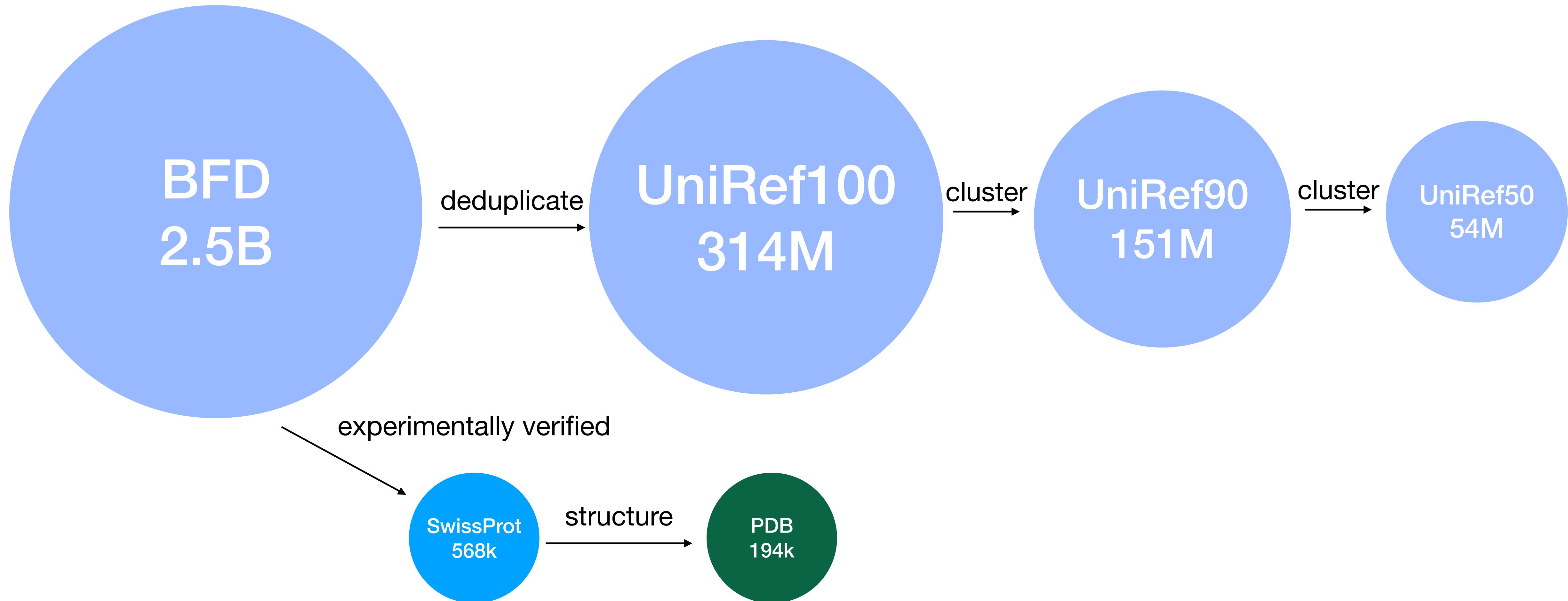
We have access to large protein databases



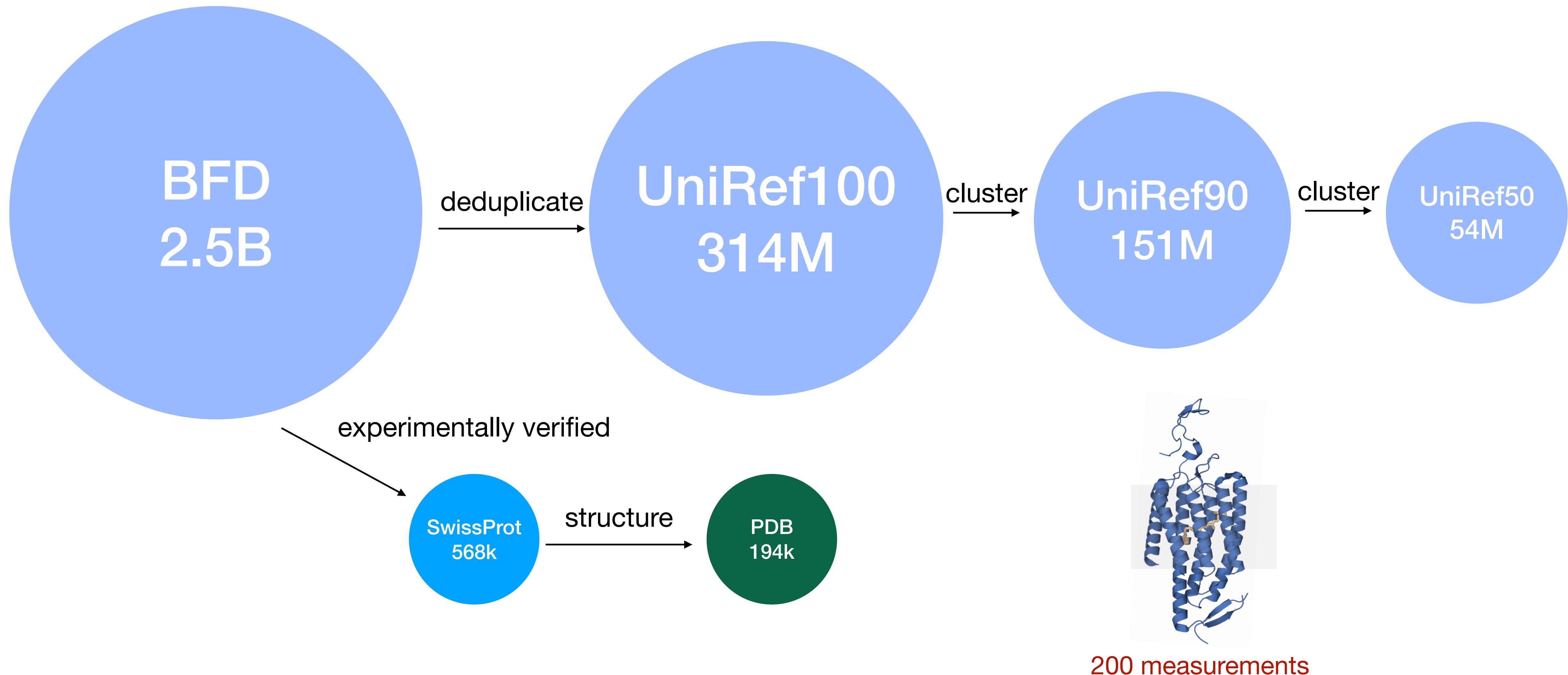
We have access to large protein databases



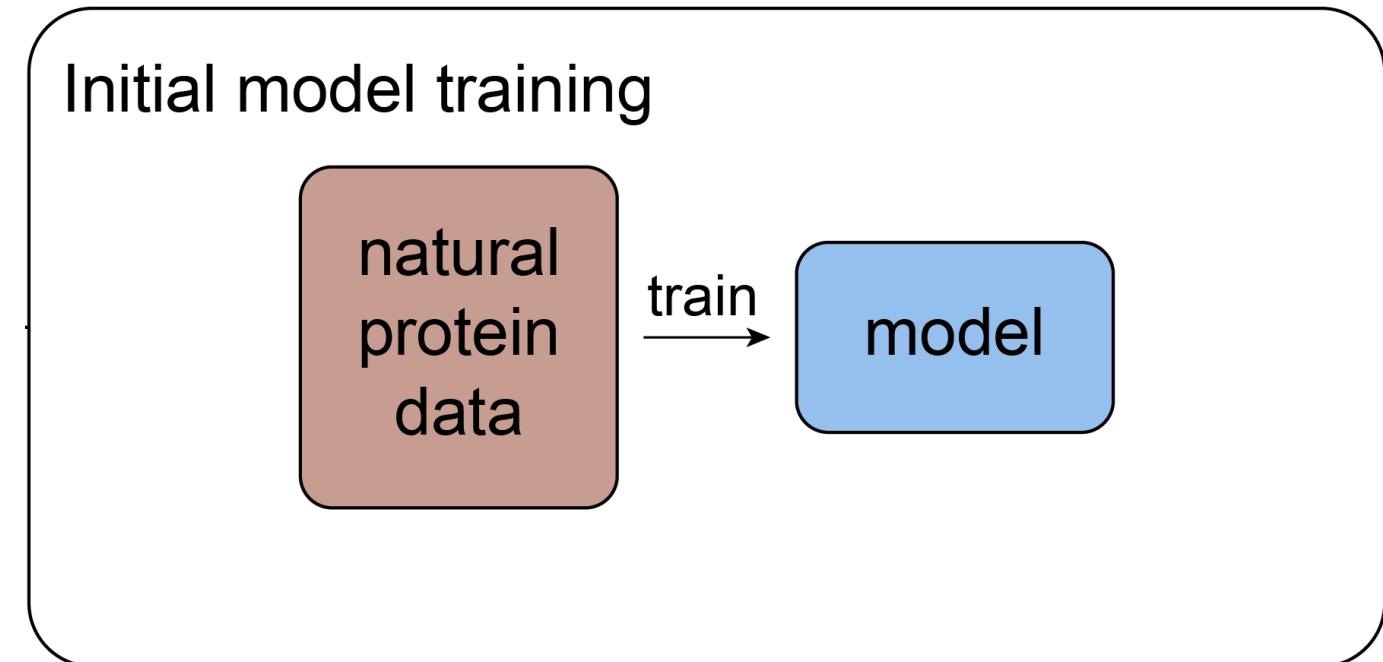
We have access to large protein databases



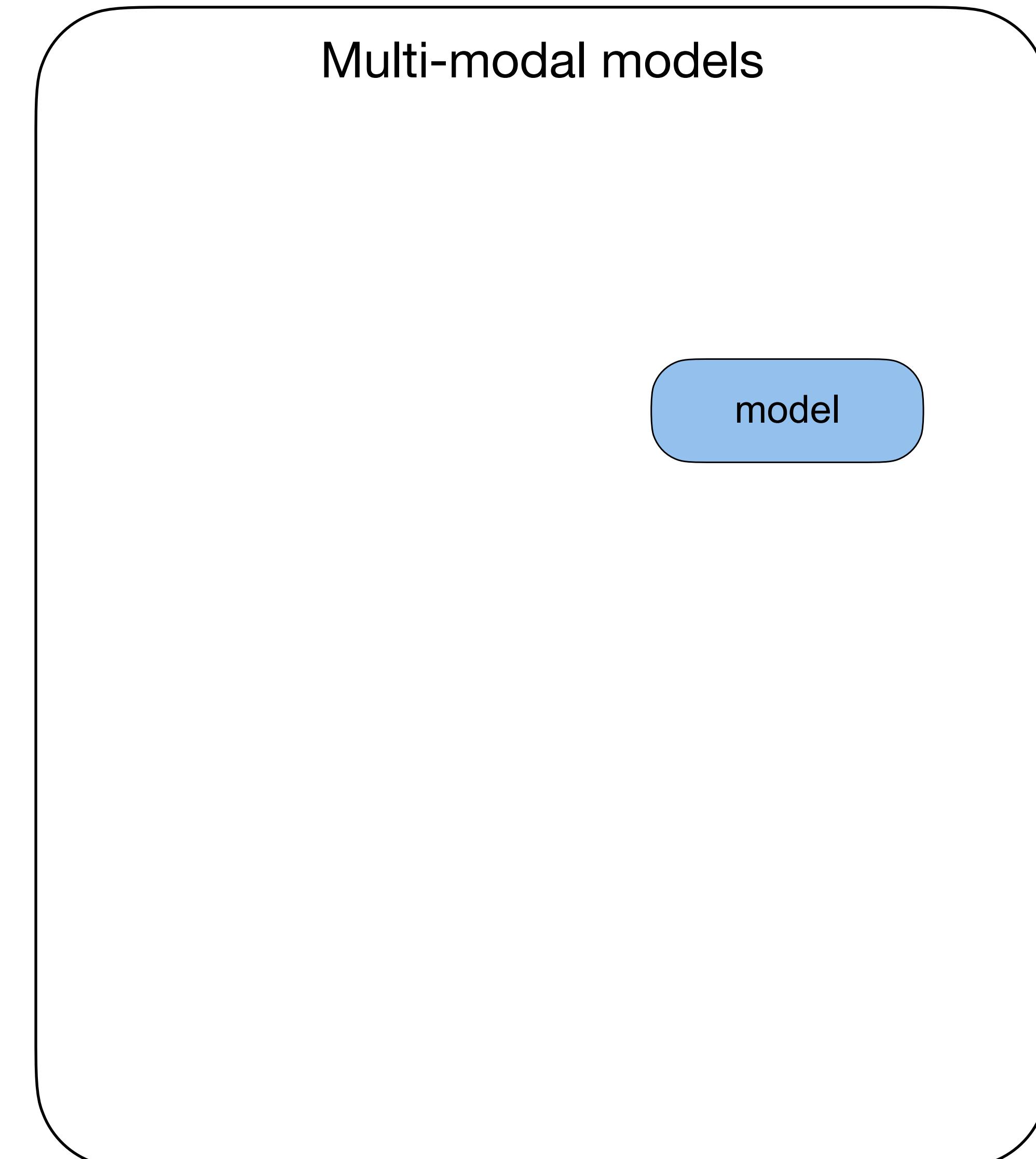
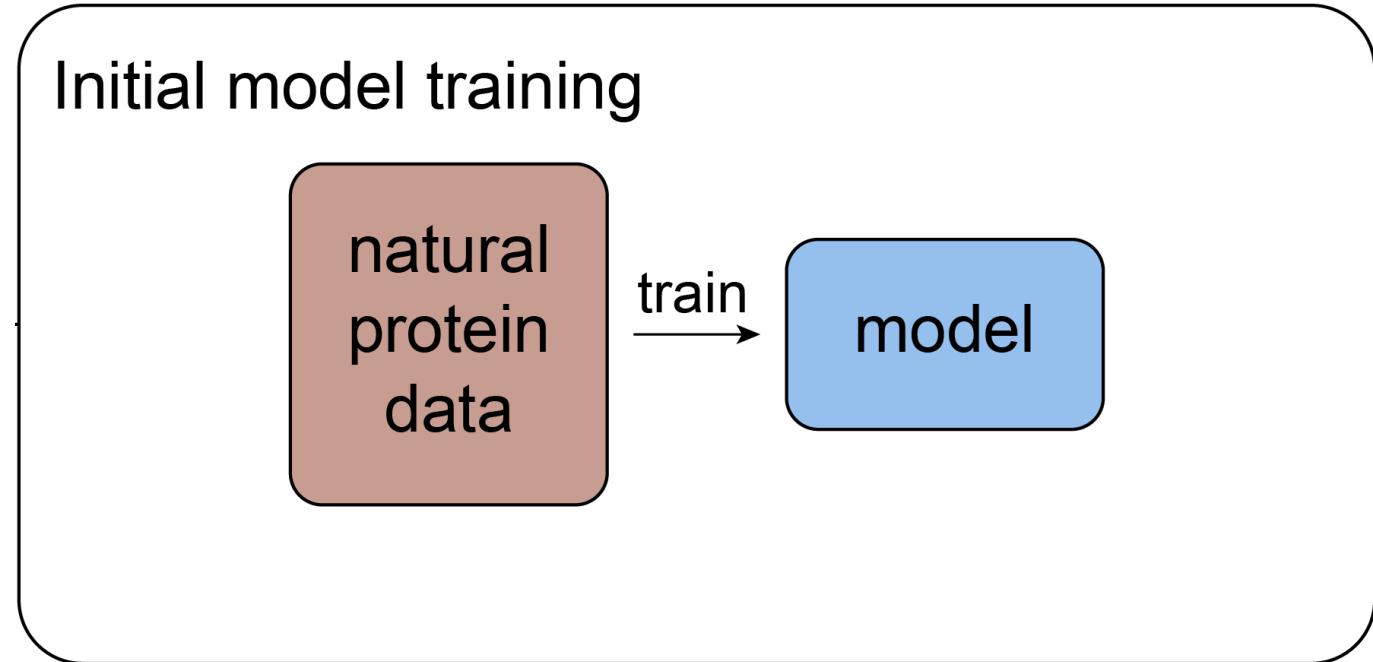
We have access to large protein databases



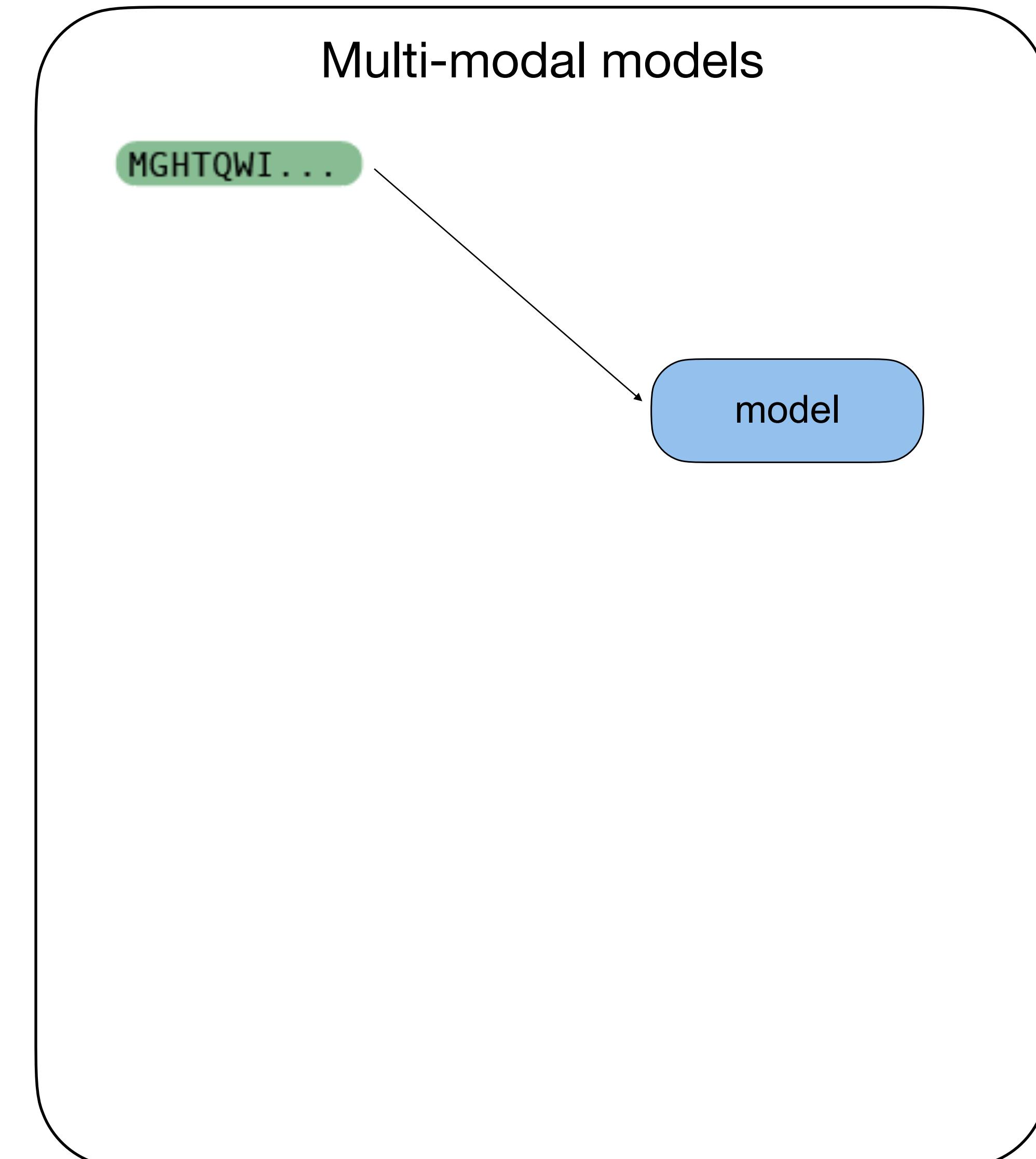
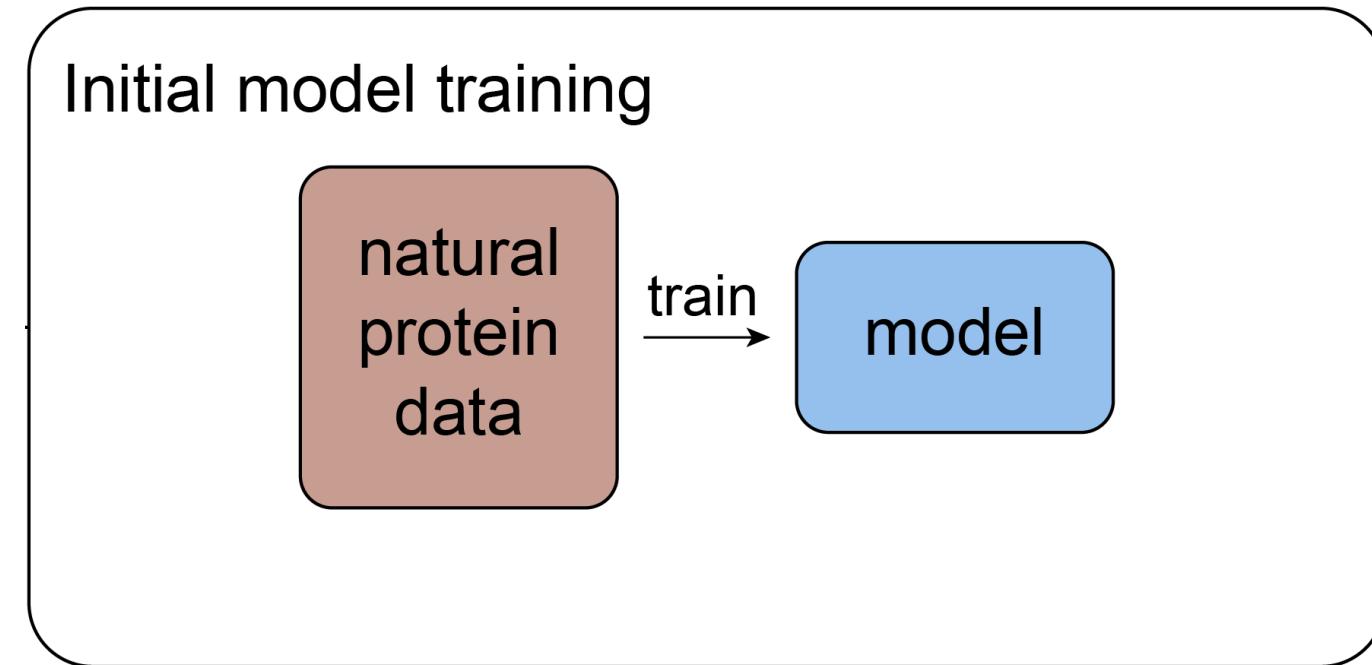
Use multiple data modalities to discover and design proteins



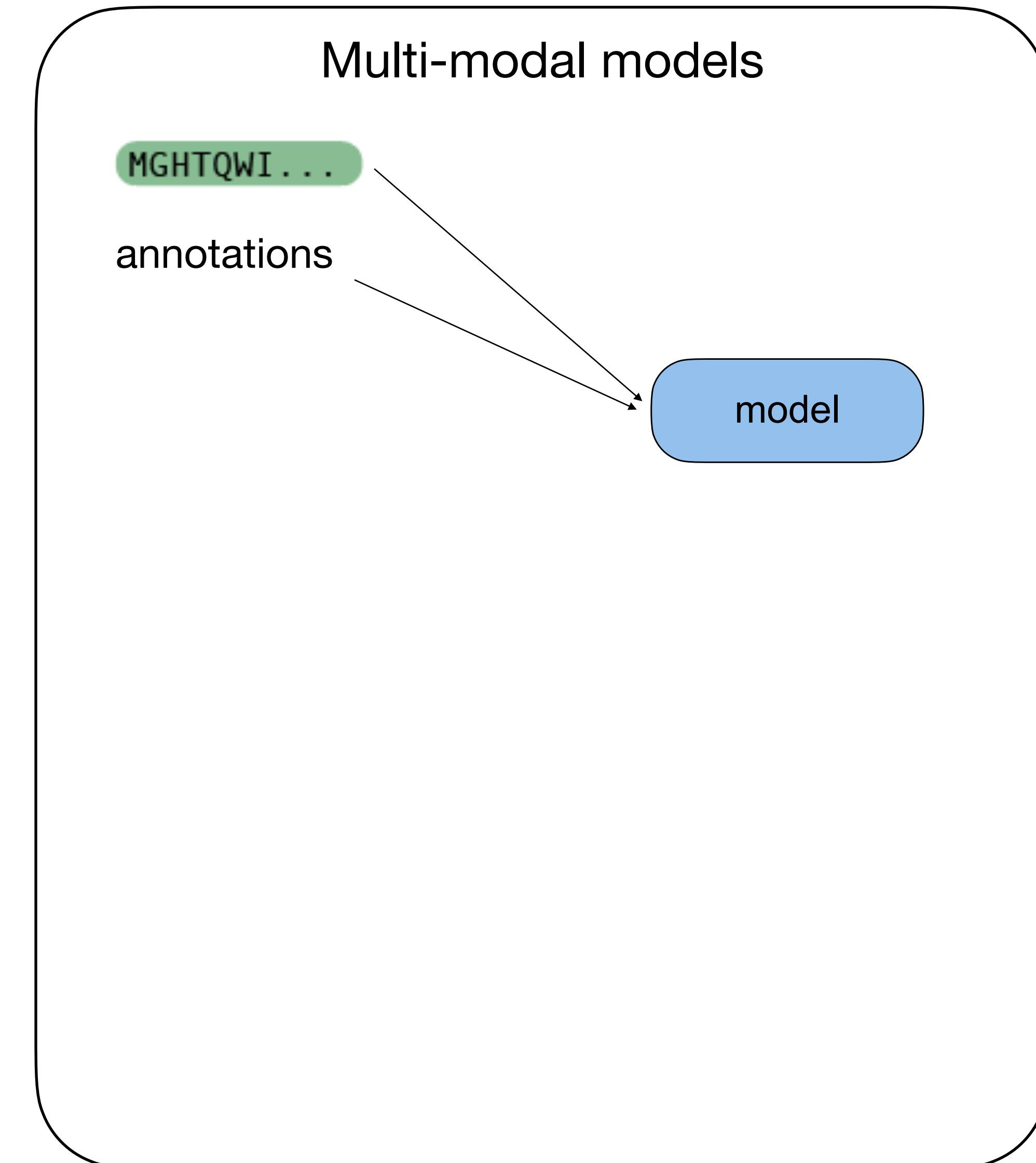
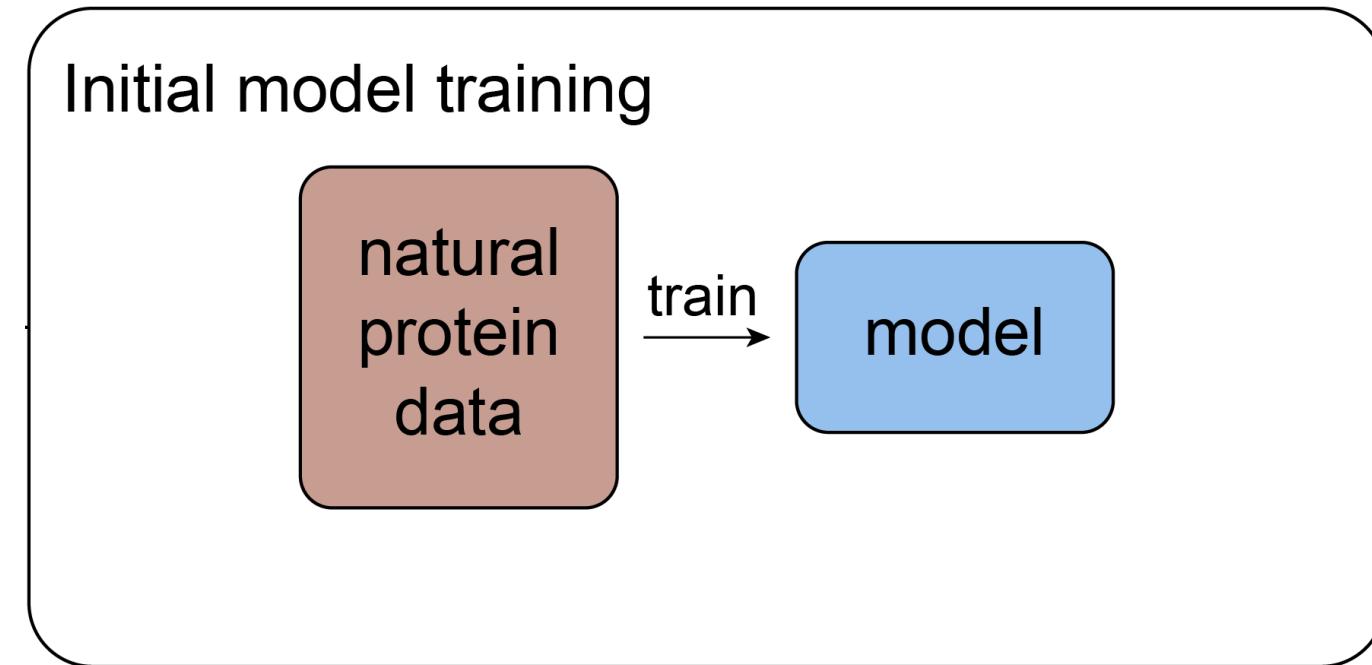
Use multiple data modalities to discover and design proteins



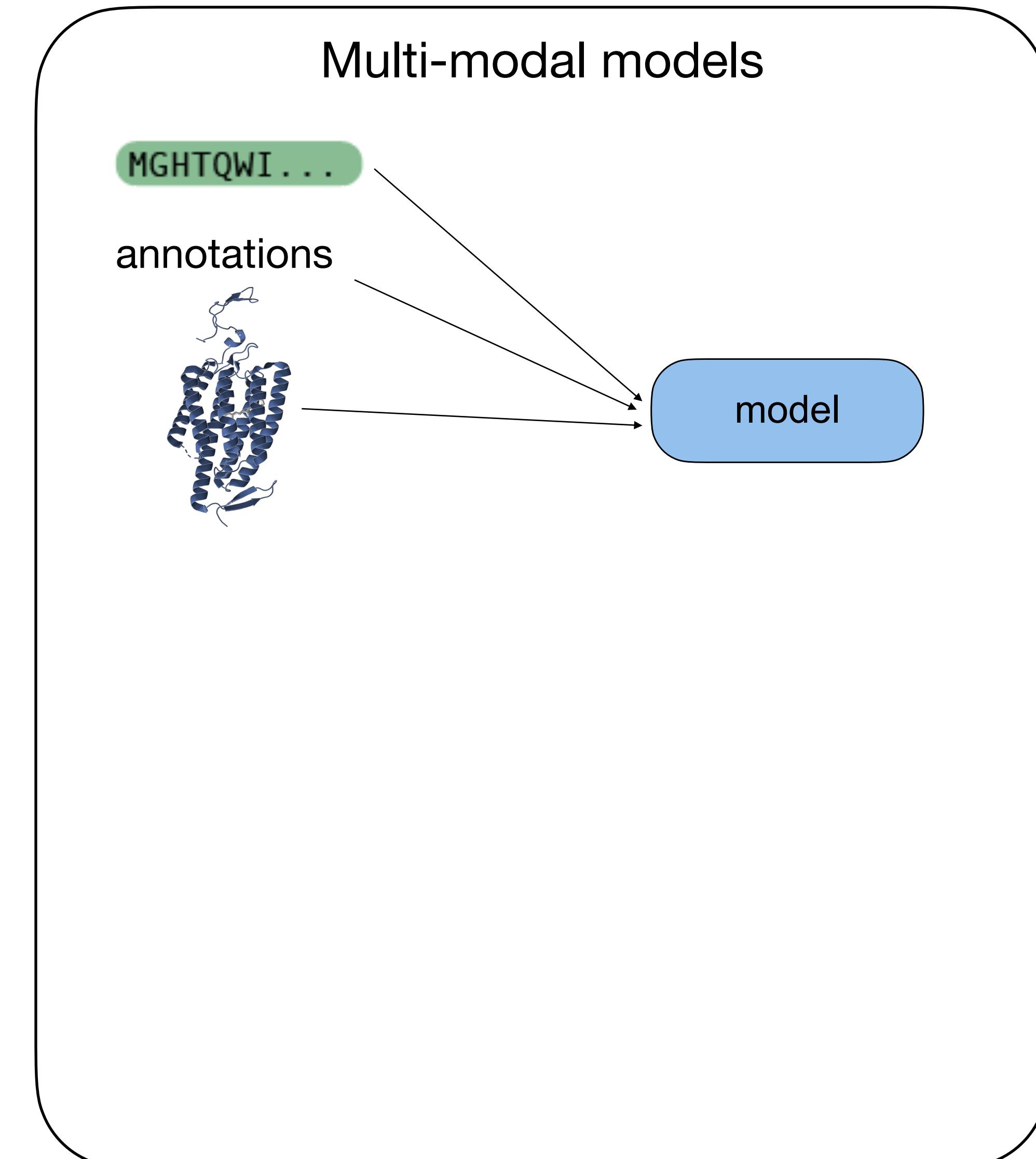
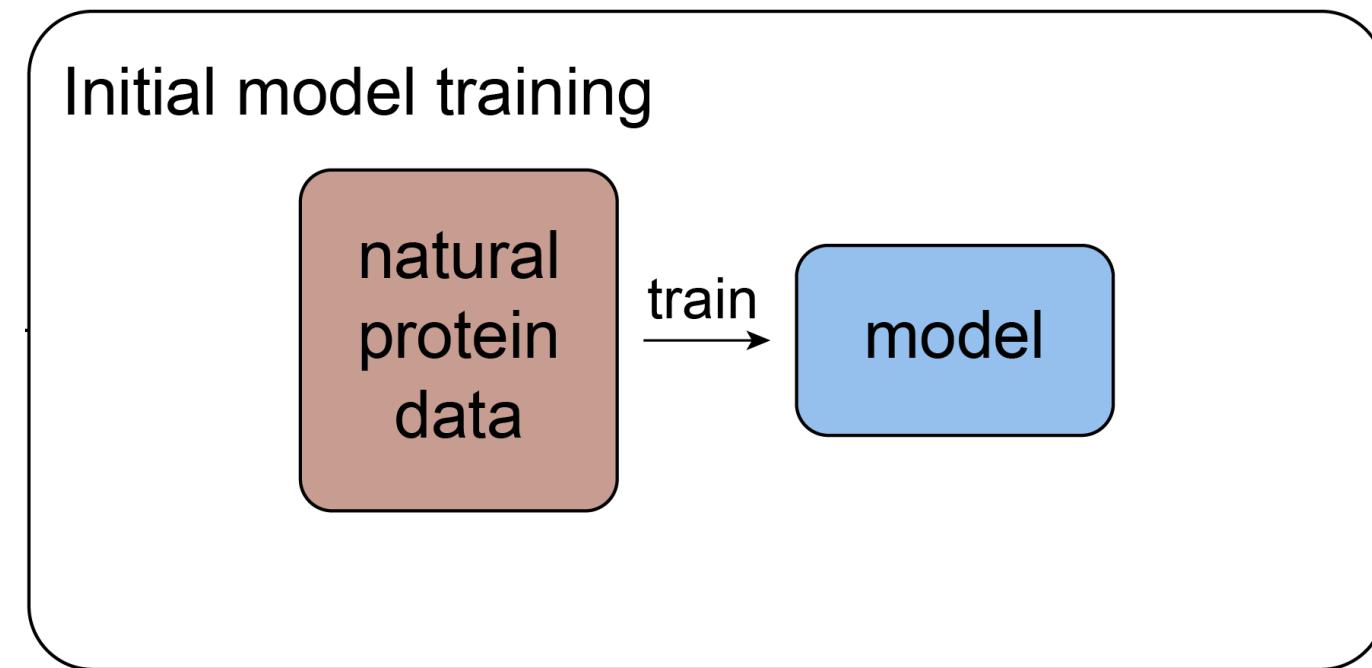
Use multiple data modalities to discover and design proteins



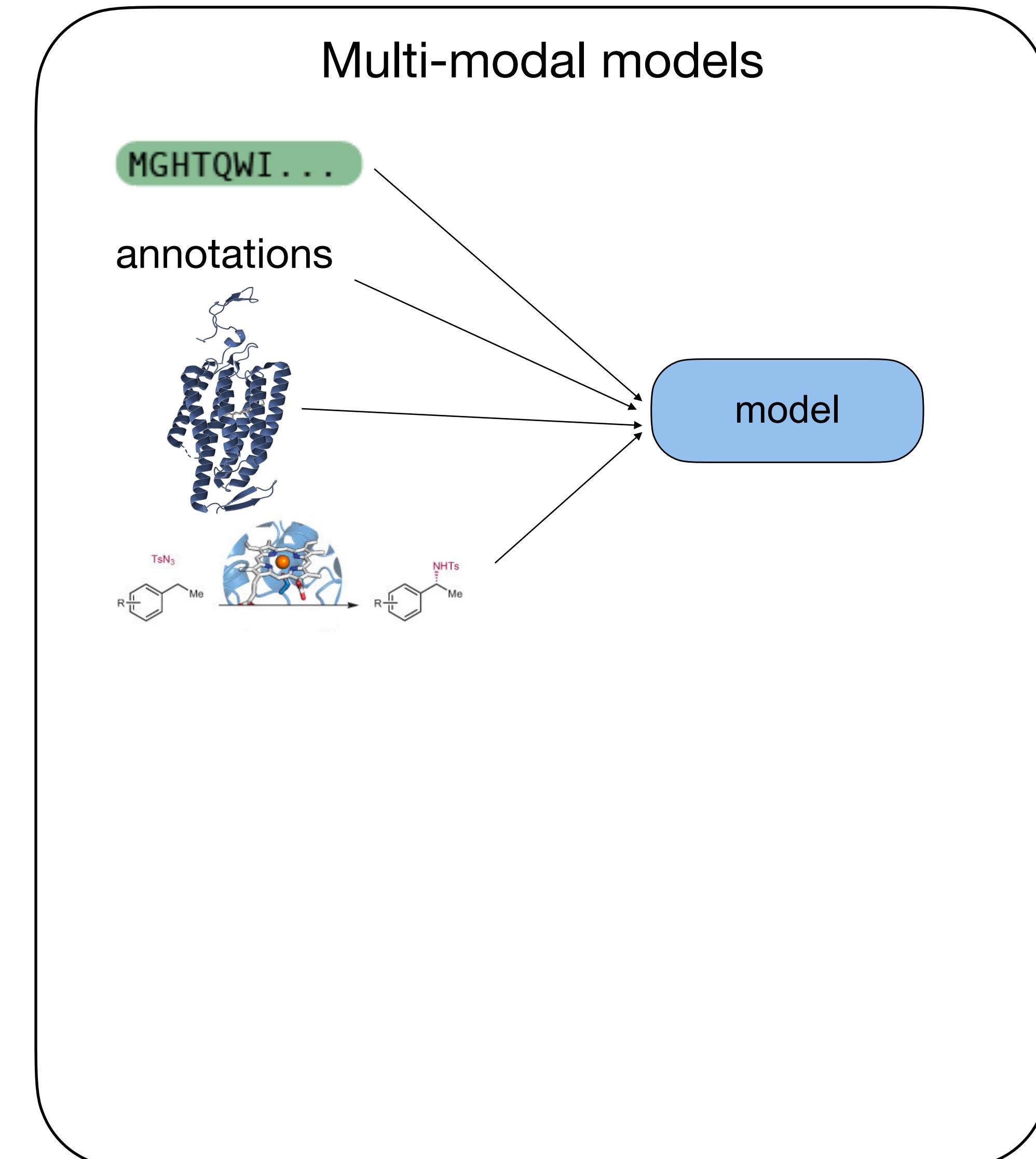
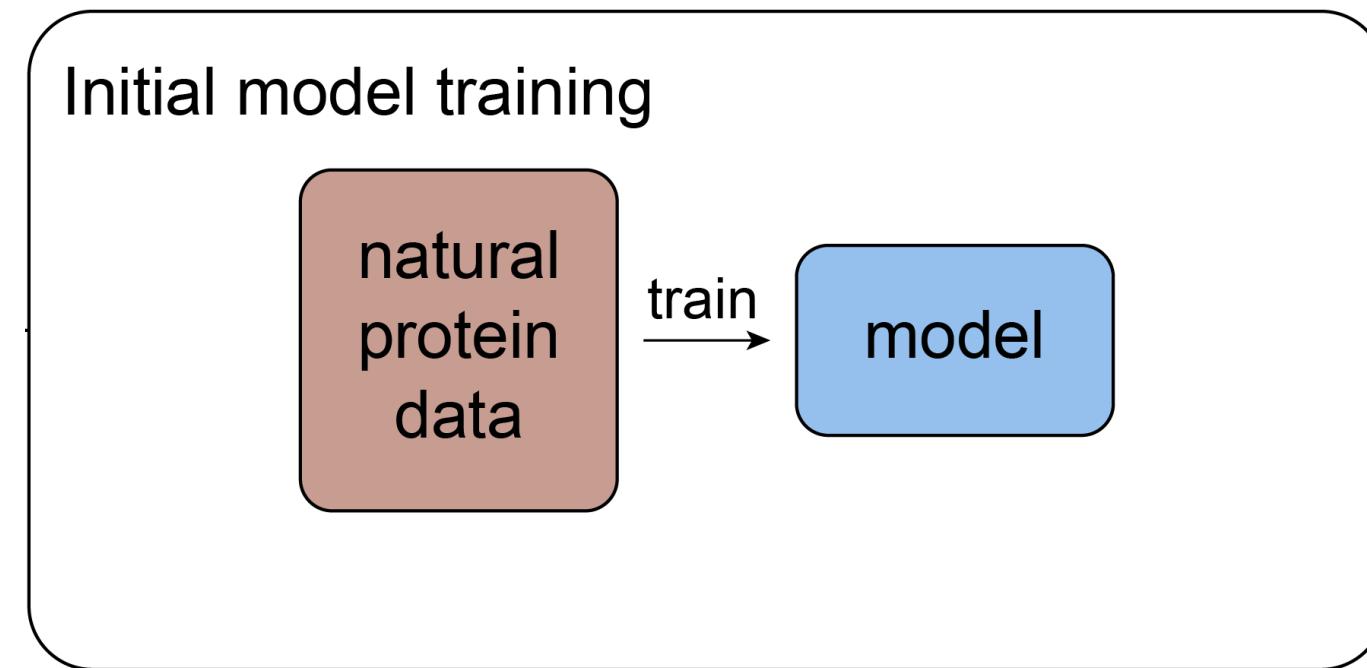
Use multiple data modalities to discover and design proteins



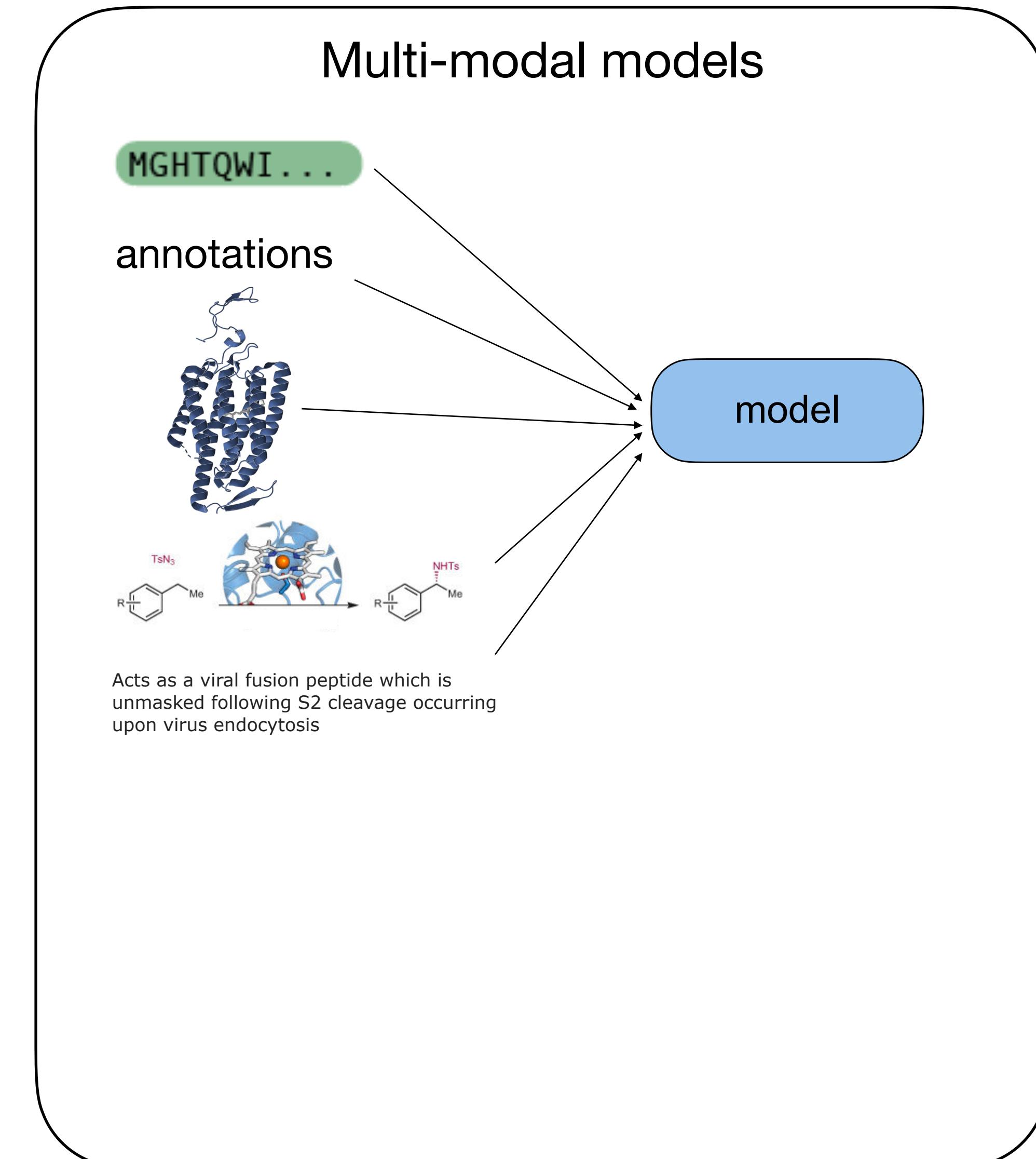
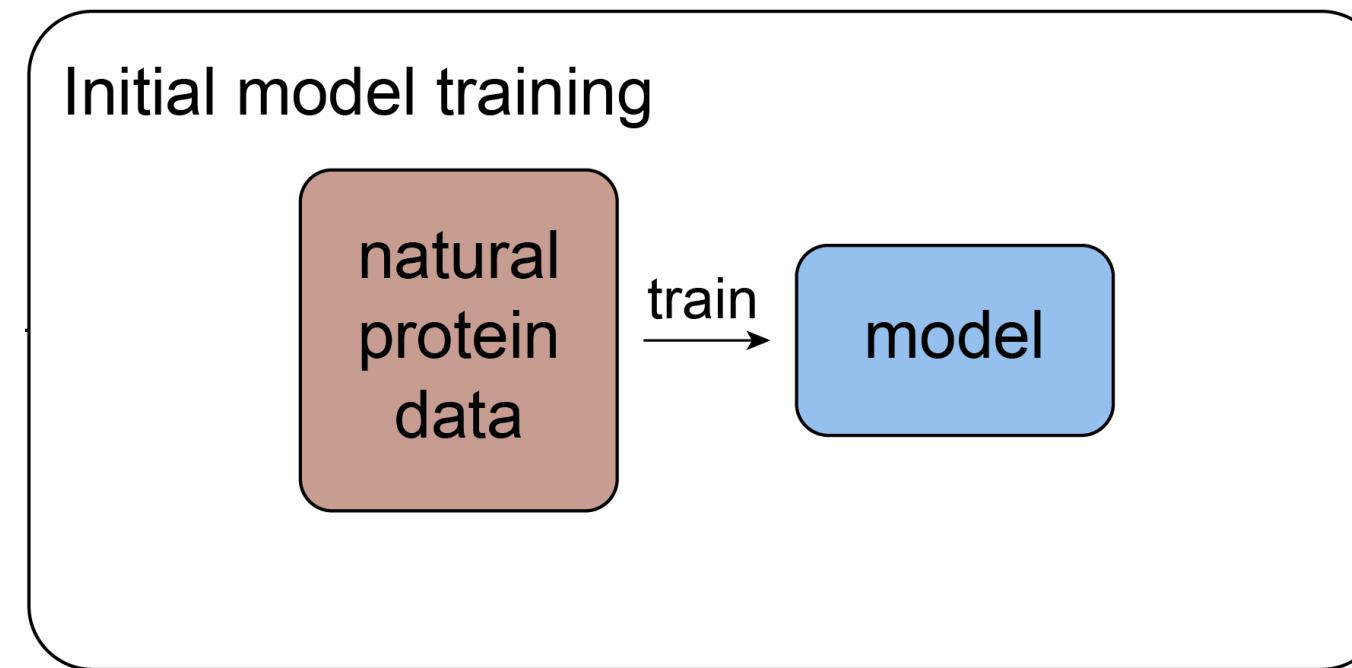
Use multiple data modalities to discover and design proteins



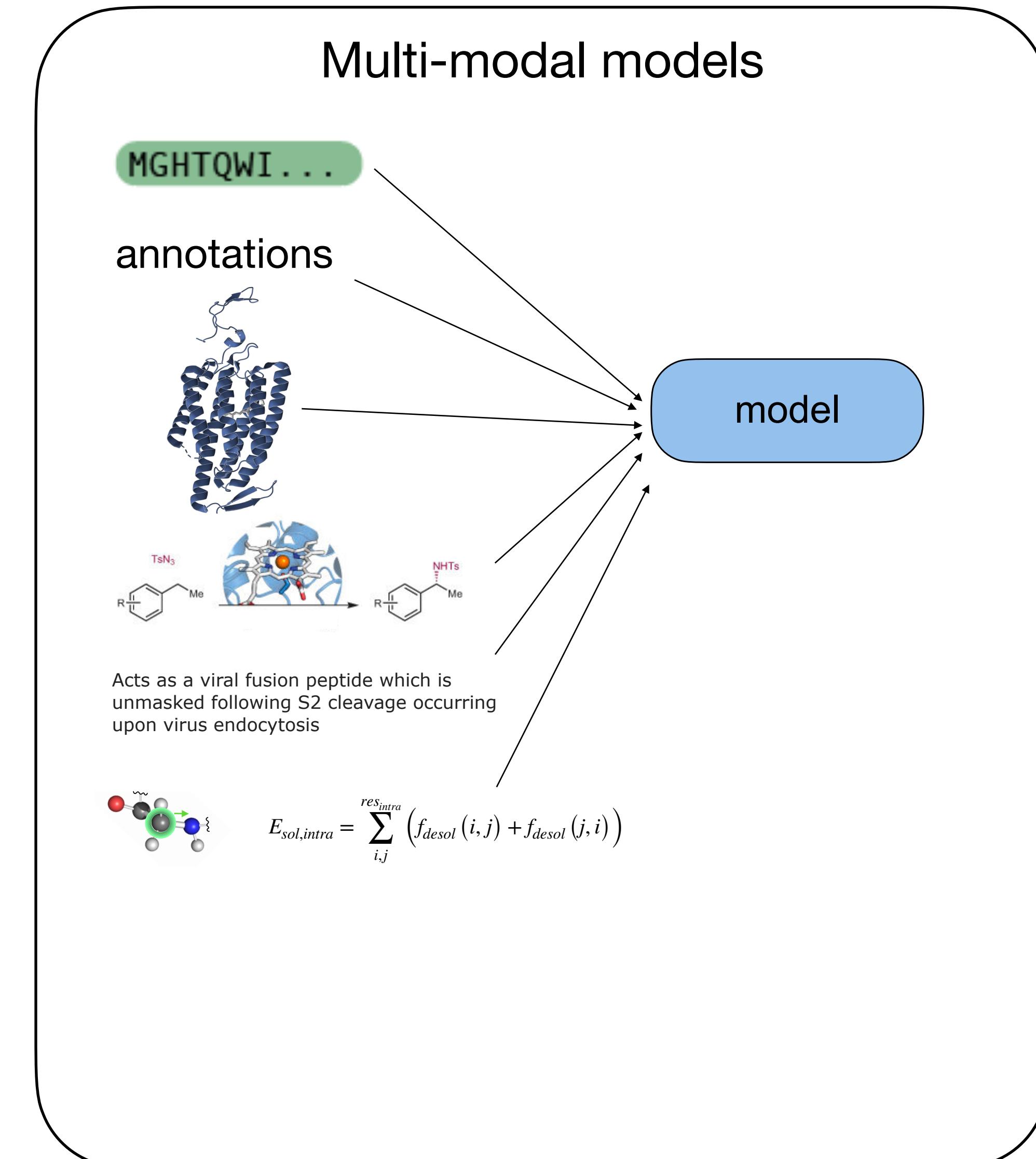
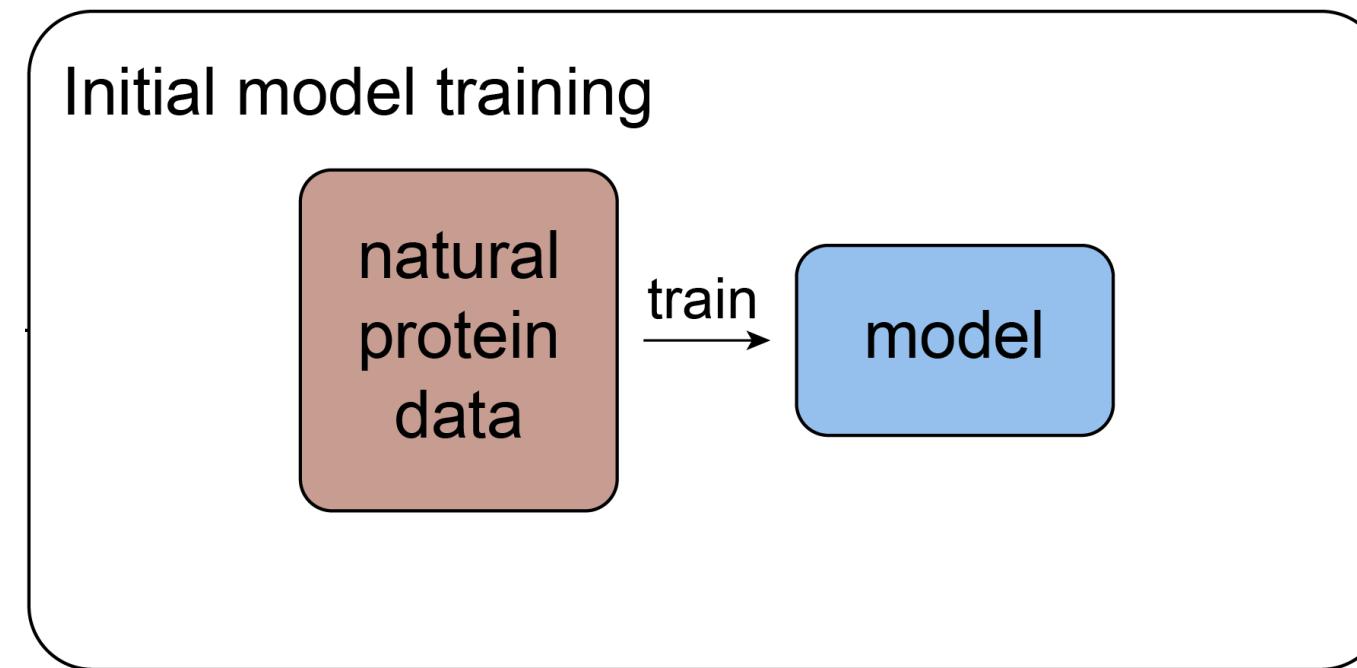
Use multiple data modalities to discover and design proteins



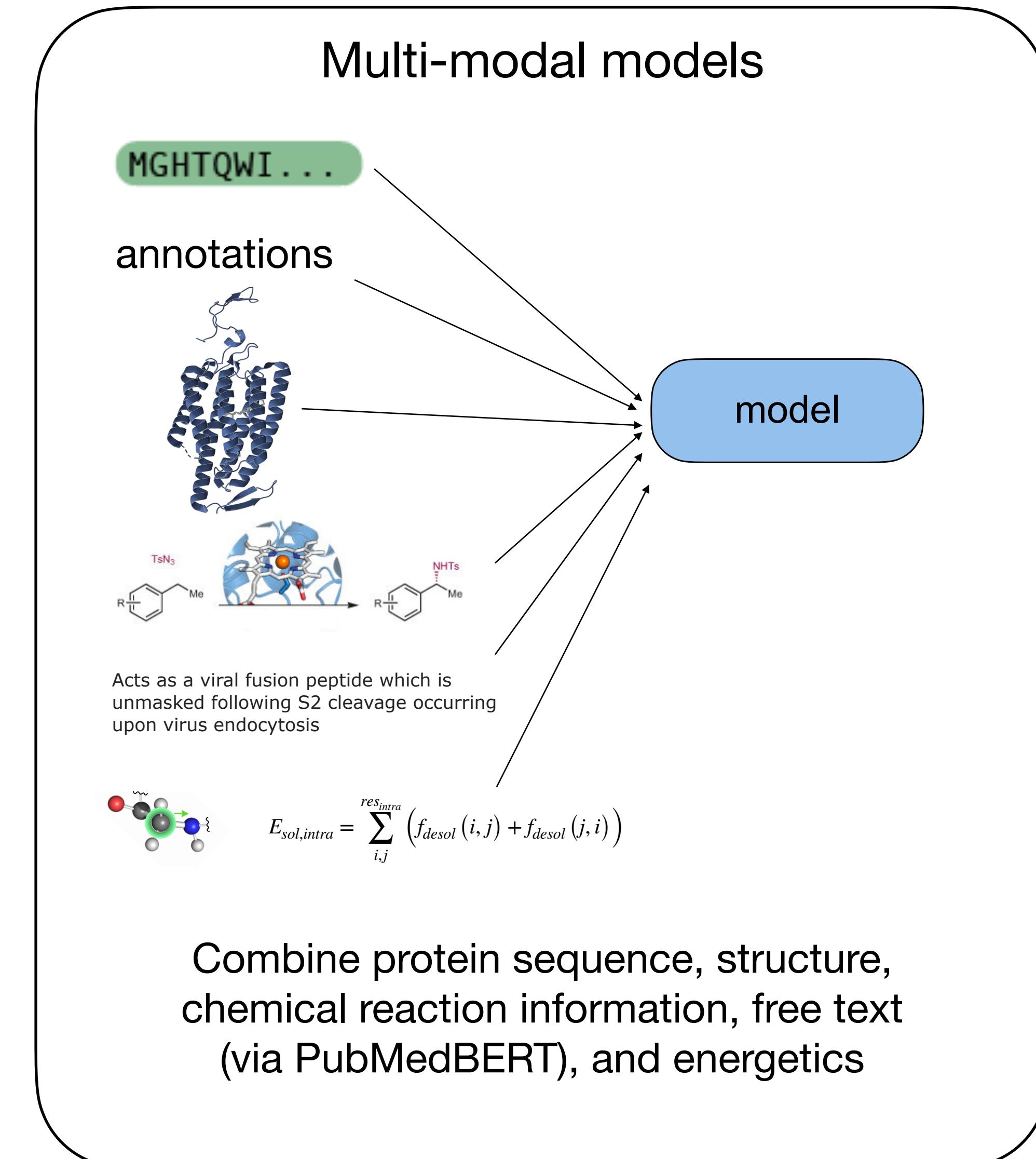
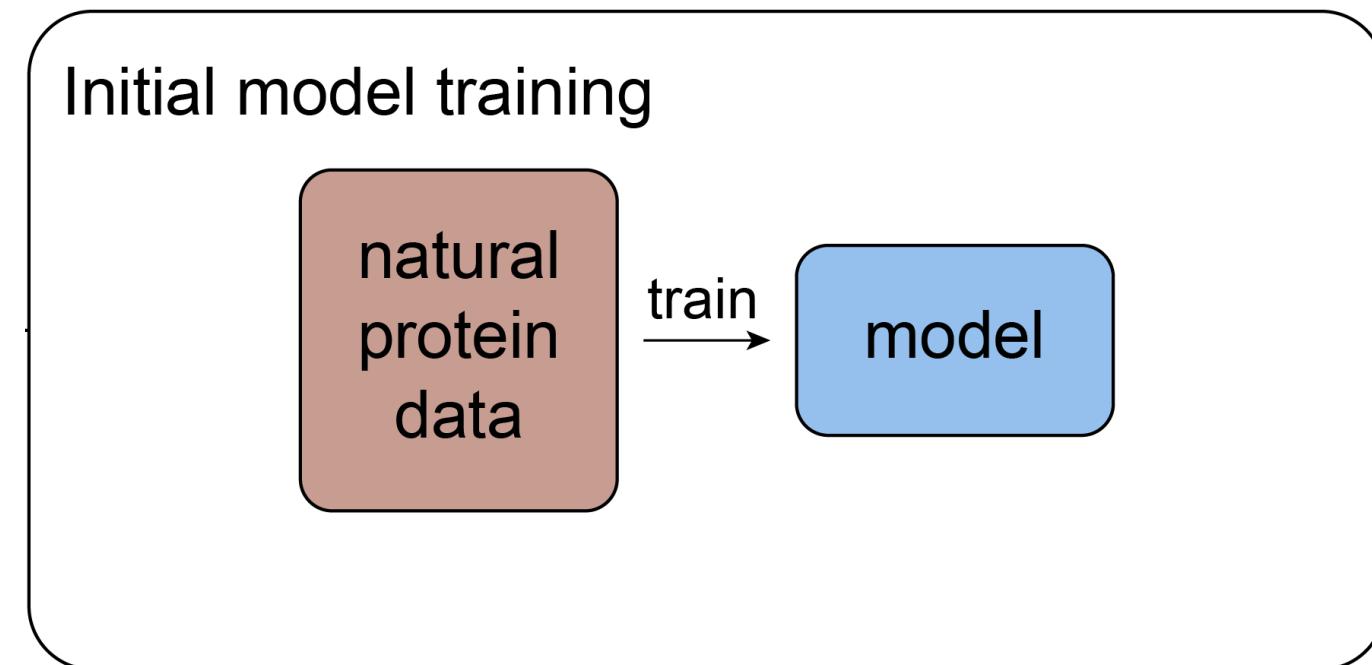
Use multiple data modalities to discover and design proteins



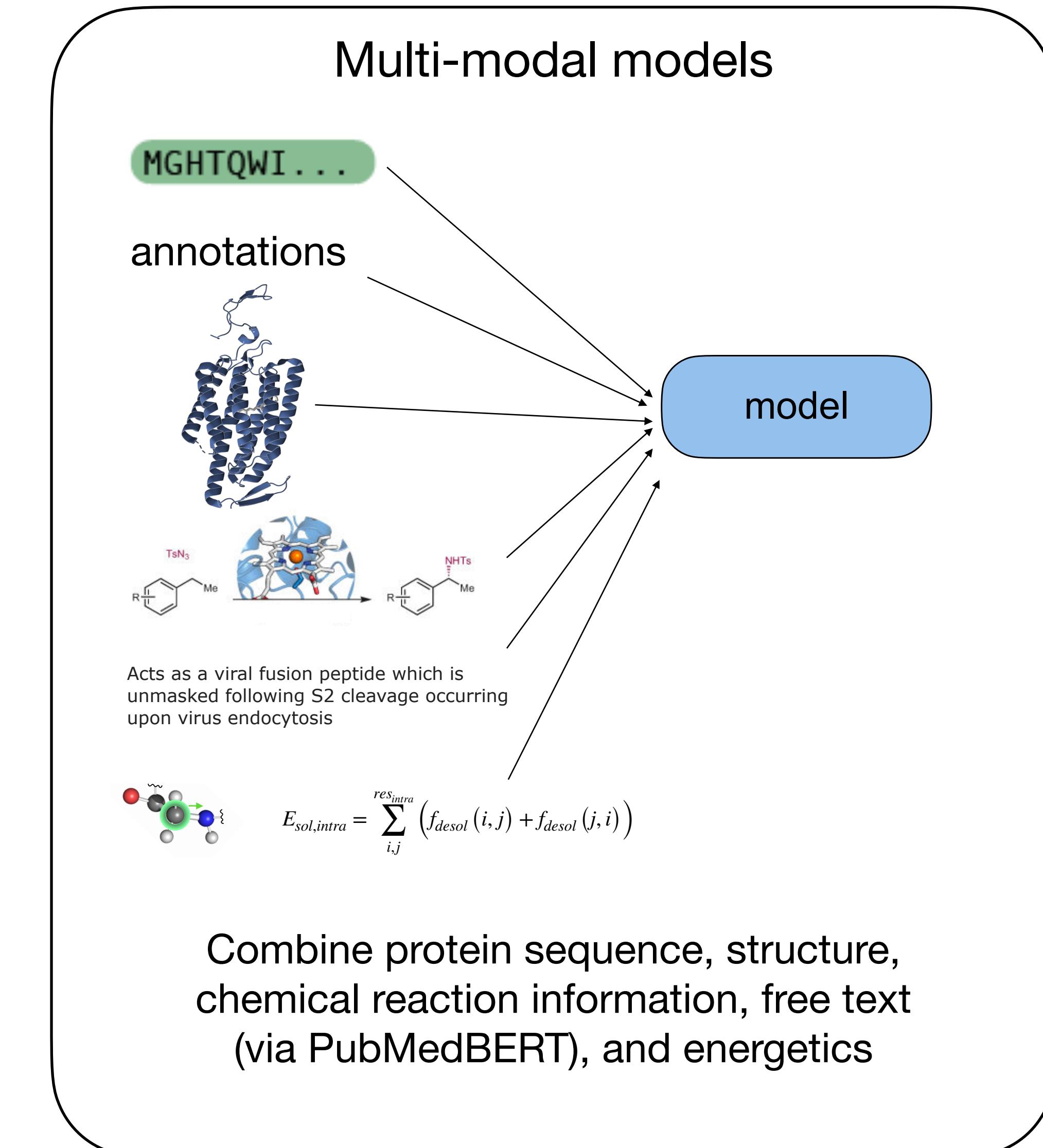
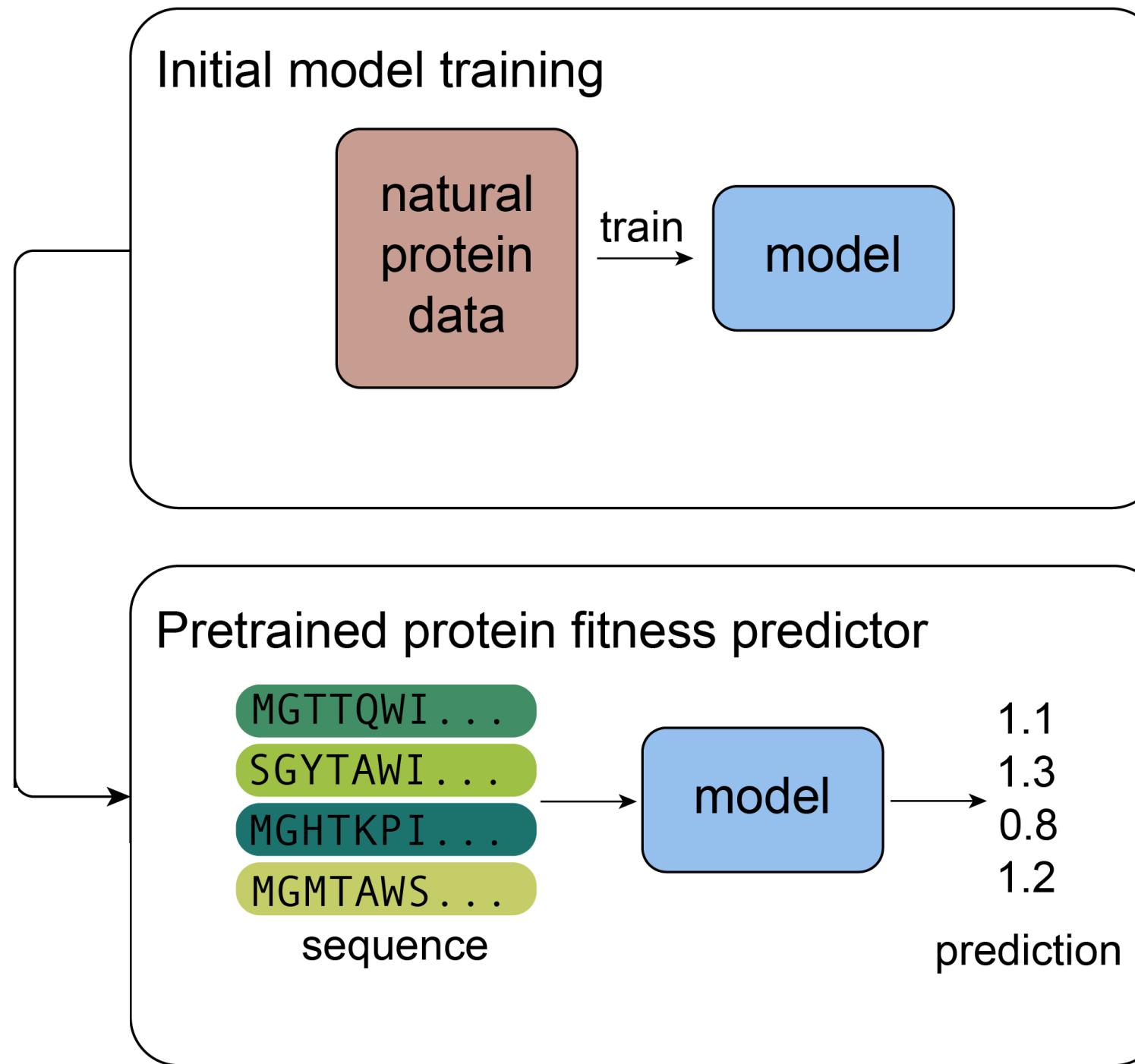
Use multiple data modalities to discover and design proteins



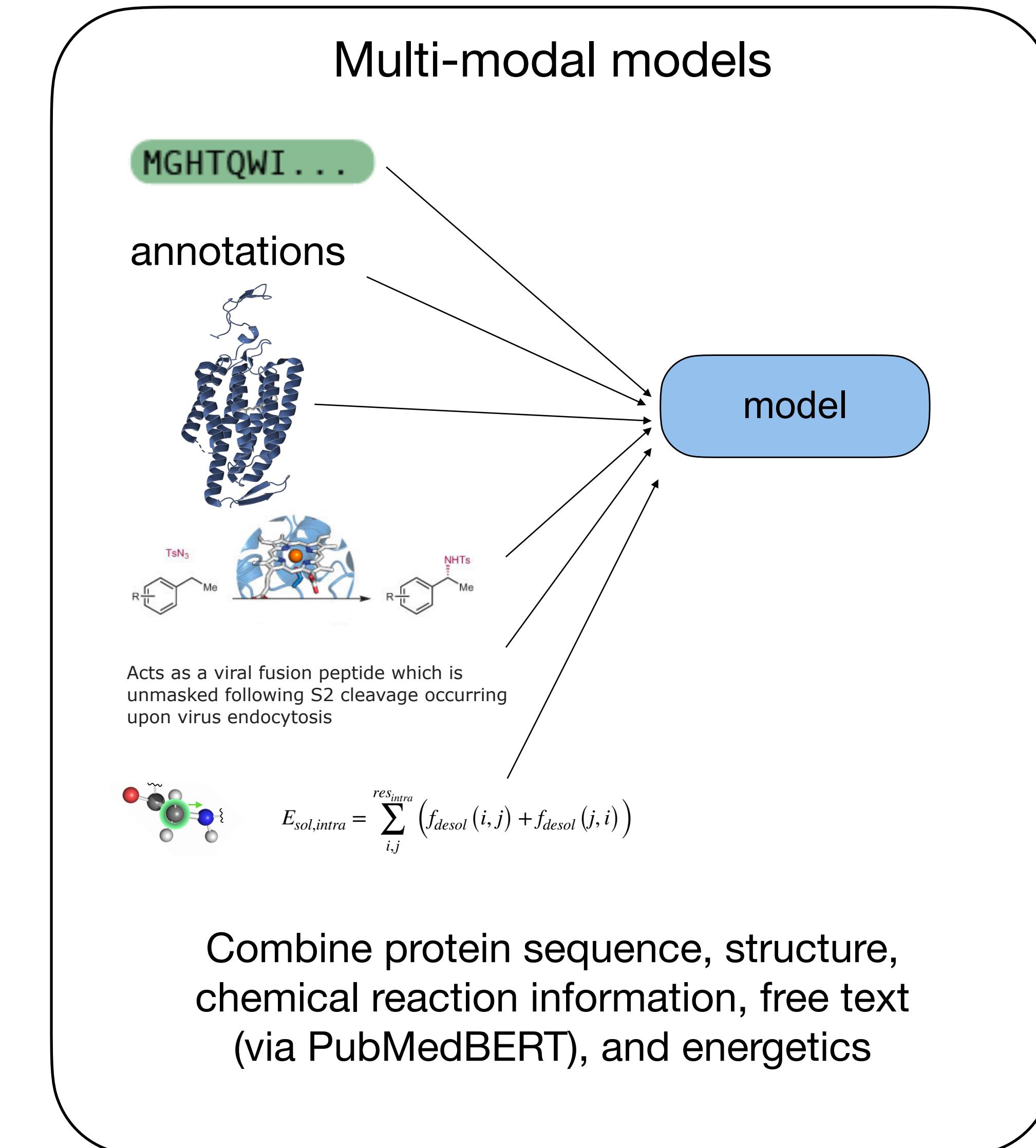
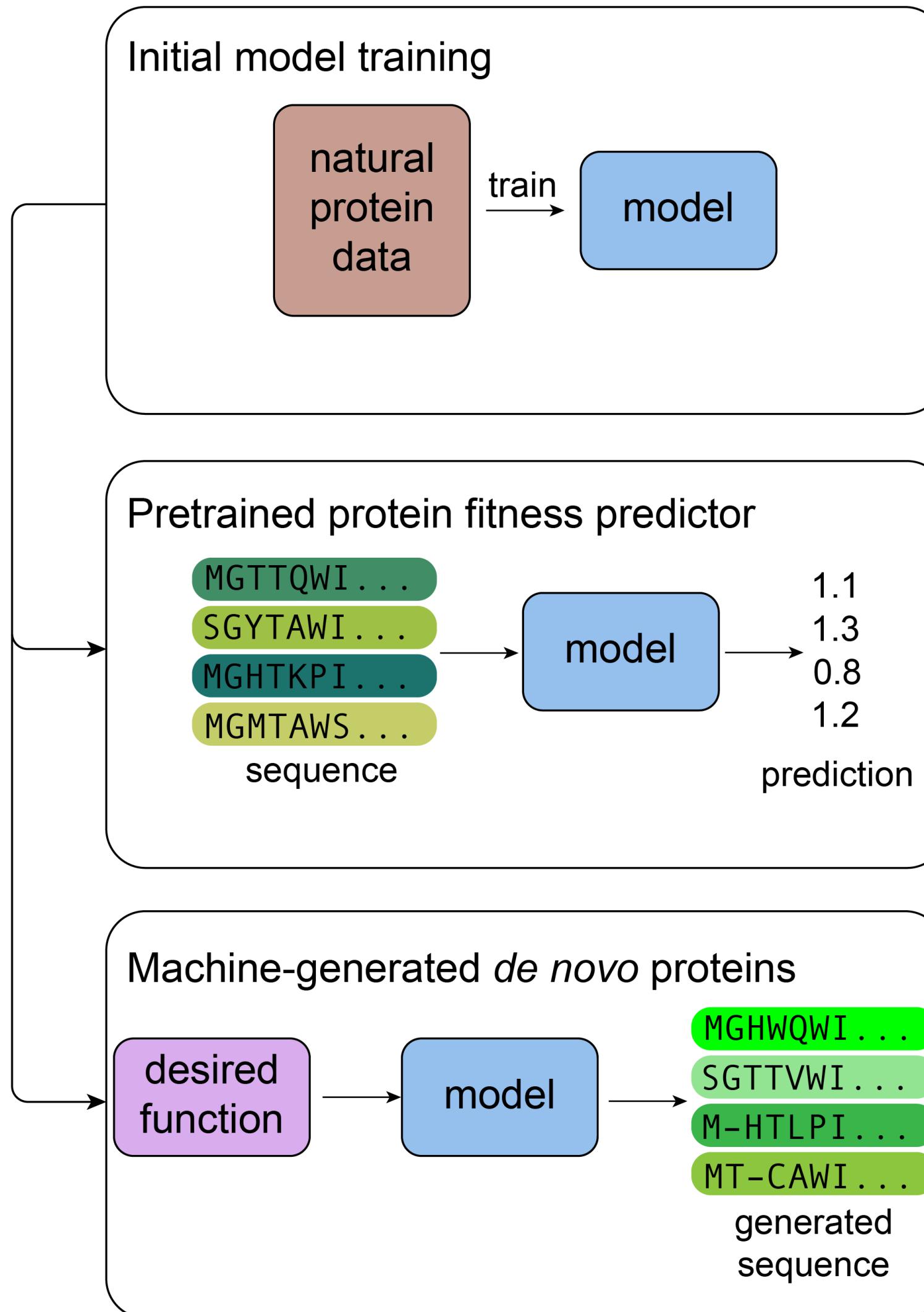
Use multiple data modalities to discover and design proteins



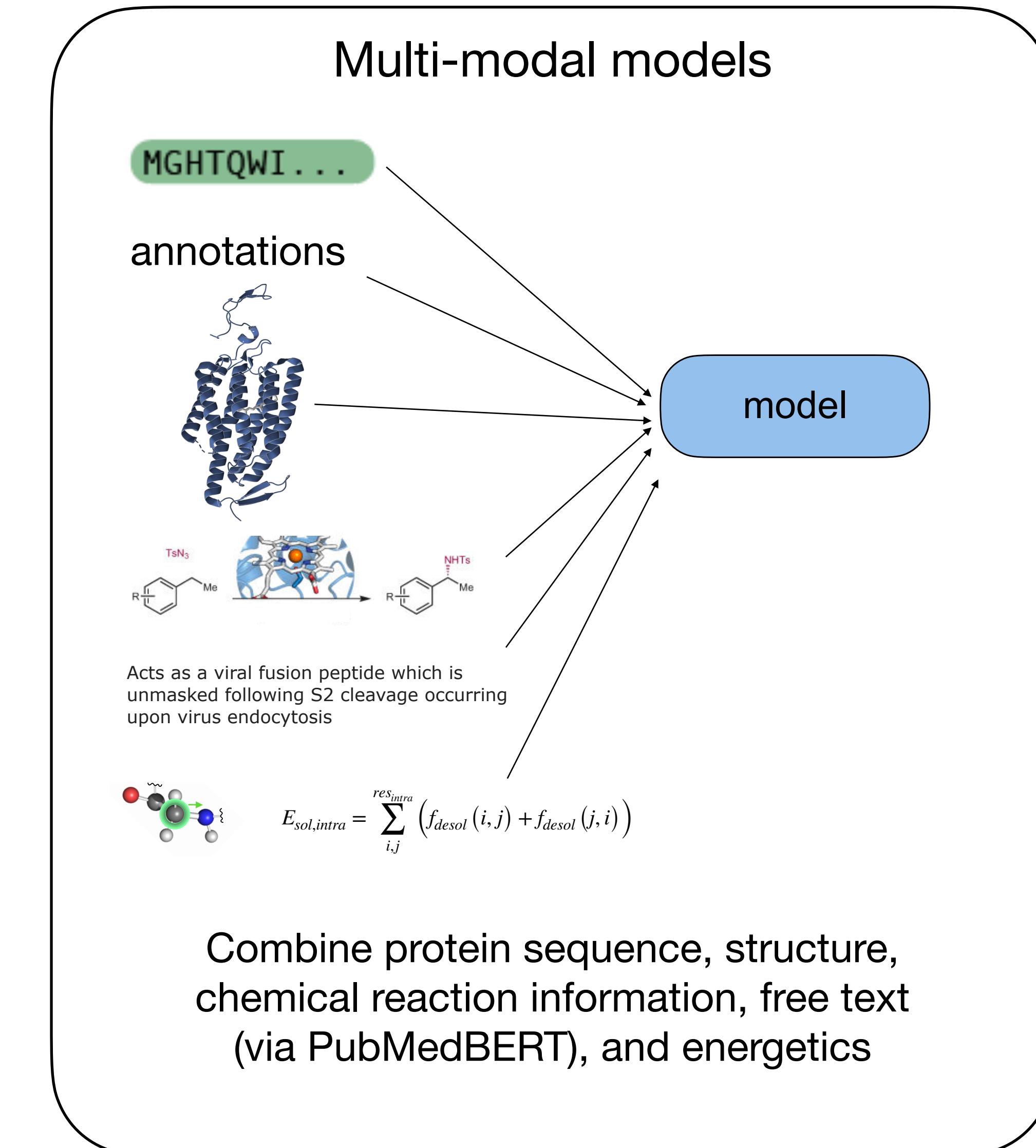
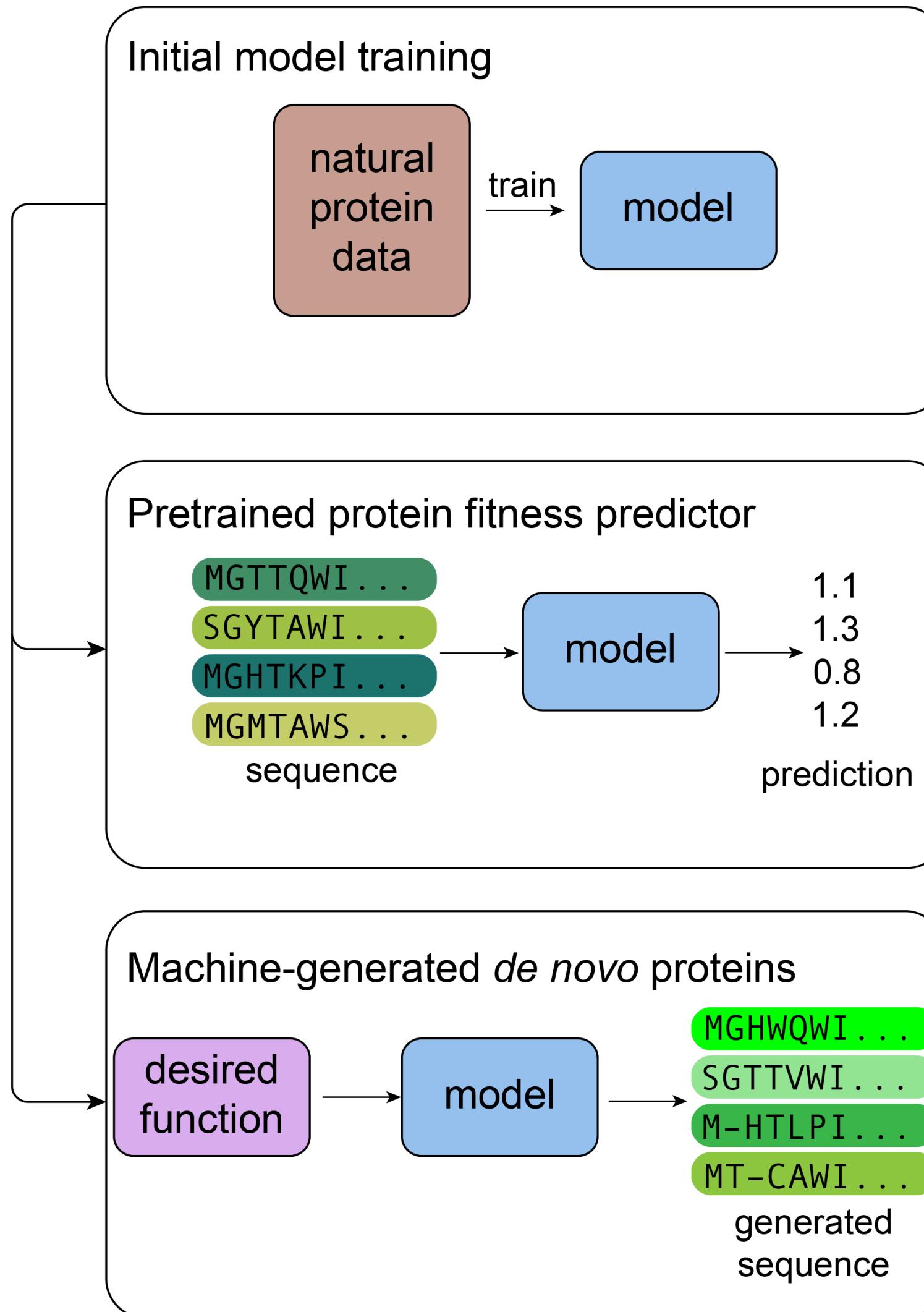
Use multiple data modalities to discover and design proteins



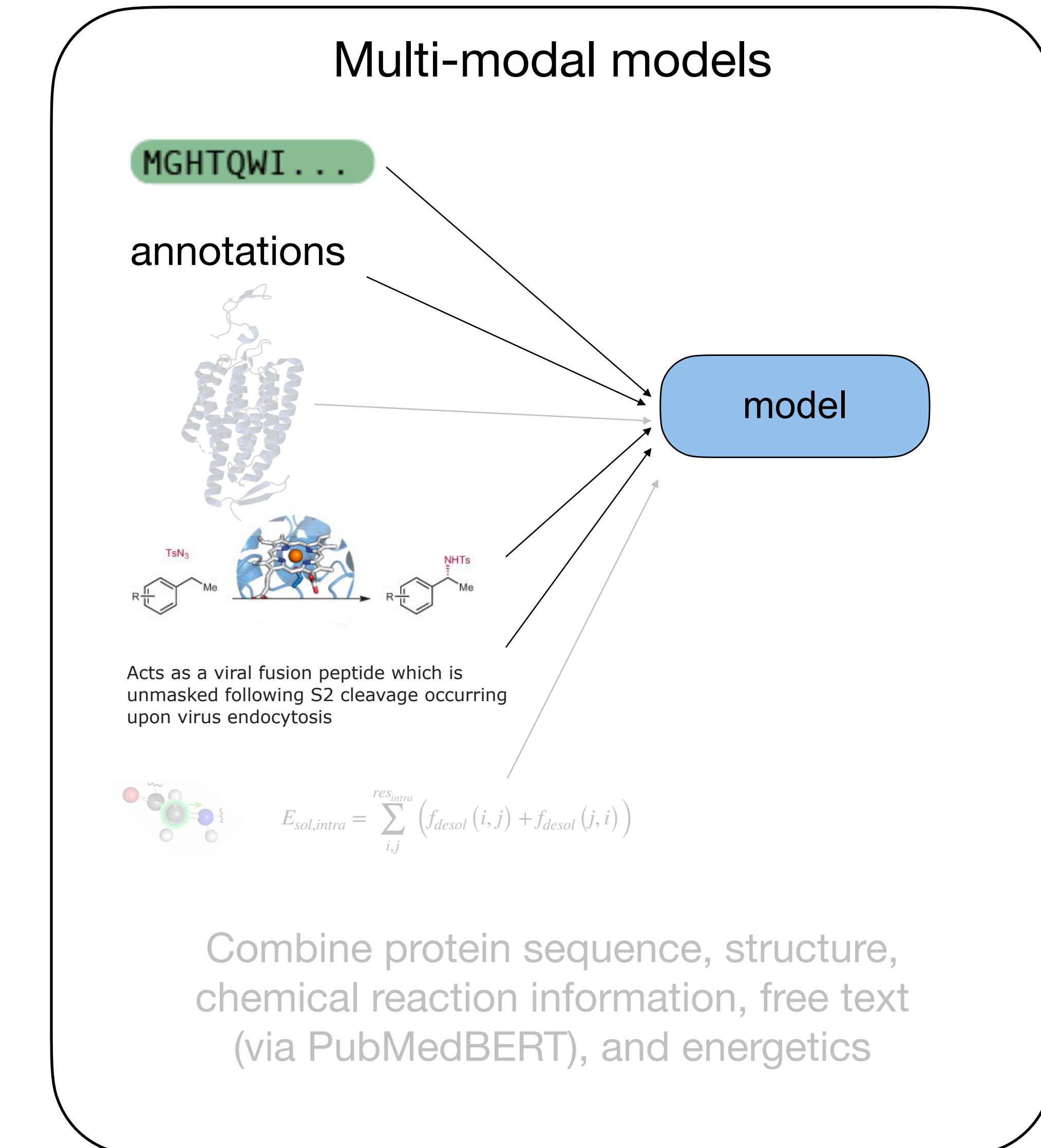
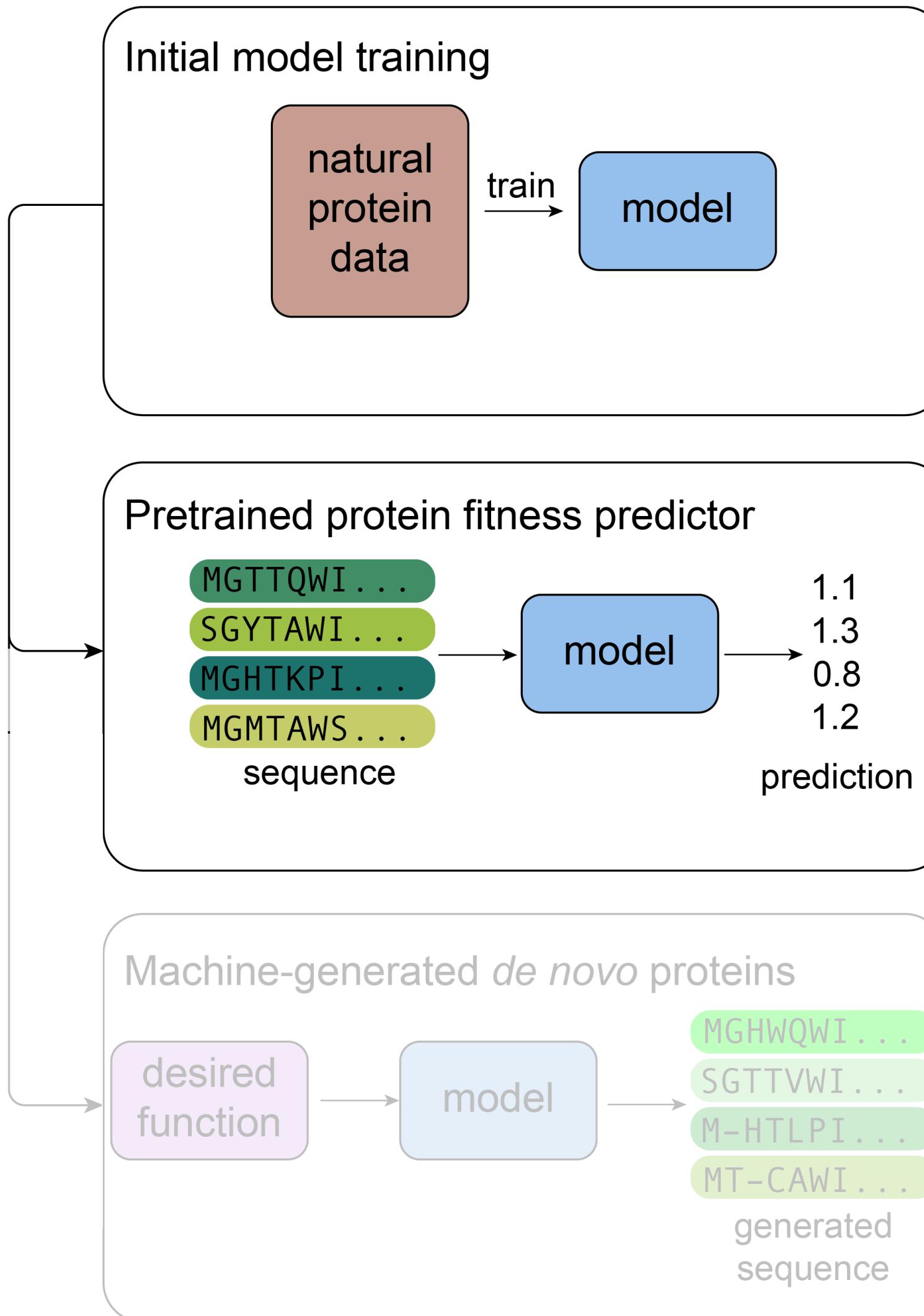
Use multiple data modalities to discover and design proteins



Use multiple data modalities to discover and design proteins

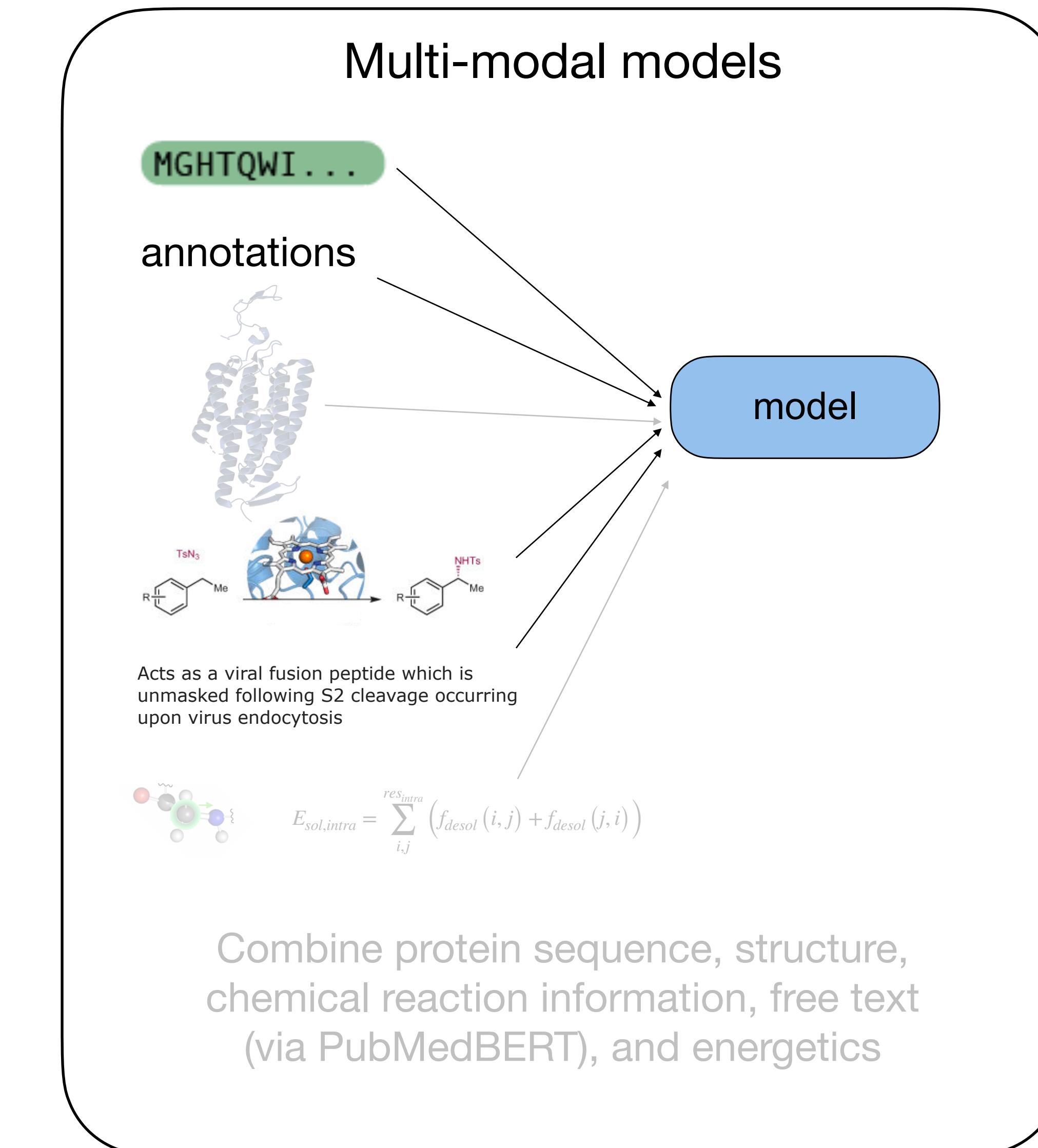
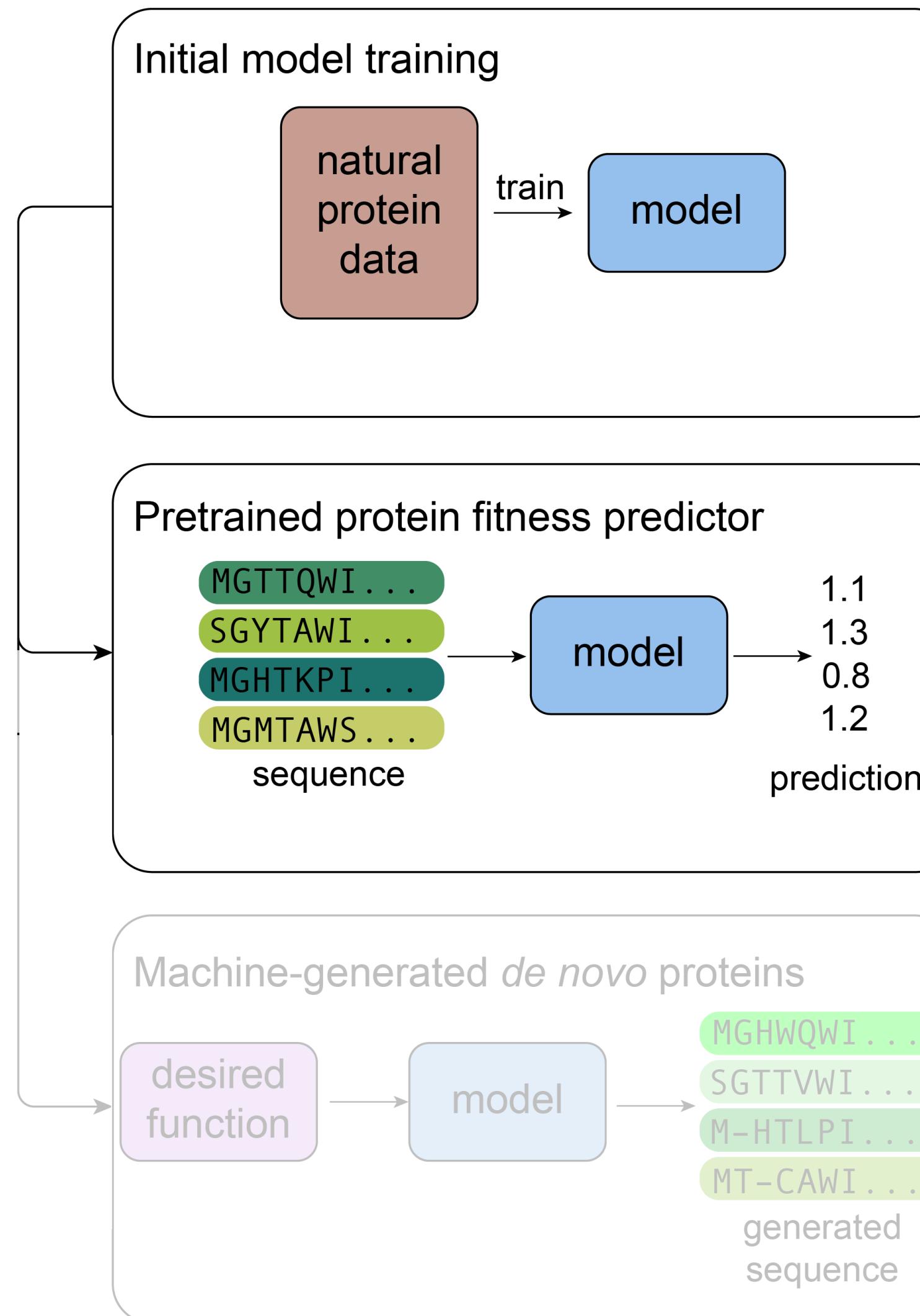


Use multiple data modalities to discover and design proteins

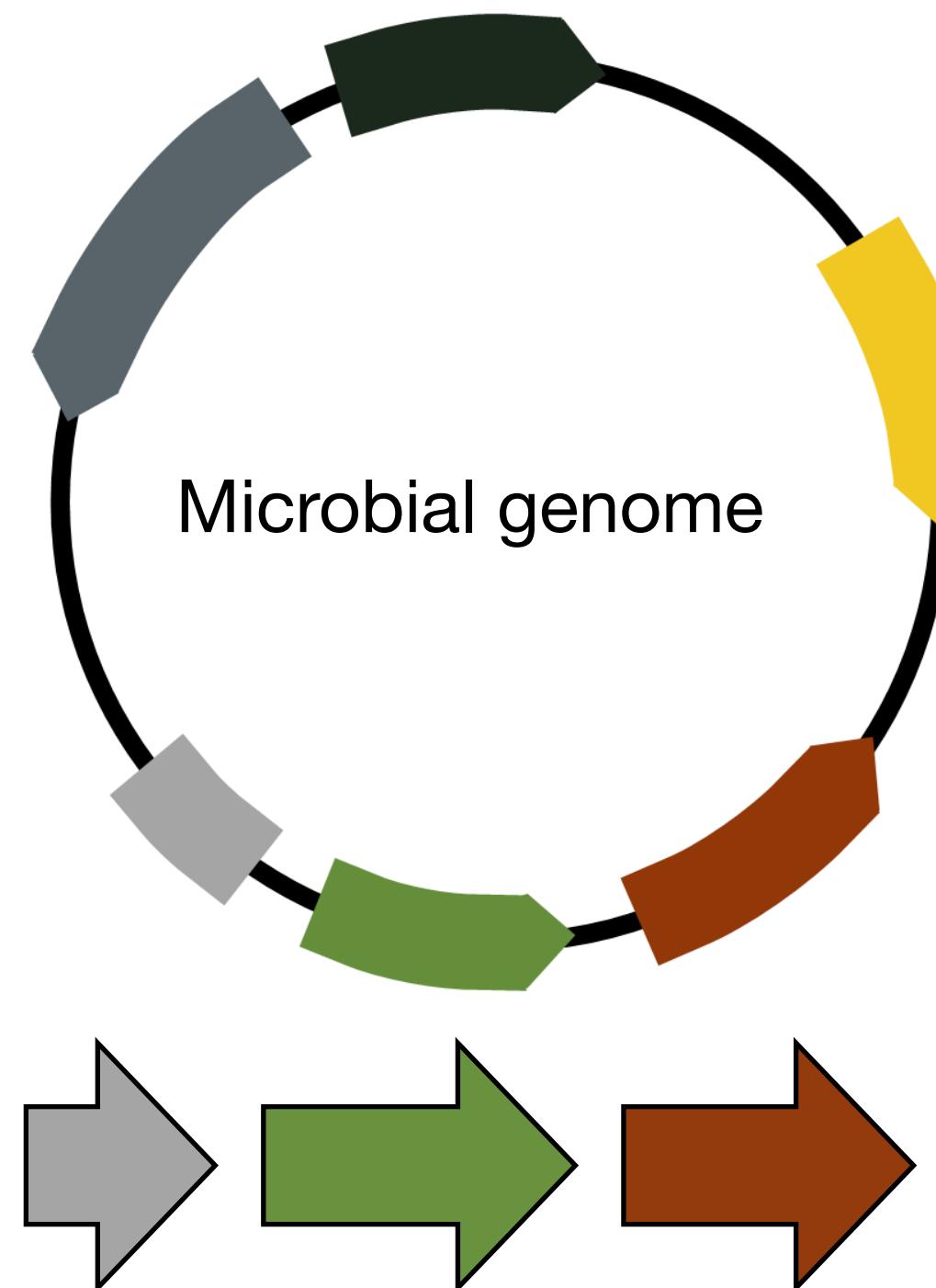


clusters of

Use multiple data modalities to discover and design proteins

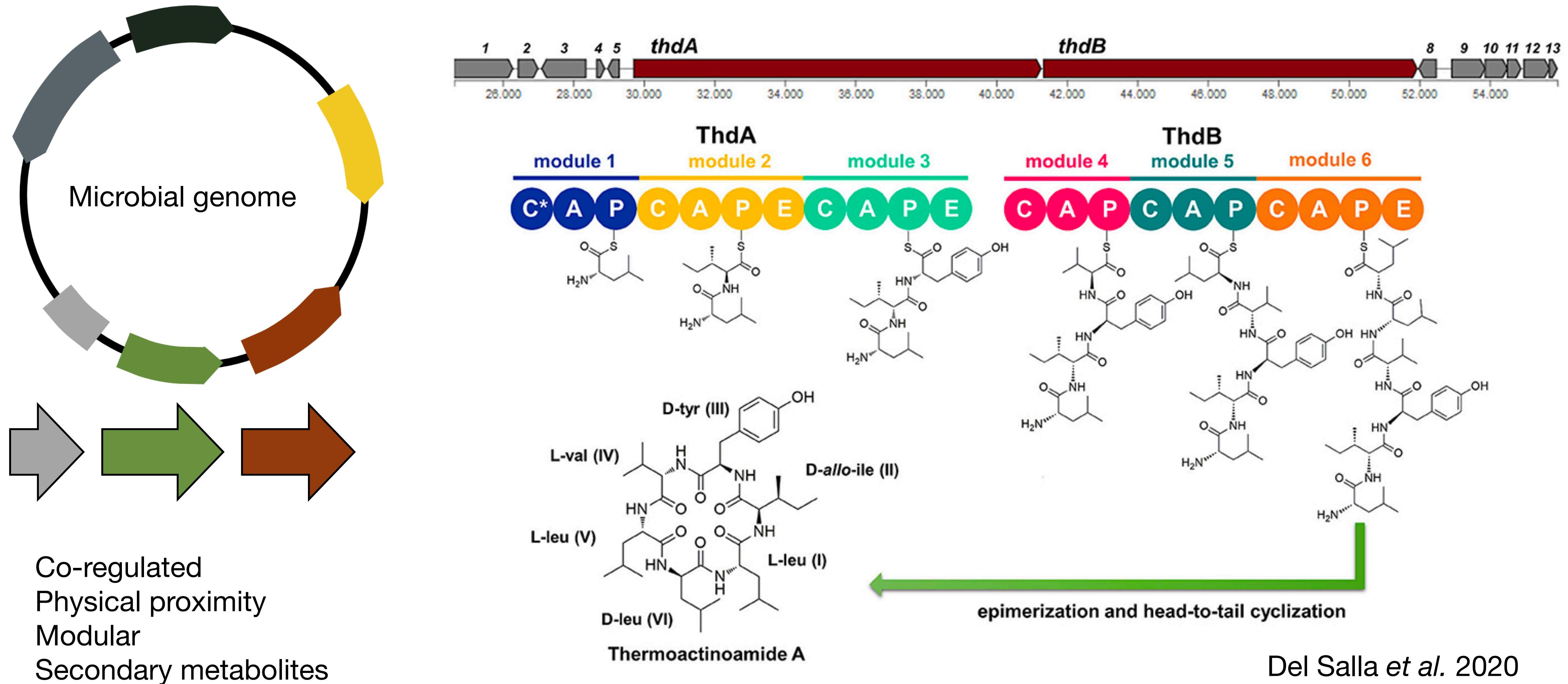


Biosynthetic gene clusters encode natural products



Co-regulated
Physical proximity
Modular
Secondary metabolites

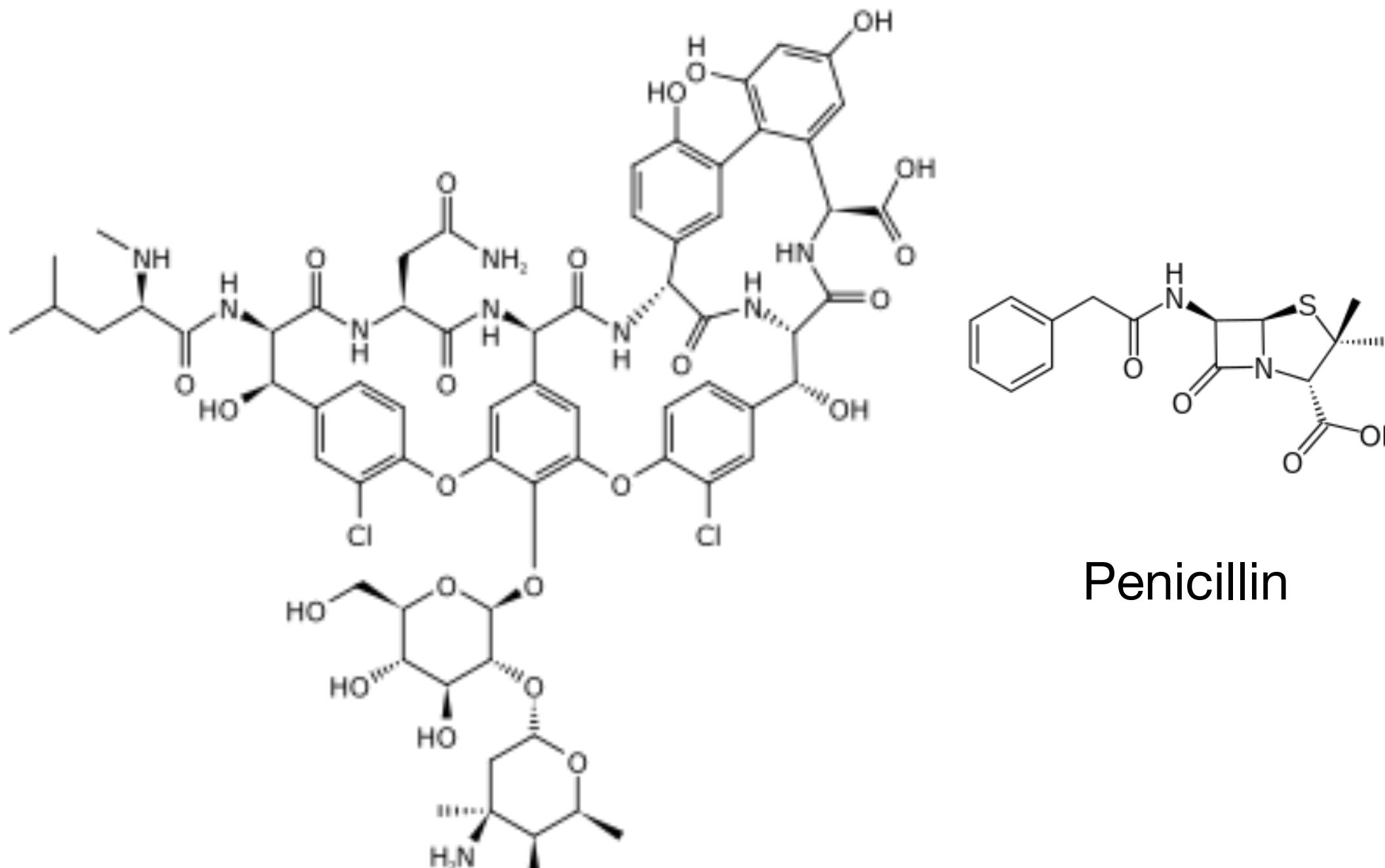
Biosynthetic gene clusters encode natural products



Natural products are sources of pharmaceuticals

Natural products are sources of pharmaceuticals

Antibiotics

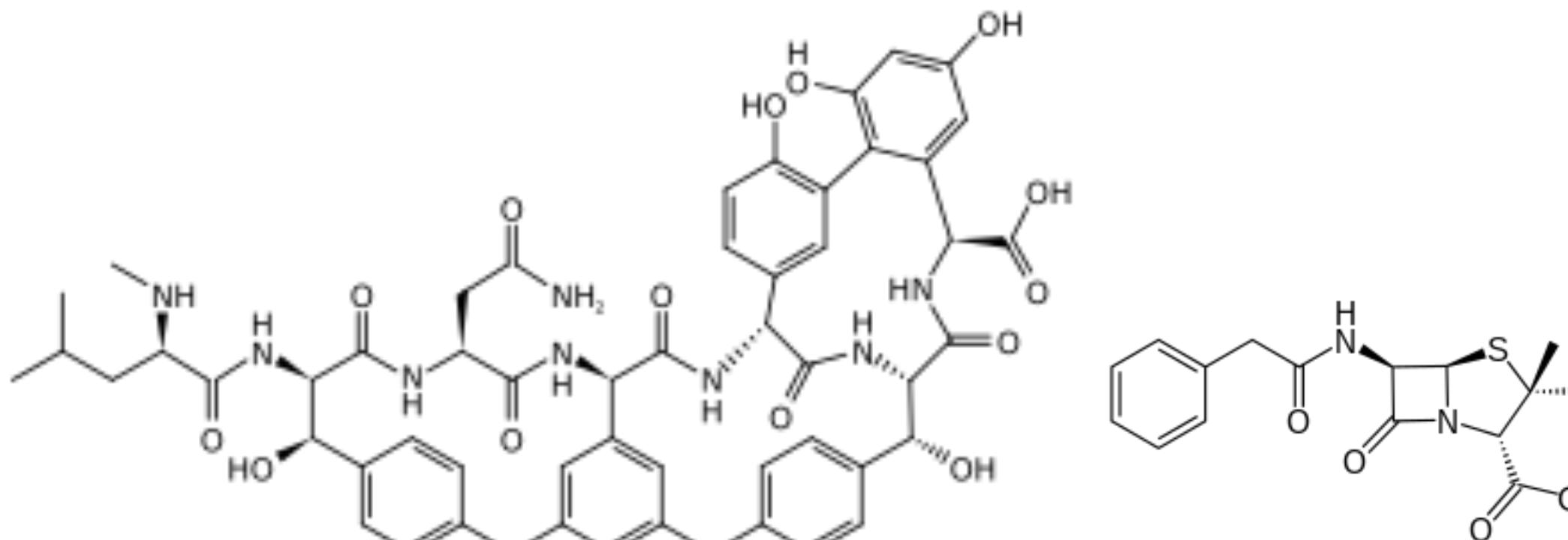


Penicillin

Vancomycin

Natural products are sources of pharmaceuticals

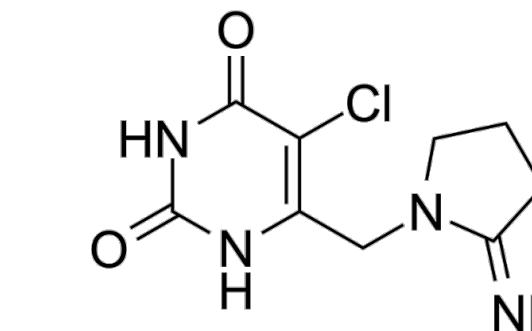
Antibiotics



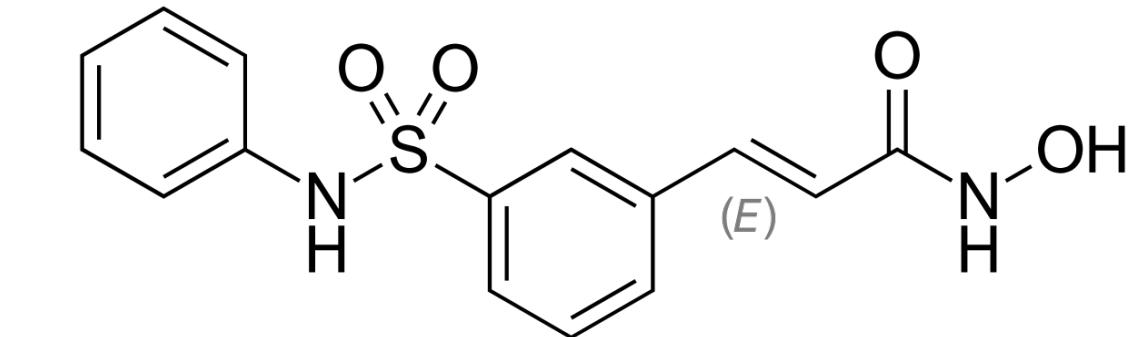
Vancomycin

Penicillin

Anti-Cancer



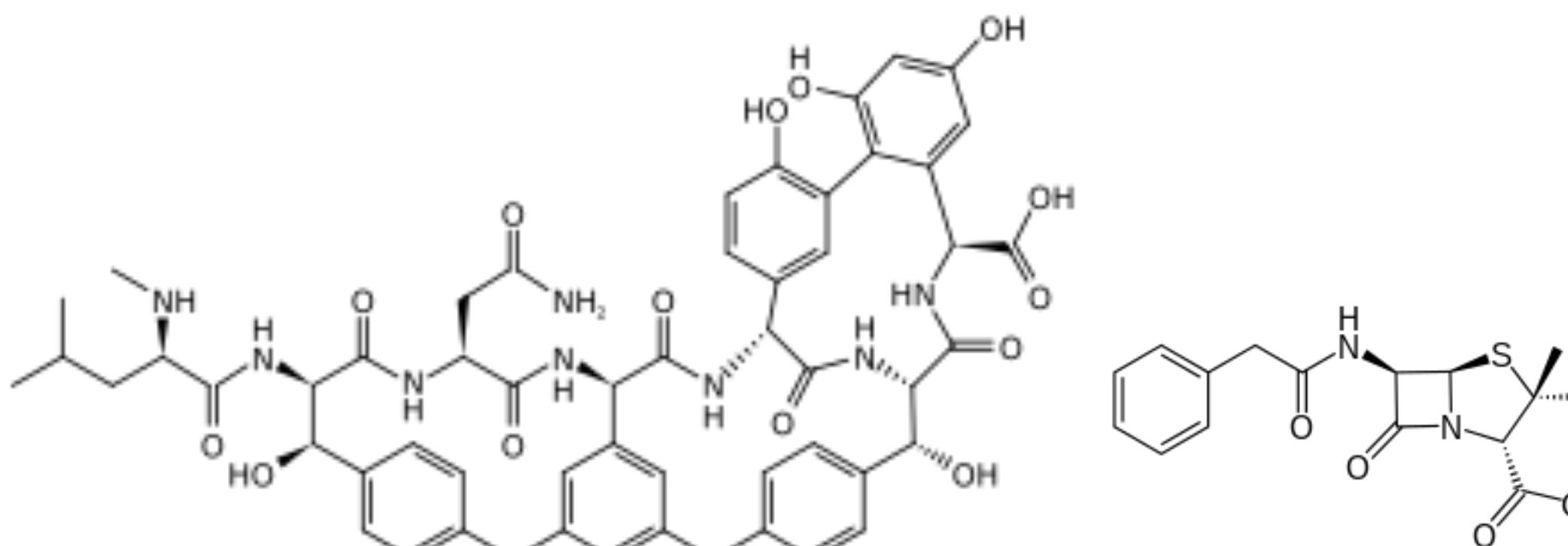
Tipiracil



Belinostat

Natural products are sources of pharmaceuticals

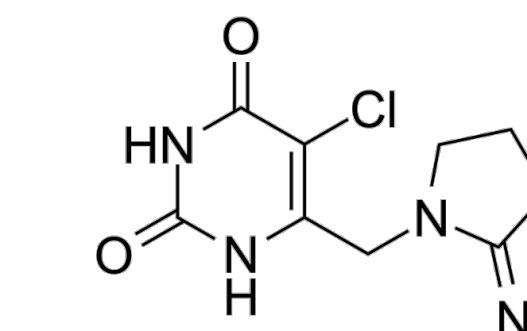
Antibiotics



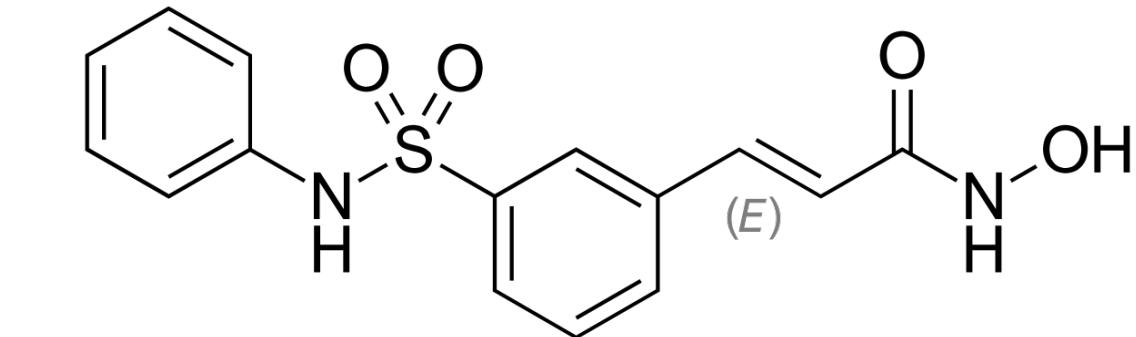
Vancomycin

Penicillin

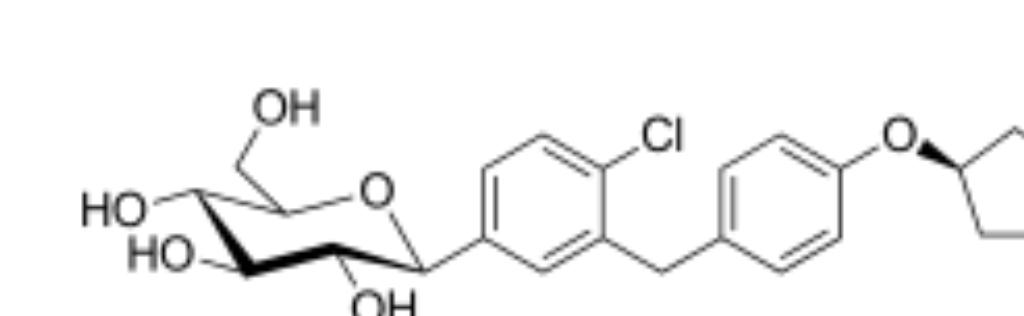
Anti-Cancer



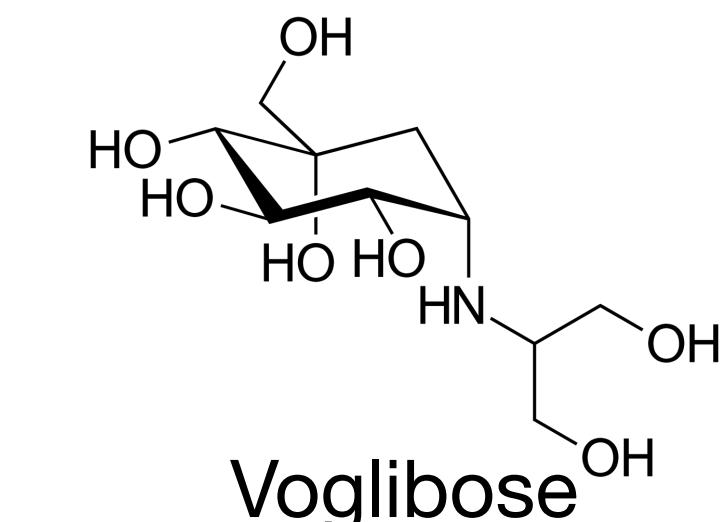
Tipiracil



Belinostat

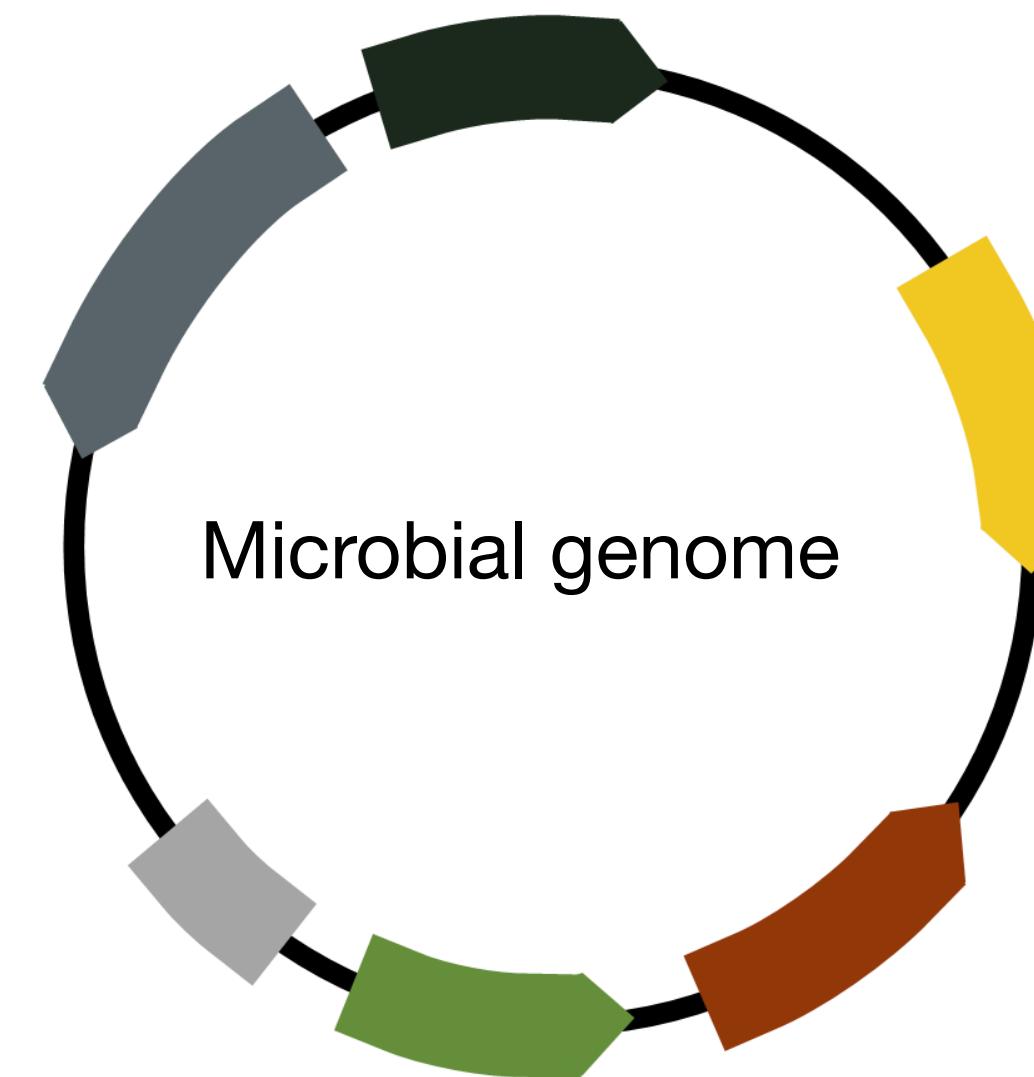


Darunavir



Voglibose

BGCs are difficult to discover and characterize



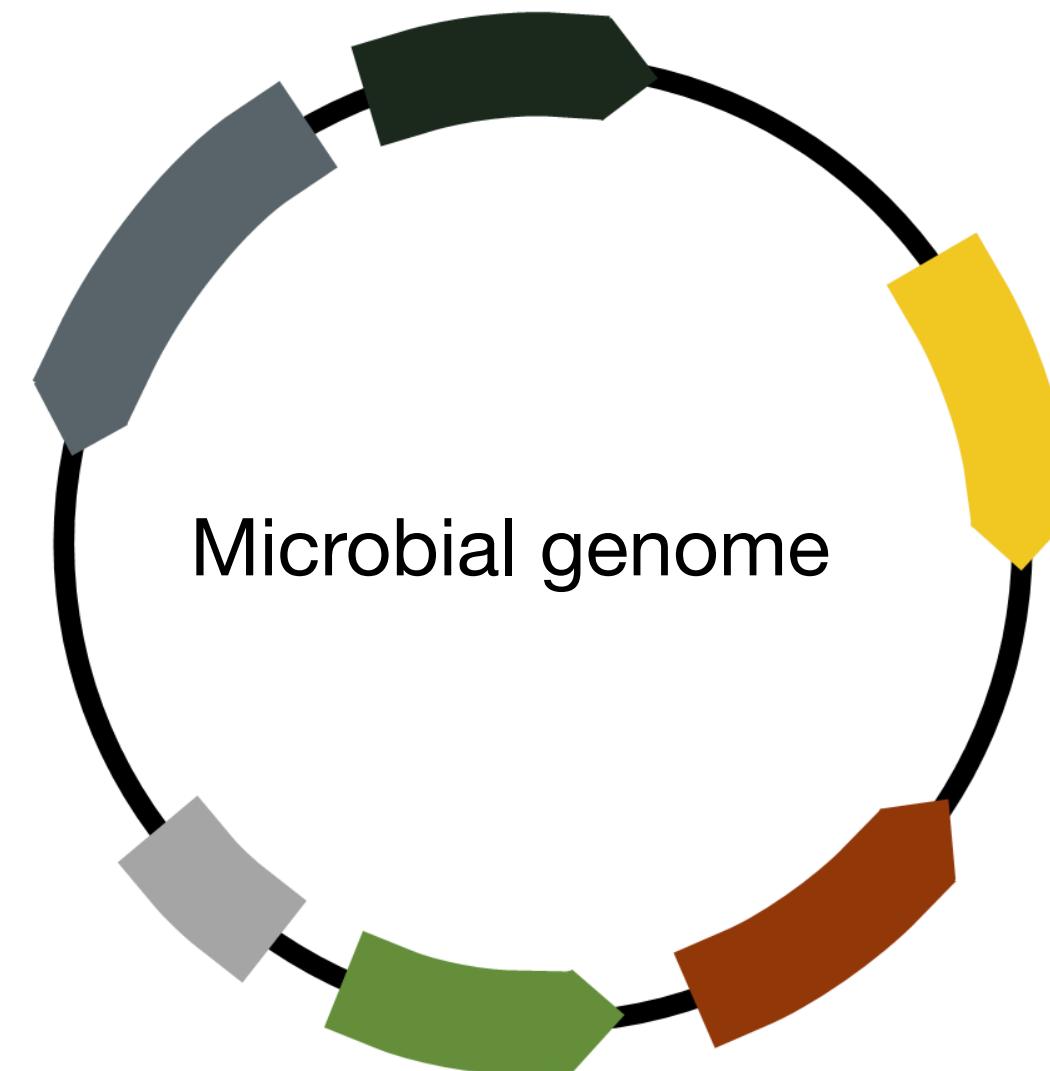
Questions

Which genes?
Where are the boundaries?
What does it make?

Challenges

Unculturable
Map product to genes?

BGCs are difficult to discover and characterize



Questions

Which genes?
Where are the boundaries?
What does it make?

Challenges

Unculturable
Map product to genes?

Only ~1400 “BGCs of known function” in MIBiG

Many methods pretend proteins are language

MFTGNDAGH

Many methods pretend proteins are language

MFTGNDAGH

Treat amino acids like tokens

Many methods pretend proteins are language

MFTGNDAGH

Treat amino acids like tokens

Use tools originally developed for language

Many methods pretend proteins are language

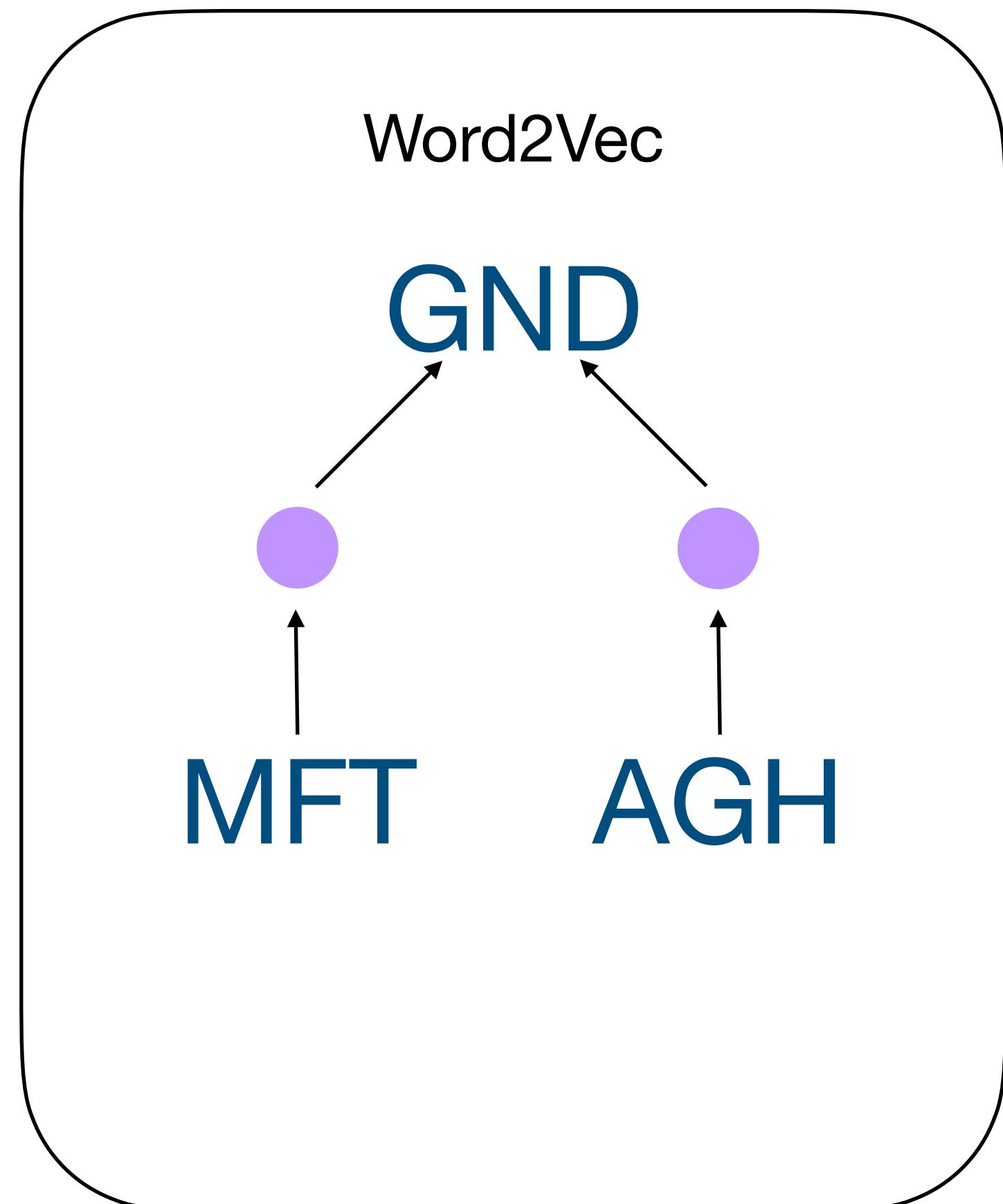
Word2Vec

MFT

AGH

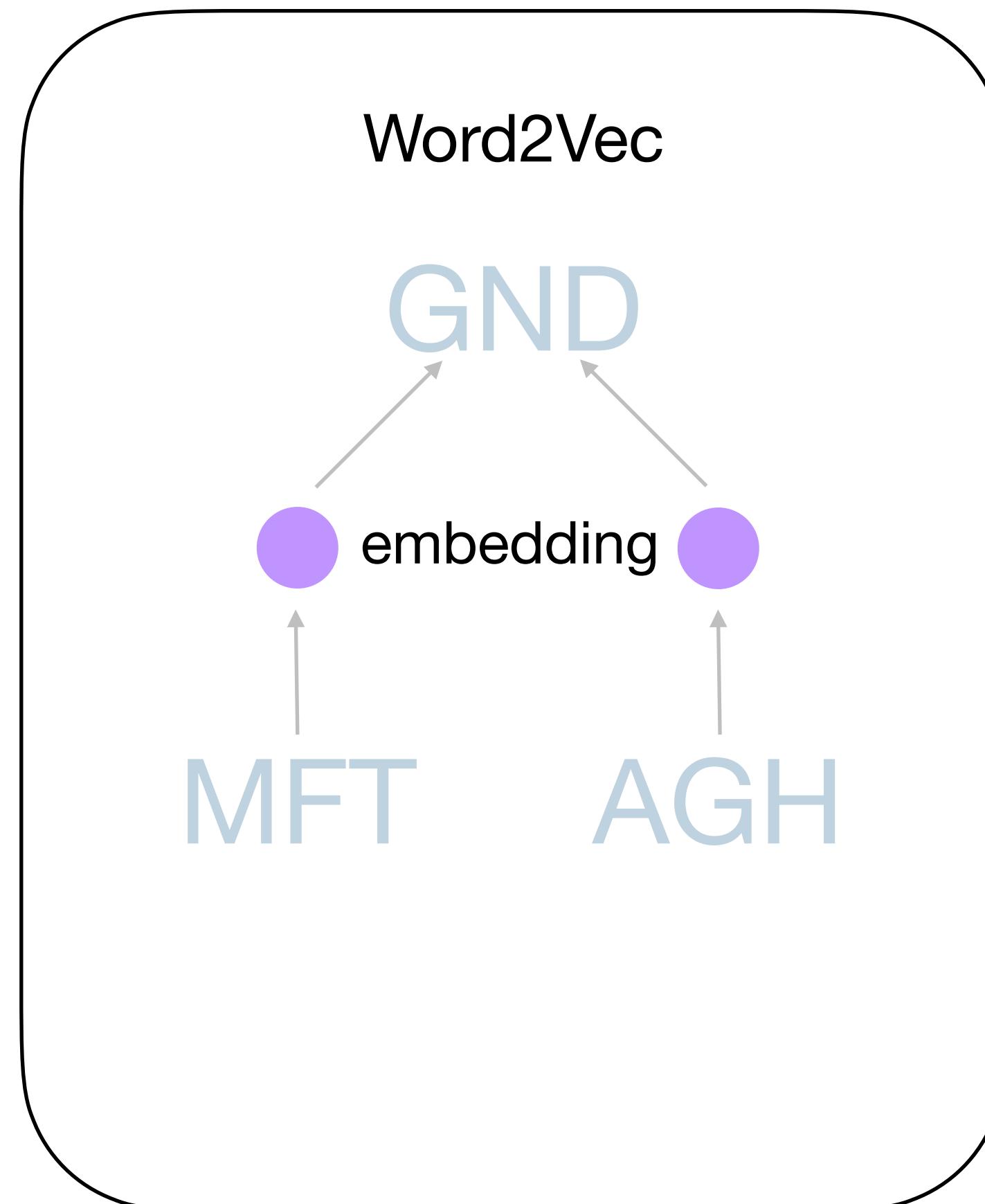
MFTGNDAGH

Many methods pretend proteins are language



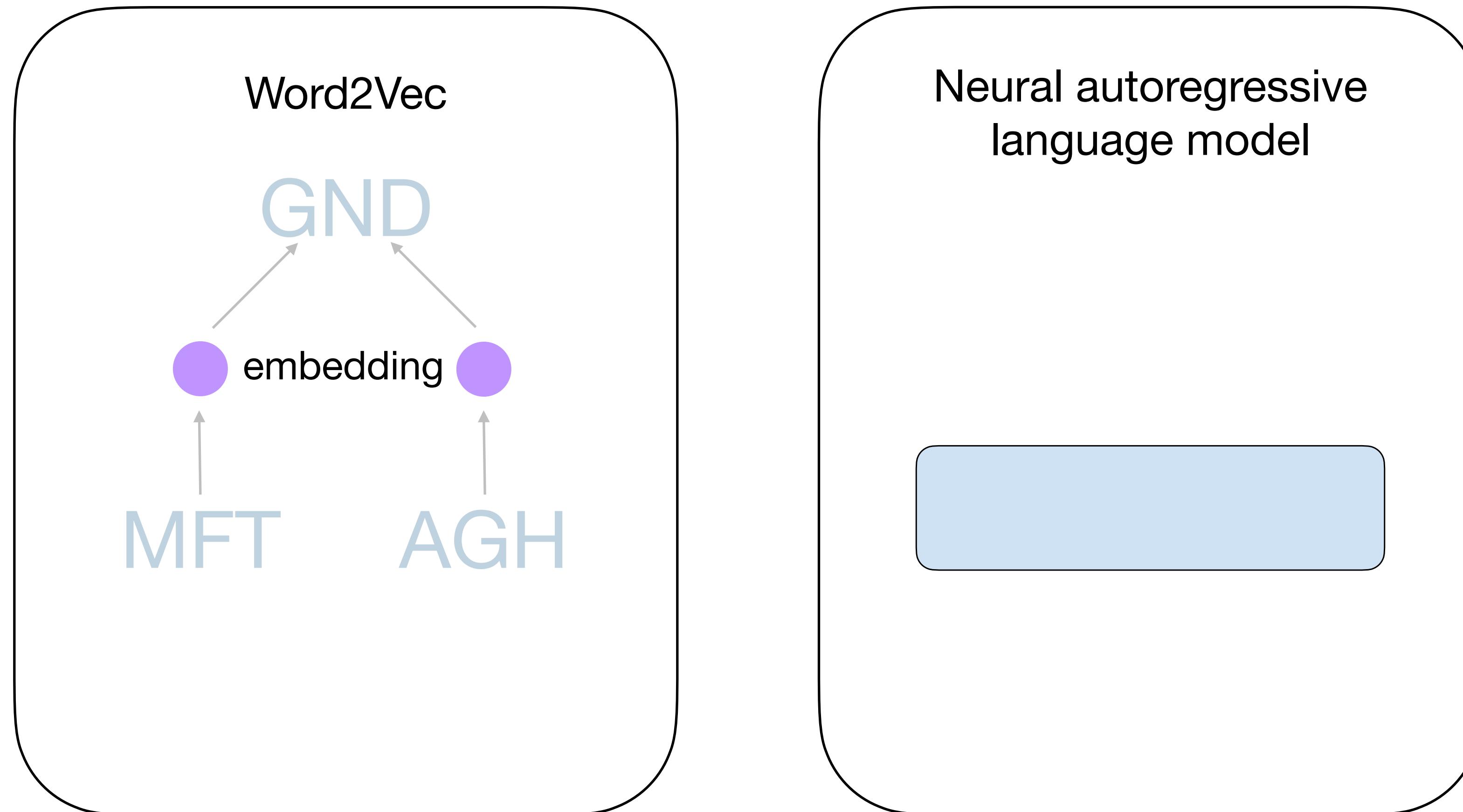
MFTGNDAGH

Many methods pretend proteins are language



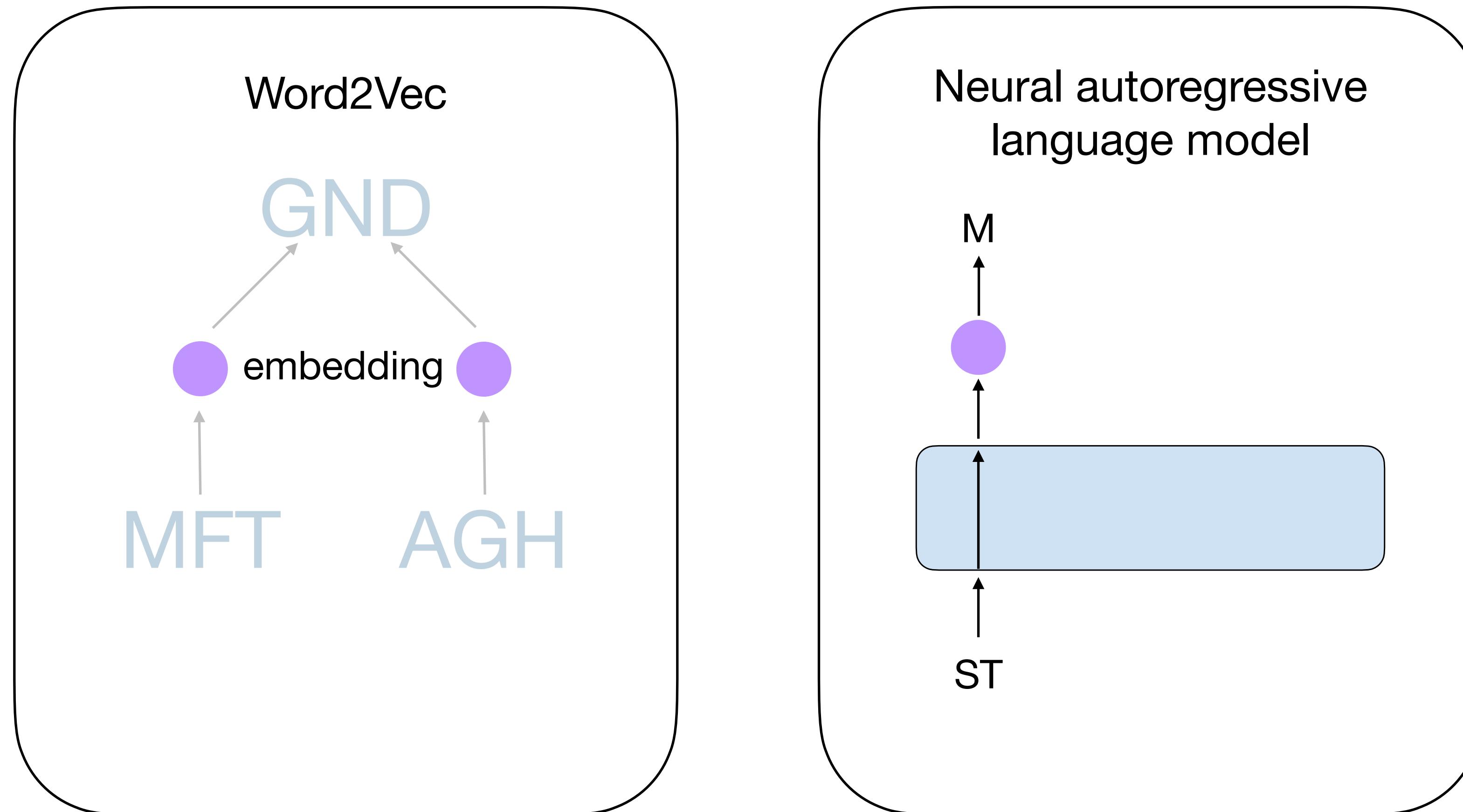
MFTGNDAGH

Many methods pretend proteins are language



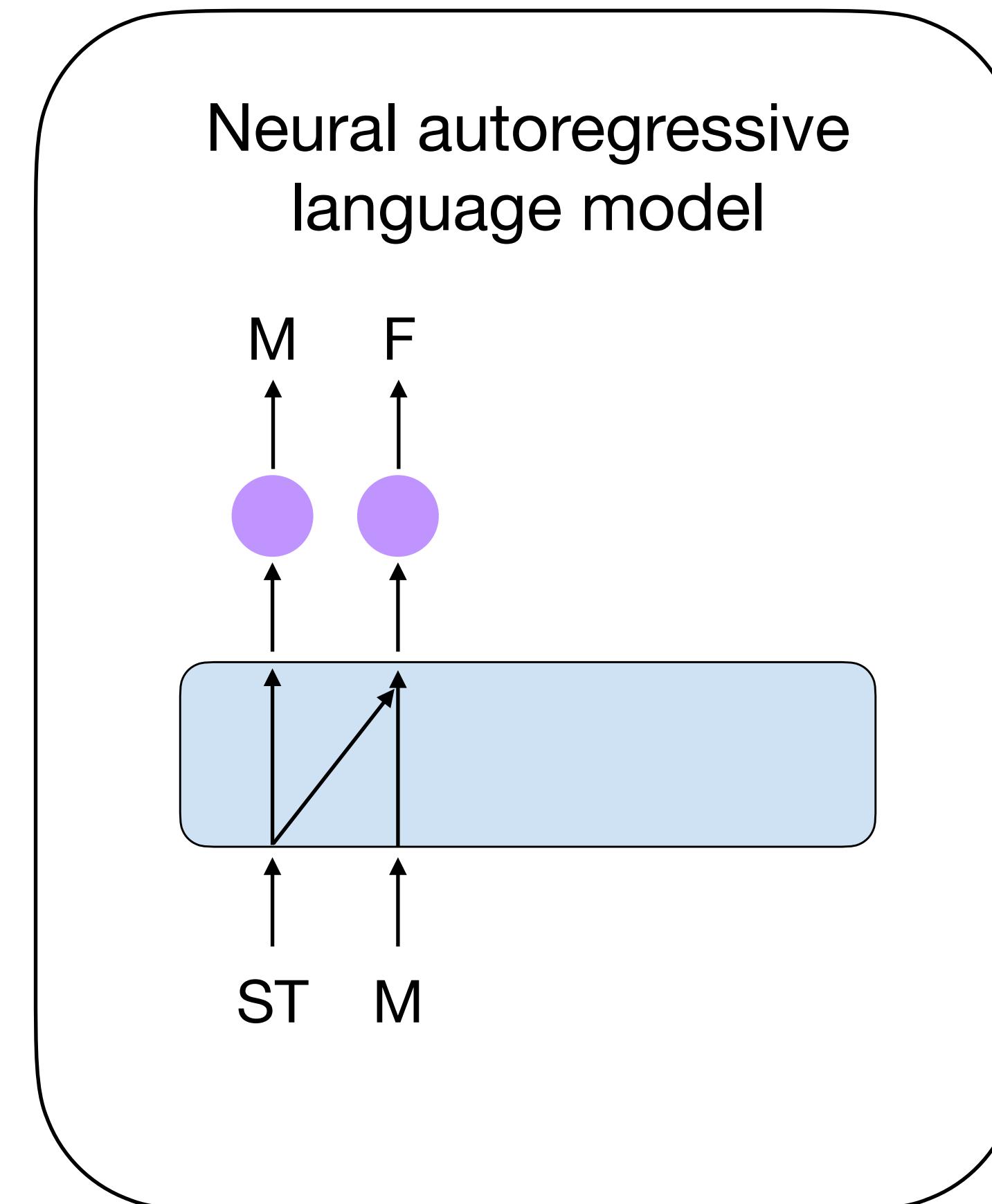
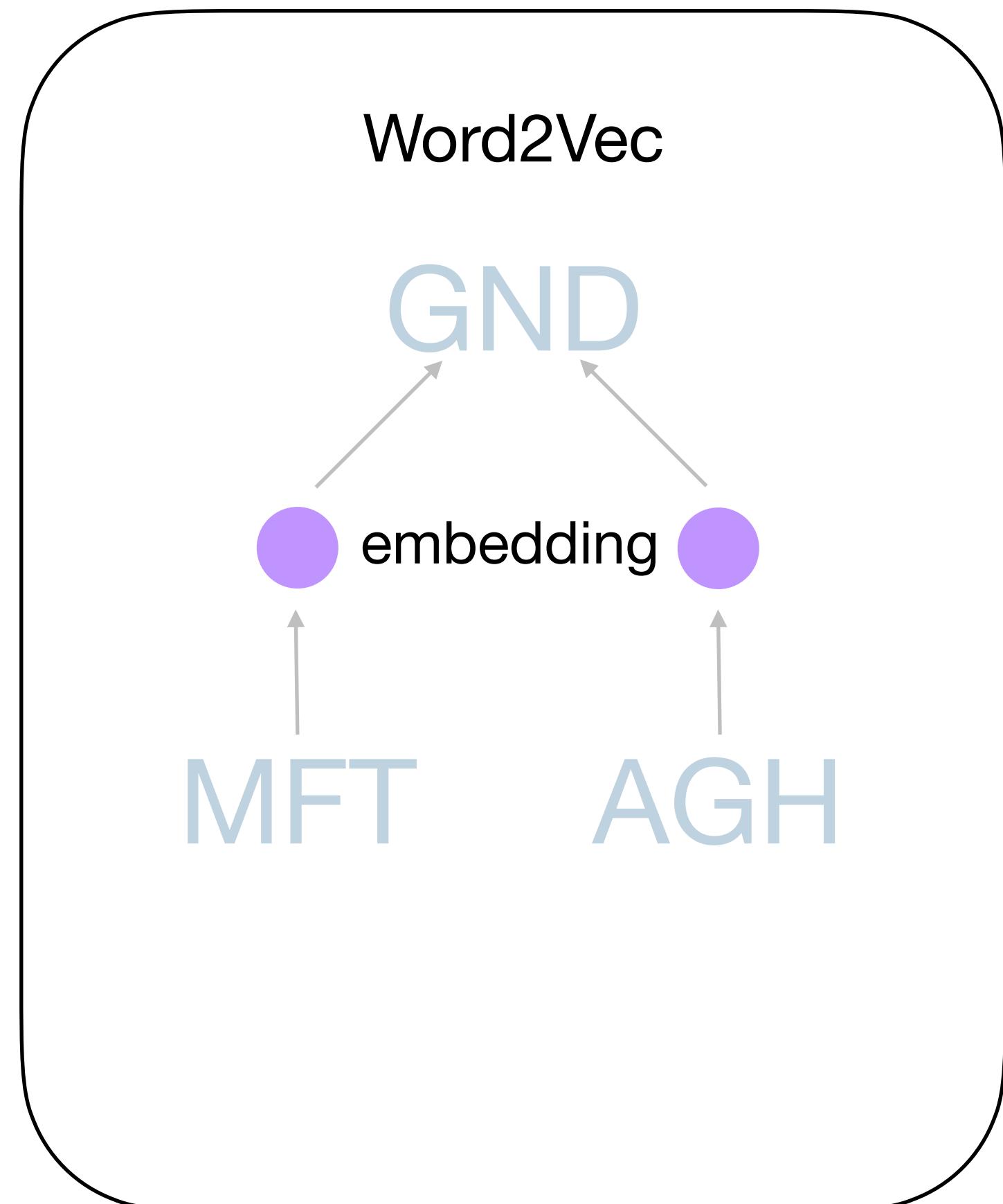
MFTGNDAGH

Many methods pretend proteins are language



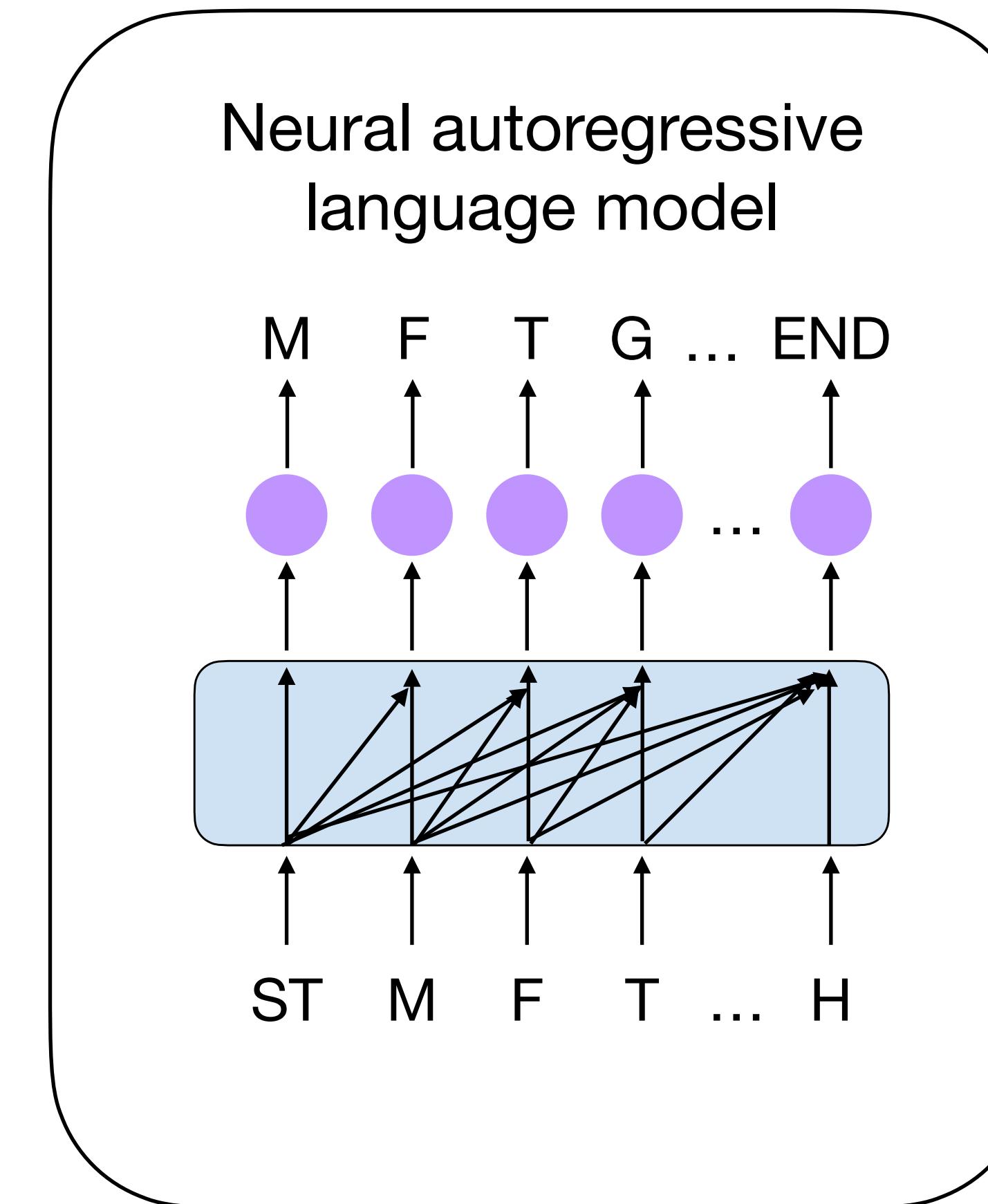
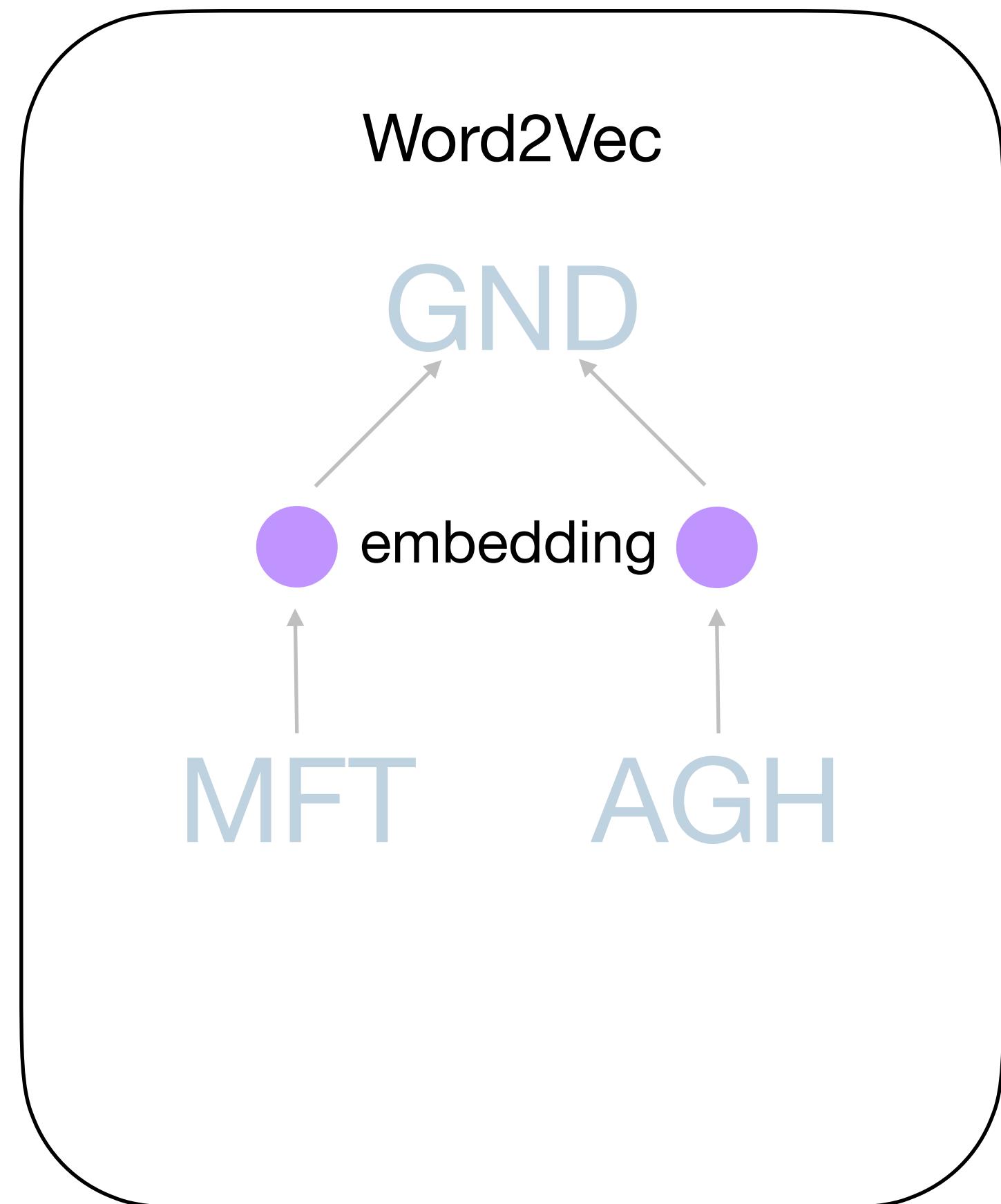
MFTGNDAGH

Many methods pretend proteins are language



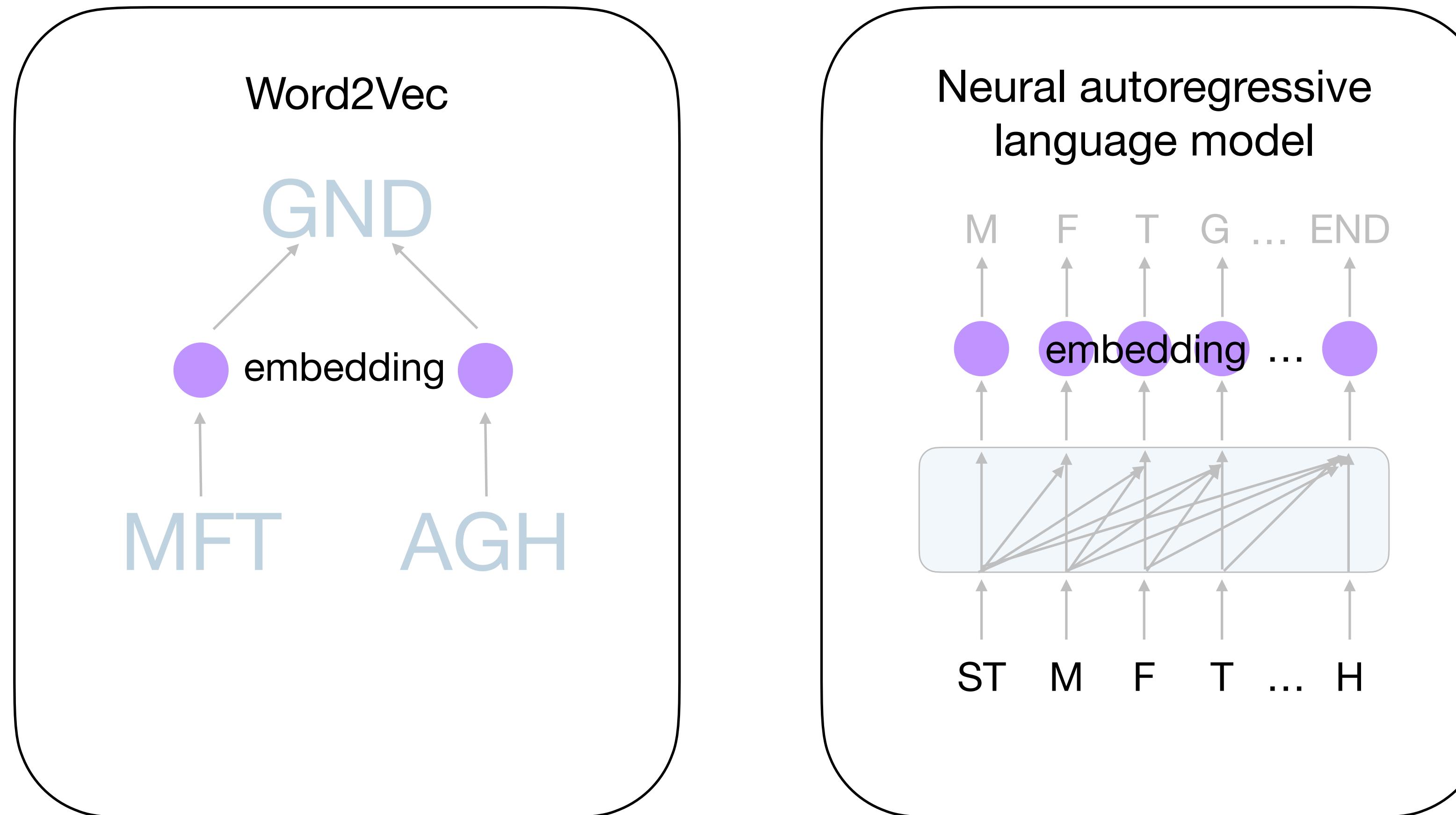
MFTGNDAGH

Many methods pretend proteins are language



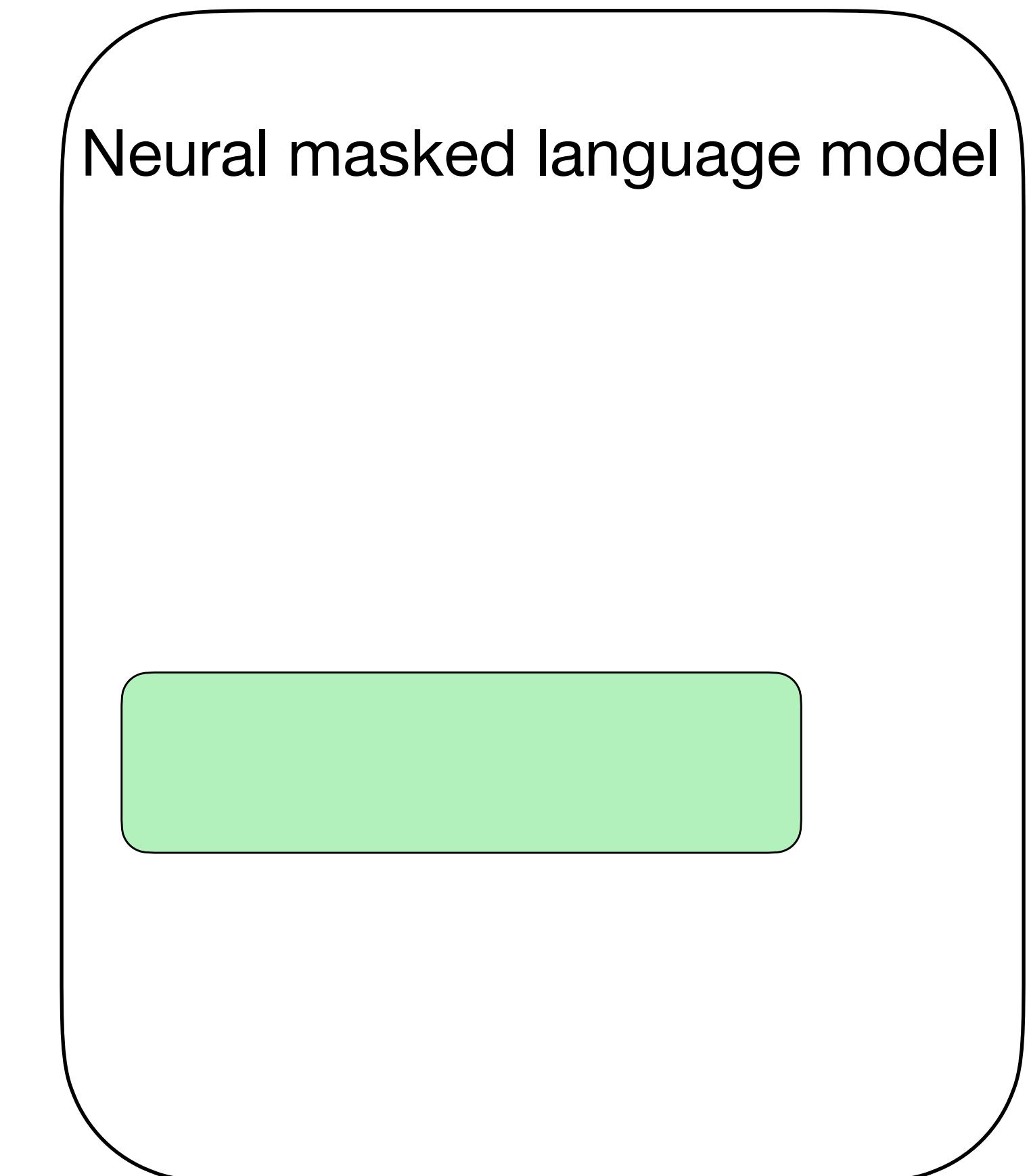
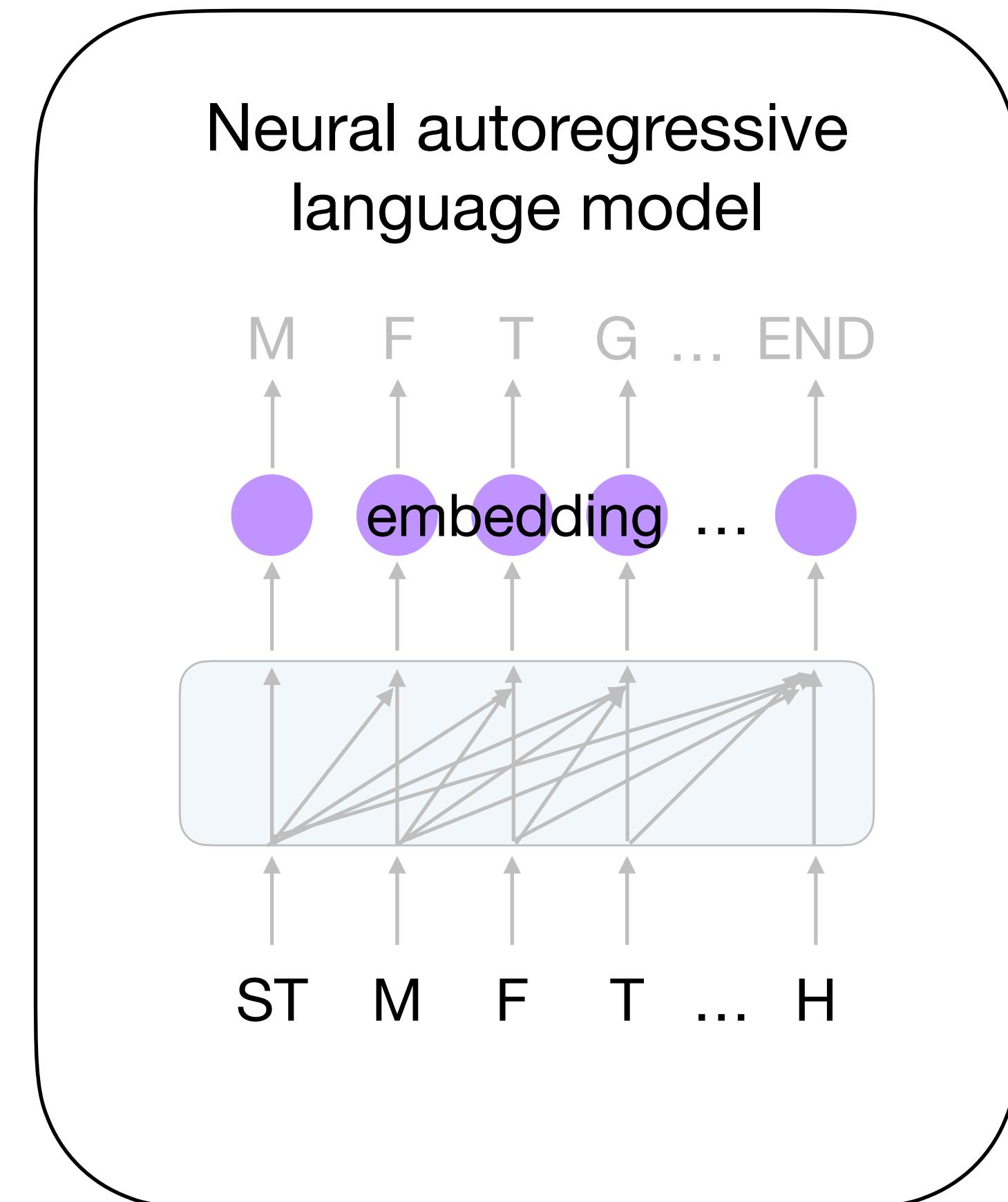
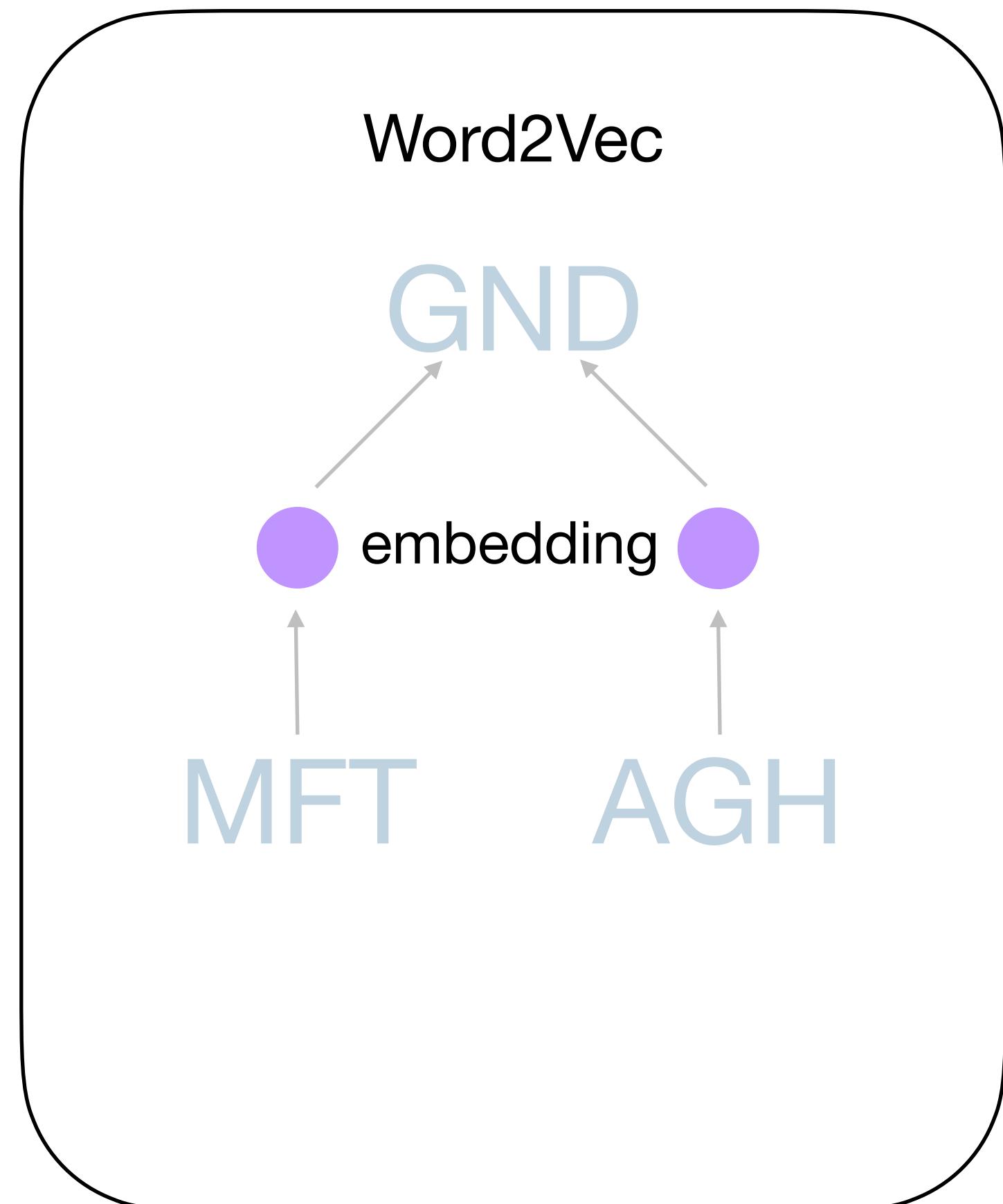
MFTGNDAGH

Many methods pretend proteins are language



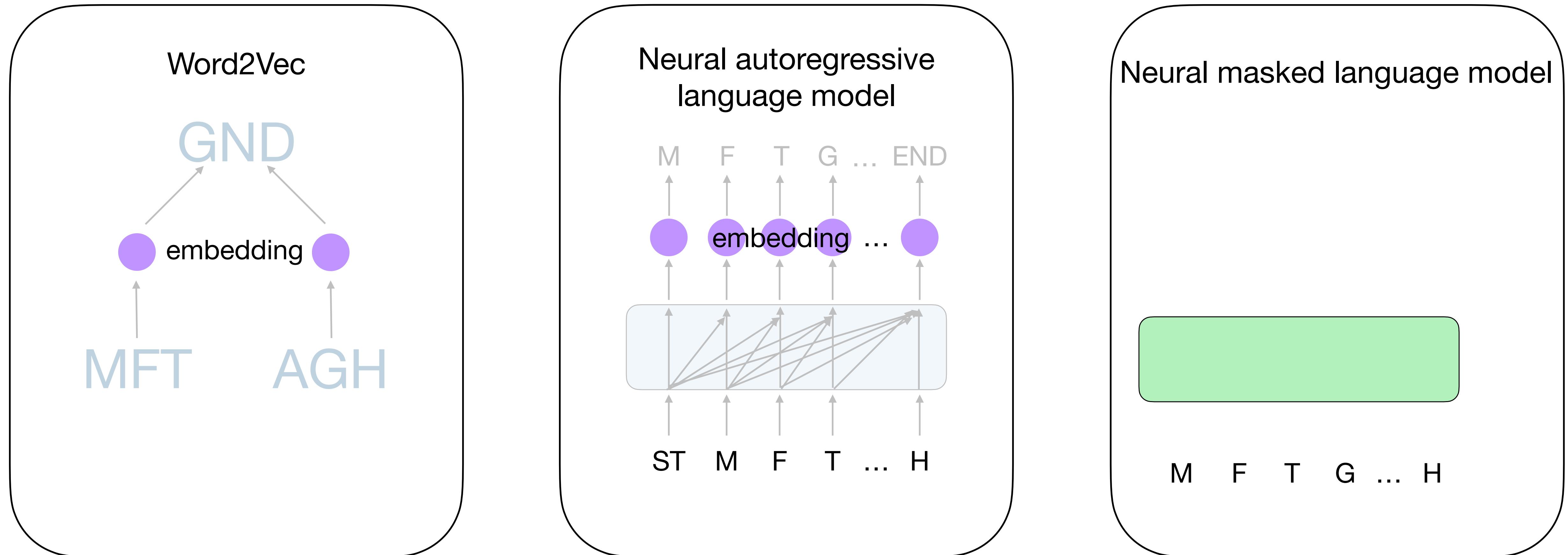
MFTGNDAGH

Many methods pretend proteins are language



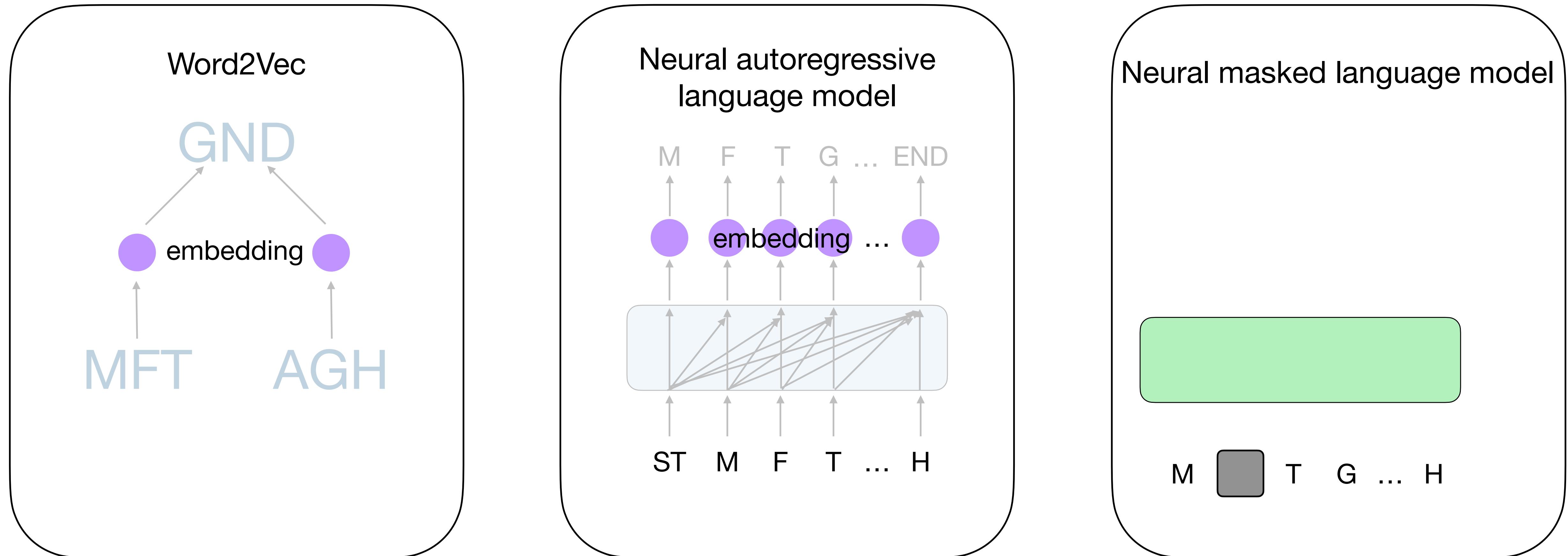
MFTGNDAGH

Many methods pretend proteins are language

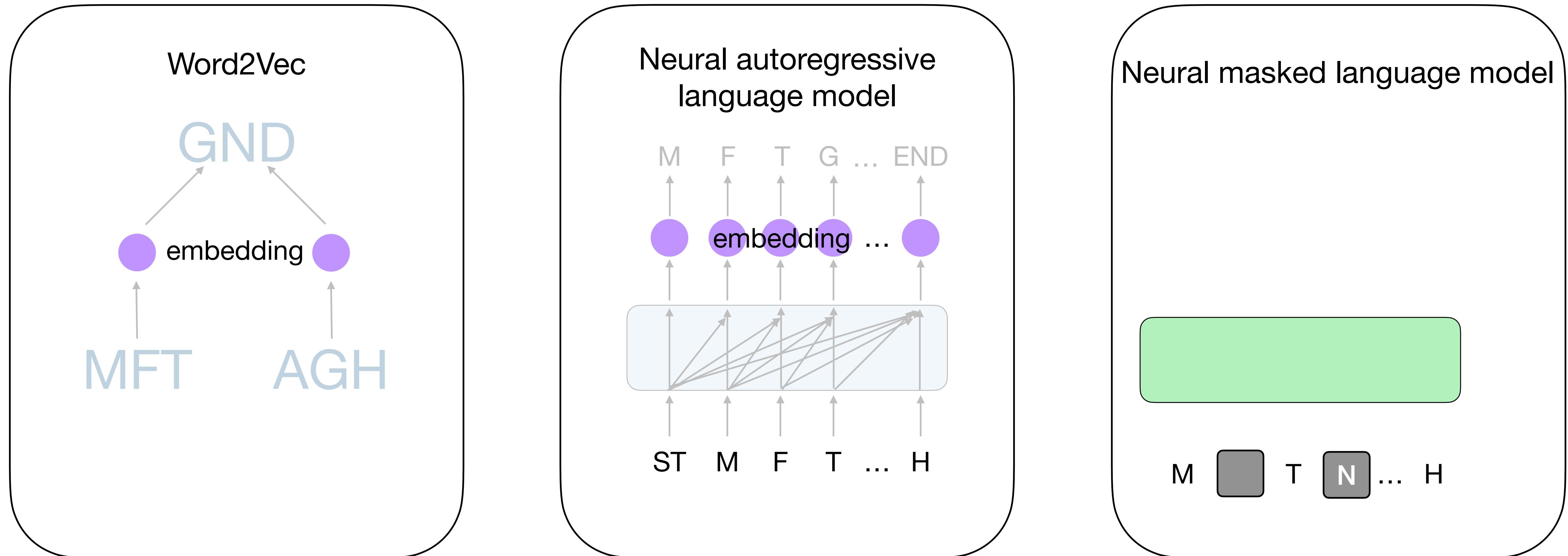


MFTGNDAGH

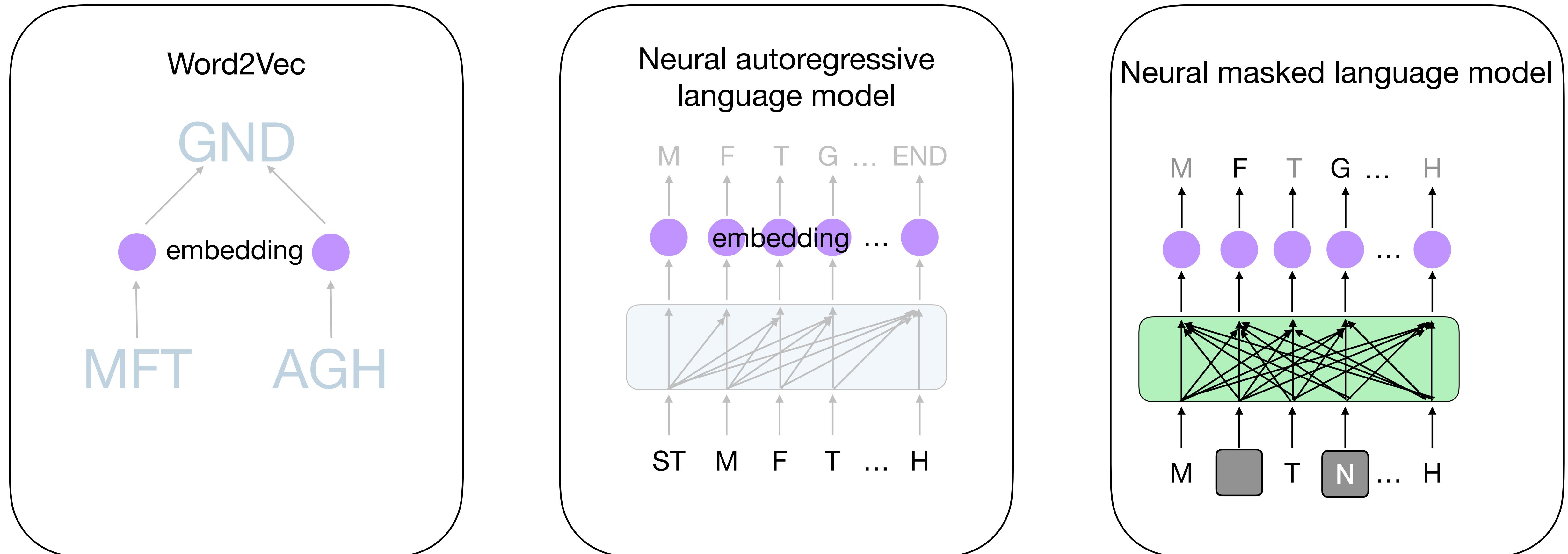
Many methods pretend proteins are language



Many methods pretend proteins are language

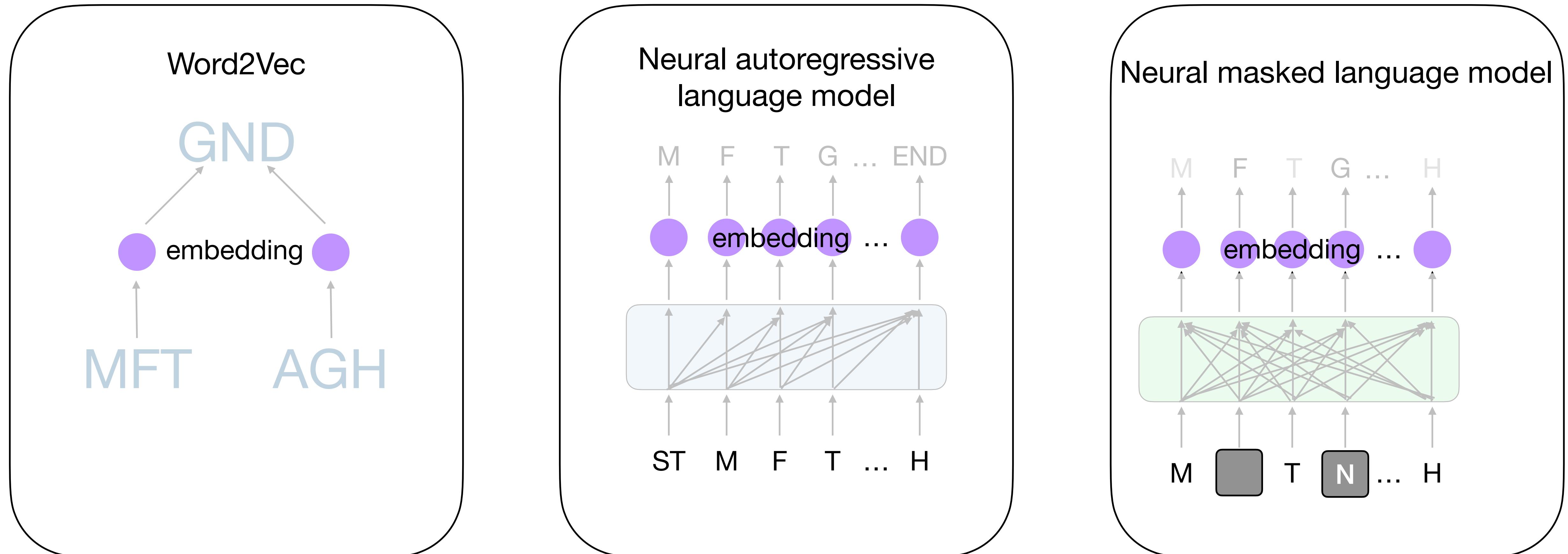


Many methods pretend proteins are language



MFTGNDAGH

Many methods pretend proteins are language



MFTGNDAGH

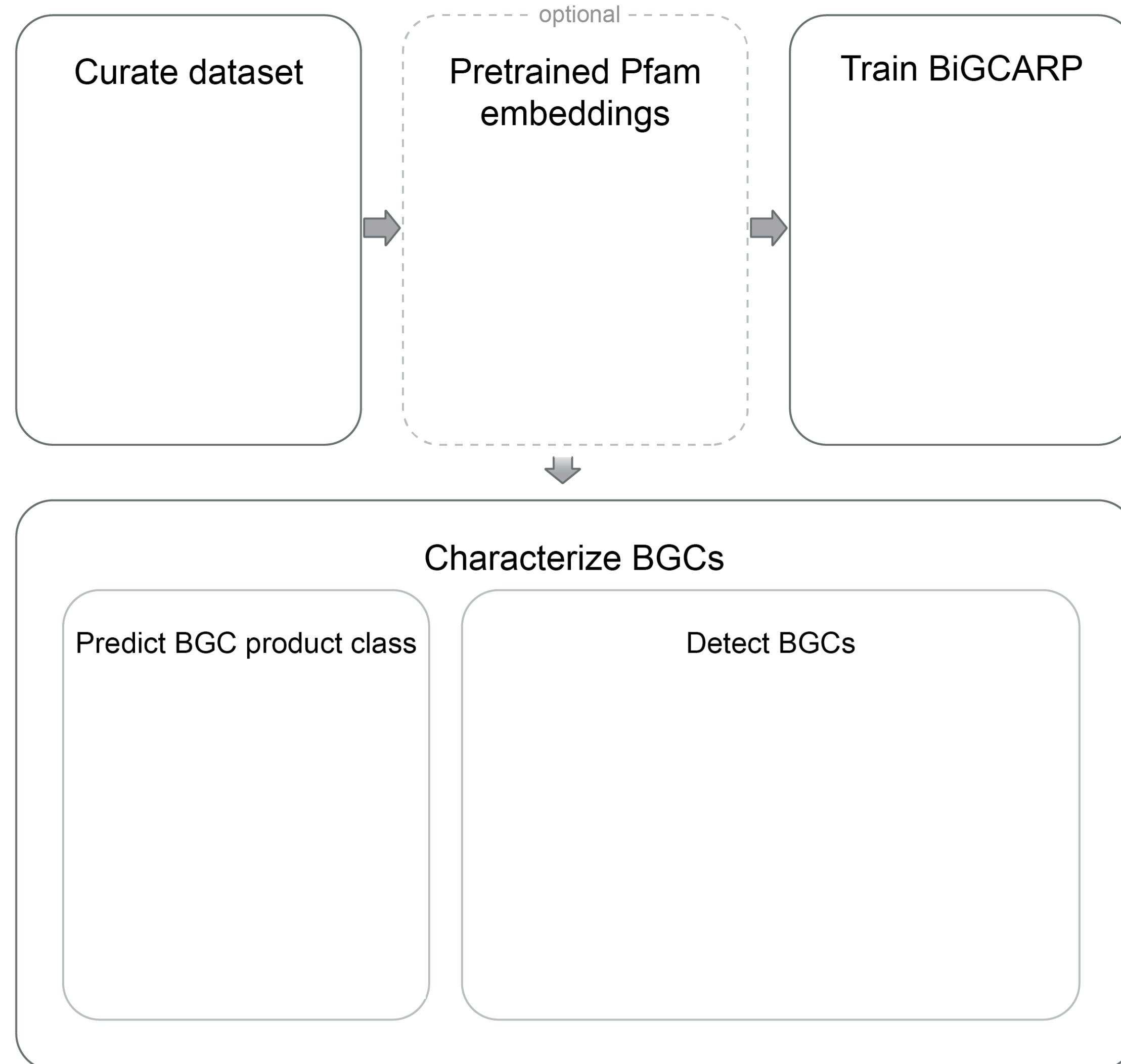
BiGCARP is a self-supervised BGC model

BiGCARP is a self-supervised BGC model

Biosynthetic Gene cluster Convolutional Autoencoding Representations of Proteins

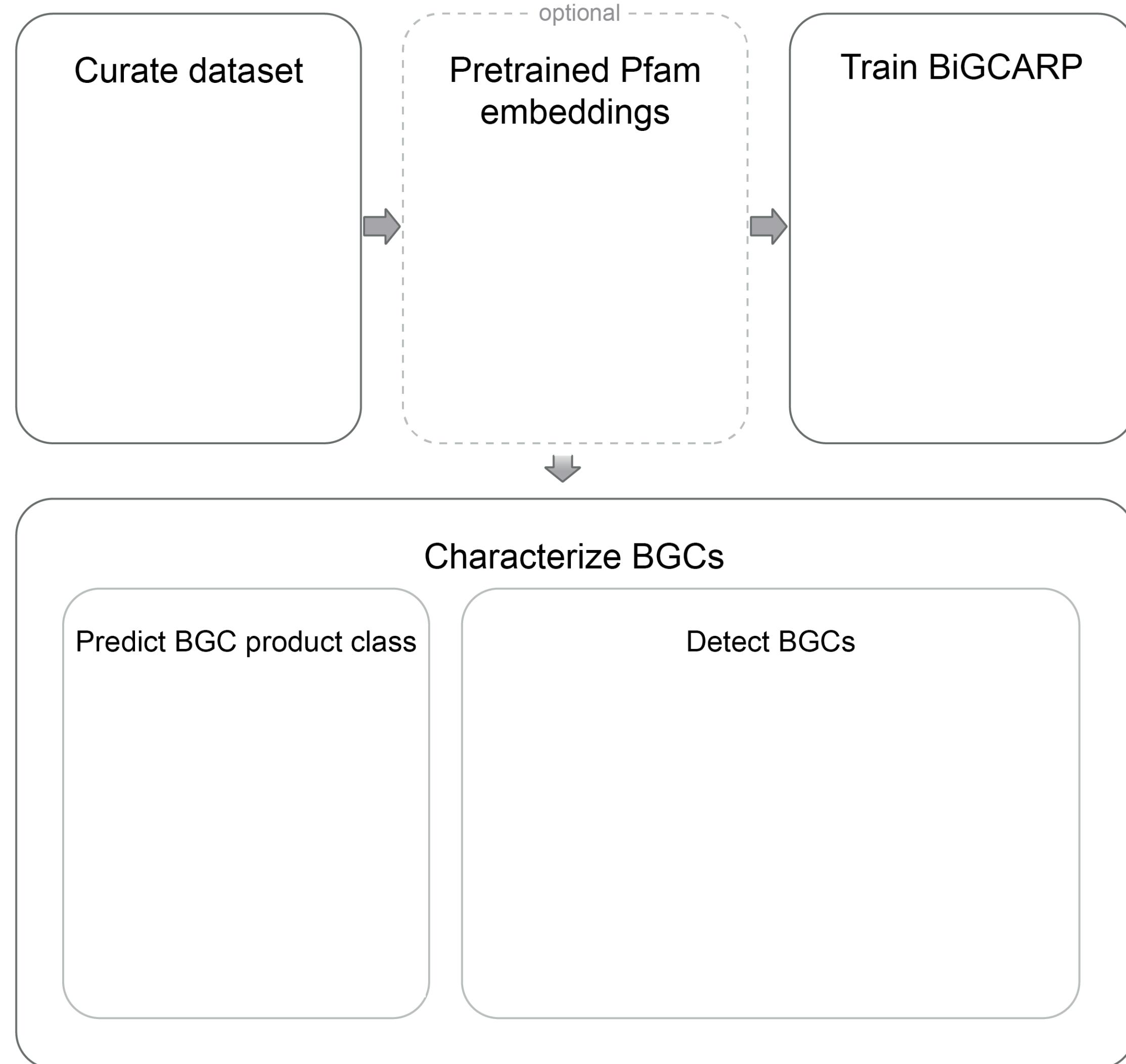
BiGCARP is a self-supervised BGC model

Biosynthetic Gene cluster Convolutional Autoencoding Representations of Proteins



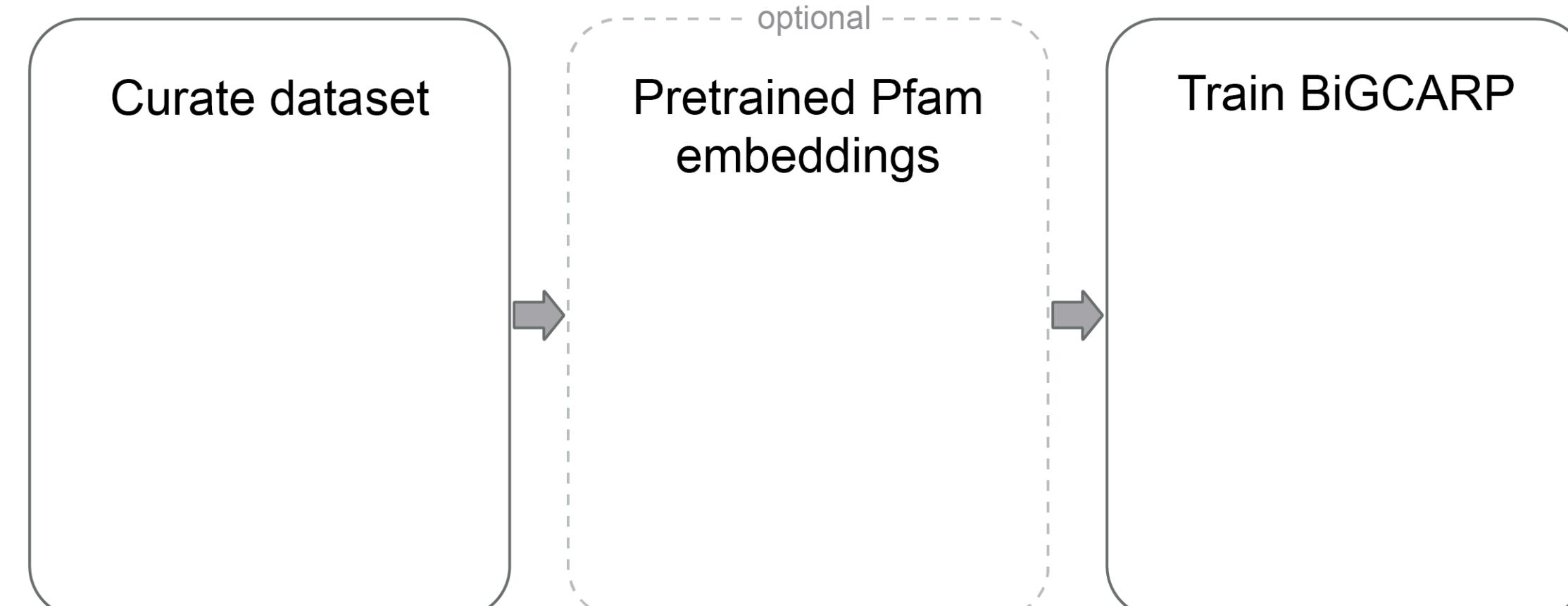
BiGCARP is a self-supervised BGC model

Biosynthetic Gene cluster Convolutional Autoencoding Representations of Proteins

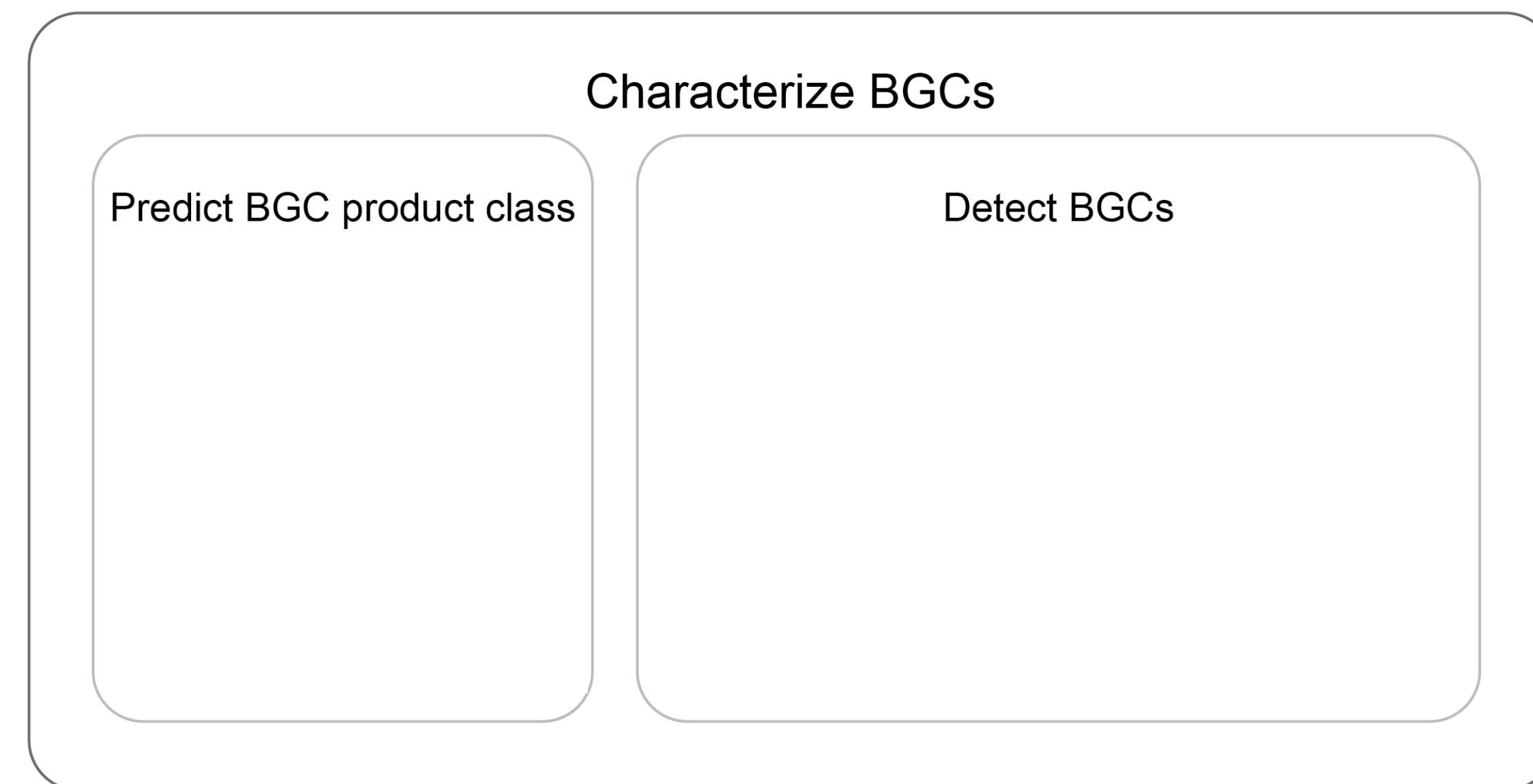


BiGCARP is a self-supervised BGC model

Biosynthetic Gene cluster Convolutional Autoencoding Representations of Proteins

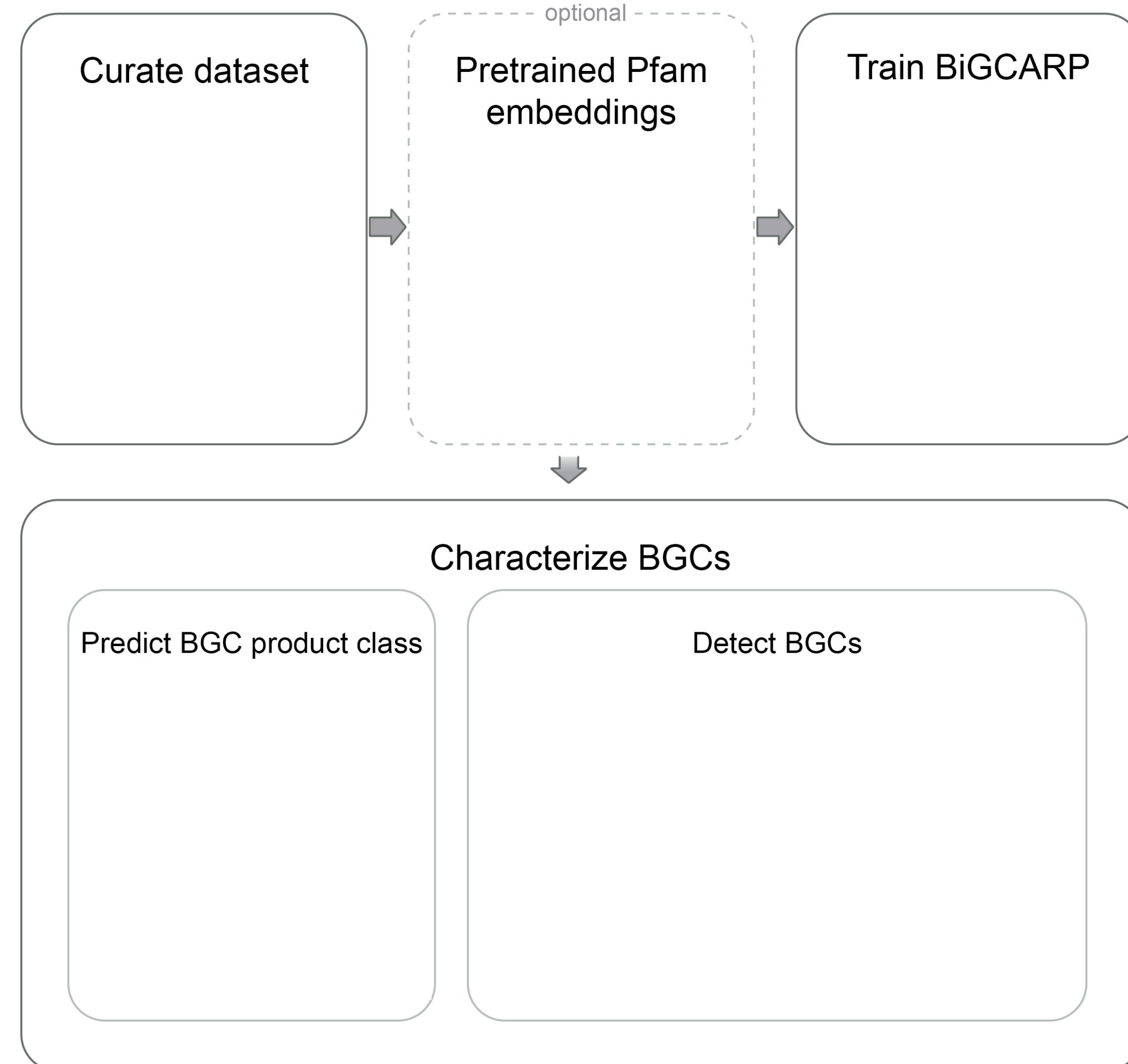


Previous data-mining efforts find **positive** examples of BGCs



BiGCARP is a self-supervised BGC model

Biosynthetic Gene cluster Convolutional Autoencoding Representations of Proteins



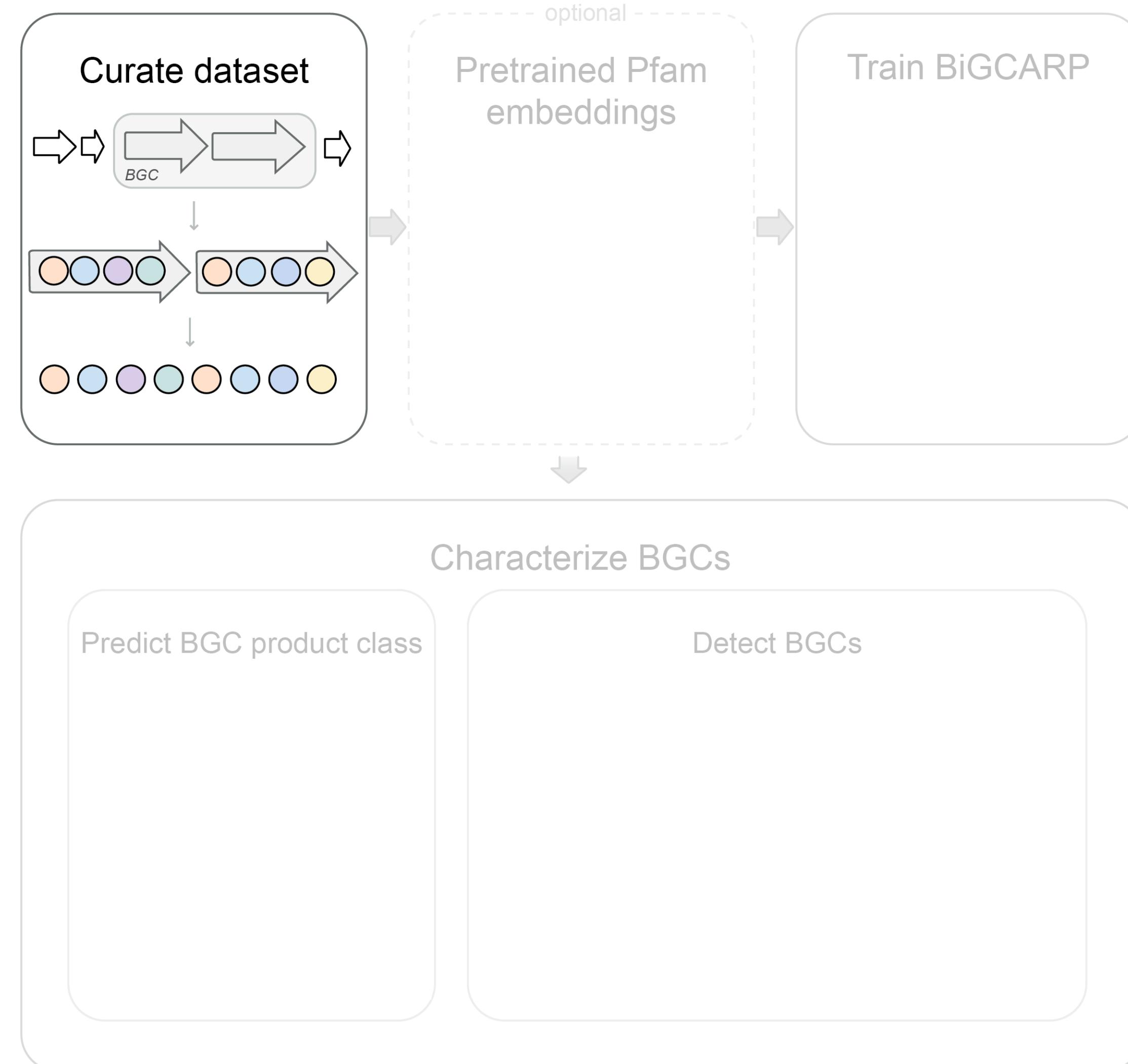
Previous data-mining efforts find **positive** examples of BGCs

True **negative** examples are difficult to find



BiGCARP is a self-supervised BGC model

Biosynthetic Gene cluster Convolutional Autoencoding Representations of Proteins



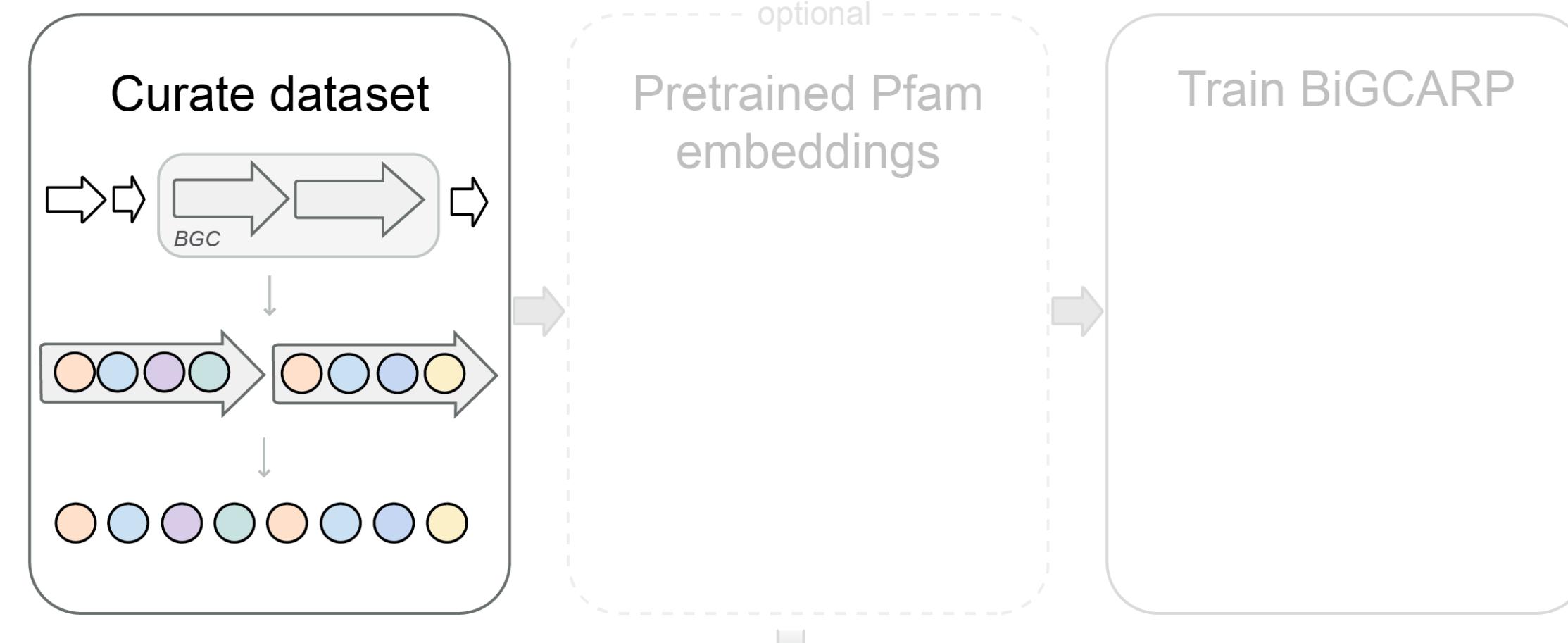
Previous data-mining efforts find **positive** examples of BGCs

True **negative** examples are difficult to find



BiGCARP is a self-supervised BGC model

Biosynthetic Gene cluster Convolutional Autoencoding Representations of Proteins



Previous data-mining efforts find **positive** examples of BGCs

True **negative** examples are difficult to find

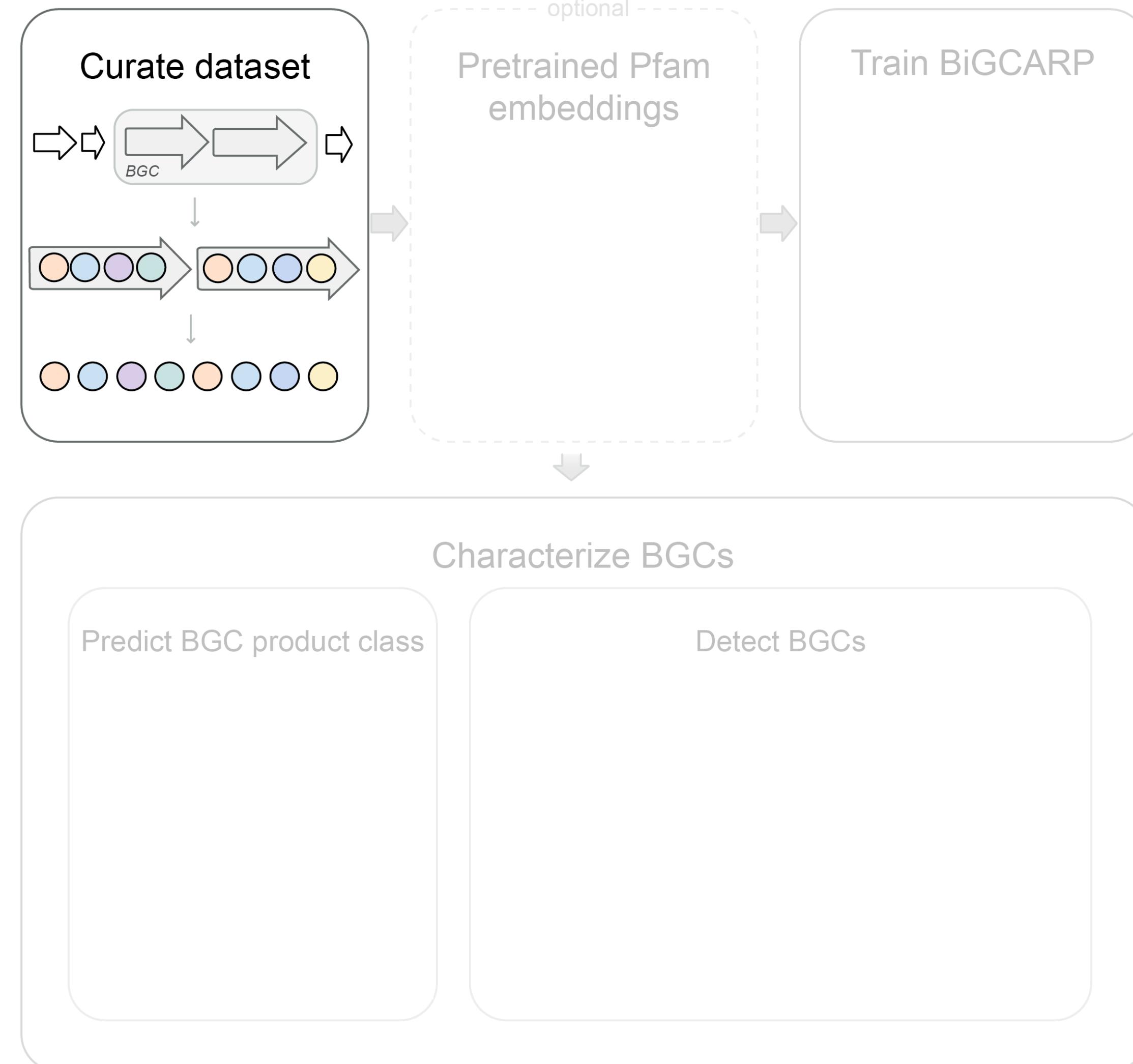
Microbial genome

...ACTGC^GGTACG...



BiGCARP is a self-supervised BGC model

Biosynthetic Gene cluster Convolutional Autoencoding Representations of Proteins

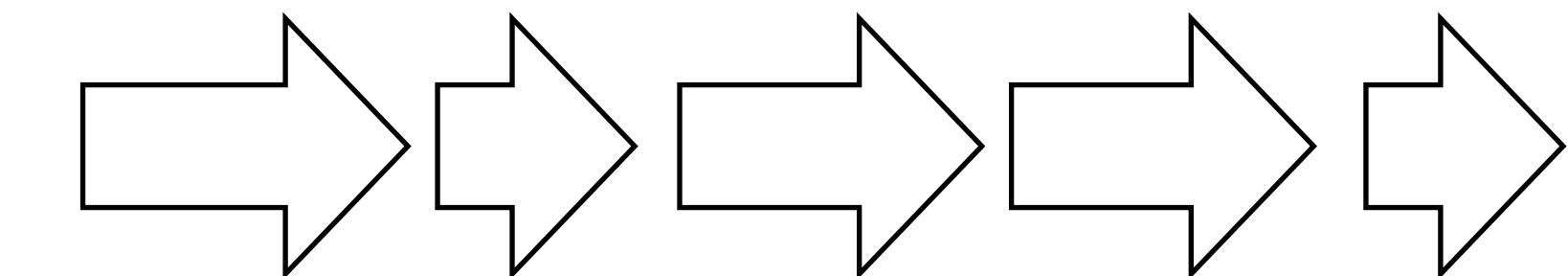


Previous data-mining efforts find **positive** examples of BGCs

True **negative** examples are difficult to find

Microbial genome

...ACTGC^GGTACG...

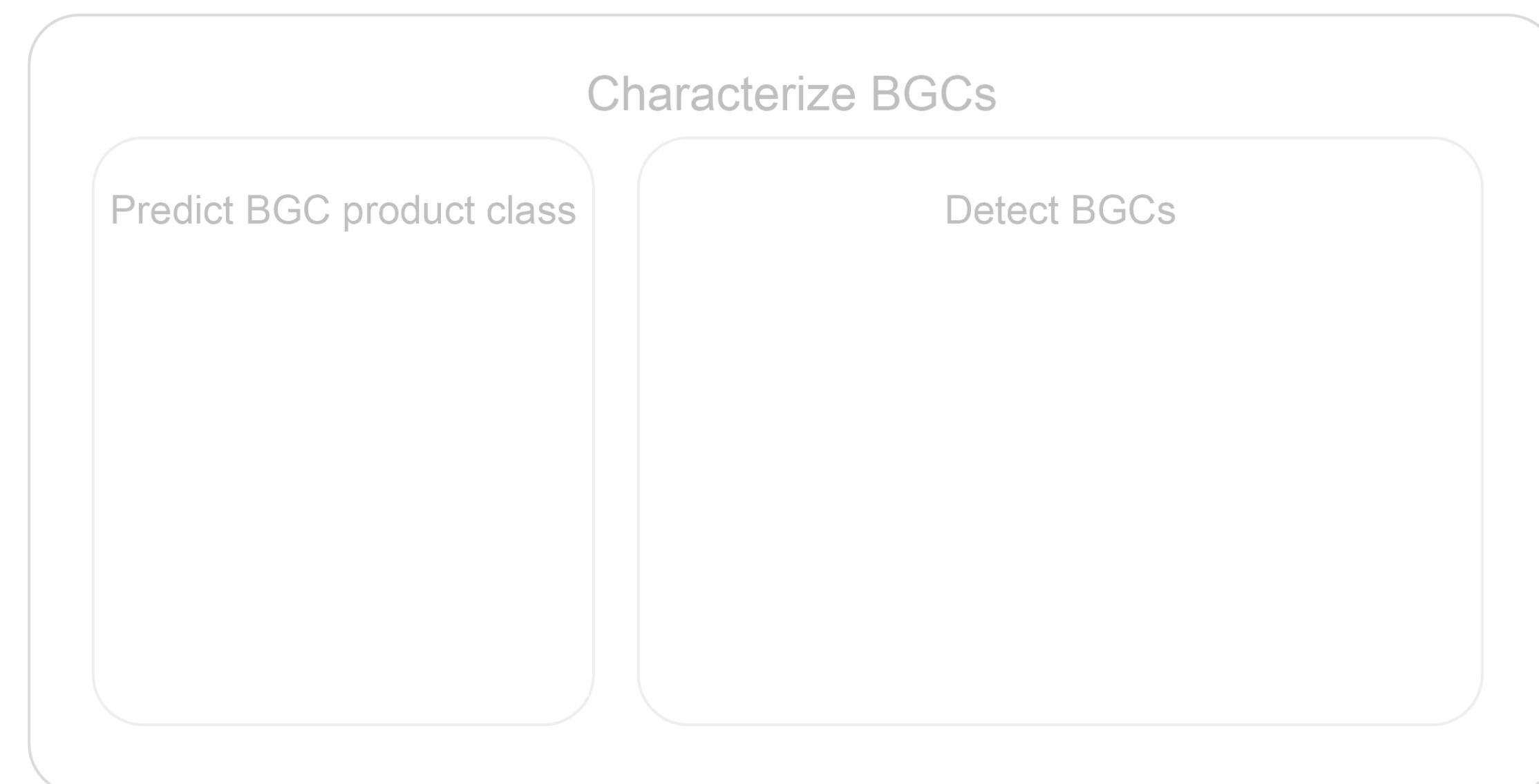
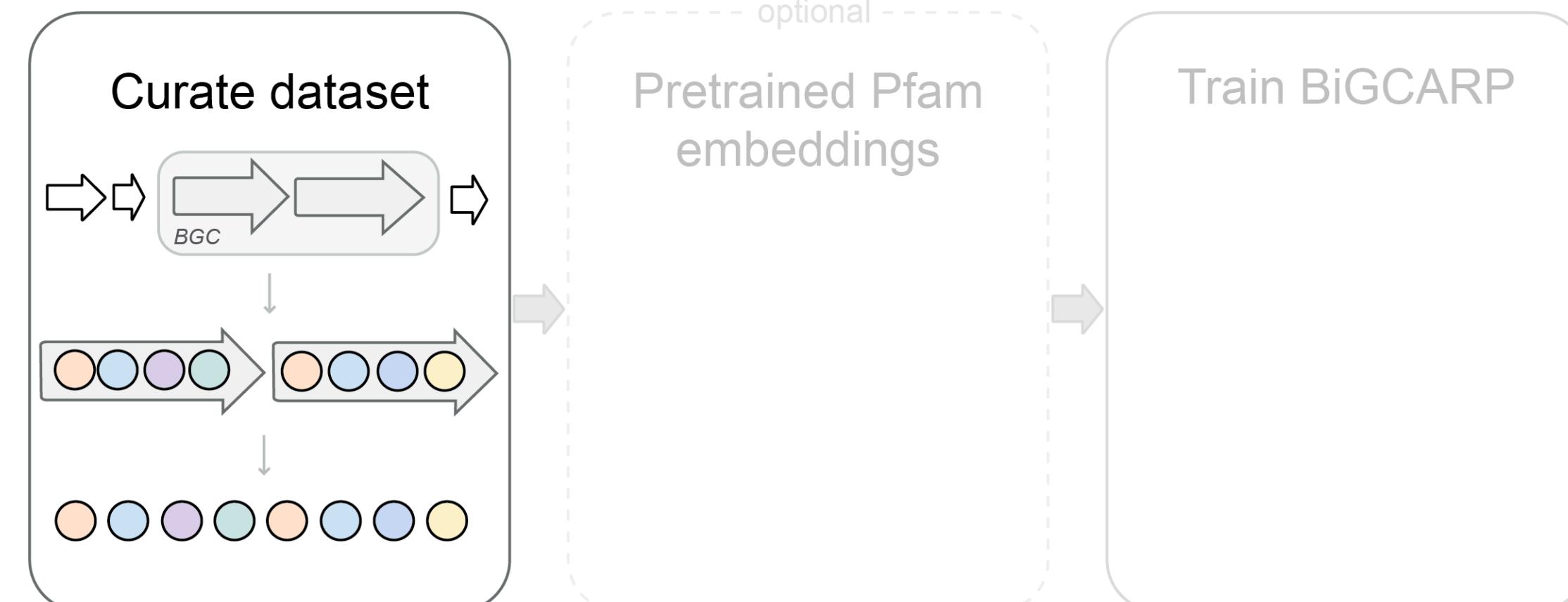


Open reading frames



BiGCARP is a self-supervised BGC model

Biosynthetic Gene cluster Convolutional Autoencoding Representations of Proteins



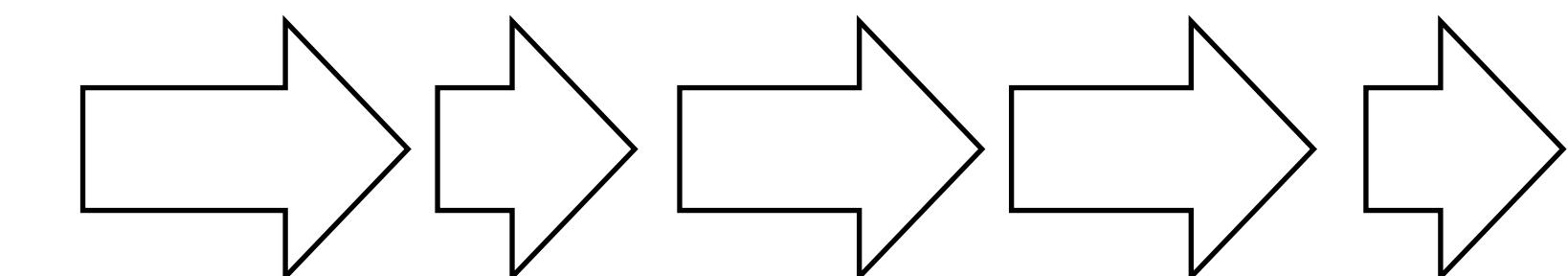
Previous data-mining efforts find **positive** examples of BGCs

True **negative** examples are difficult to find

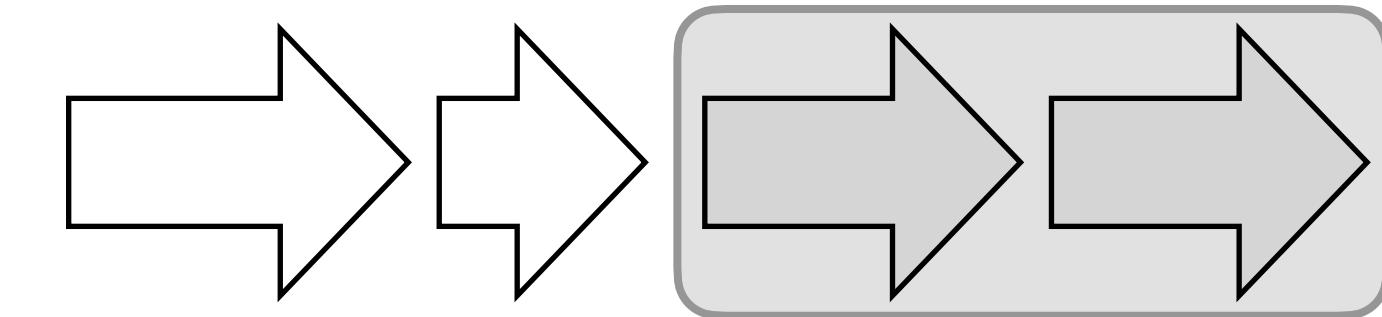
Microbial genome

...ACTGC~~G~~TTACG...

Open reading frames

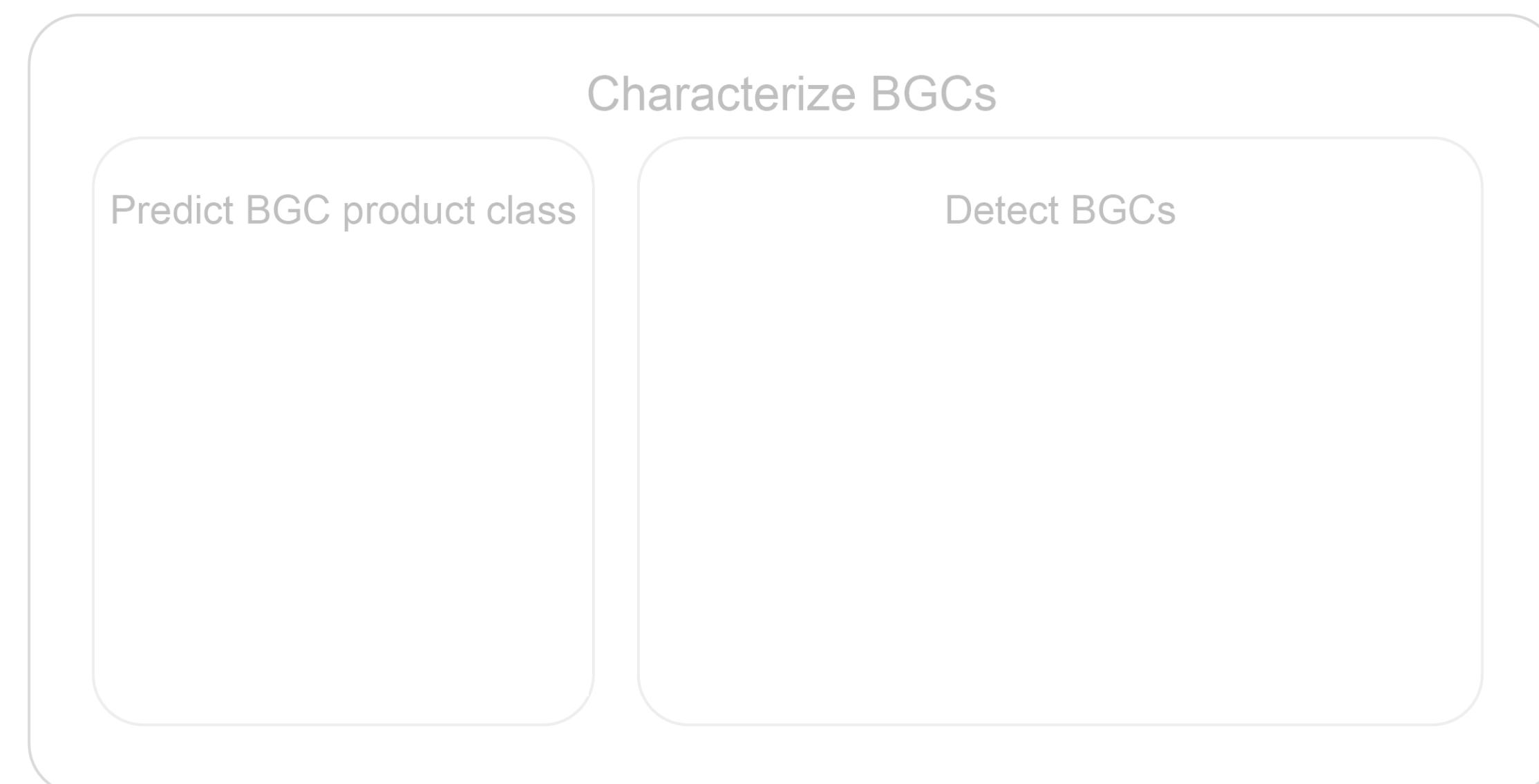
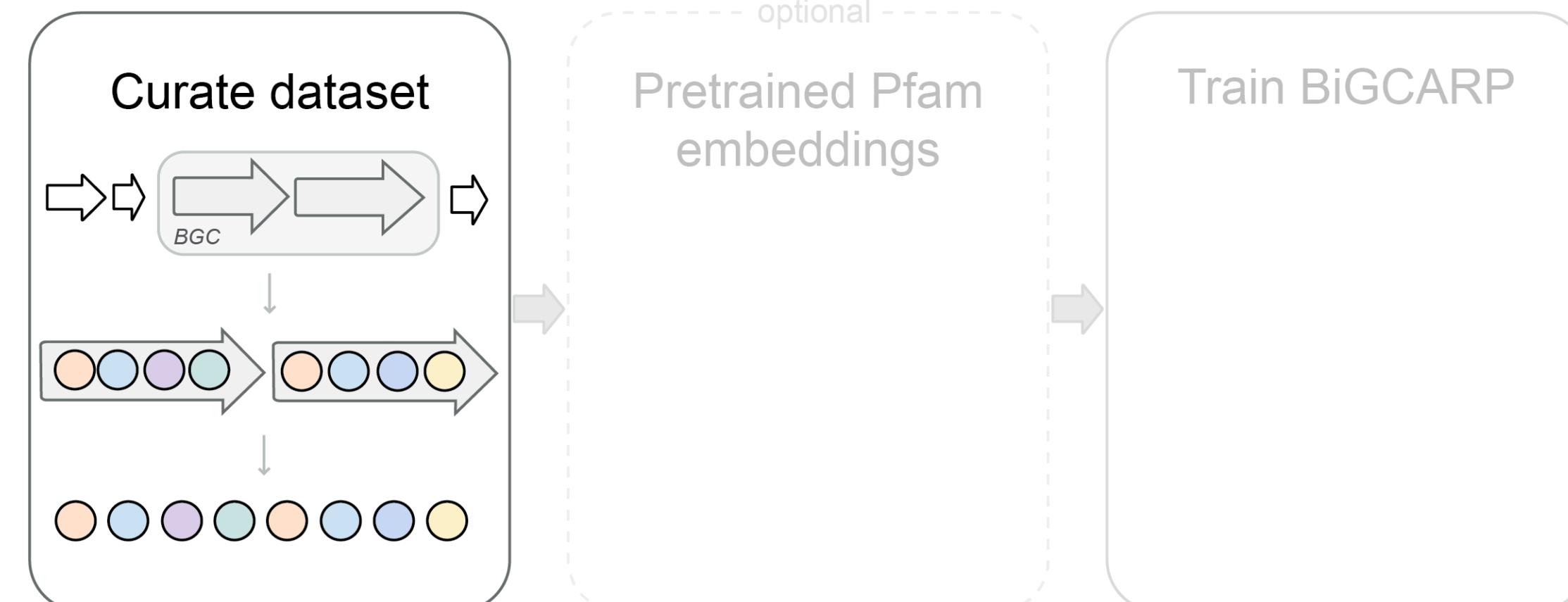


Identify BGCs with
antiSMASH



BiGCARP is a self-supervised BGC model

Biosynthetic Gene cluster Convolutional Autoencoding Representations of Proteins



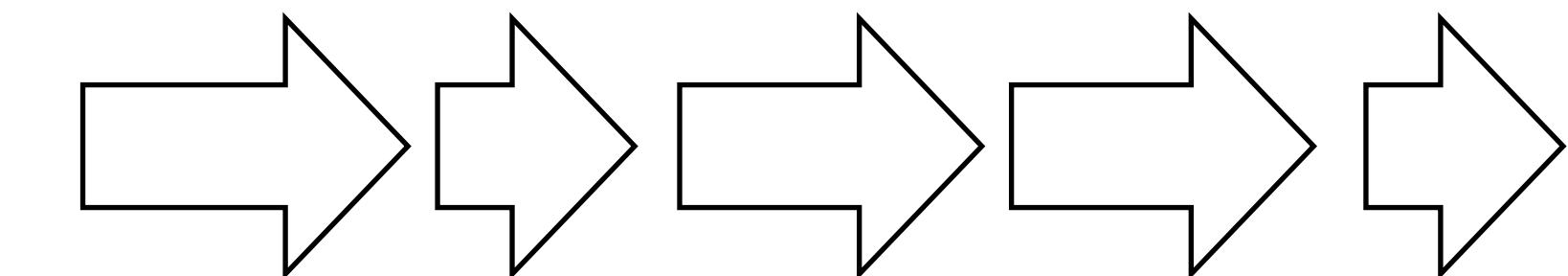
Previous data-mining efforts find **positive** examples of BGCs

True **negative** examples are difficult to find

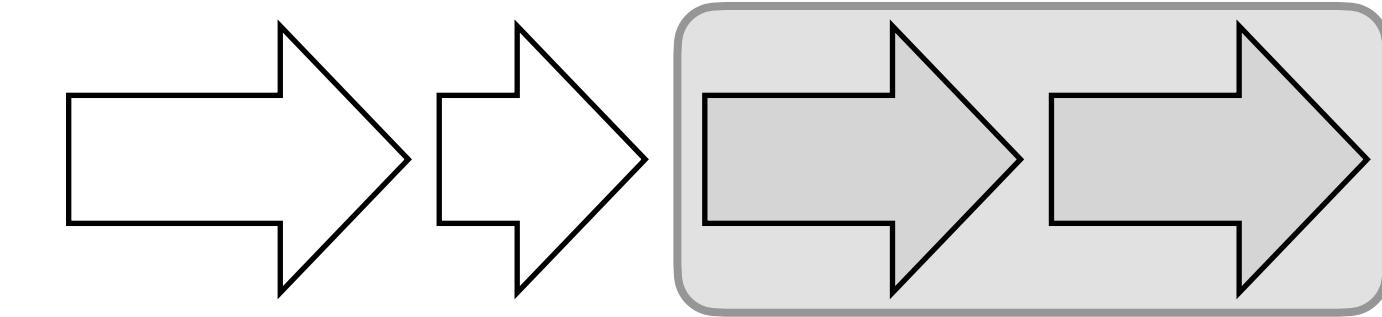
Microbial genome

...ACTGCGGTTACG...

Open reading frames



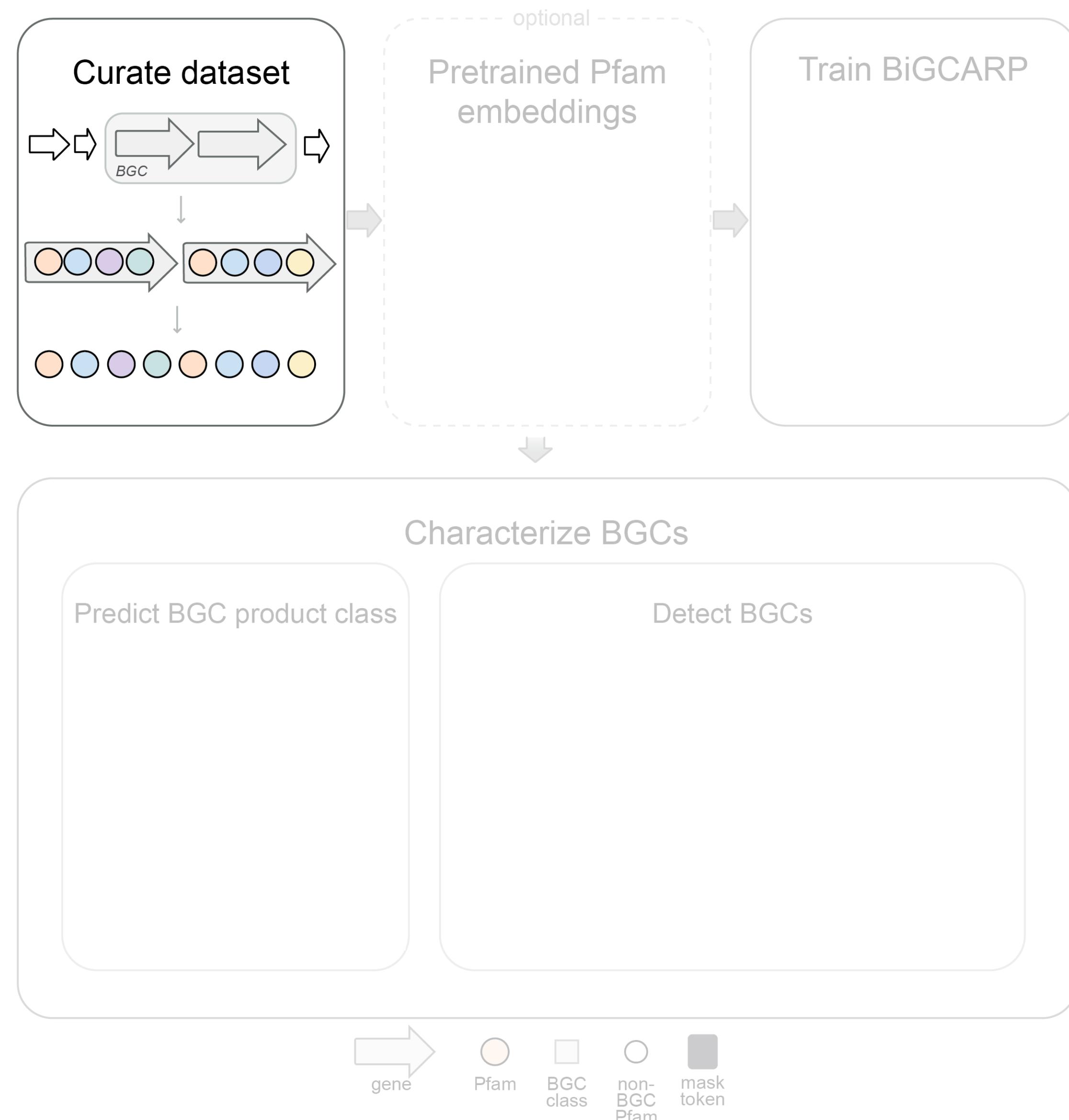
Identify BGCs with
antiSMASH



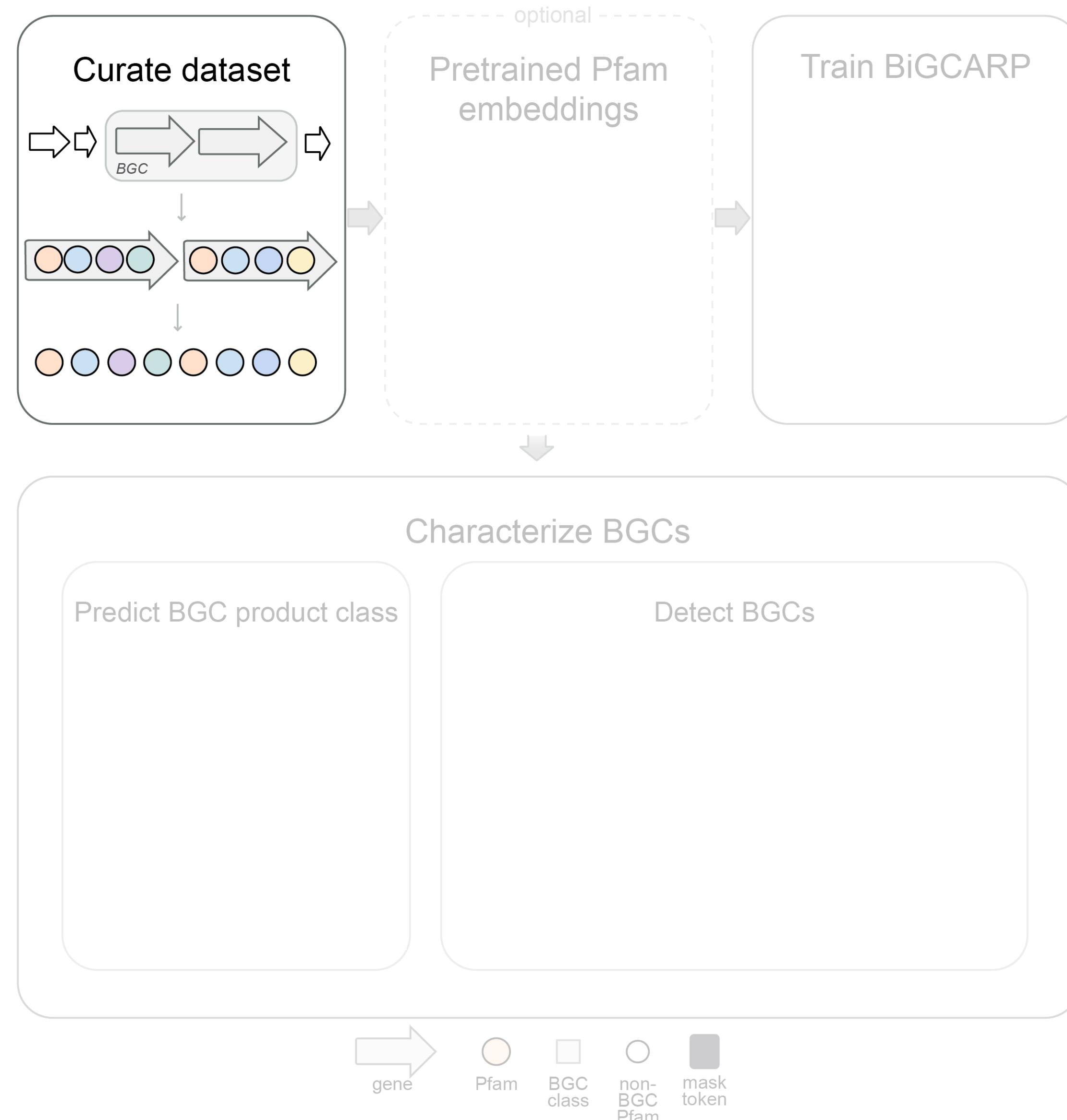
127k BGCs from 55 product classes



We represent BGCs as sequences of domains

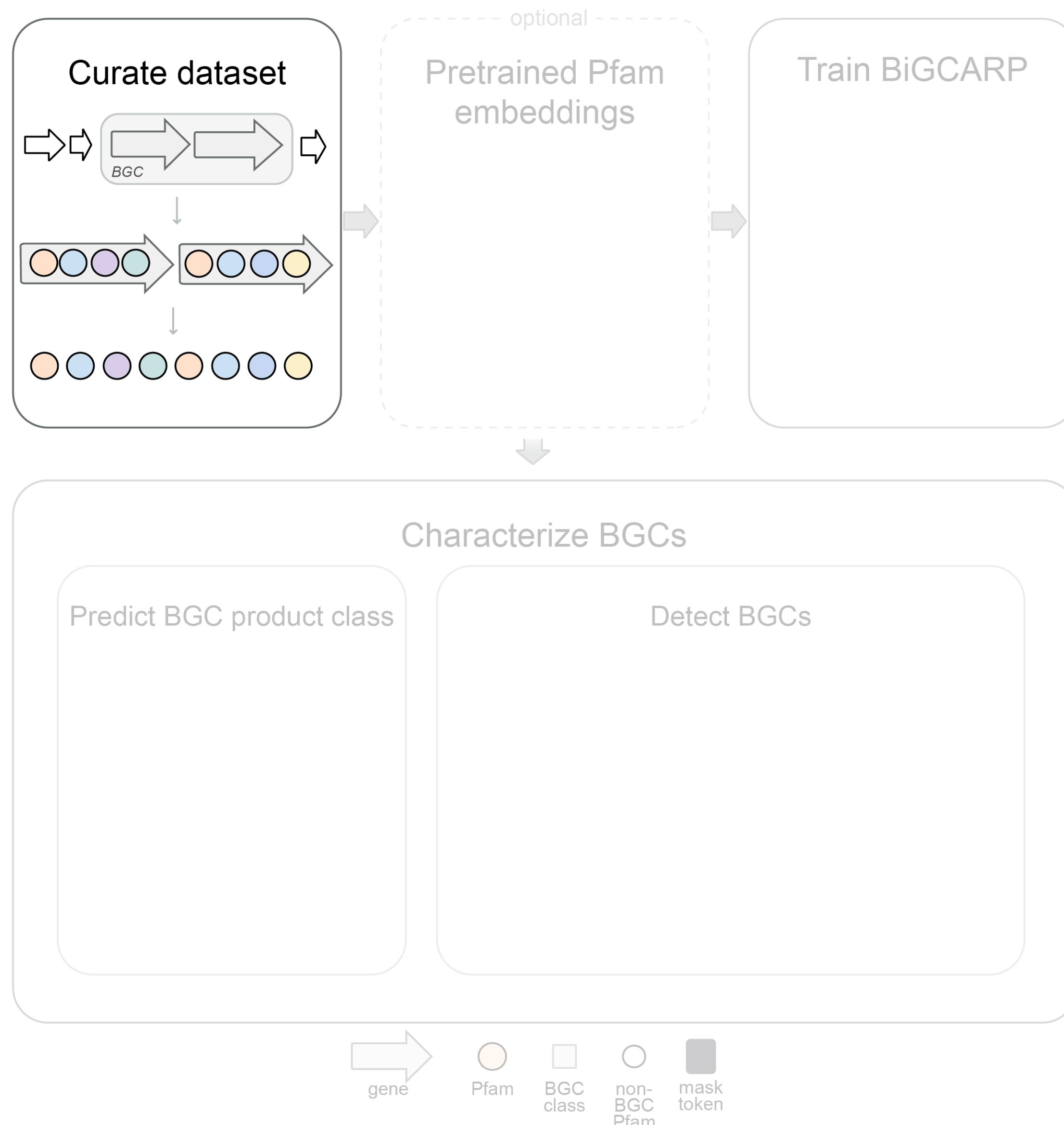


We represent BGCs as sequences of domains



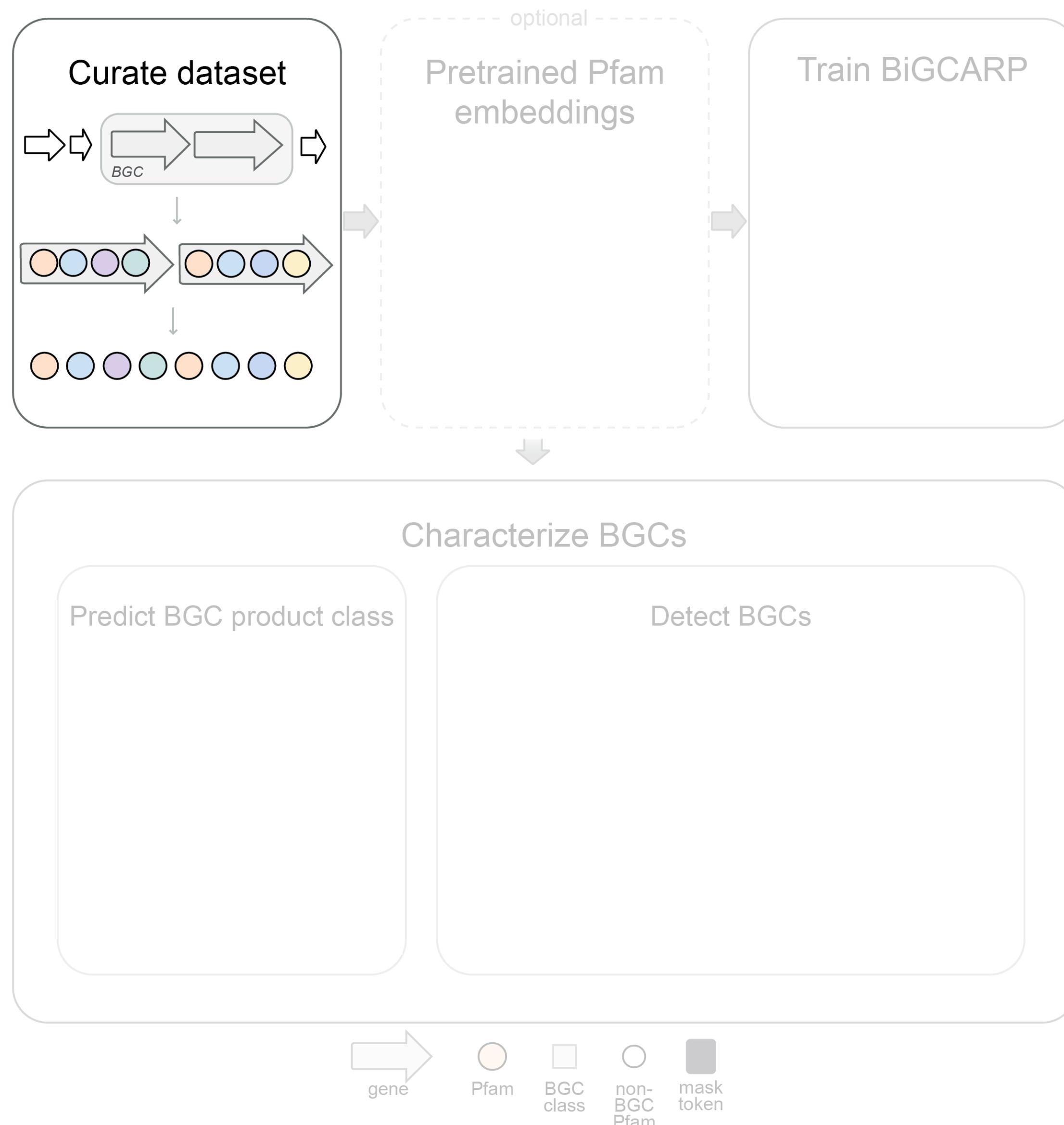
Domain: protein subunit with

We represent BGCs as sequences of domains



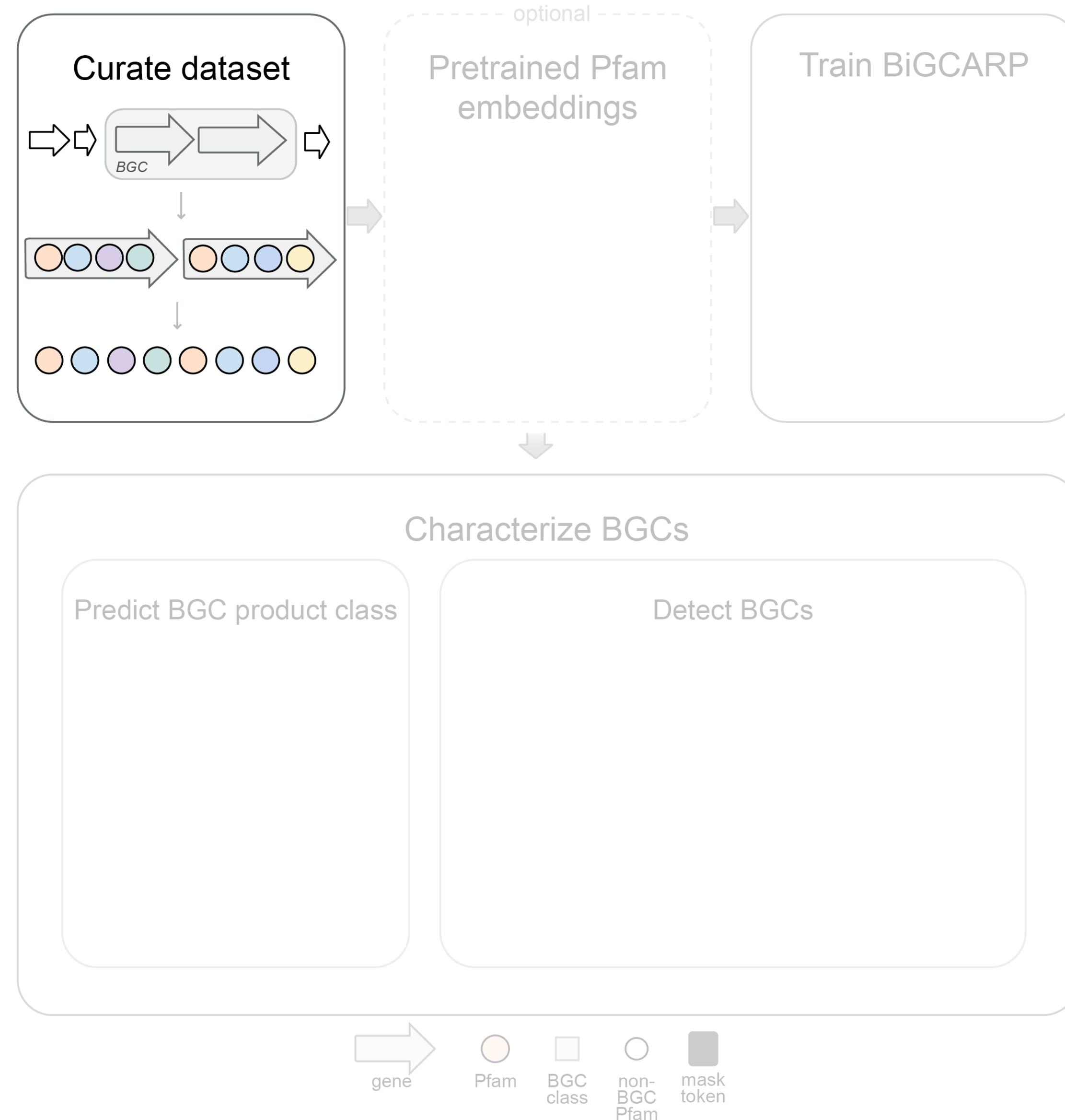
Domain: protein subunit with **compact structure**

We represent BGCs as sequences of domains

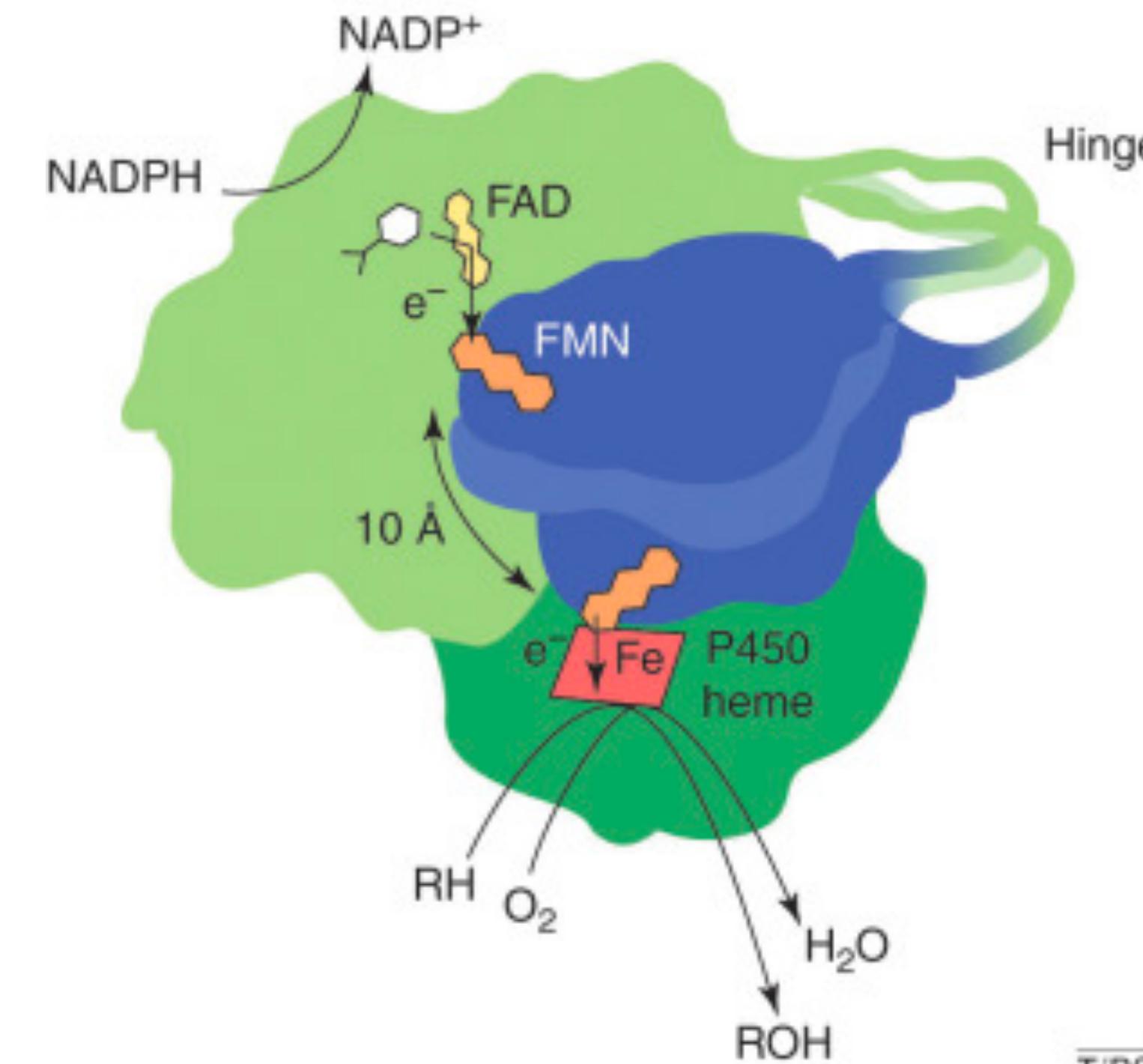


Domain: protein subunit with **compact structure**
modular function

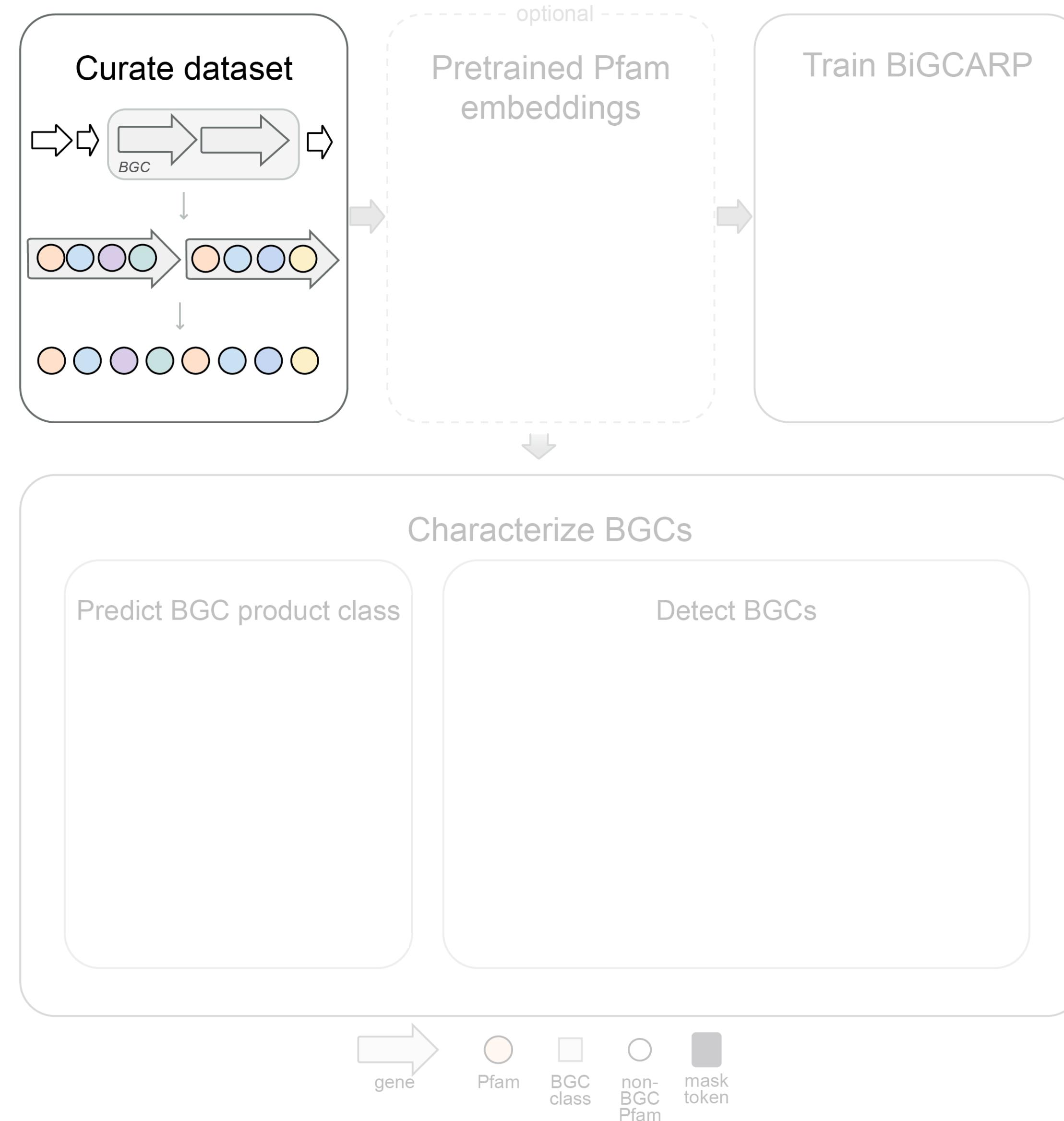
We represent BGCs as sequences of domains



Domain: protein subunit with **compact structure**
modular function

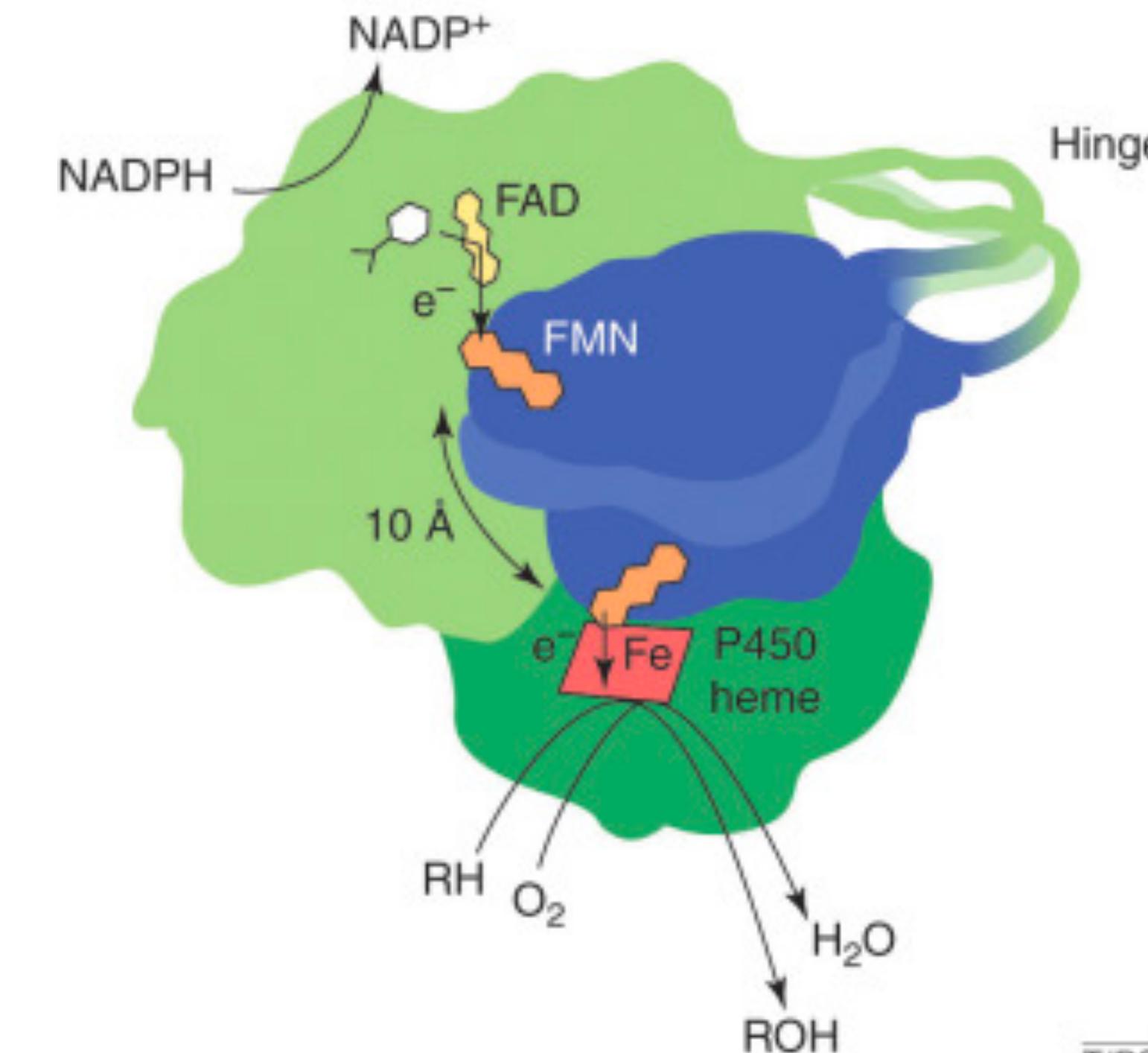


We represent BGCs as sequences of domains

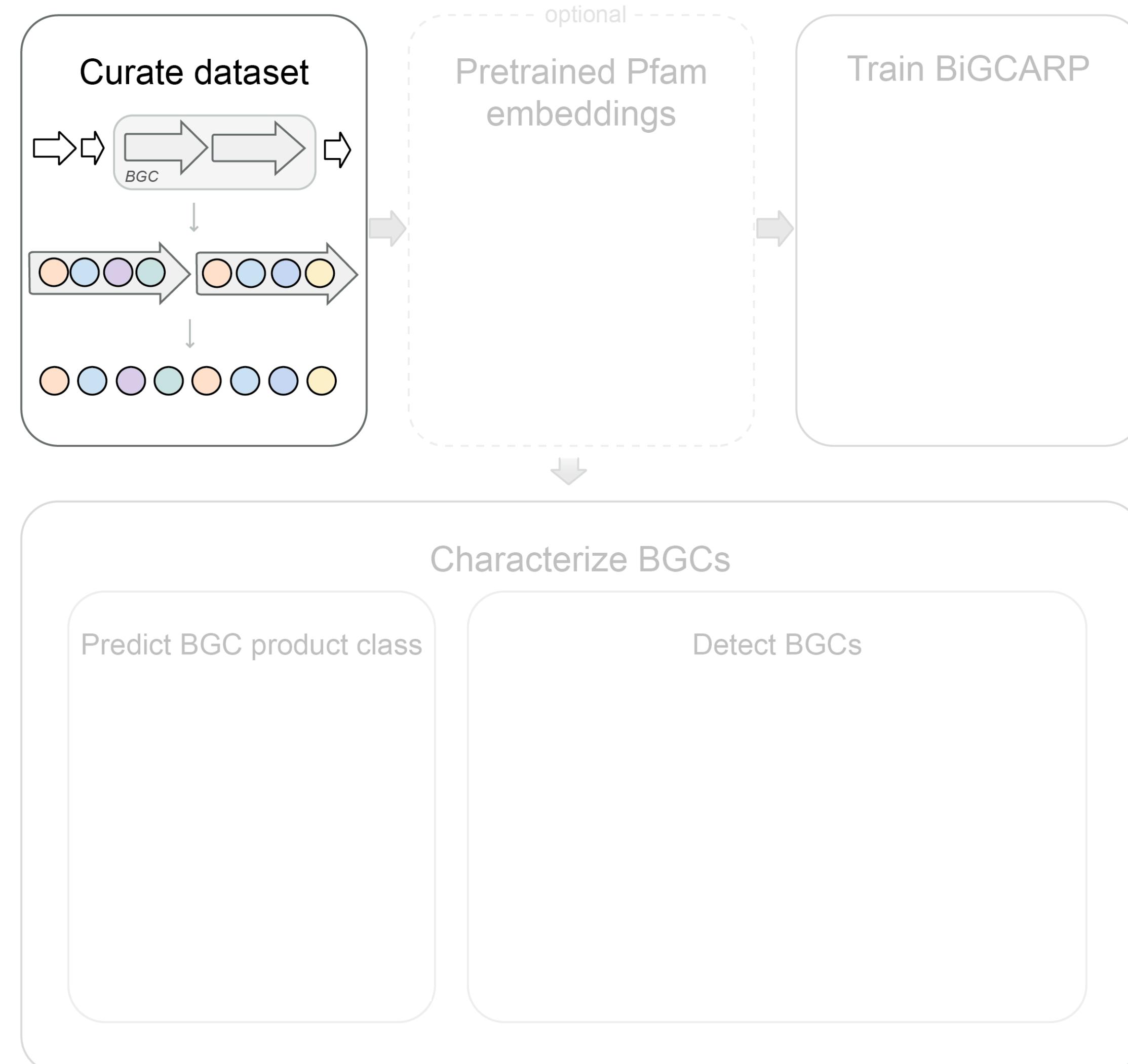


Domain: protein subunit with **compact structure**
modular function

Compromise between granularity,
vocab size, and sequence length



We represent BGCs as sequences of domains

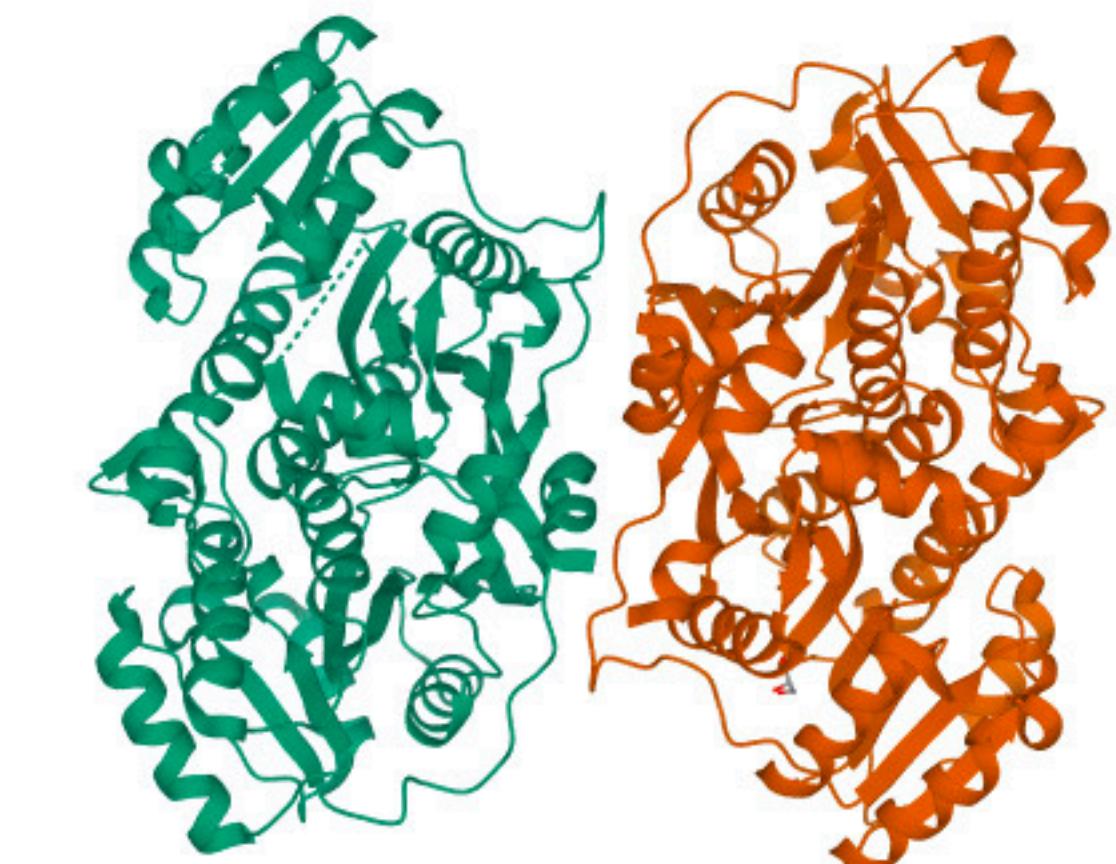


Domain: protein subunit with **compact structure**
modular function

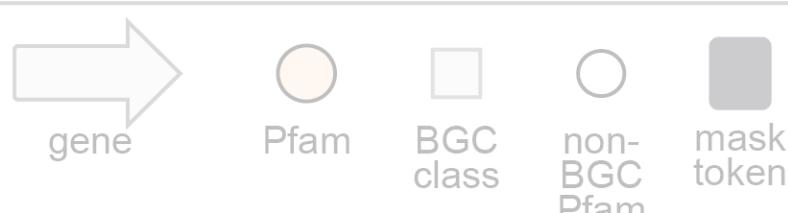
Compromise between granularity,
vocab size, and sequence length



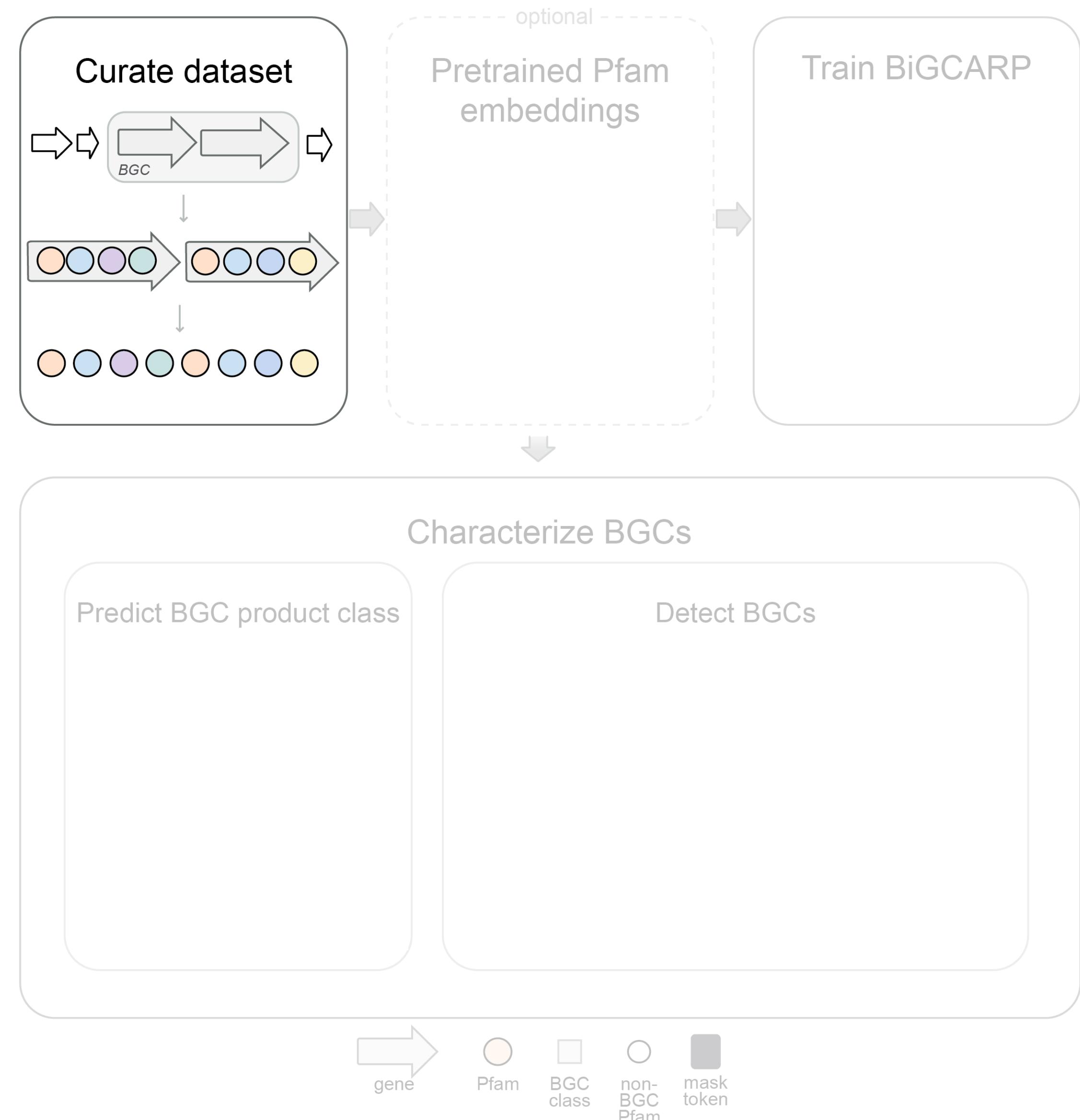
PF02353
Mycolic acid cyclopropane synthetase



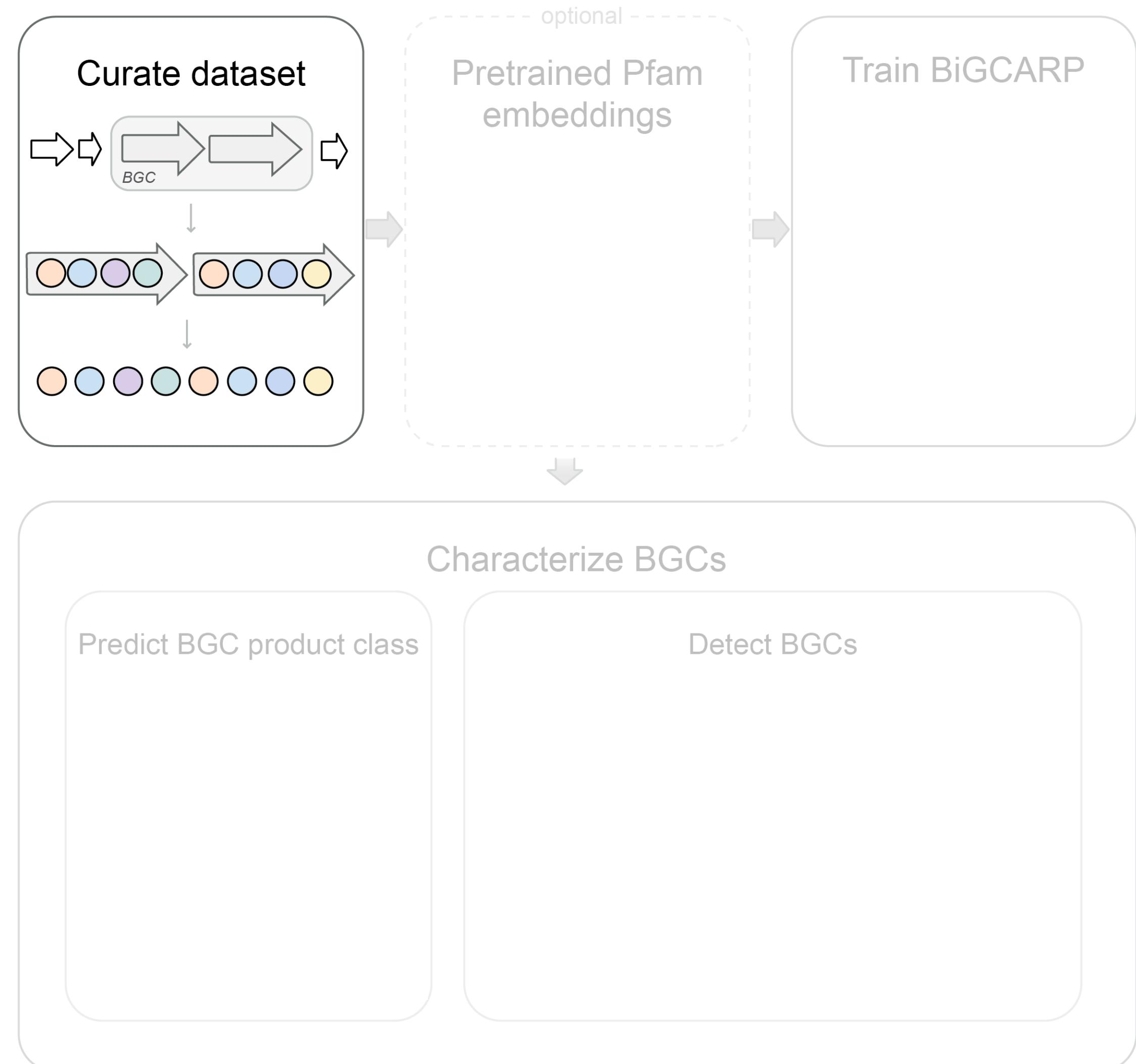
PF08659
ketoreductase domain



We represent BGCs as sequences of domains



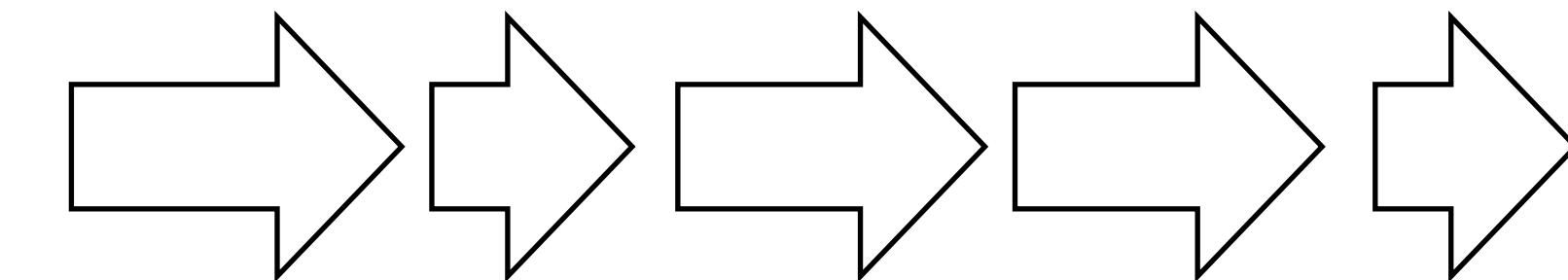
We represent BGCs as sequences of domains



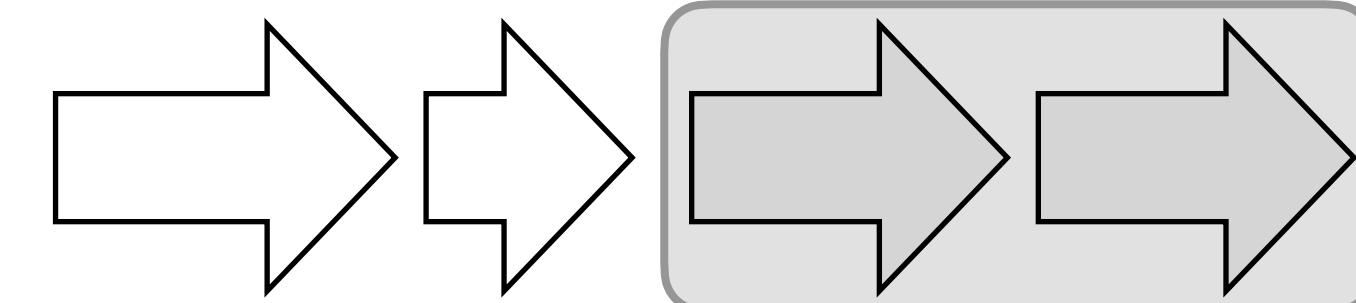
Microbial genome

...ACTGCGGTTACG...

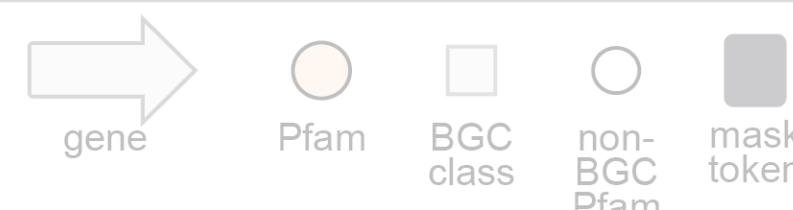
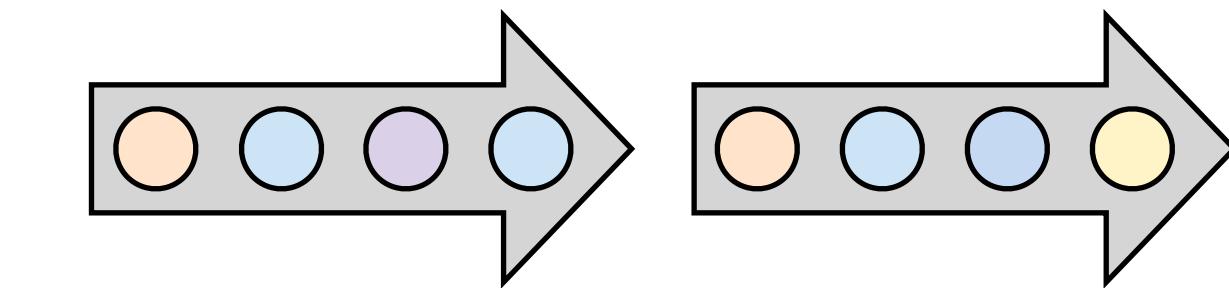
Open reading frames



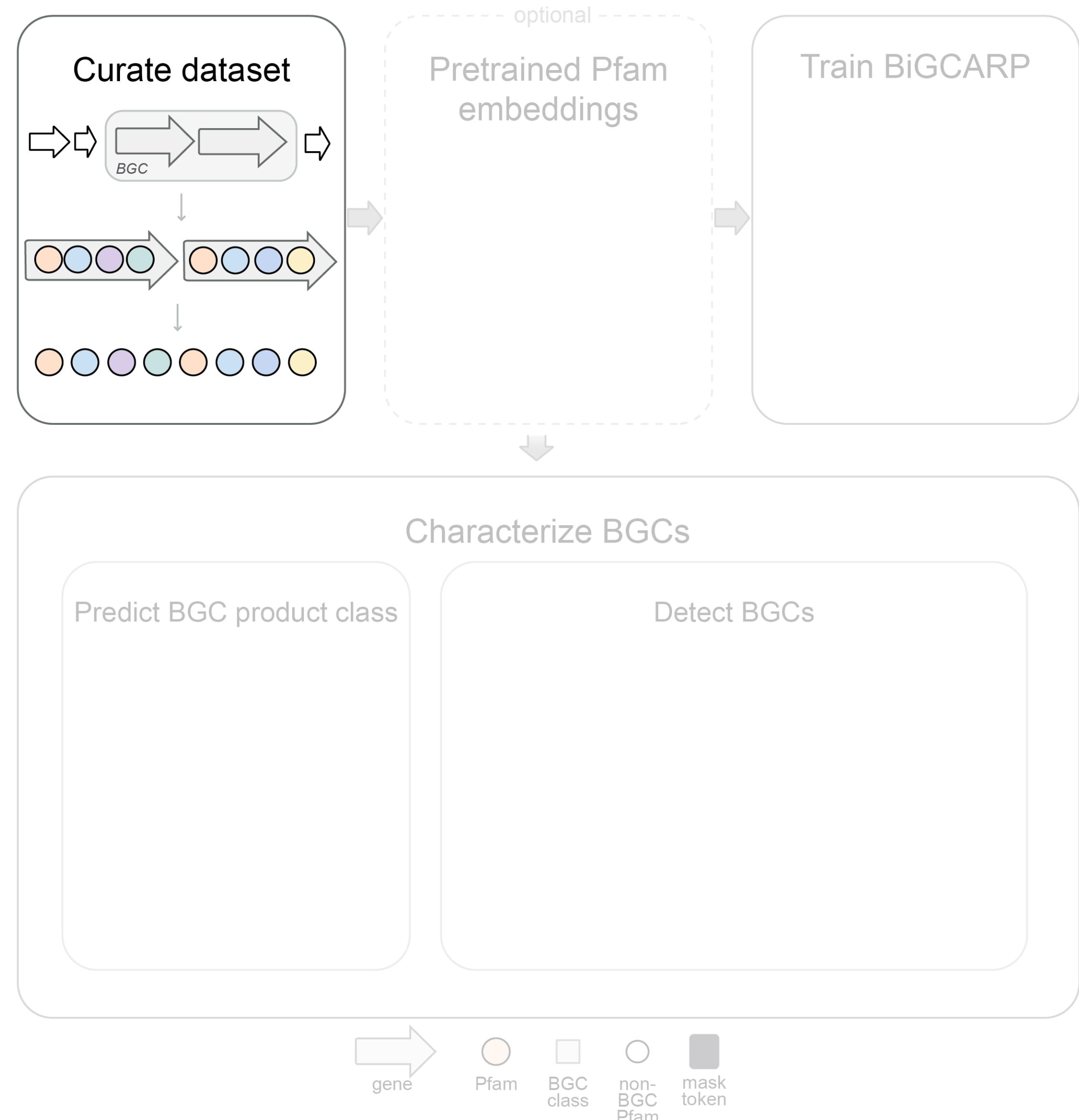
Identify BGCs with
antiSMASH



Identify Pfam domains



We represent BGCs as sequences of domains



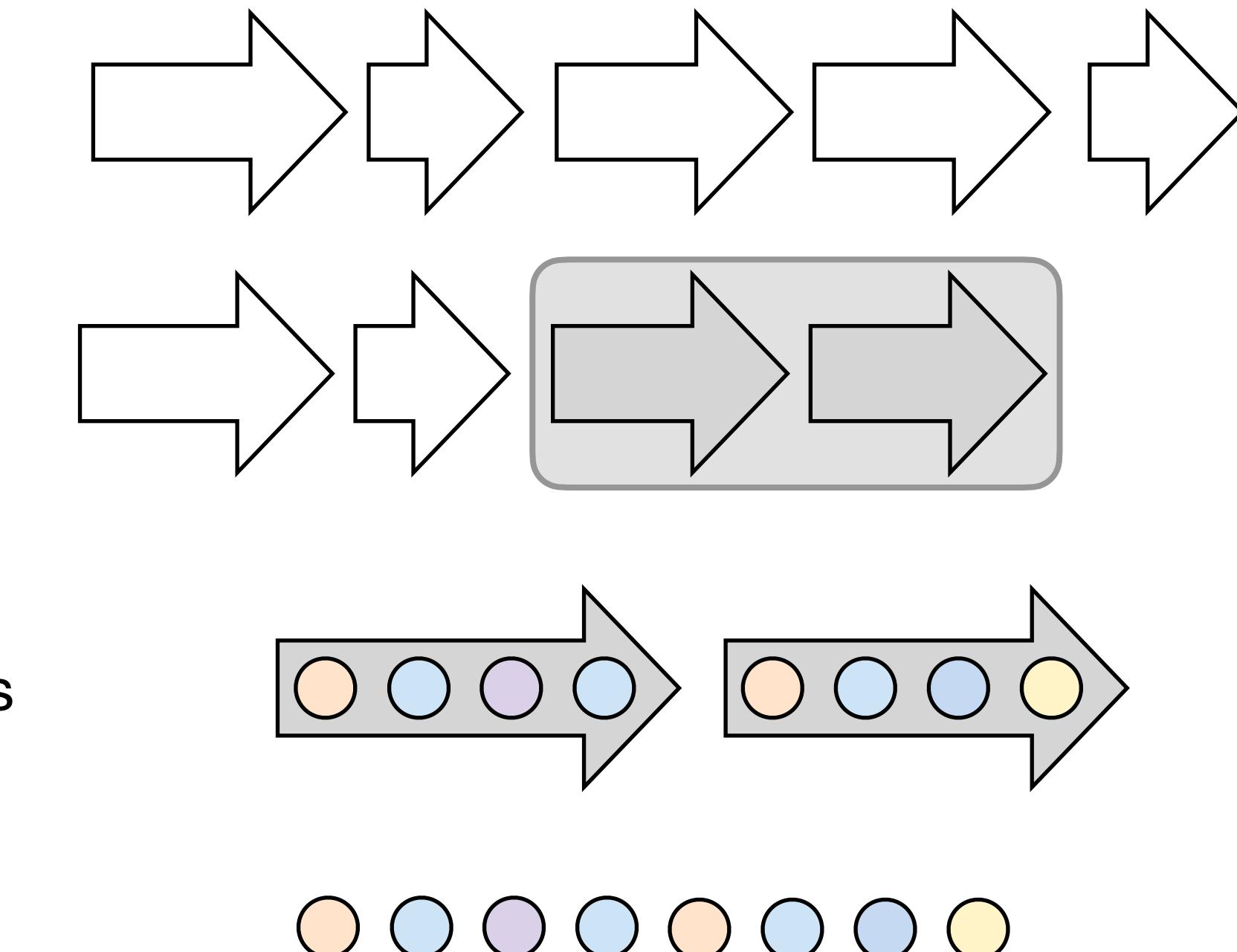
Microbial genome

Open reading frames

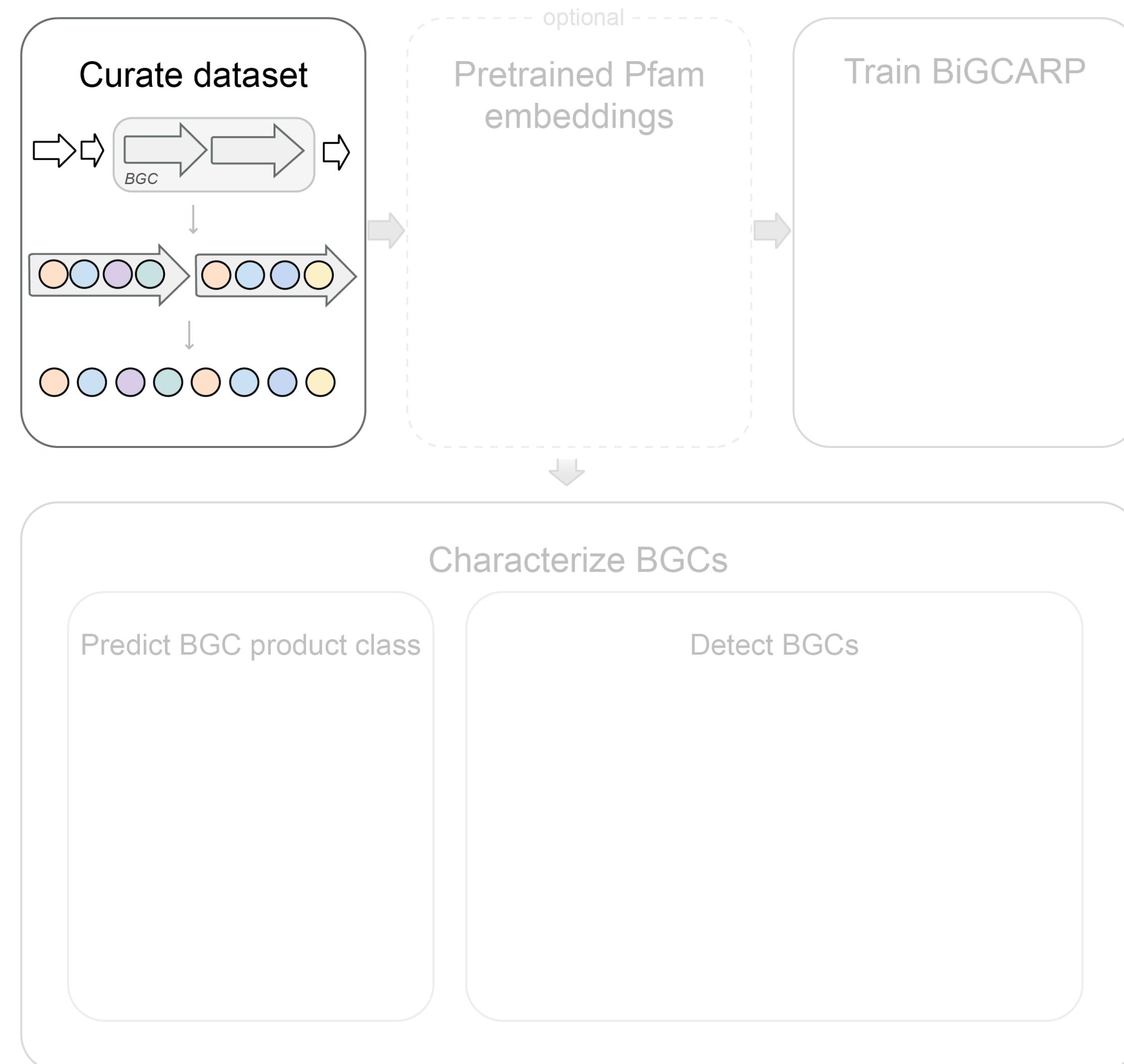
Identify BGCs with
antiSMASH

Identify Pfam domains

...ACTGCGGTTACG...



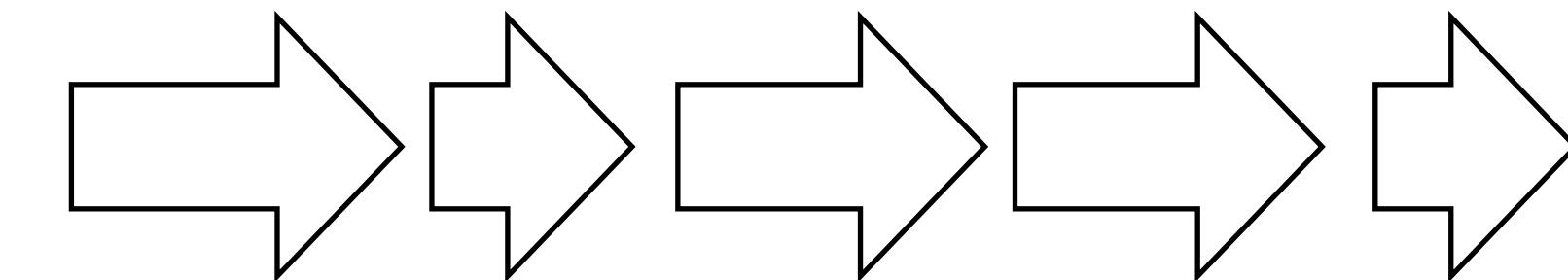
We represent BGCs as sequences of domains



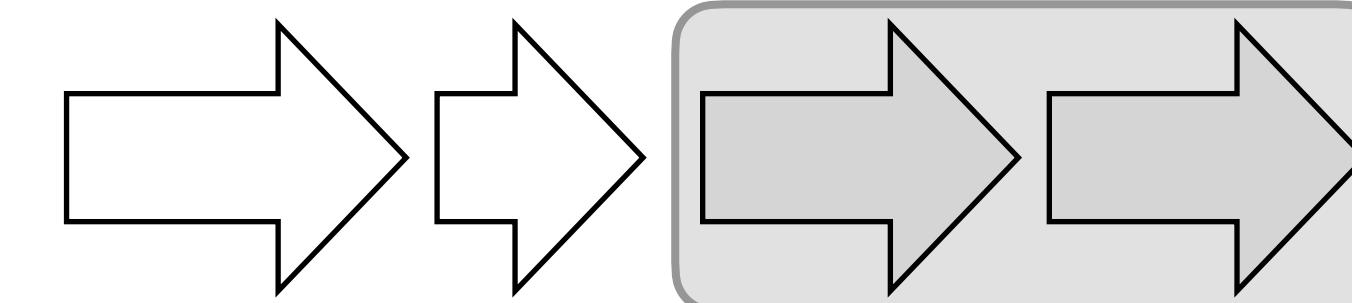
Microbial genome

...ACTGCGGTTACG...

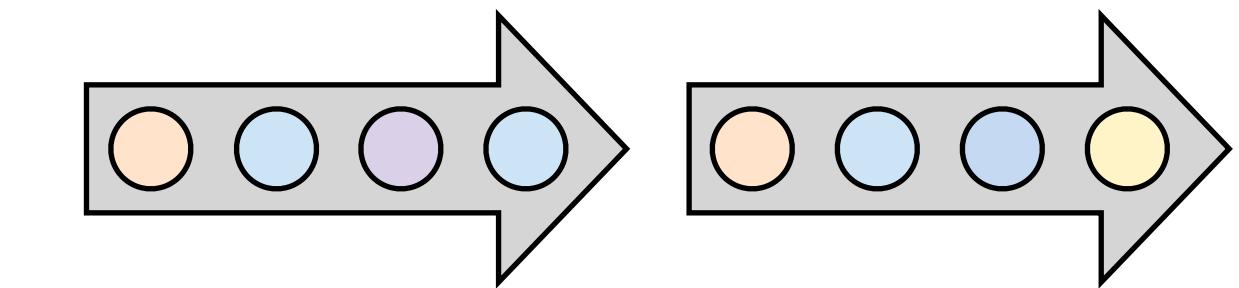
Open reading frames



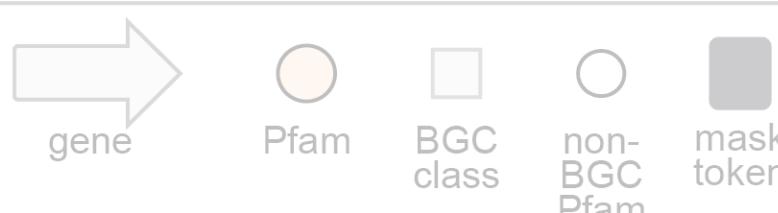
Identify BGCs with
antiSMASH



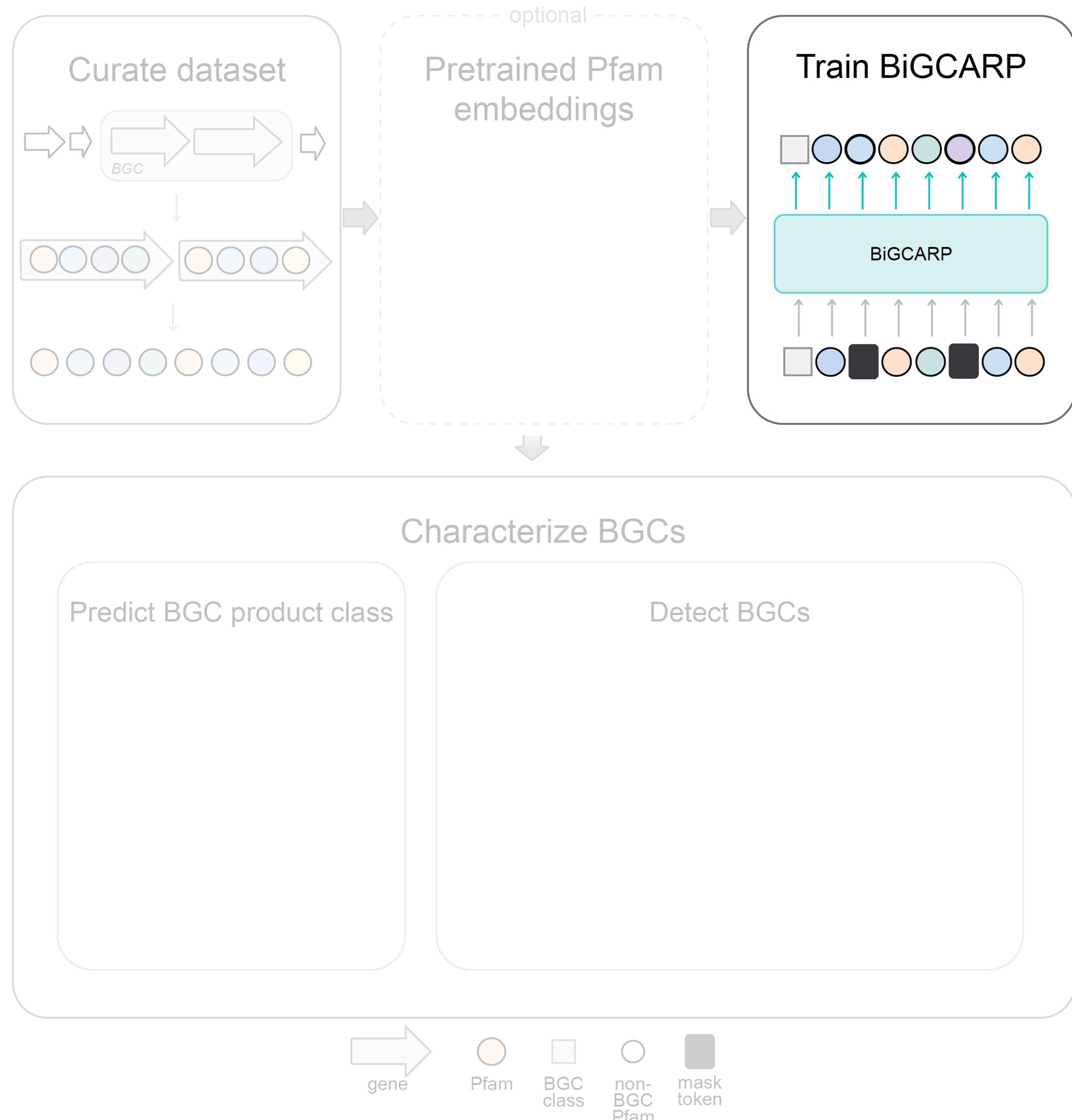
Identify Pfam domains



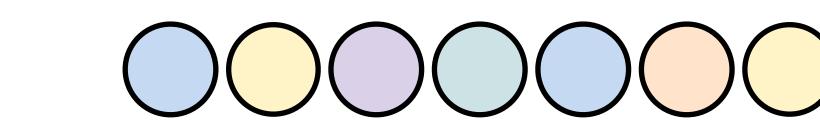
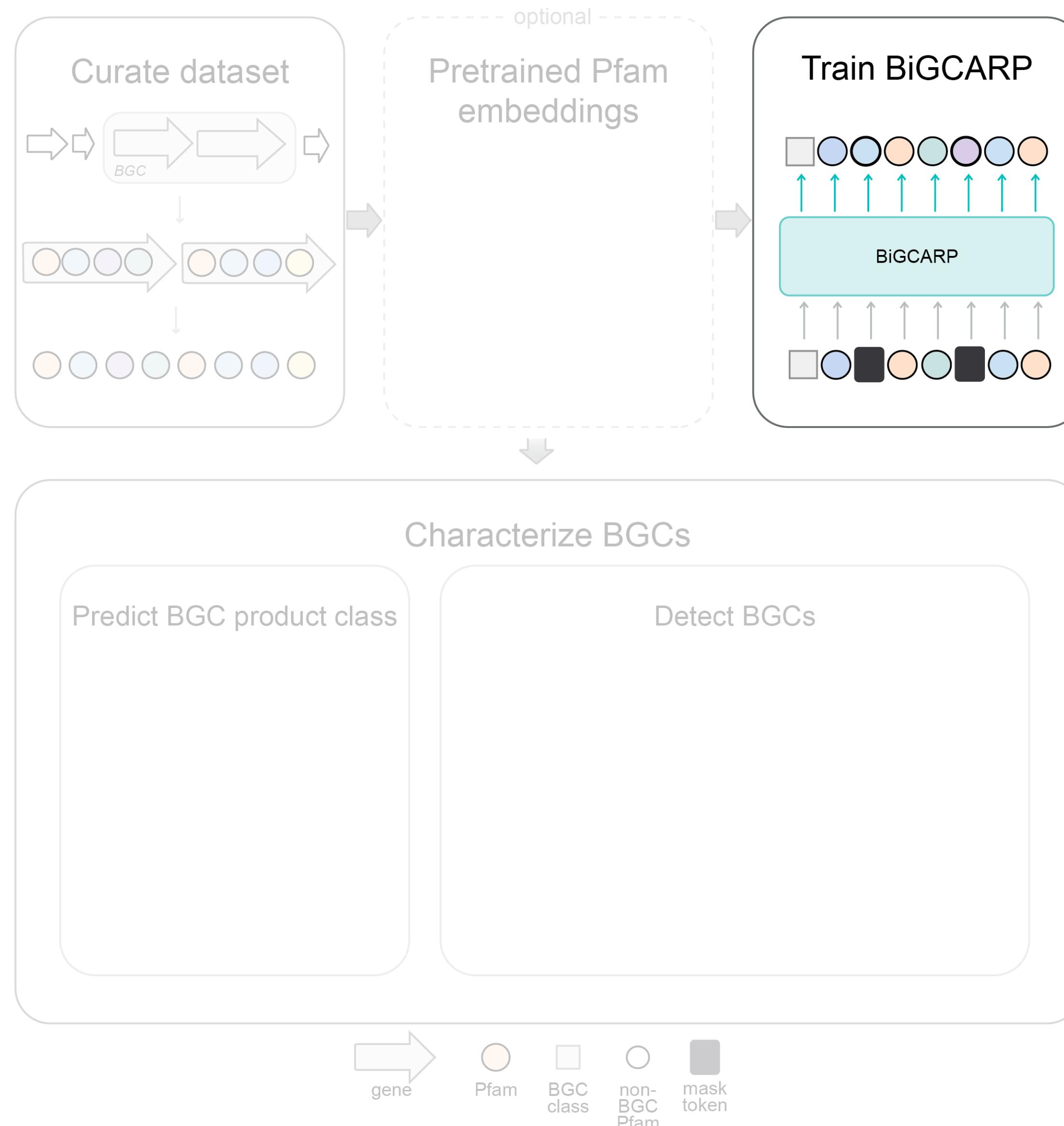
127k BGCs, 19.5k unique domain, 55 product classes



BiGCARP is trained to reconstruct corrupted BGCs

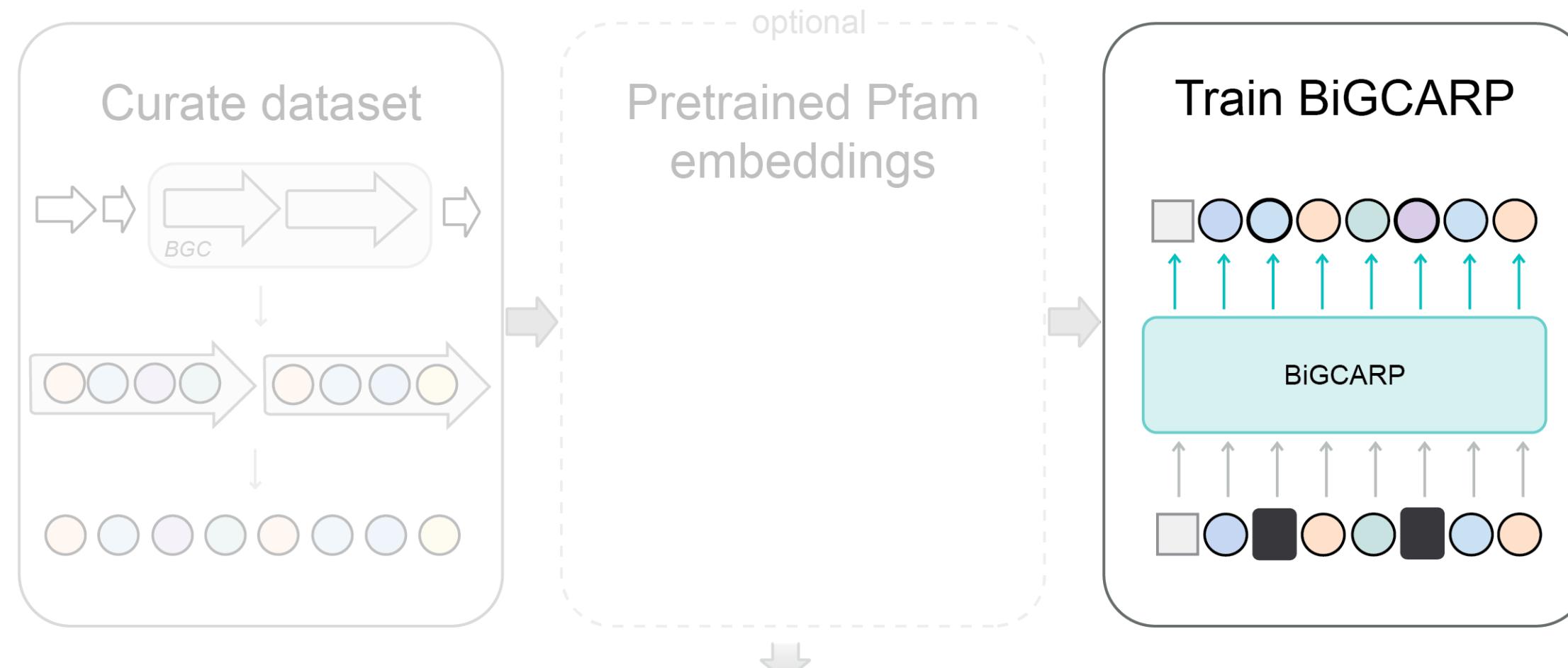


BiGCARP is trained to reconstruct corrupted BGCs

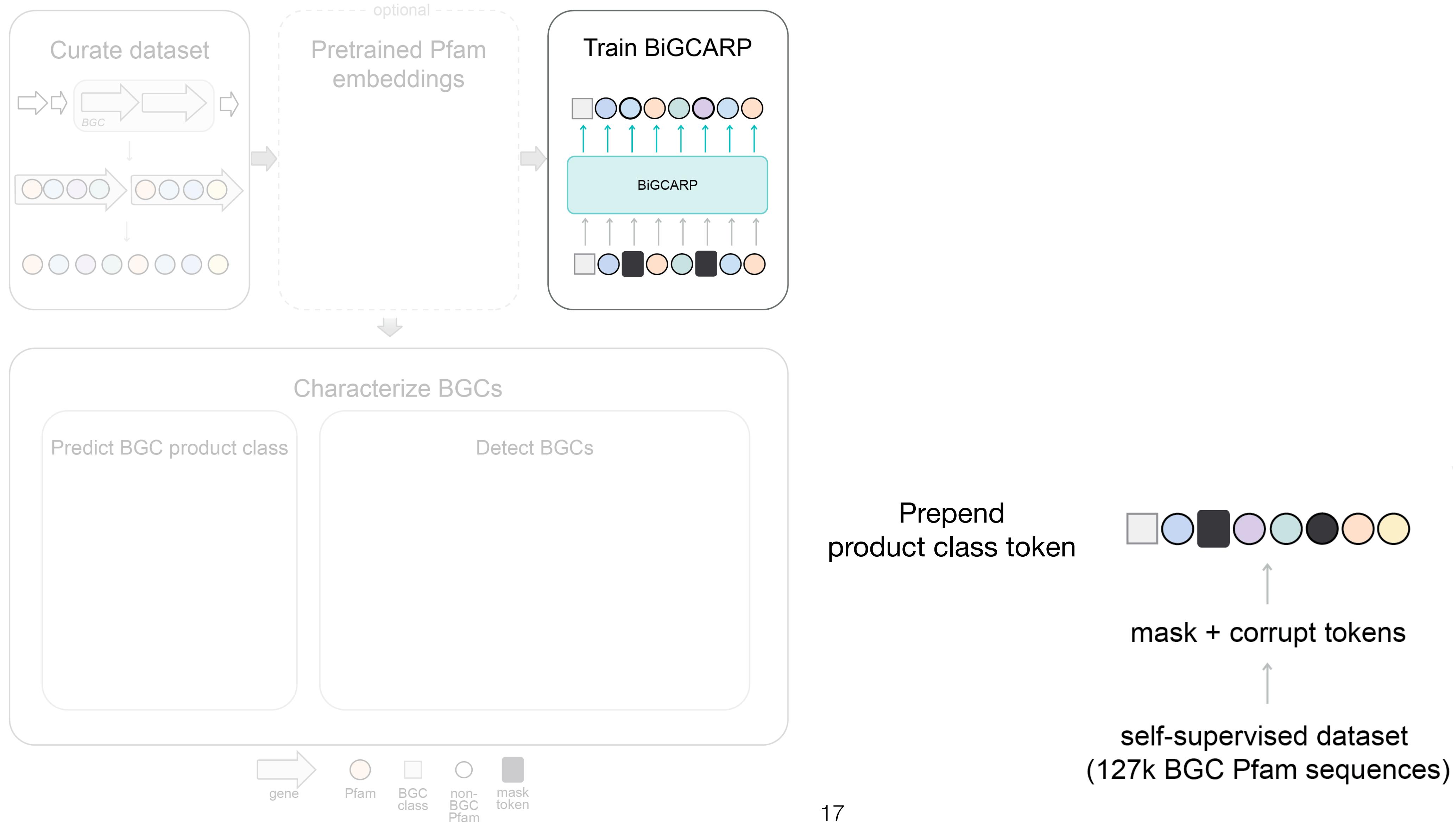


self-supervised dataset
(127k BGC Pfam sequences)

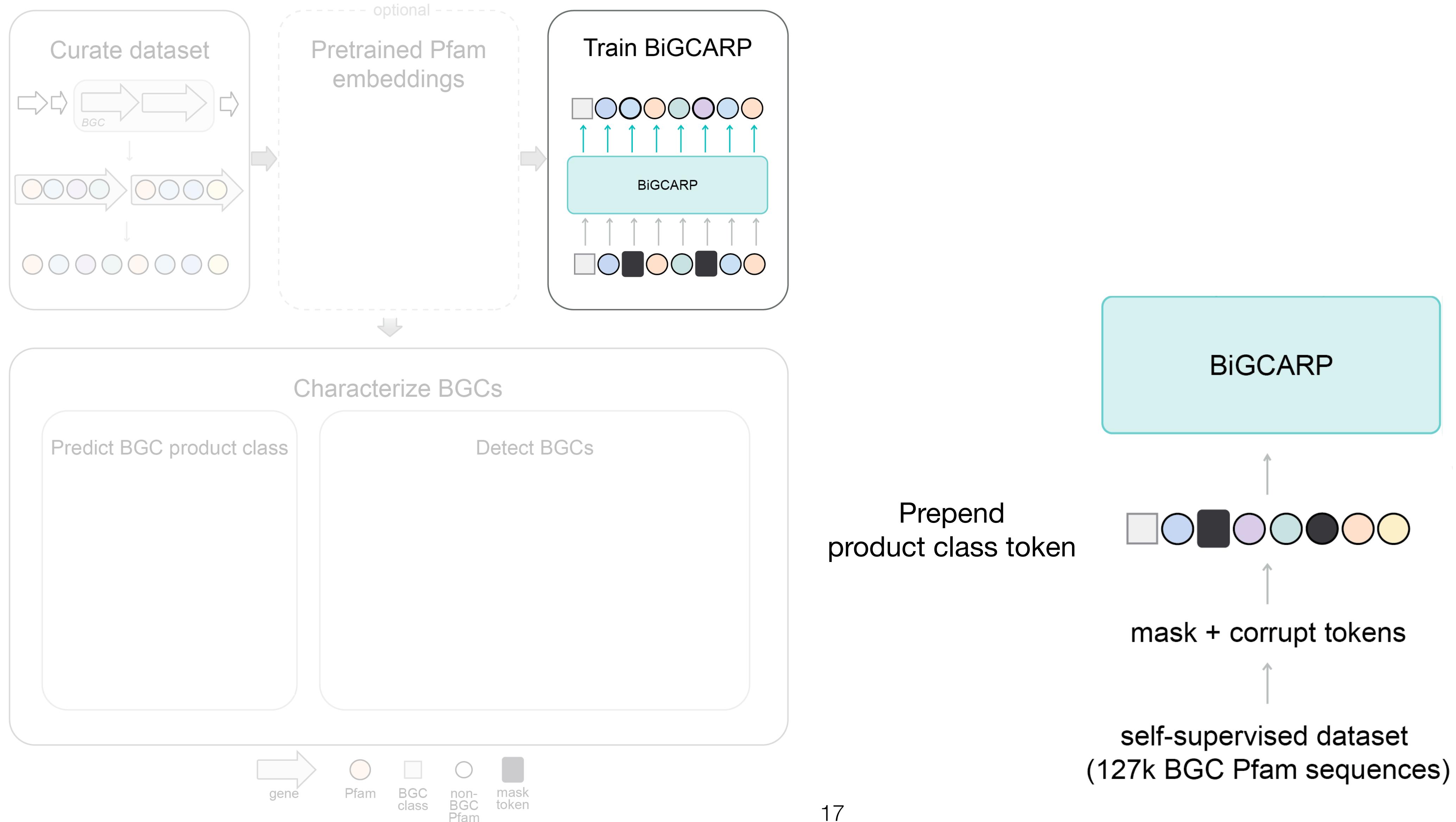
BiGCARP is trained to reconstruct corrupted BGCs



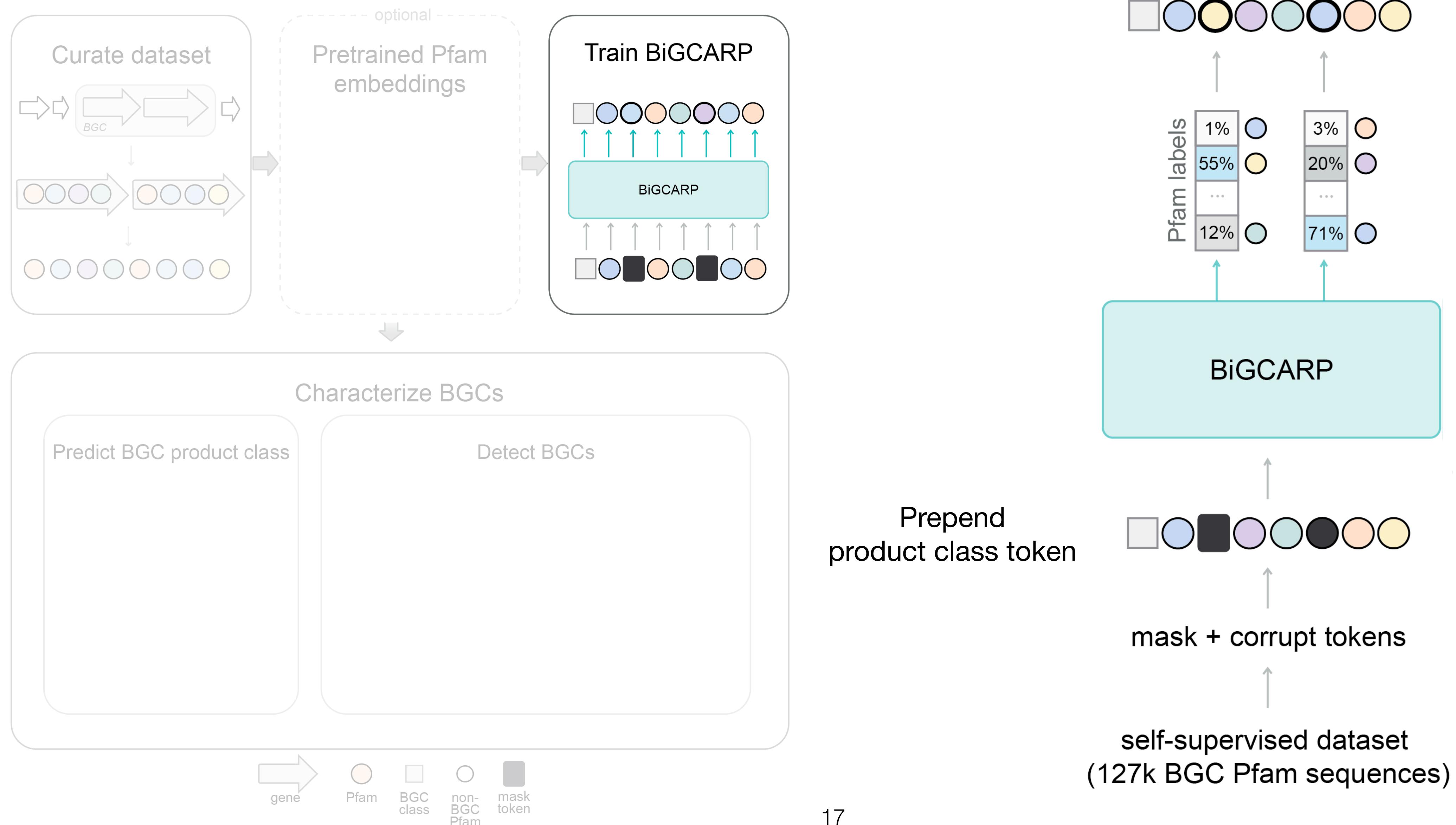
BiGCARP is trained to reconstruct corrupted BGCs



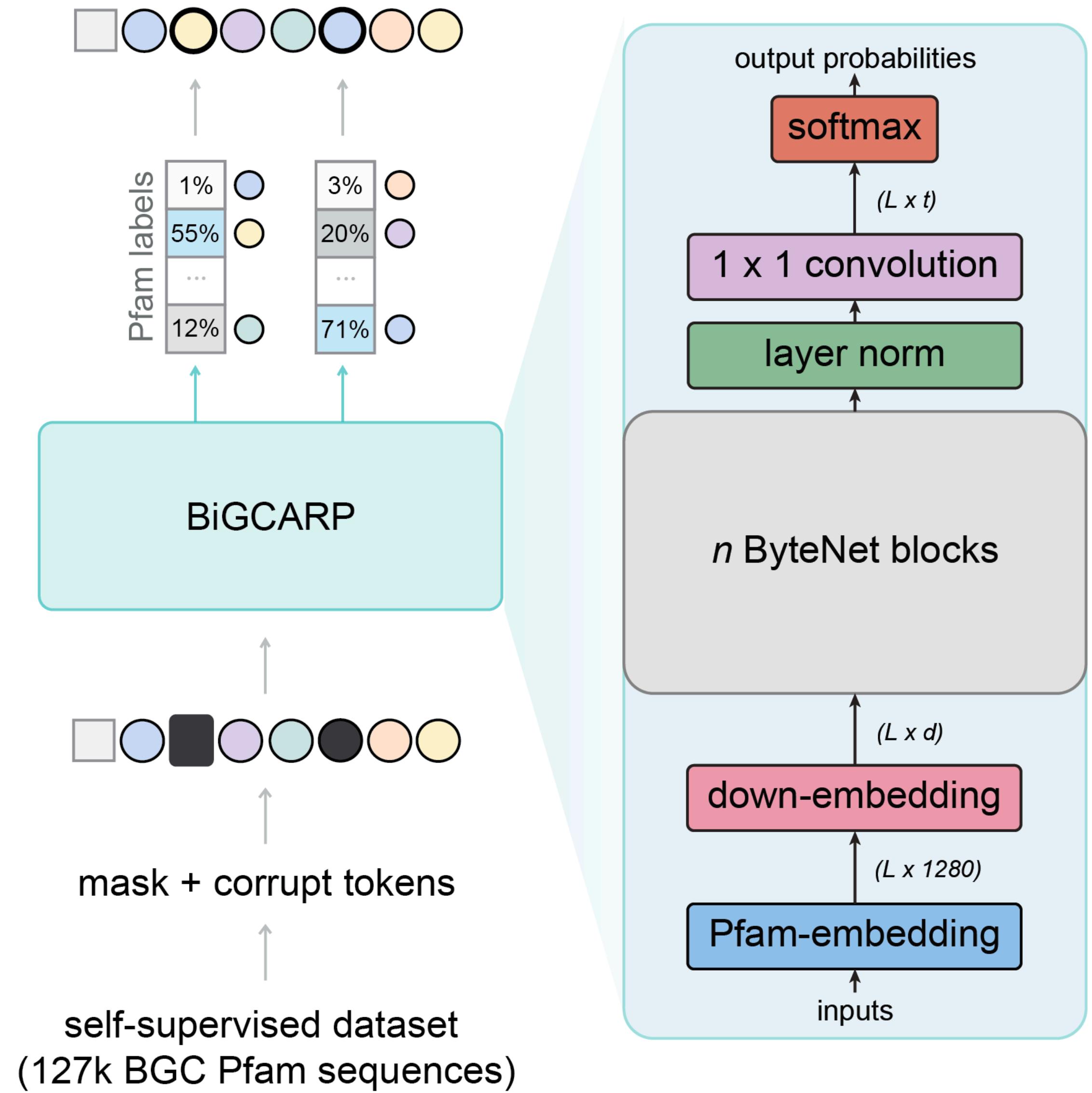
BiGCARP is trained to reconstruct corrupted BGCs



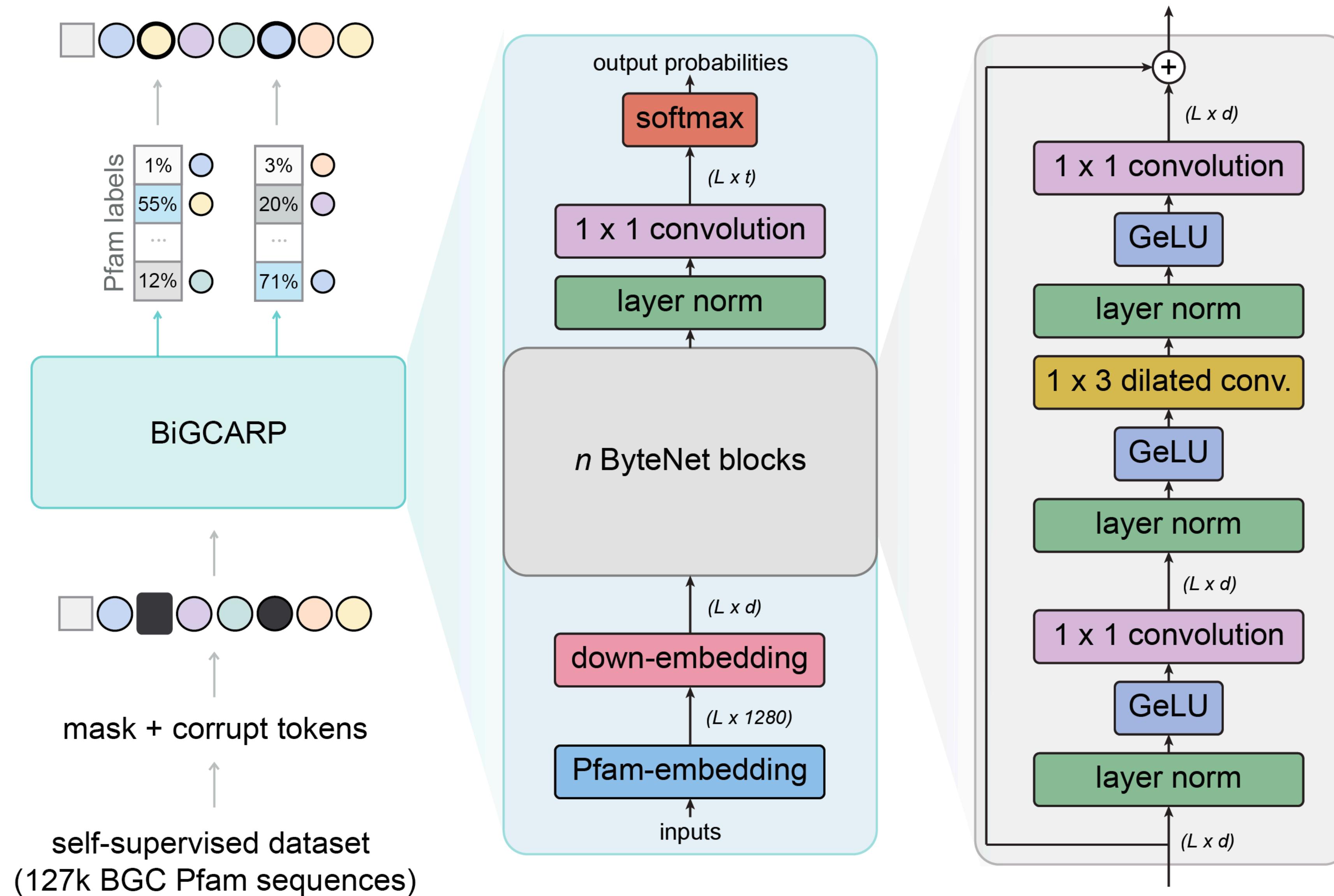
BiGCARP is trained to reconstruct corrupted BGCs



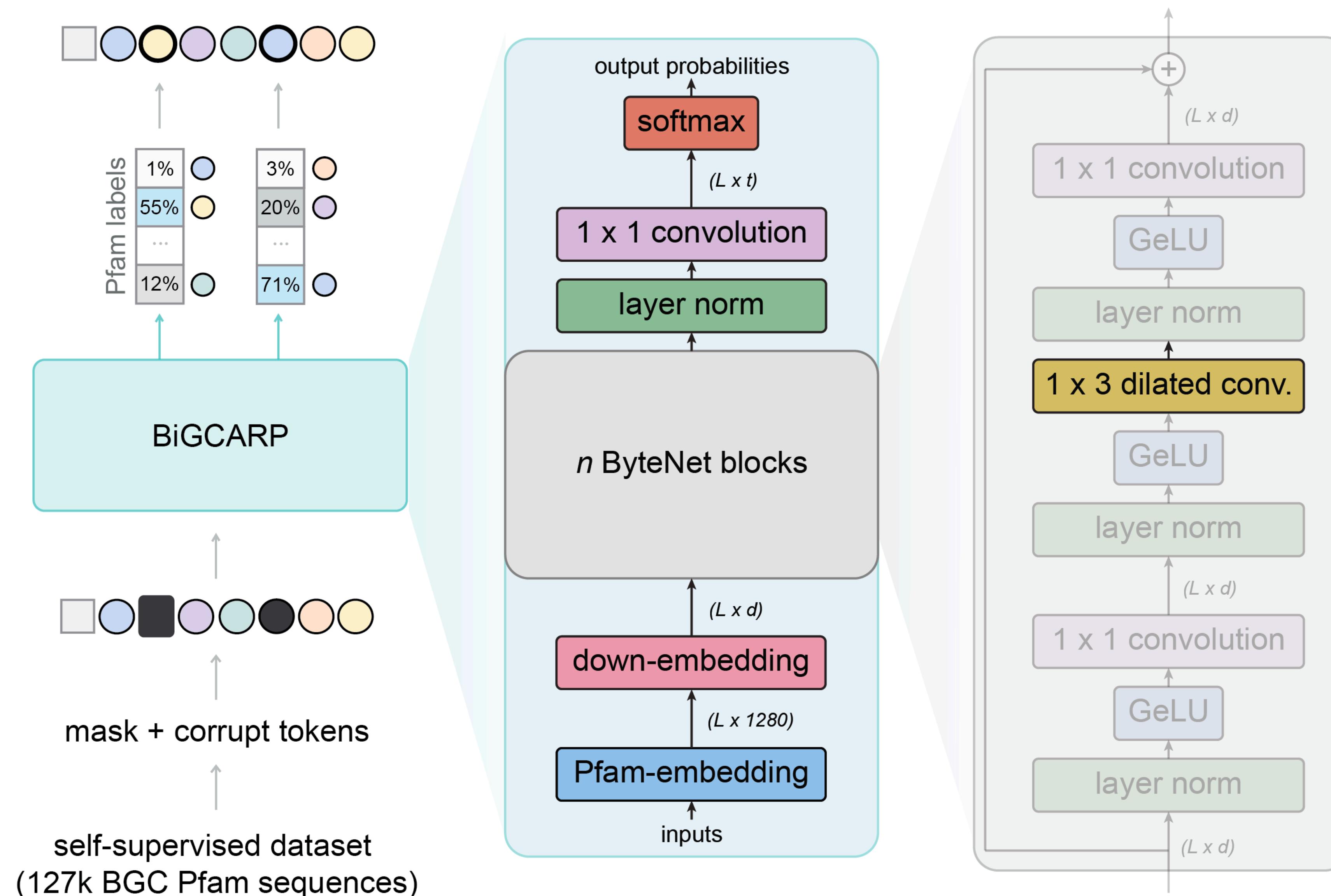
BiGCARP is trained to reconstruct corrupted BGCs



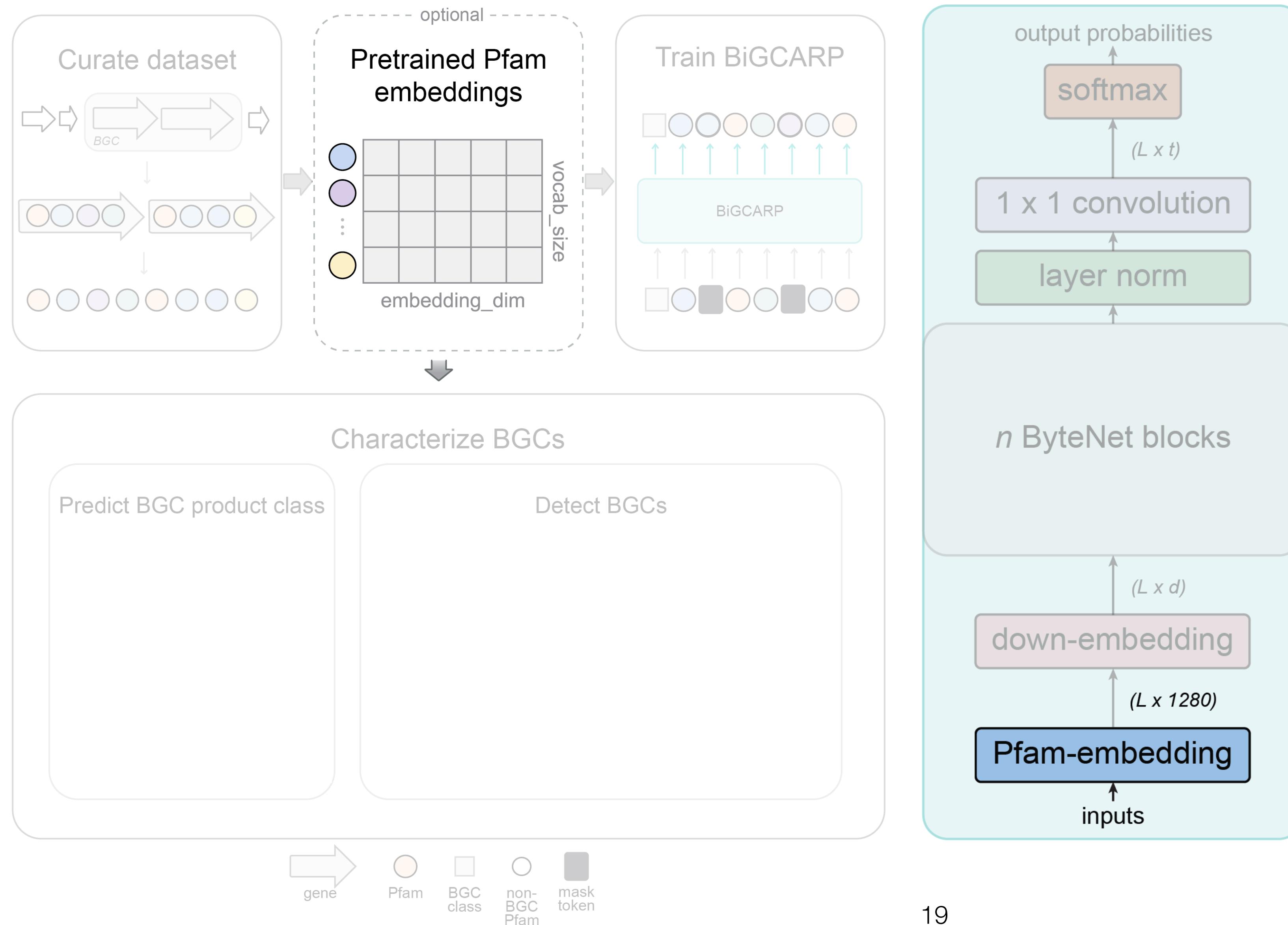
BiGCARP is trained to reconstruct corrupted BGCs



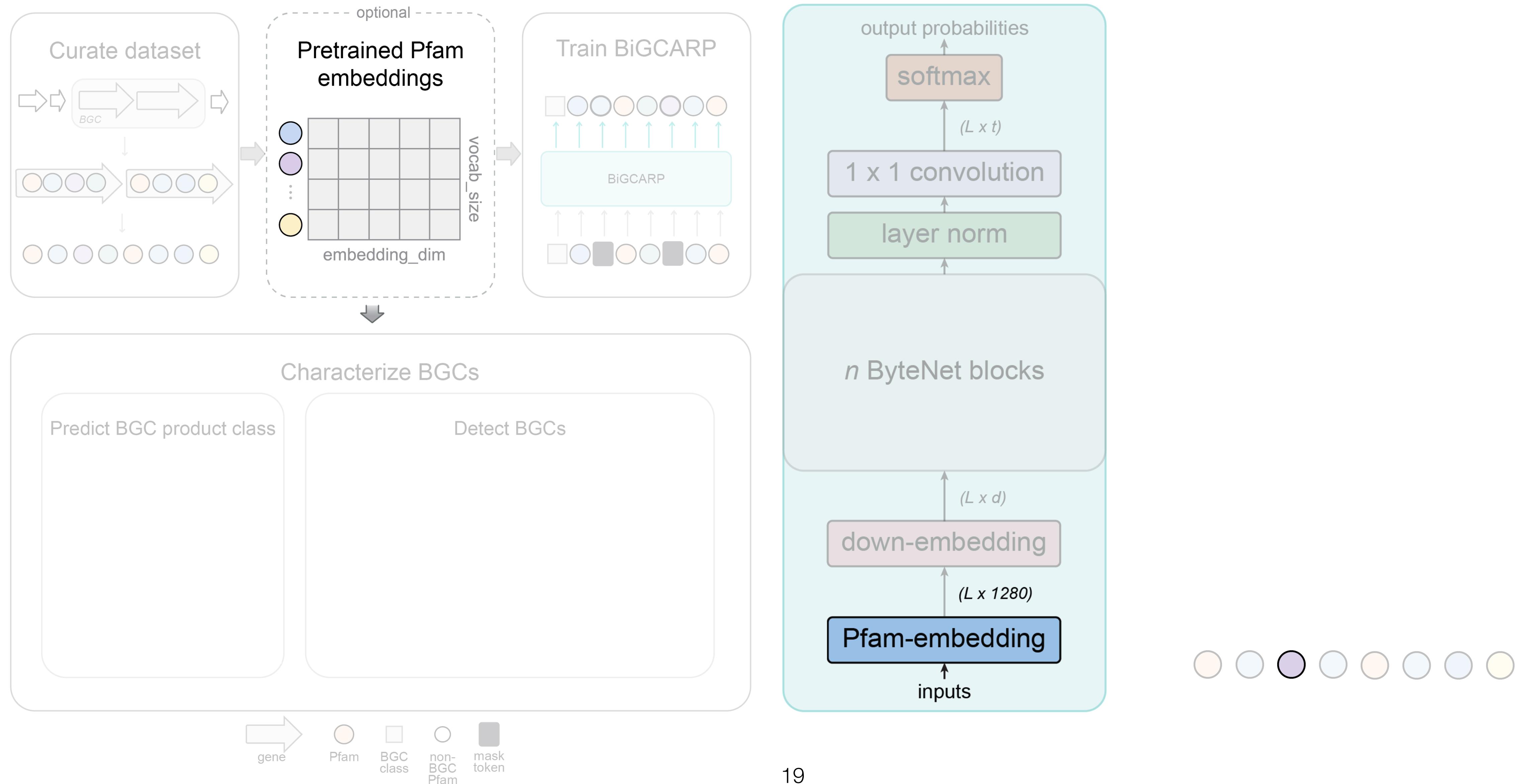
BiGCARP is trained to reconstruct corrupted BGCs



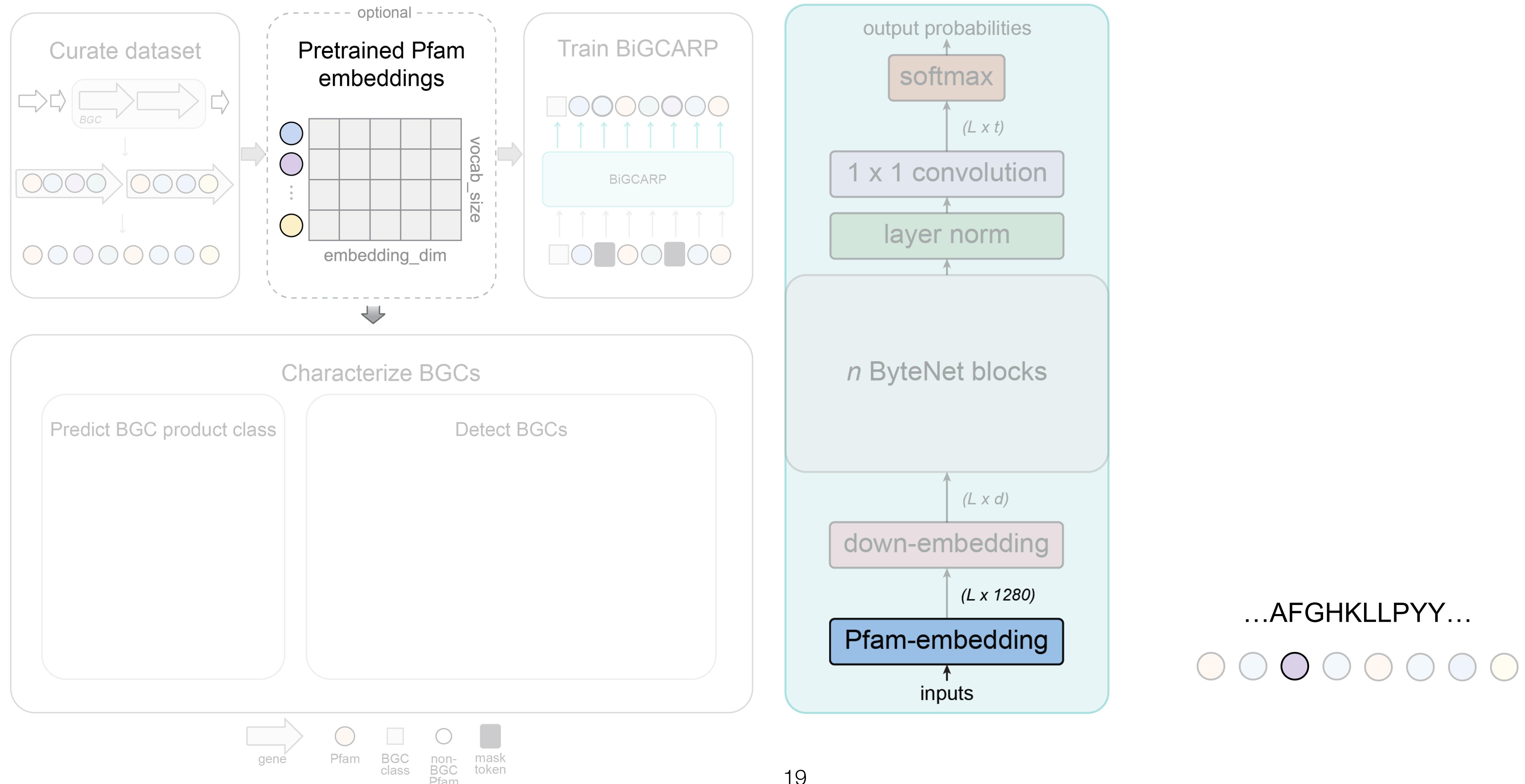
Try 3 different initial domain embeddings



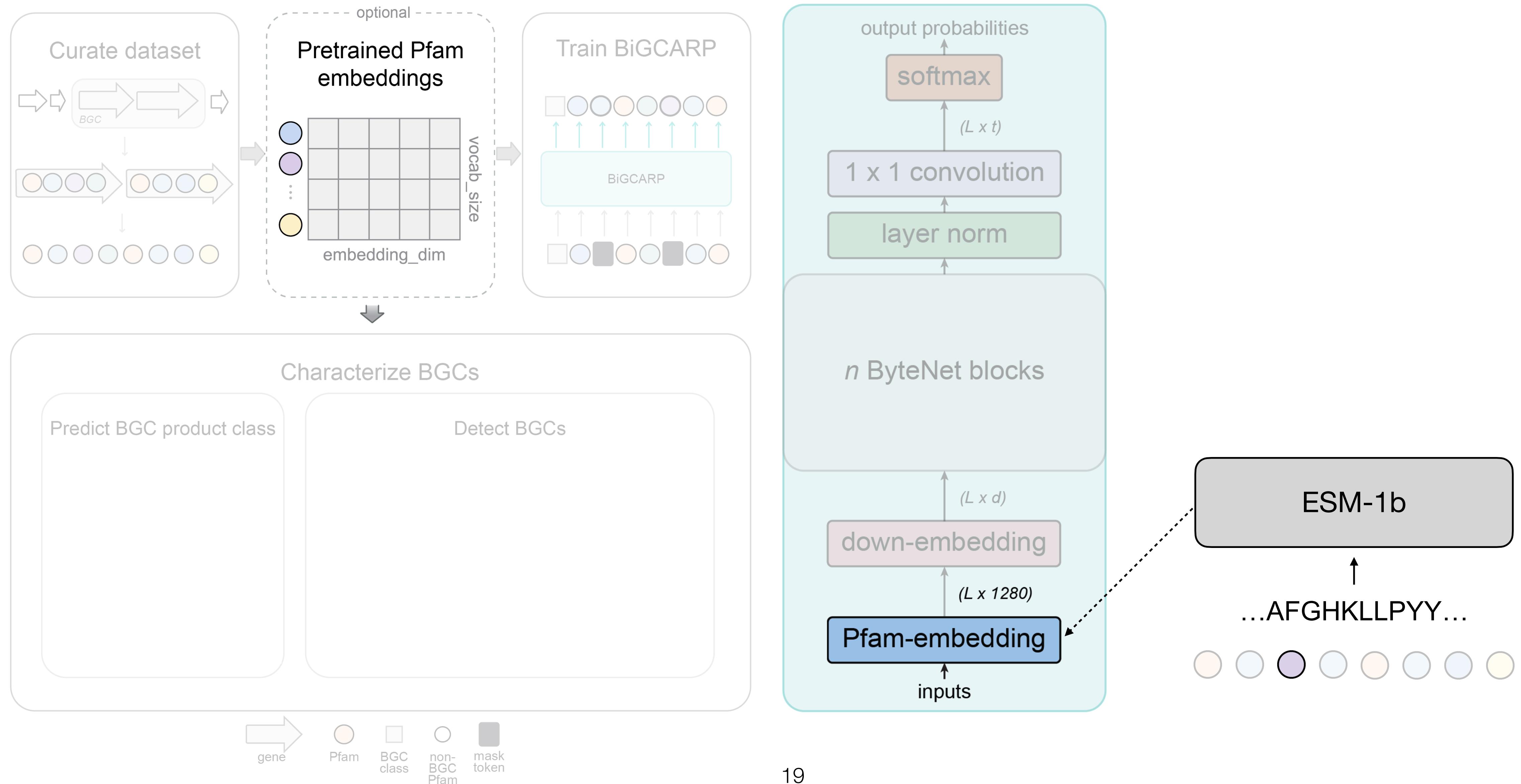
Try 3 different initial domain embeddings



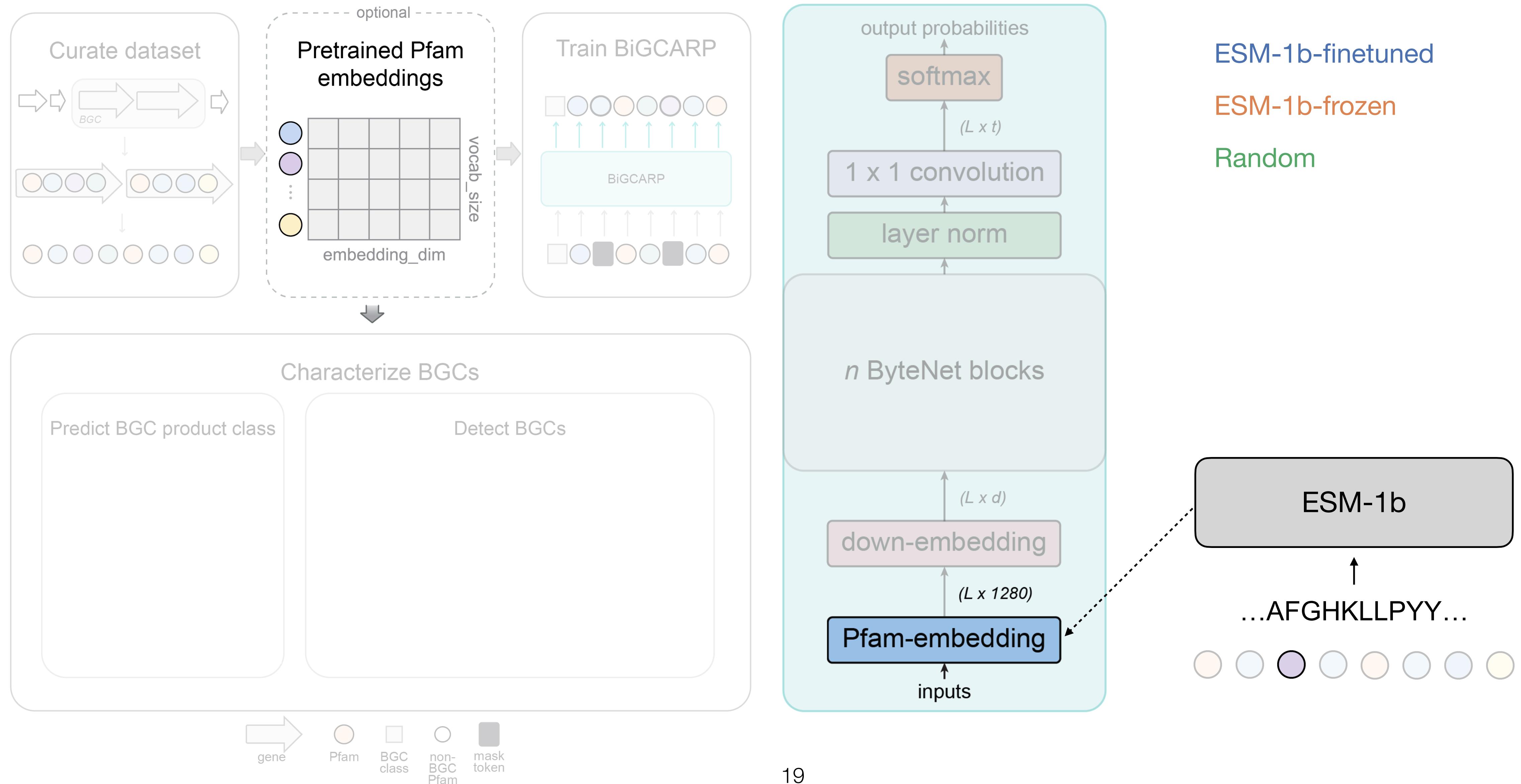
Try 3 different initial domain embeddings



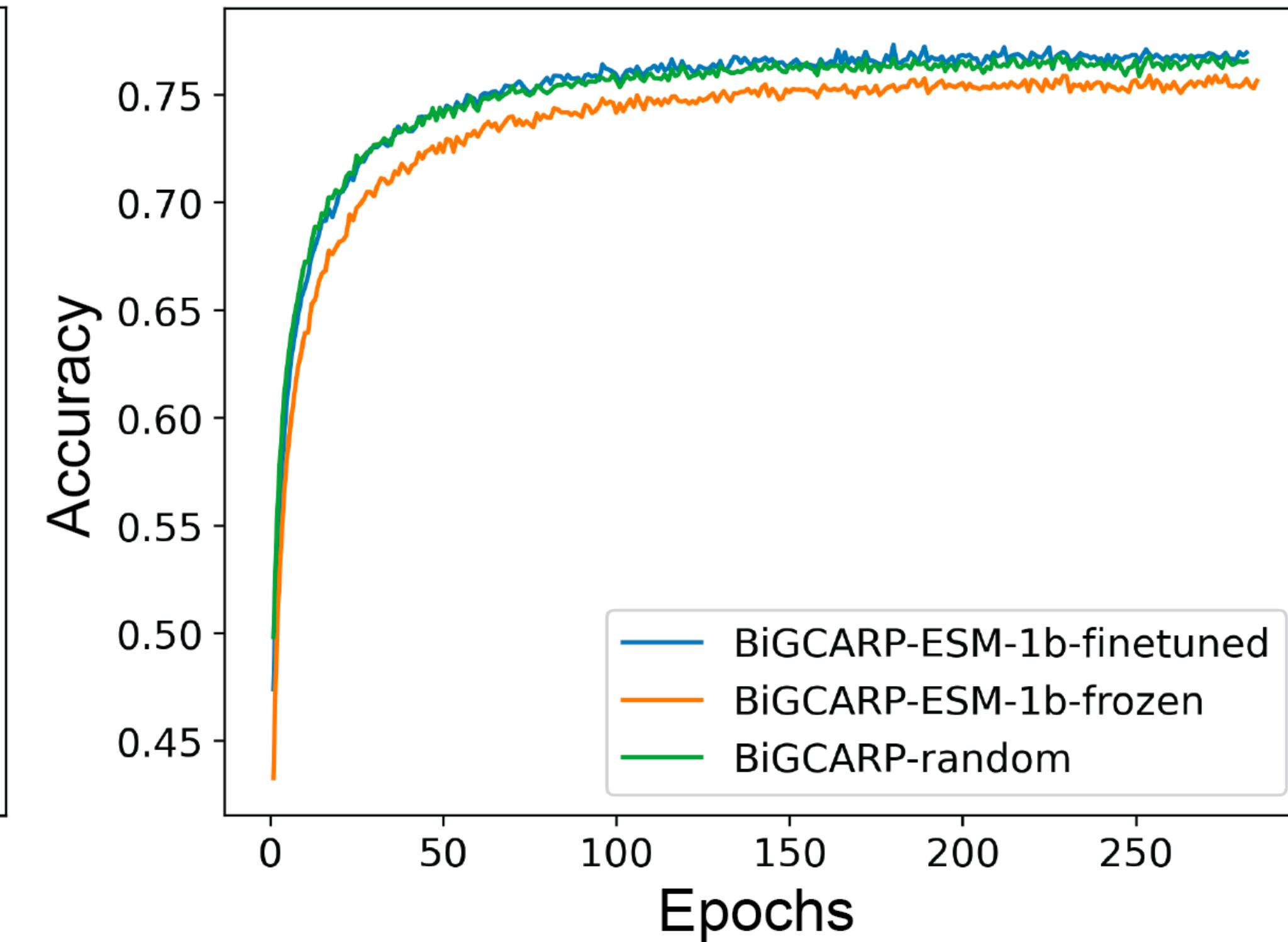
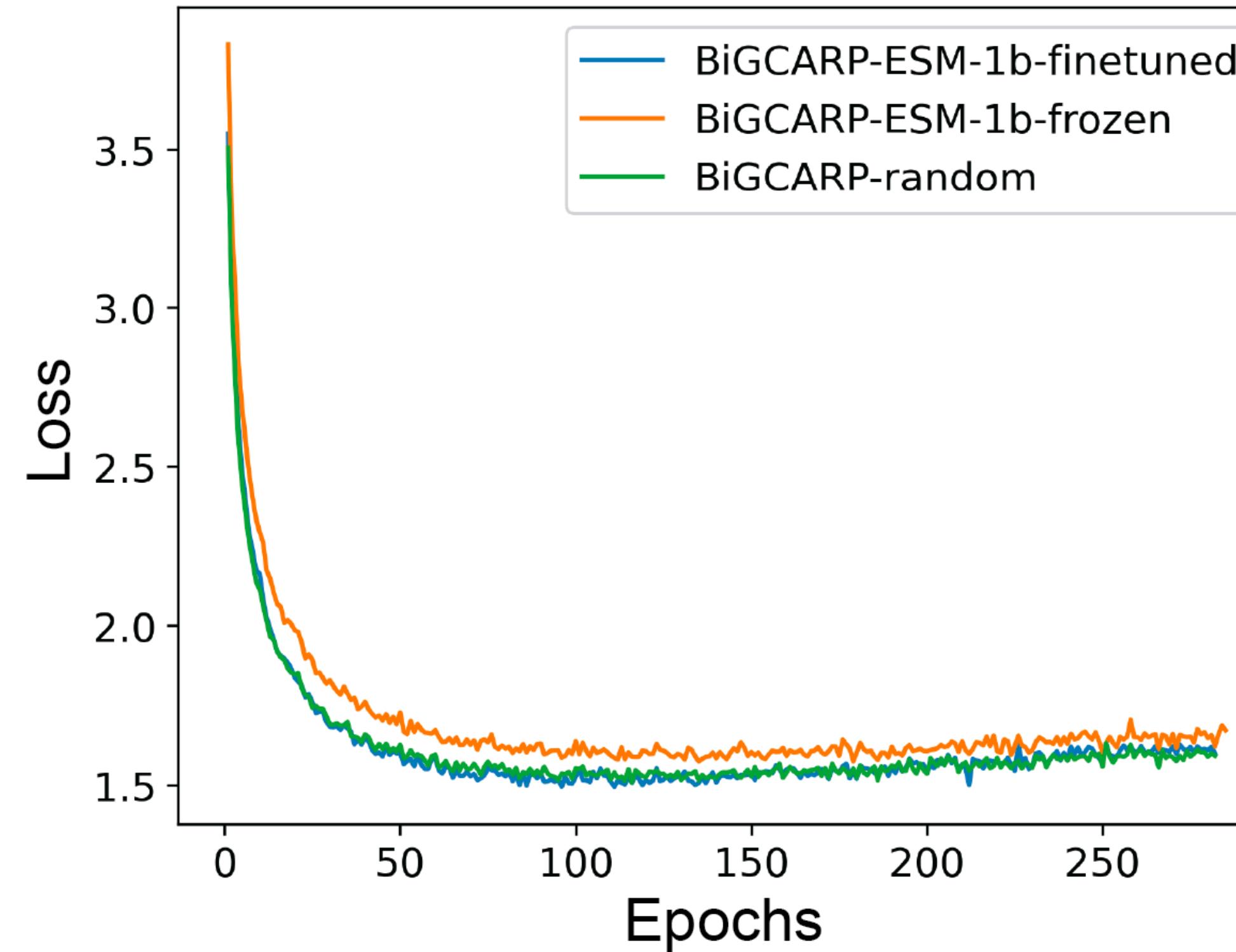
Try 3 different initial domain embeddings



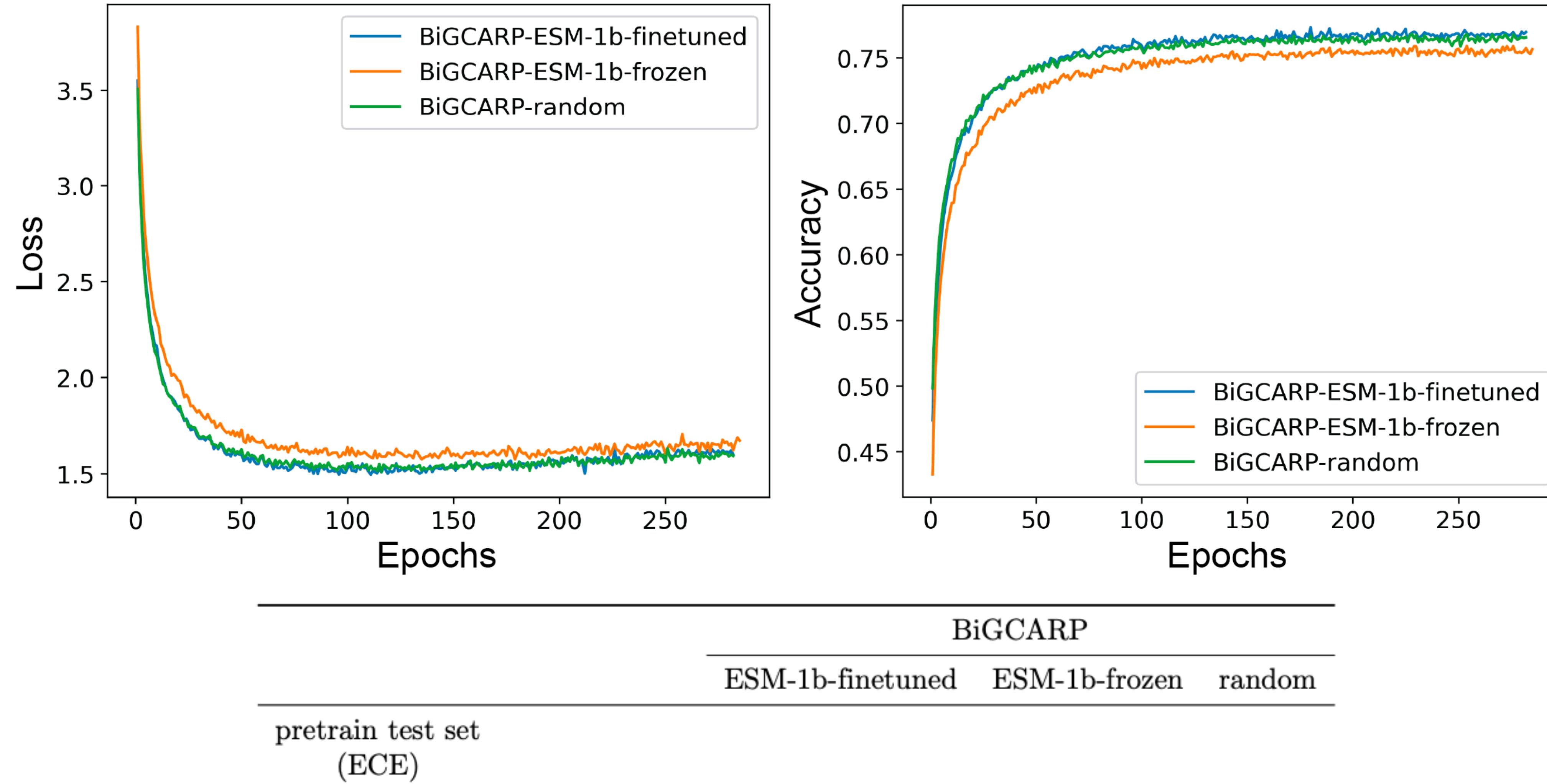
Try 3 different initial domain embeddings



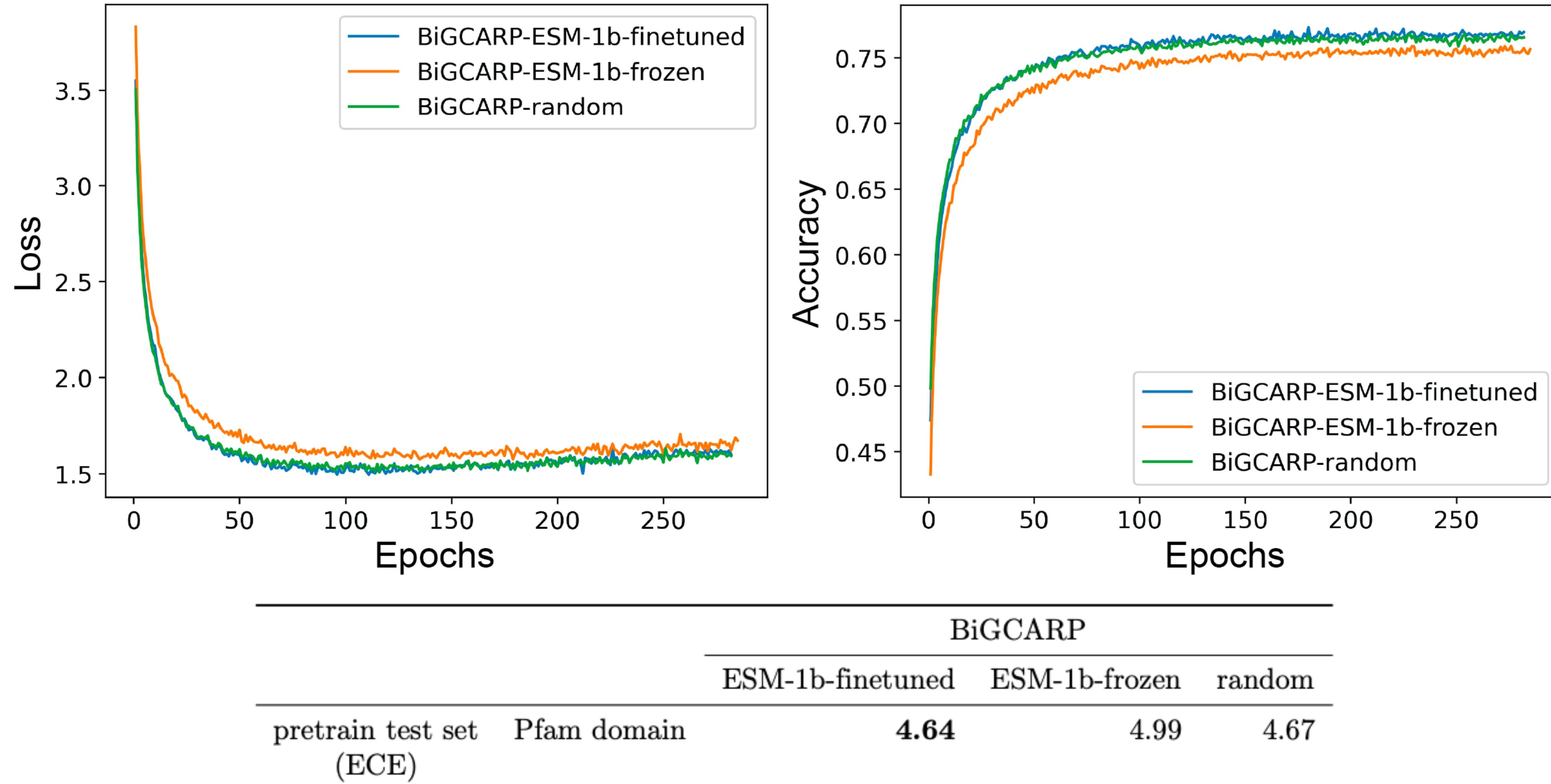
BiGCARP learns to reconstruct BGC domain sequences



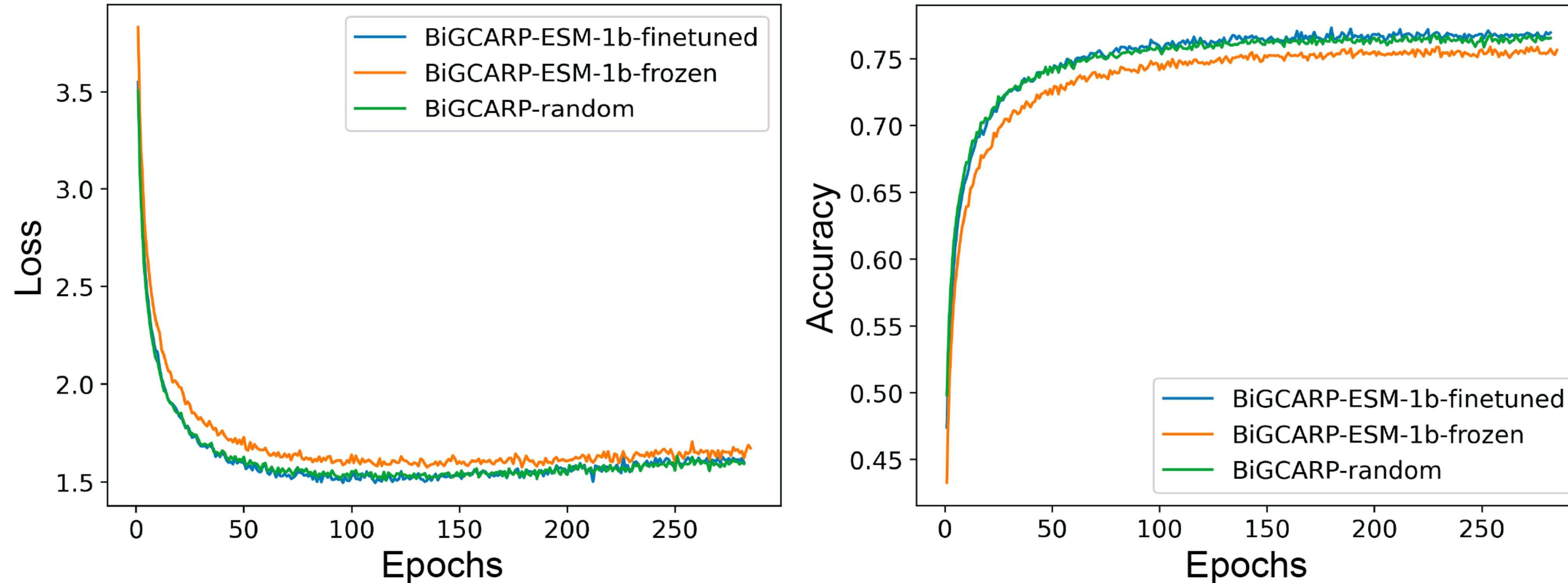
BiGCARP learns to reconstruct BGC domain sequences



BiGCARP learns to reconstruct BGC domain sequences

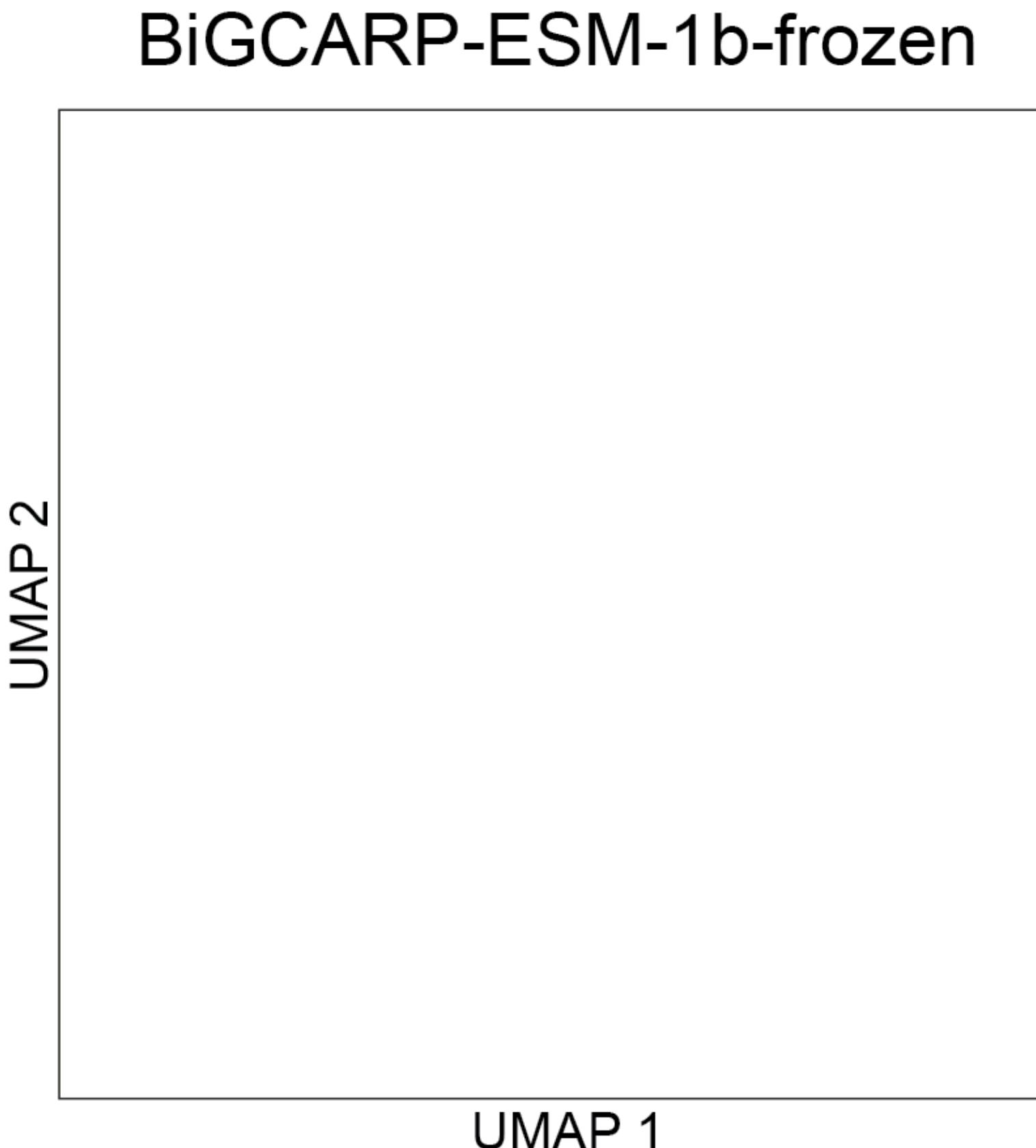


BiGCARP learns to reconstruct BGC domain sequences

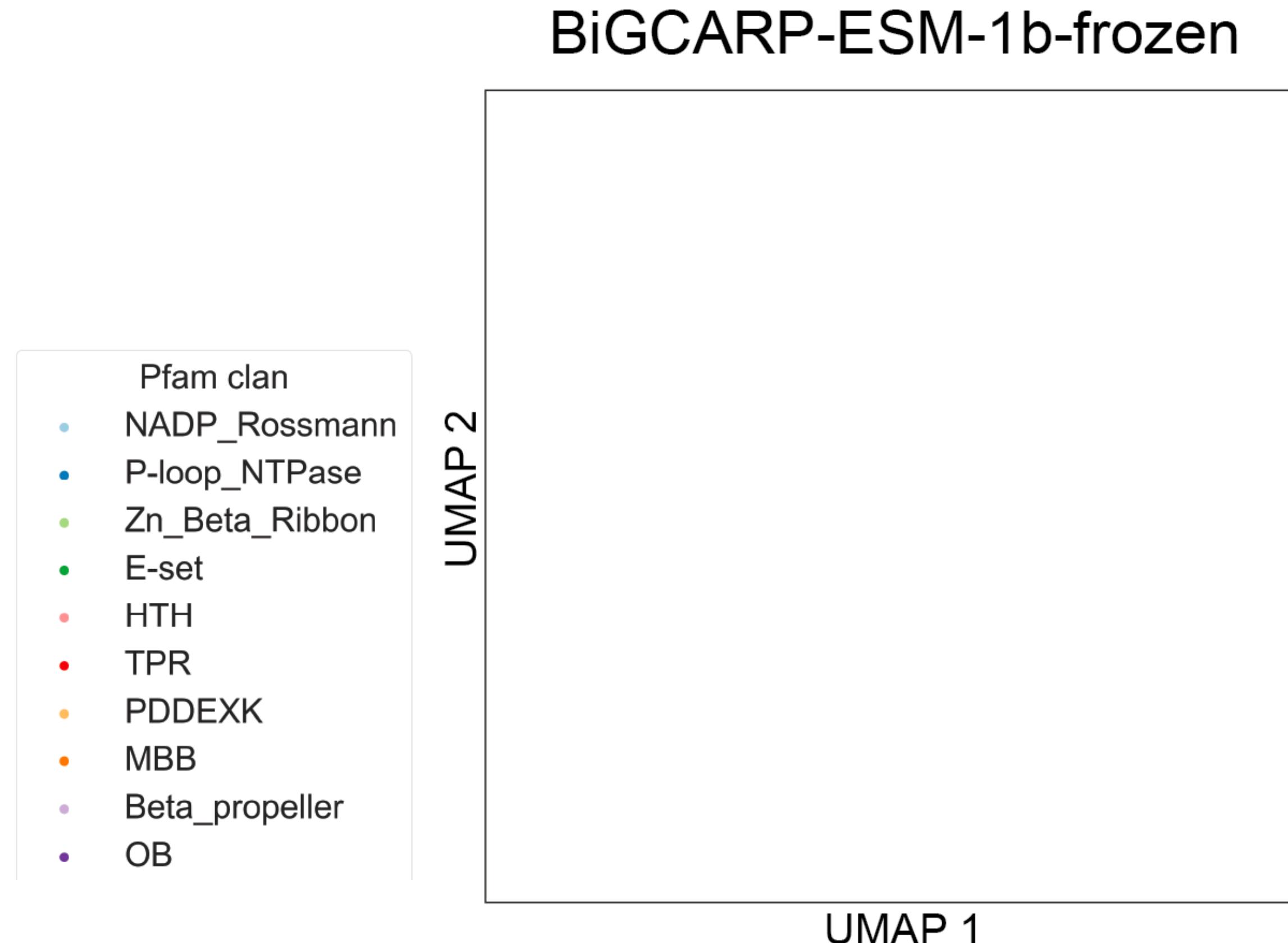


BiGCARP					
		ESM-1b-finetuned	ESM-1b-frozen	random	
pretrain	test set	Pfam domain			
	(ECE)	product class	4.64 1.50	4.99 1.46 1.50	4.67

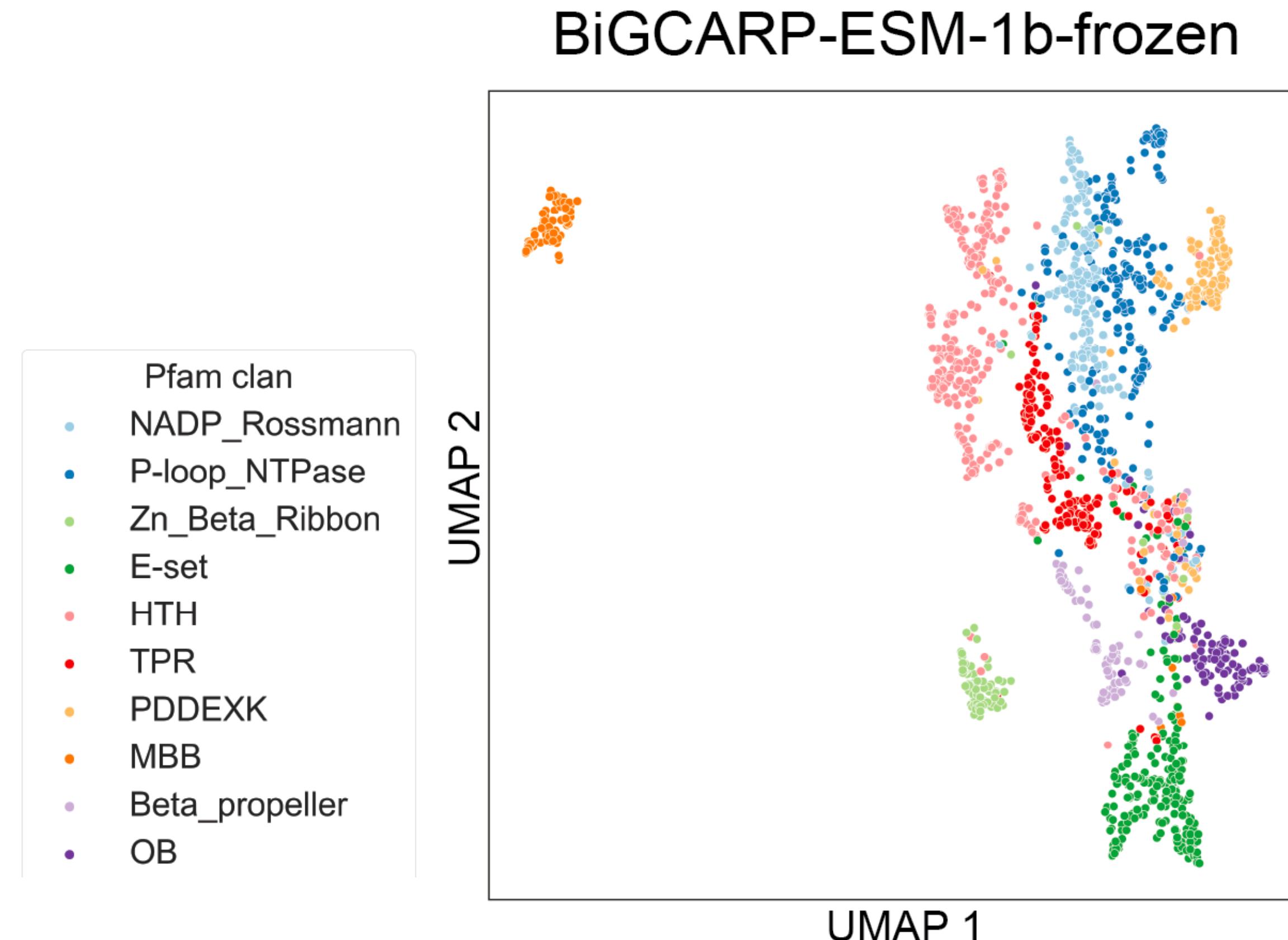
BiGCARP embeddings encode relevant Pfam information



BiGCARP embeddings encode relevant Pfam information

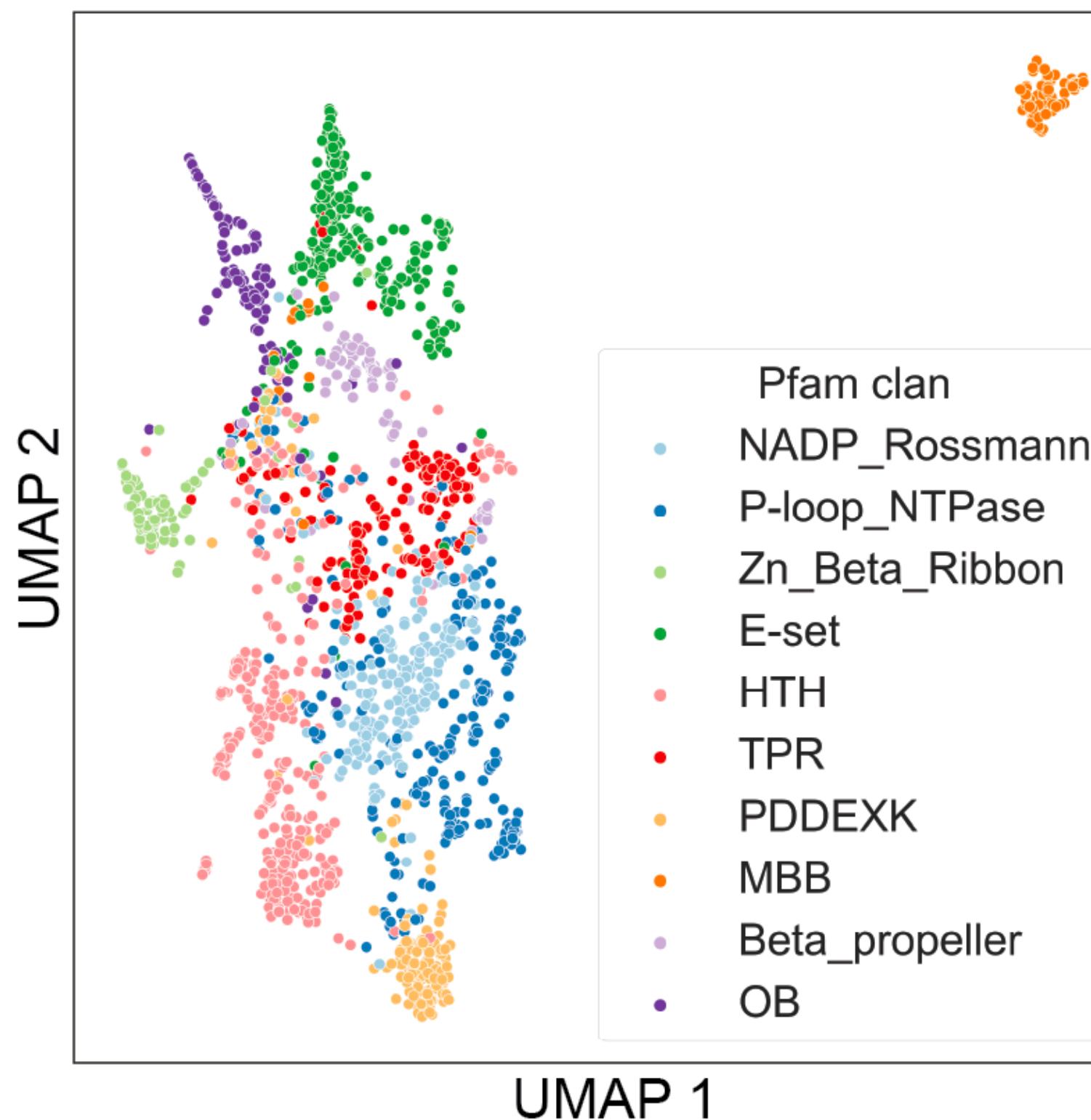


BiGCARP embeddings encode relevant Pfam information

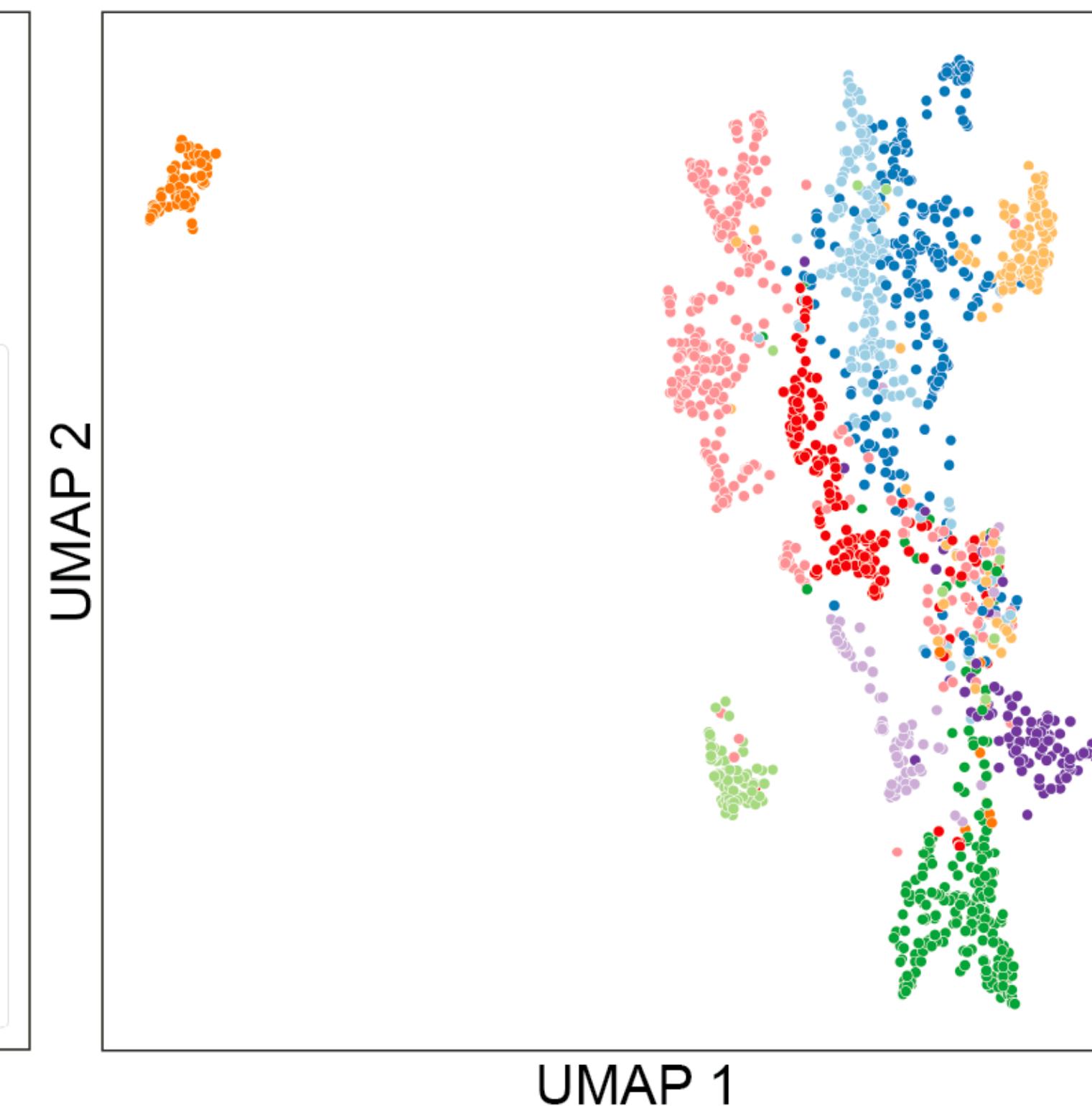


BiGCARP embeddings encode relevant Pfam information

BiGCARP-ESM-1b-finetuned

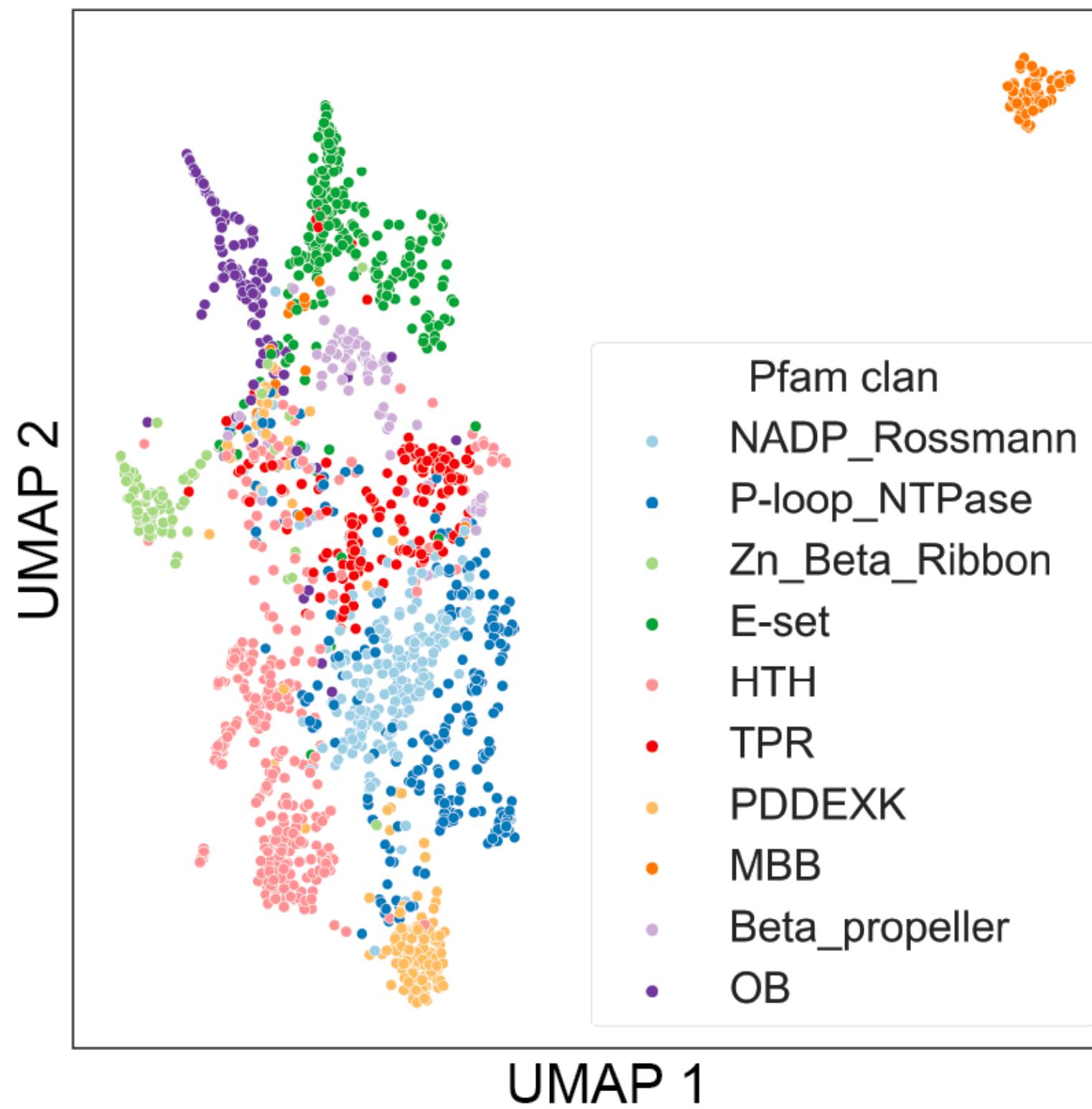


BiGCARP-ESM-1b-frozen

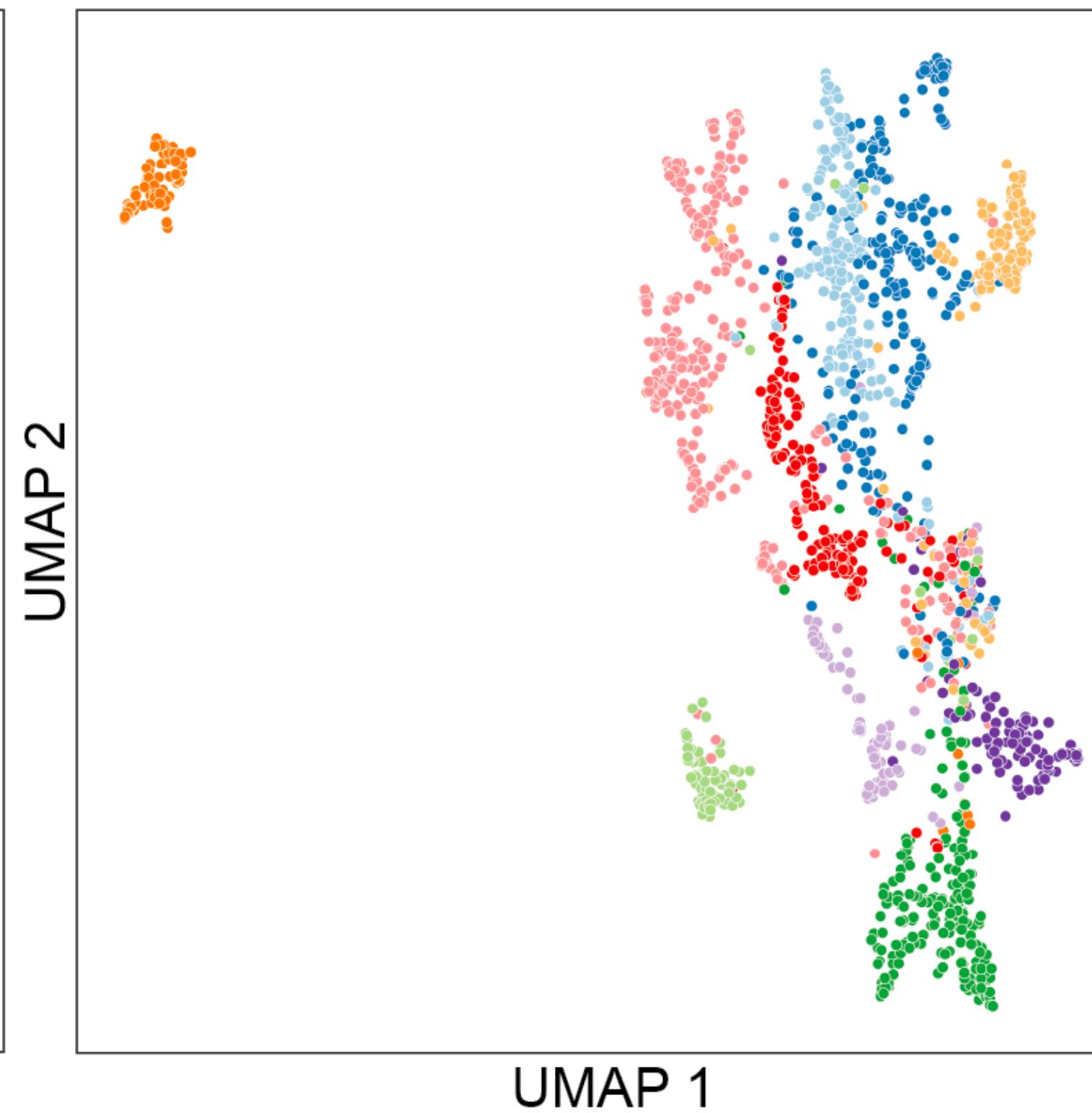


BiGCARP embeddings encode relevant Pfam information

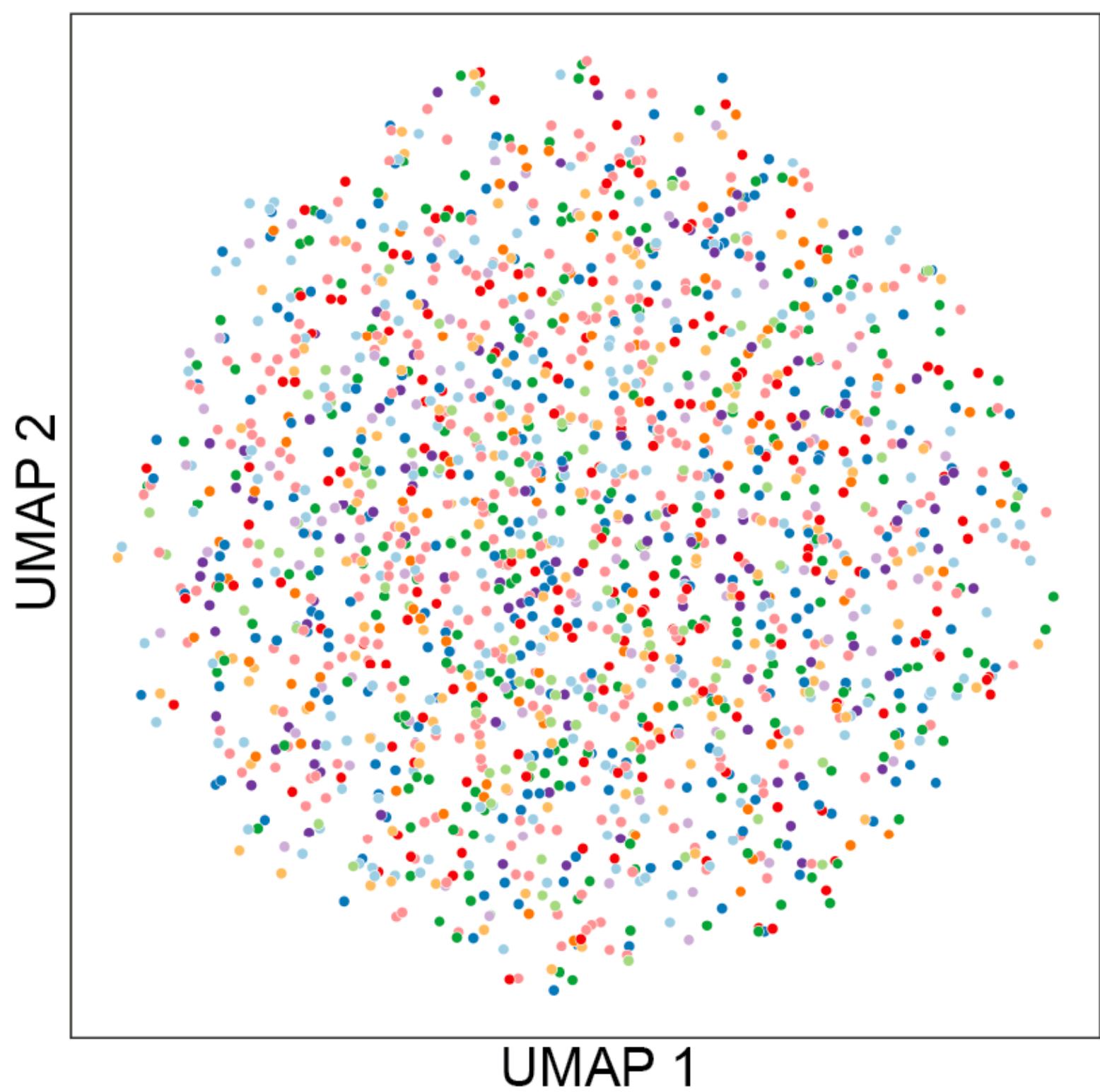
BiGCARP-ESM-1b-finetuned



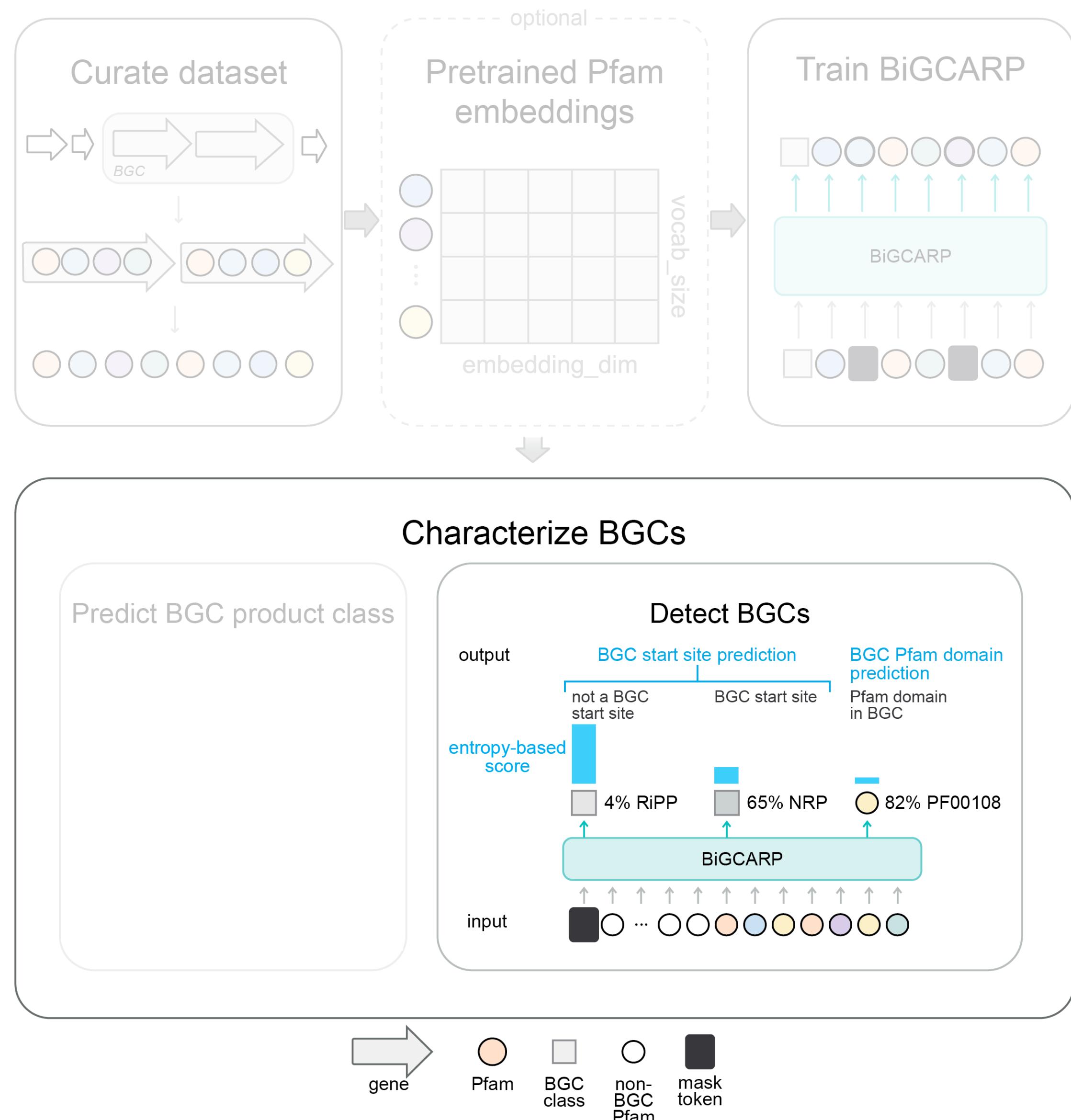
BiGCARP-ESM-1b-frozen



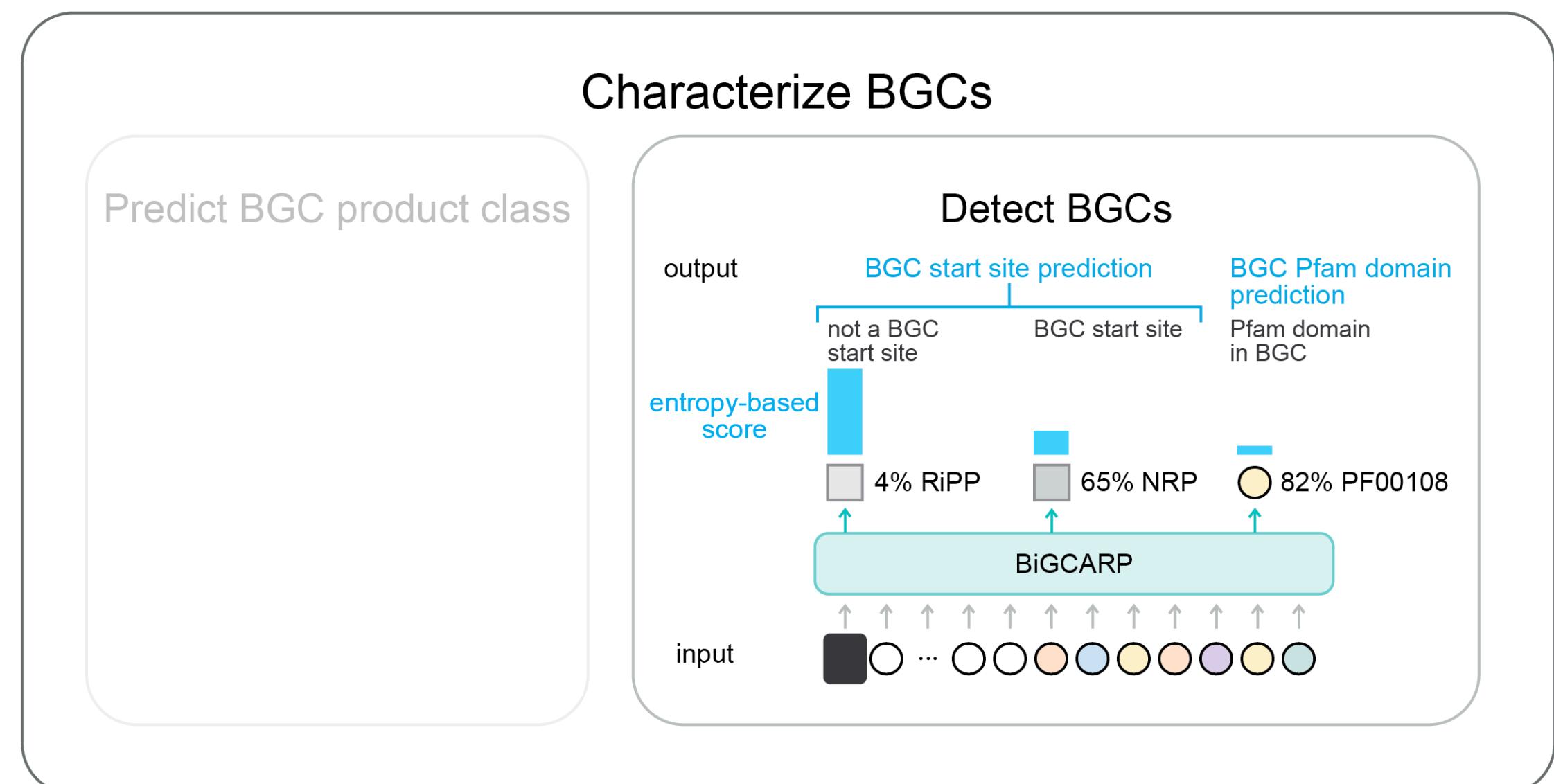
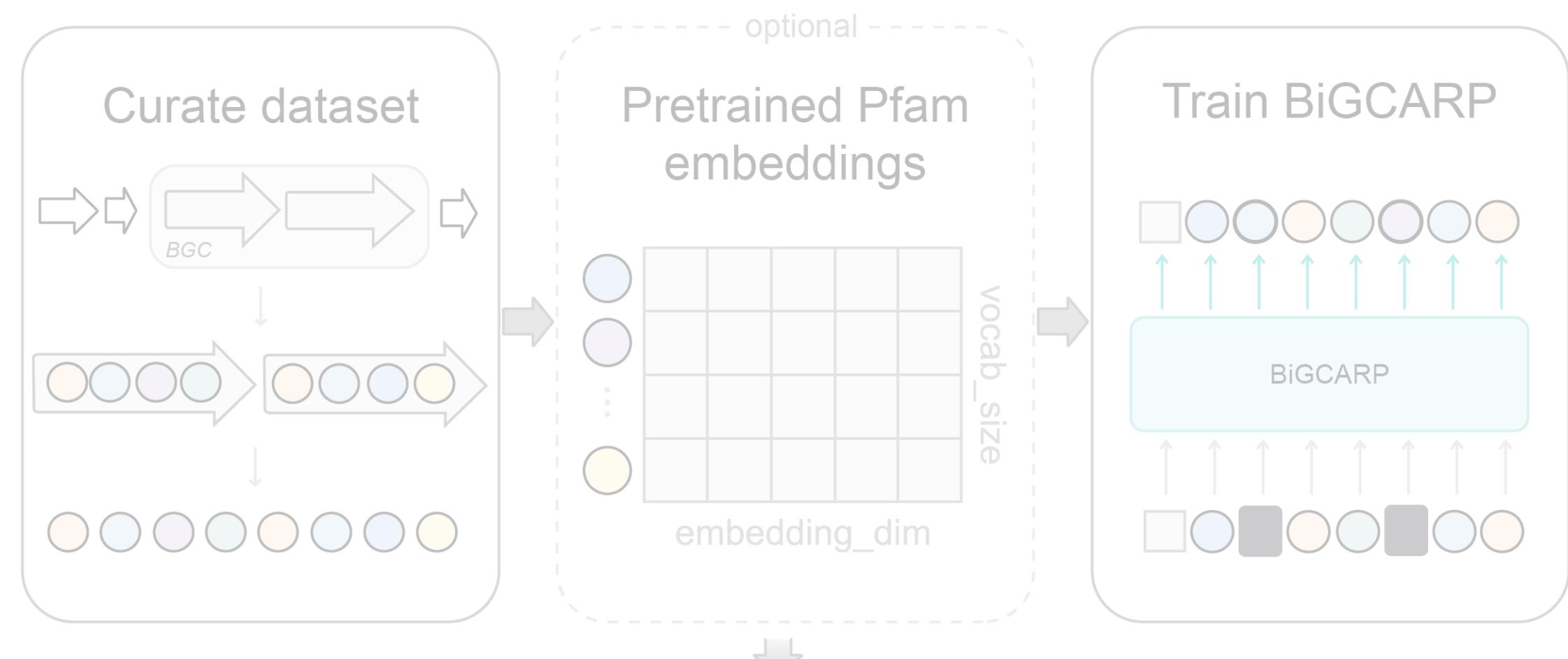
BiGCARP-random



BiGCARP is an unsupervised BGC detector

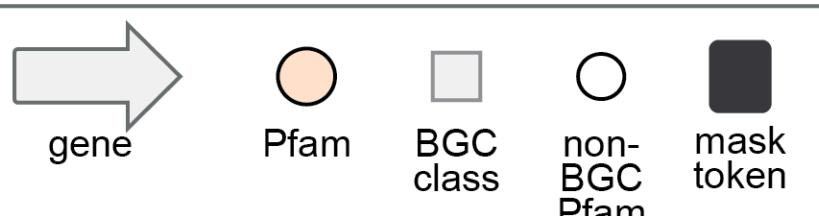


BiGCARP is an unsupervised BGC detector

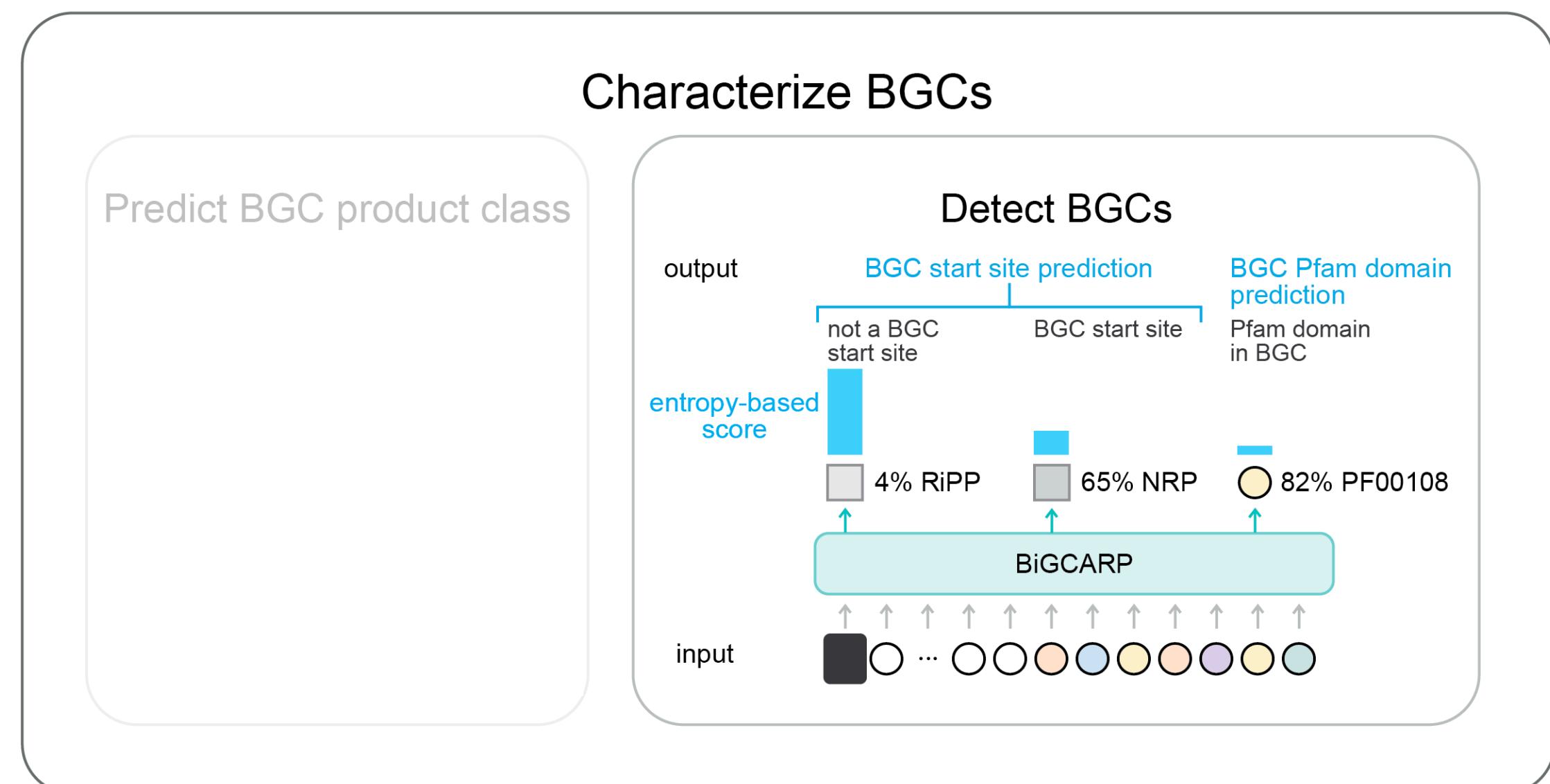
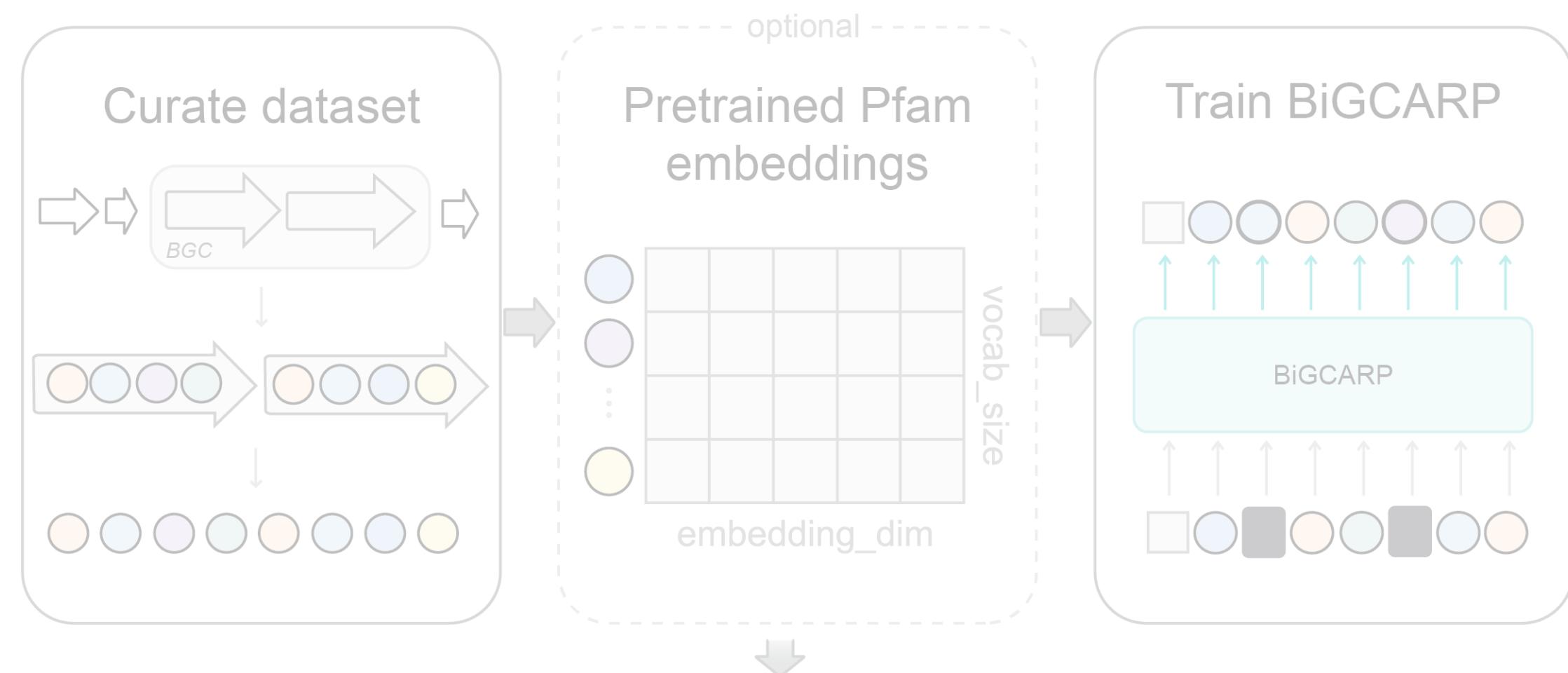


9 bacterial genomes

...ACTGCGGTTACG...



BiGCARP is an unsupervised BGC detector

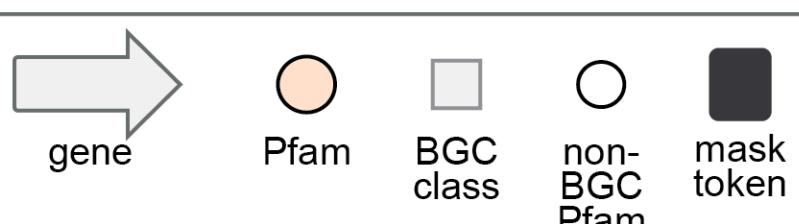


All windows of 64 domains

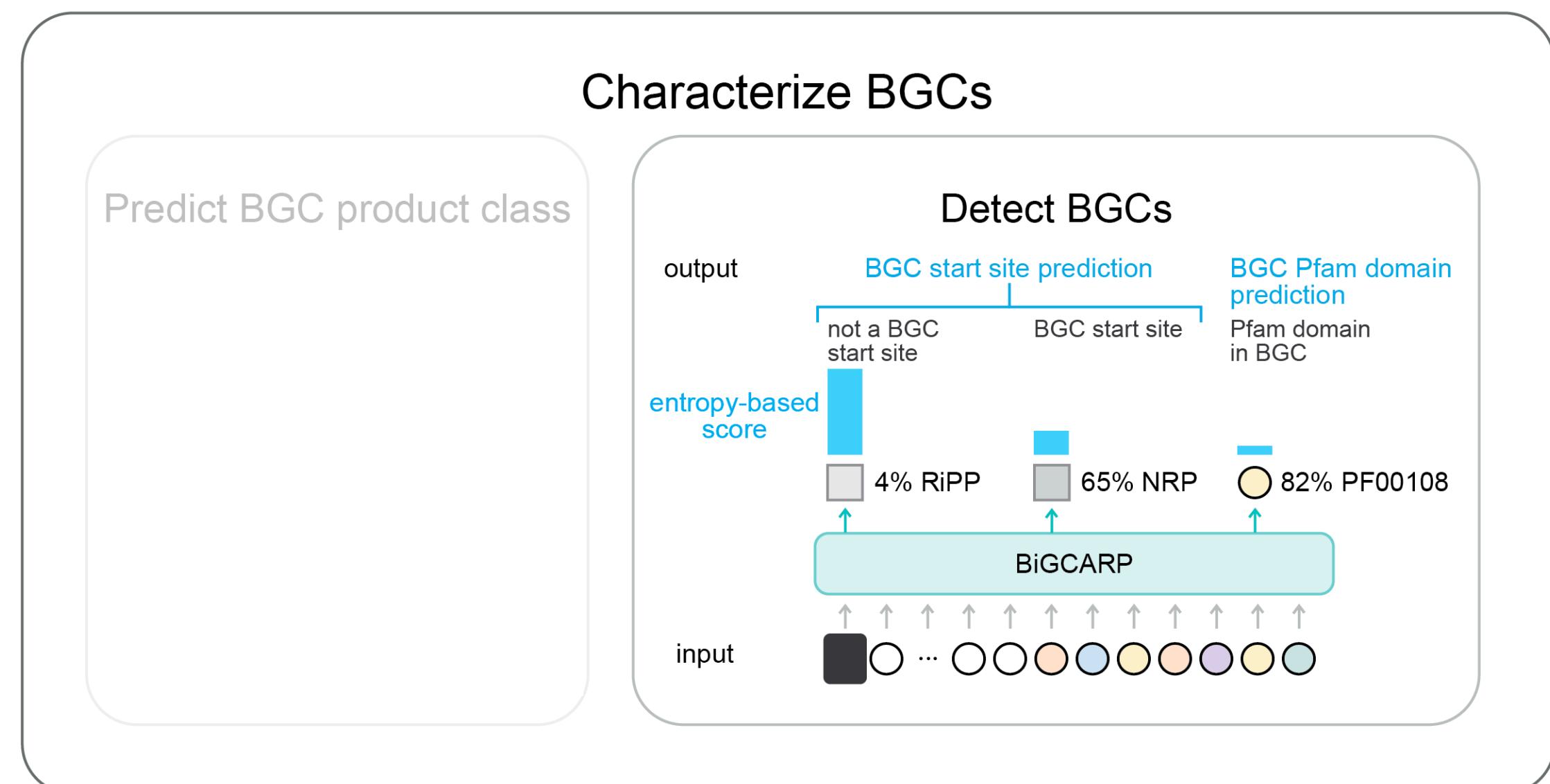
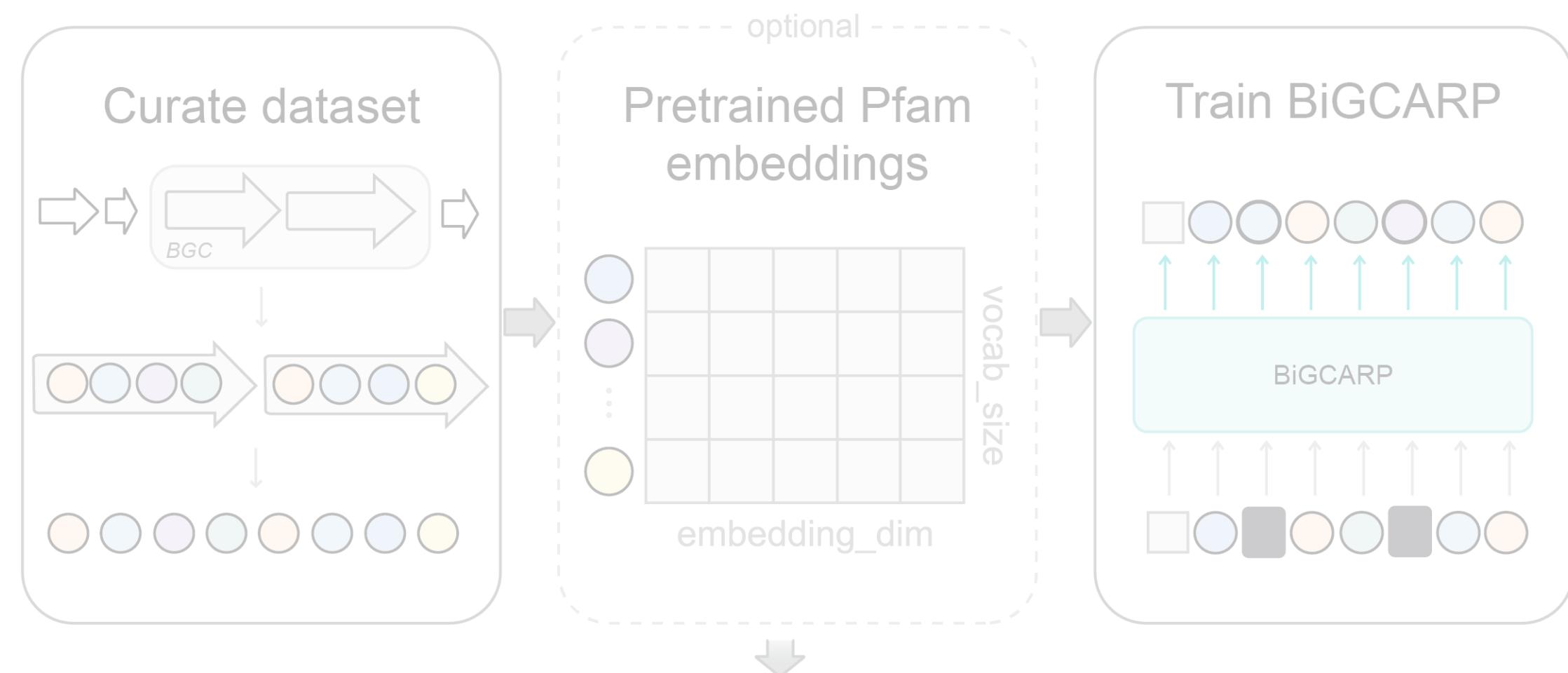
9 bacterial genomes



...ACTGCGGTTACG...



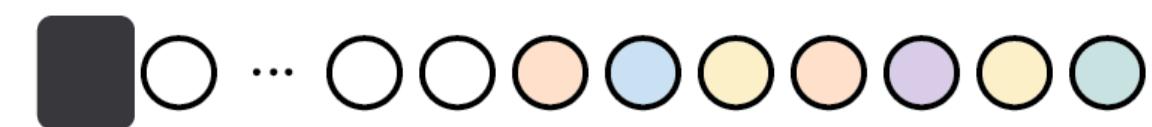
BiGCARP is an unsupervised BGC detector



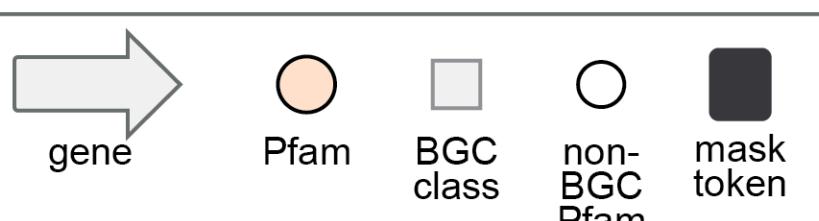
All windows of 64 domains

Append mask

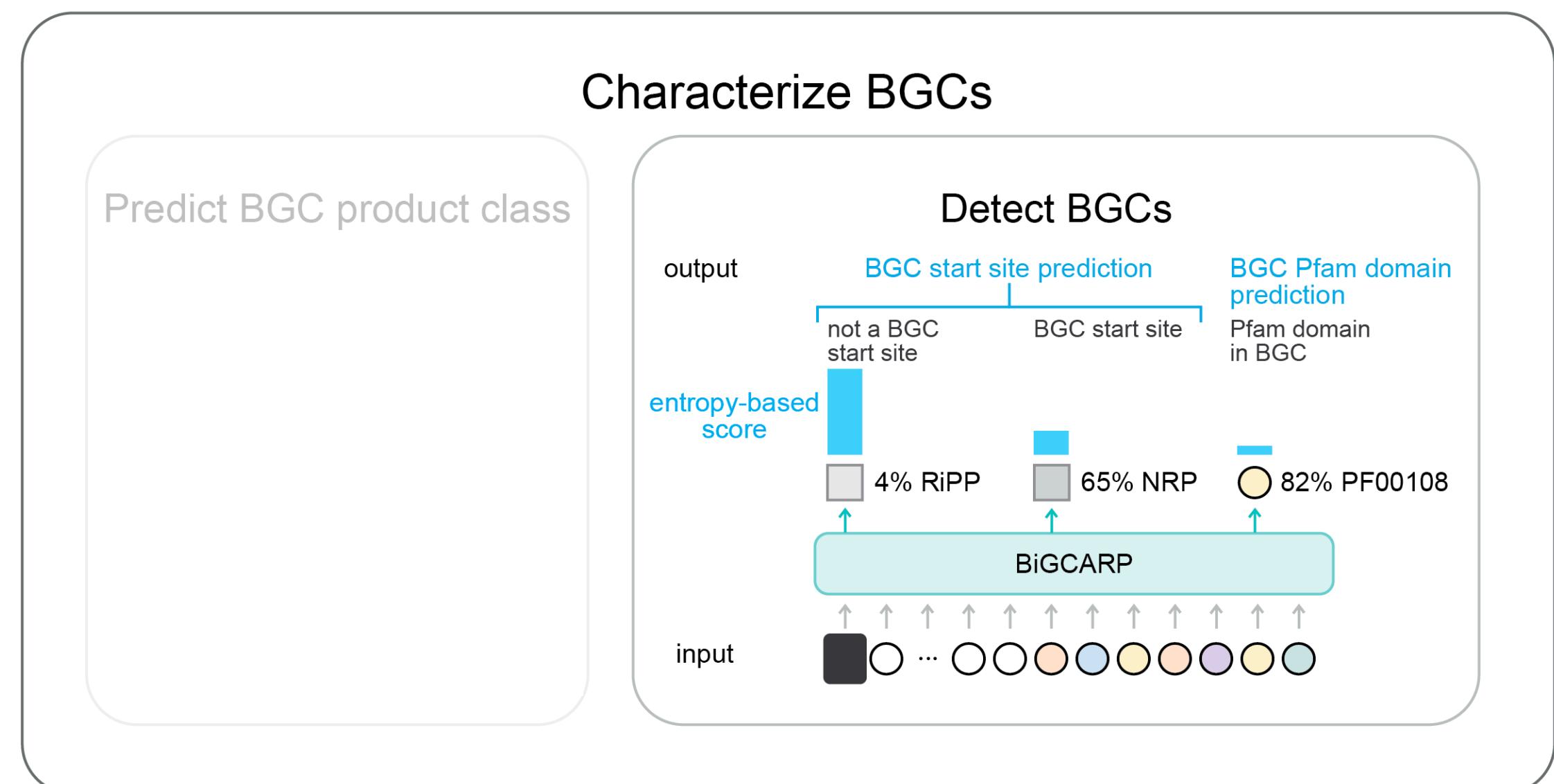
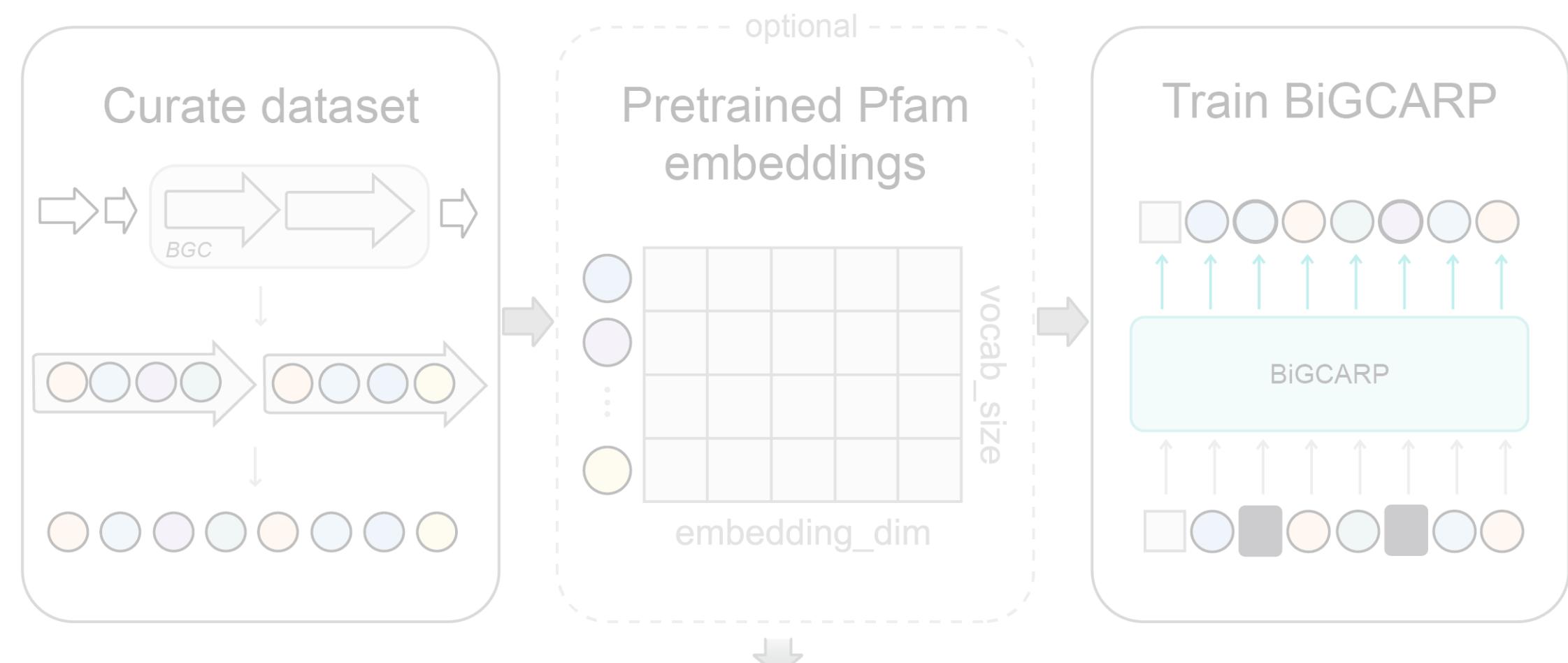
9 bacterial genomes



...ACTGCGGTTACG...



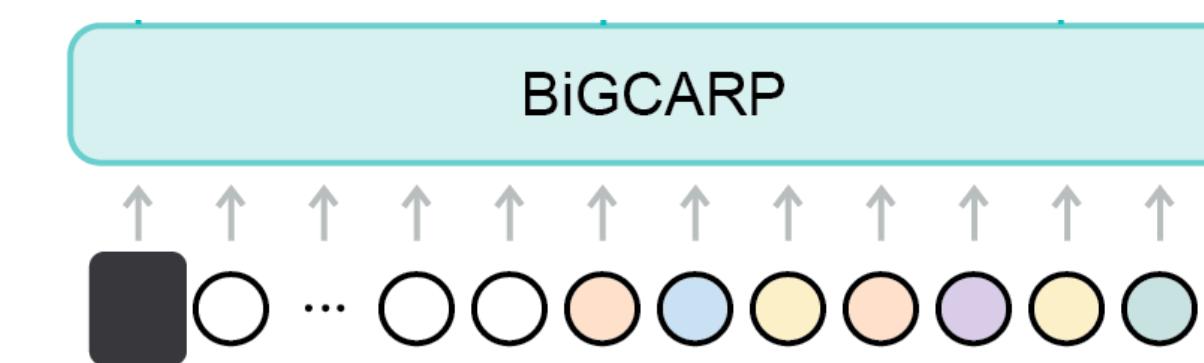
BiGCARP is an unsupervised BGC detector



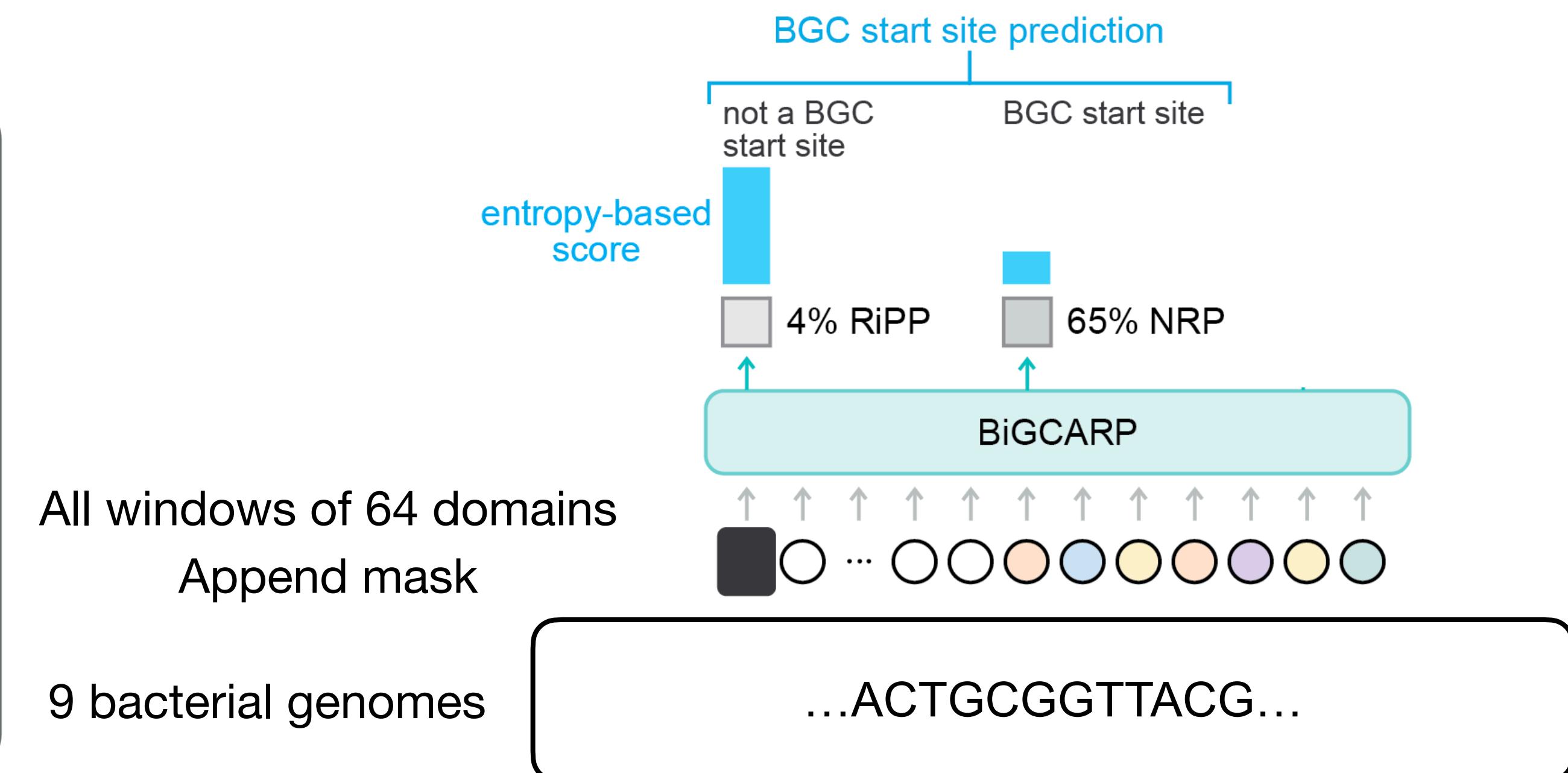
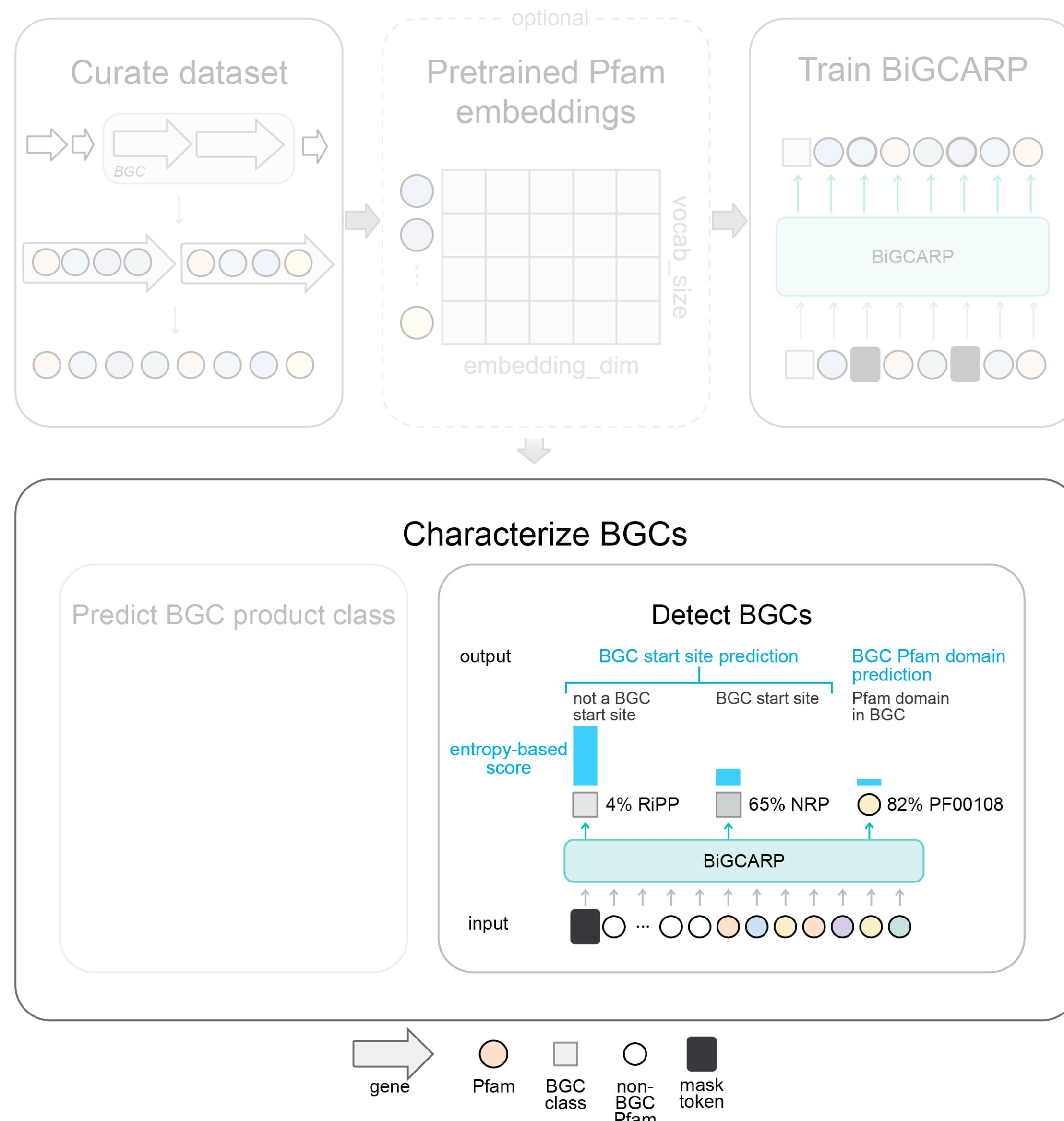
All windows of 64 domains

Append mask

9 bacterial genomes

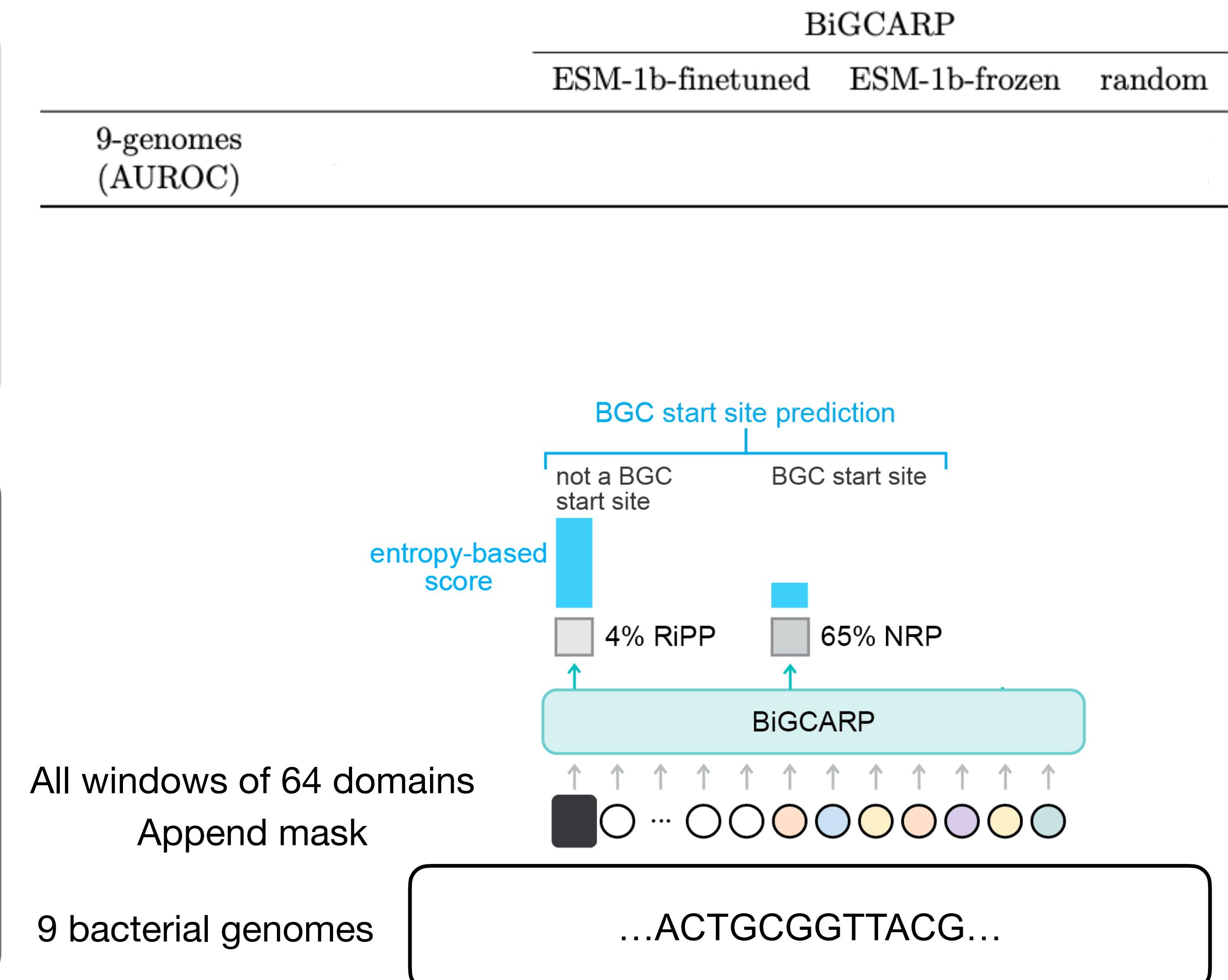
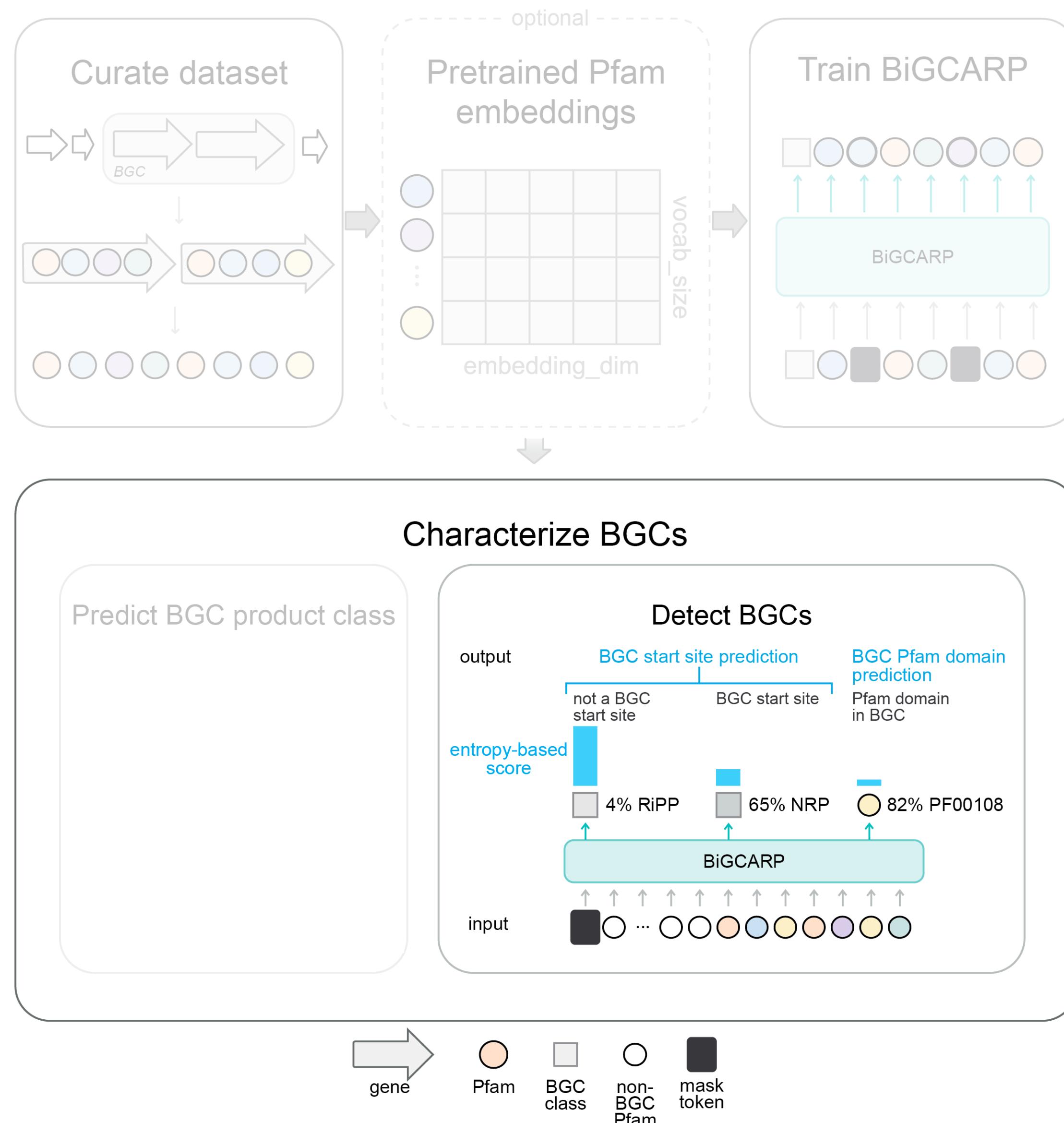


BiGCARP is an unsupervised BGC detector

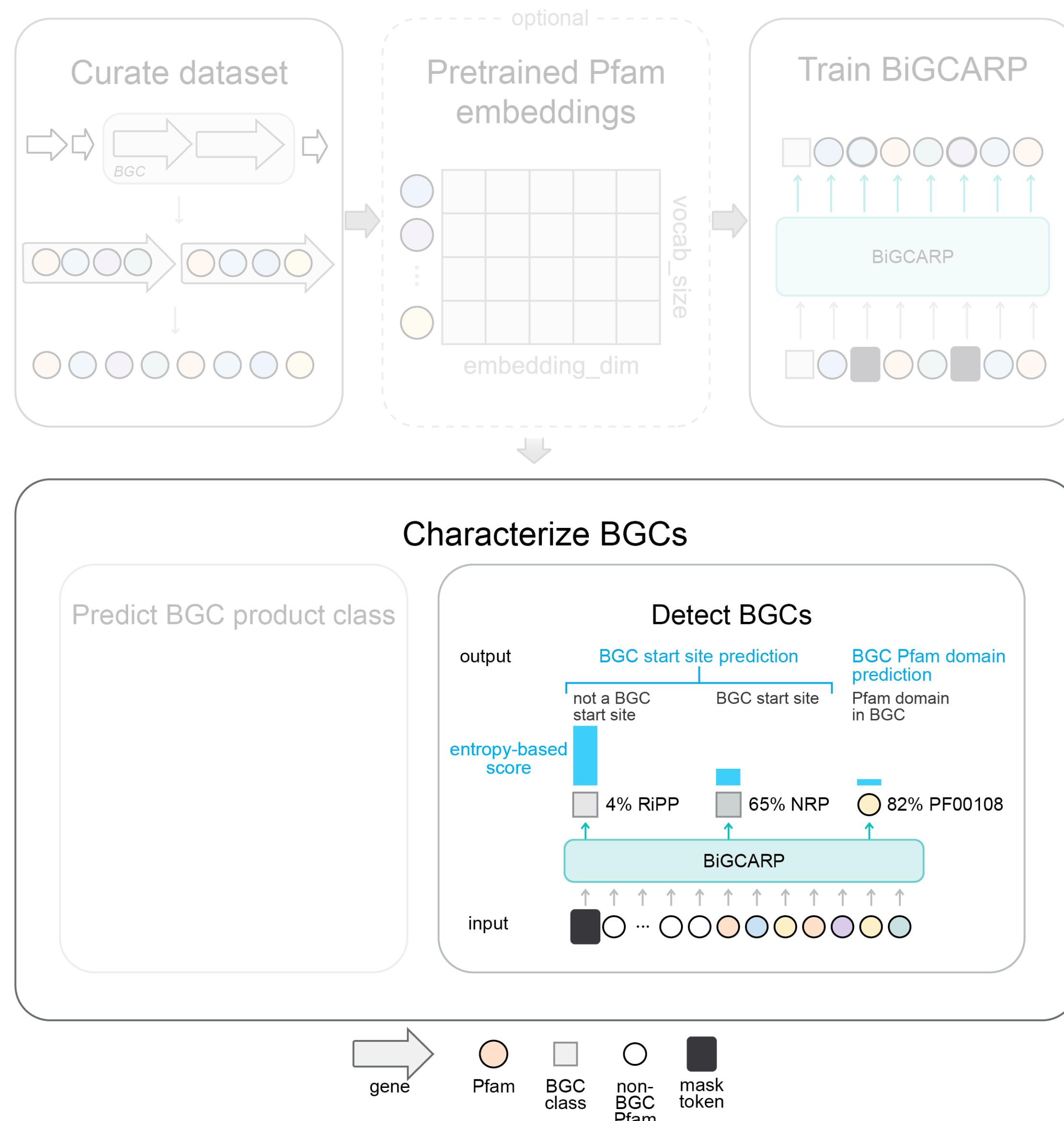


All windows of 64 domains
Append mask
9 bacterial genomes

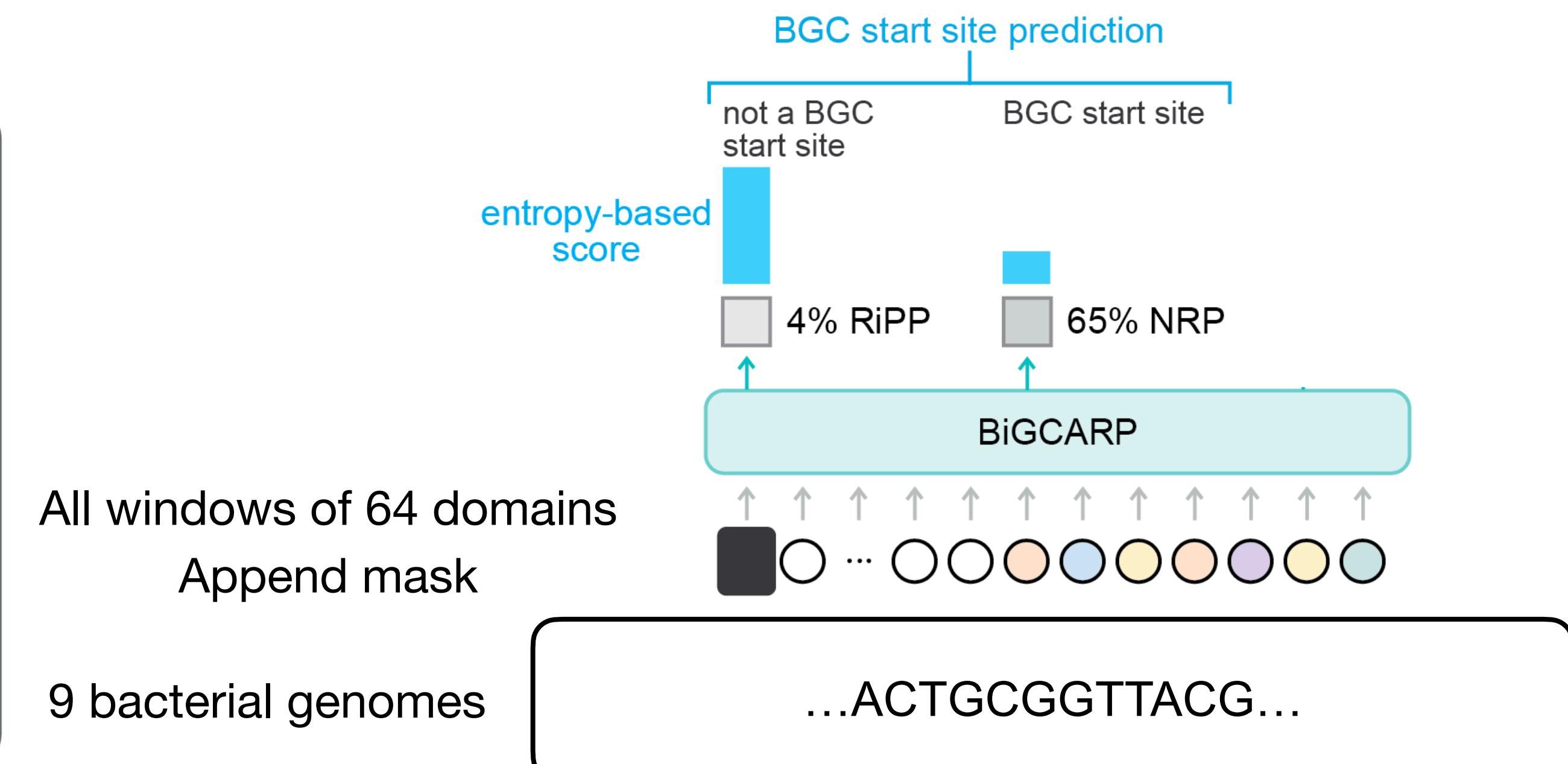
BiGCARP is an unsupervised BGC detector



BiGCARP is an unsupervised BGC detector



BiGCARP			
	ESM-1b-finetuned	ESM-1b-frozen	random
9-genomes (AUROC)	start	0.720	0.701 0.723

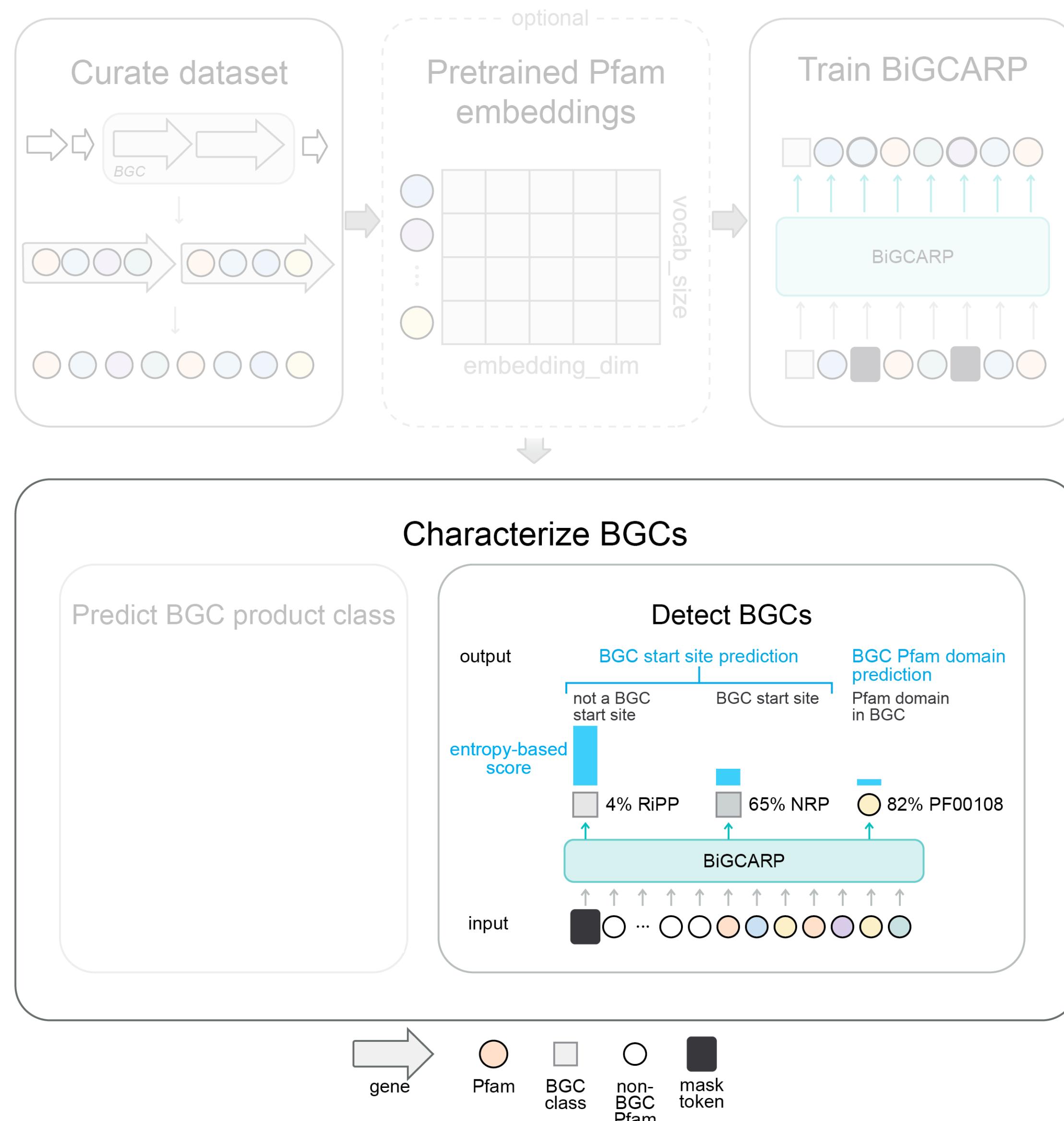


All windows of 64 domains

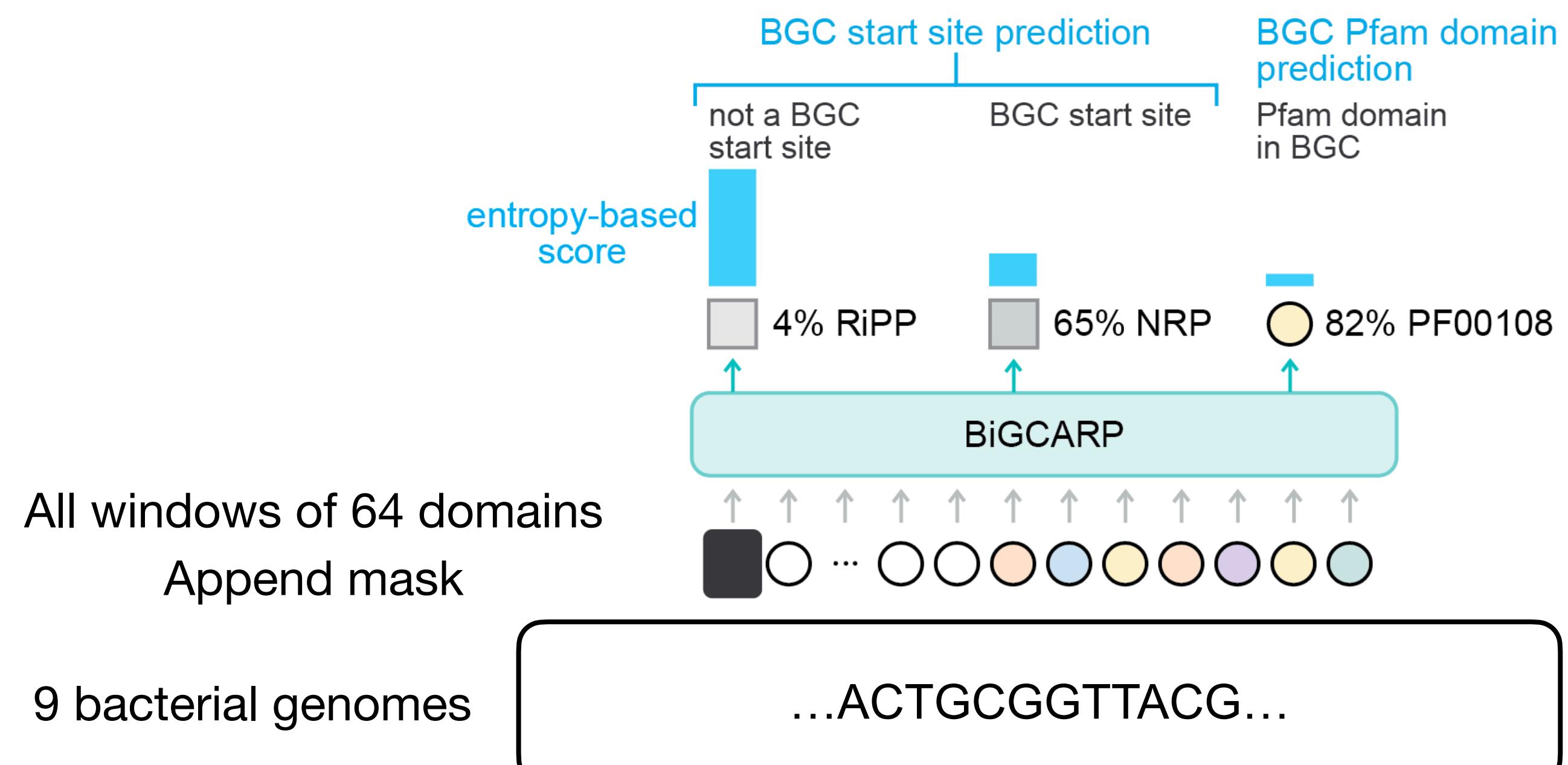
Append mask

9 bacterial genomes

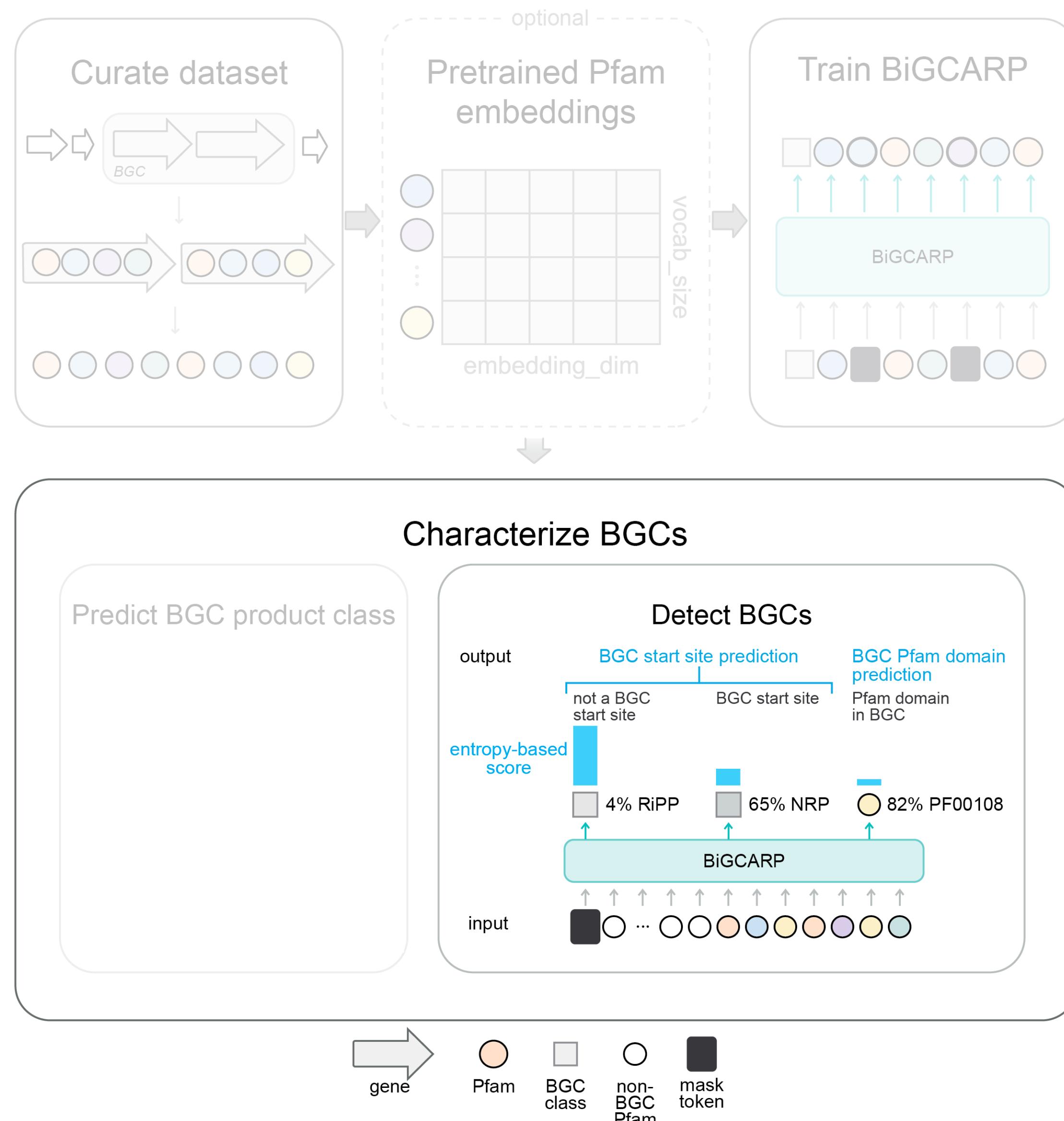
BiGCARP is an unsupervised BGC detector



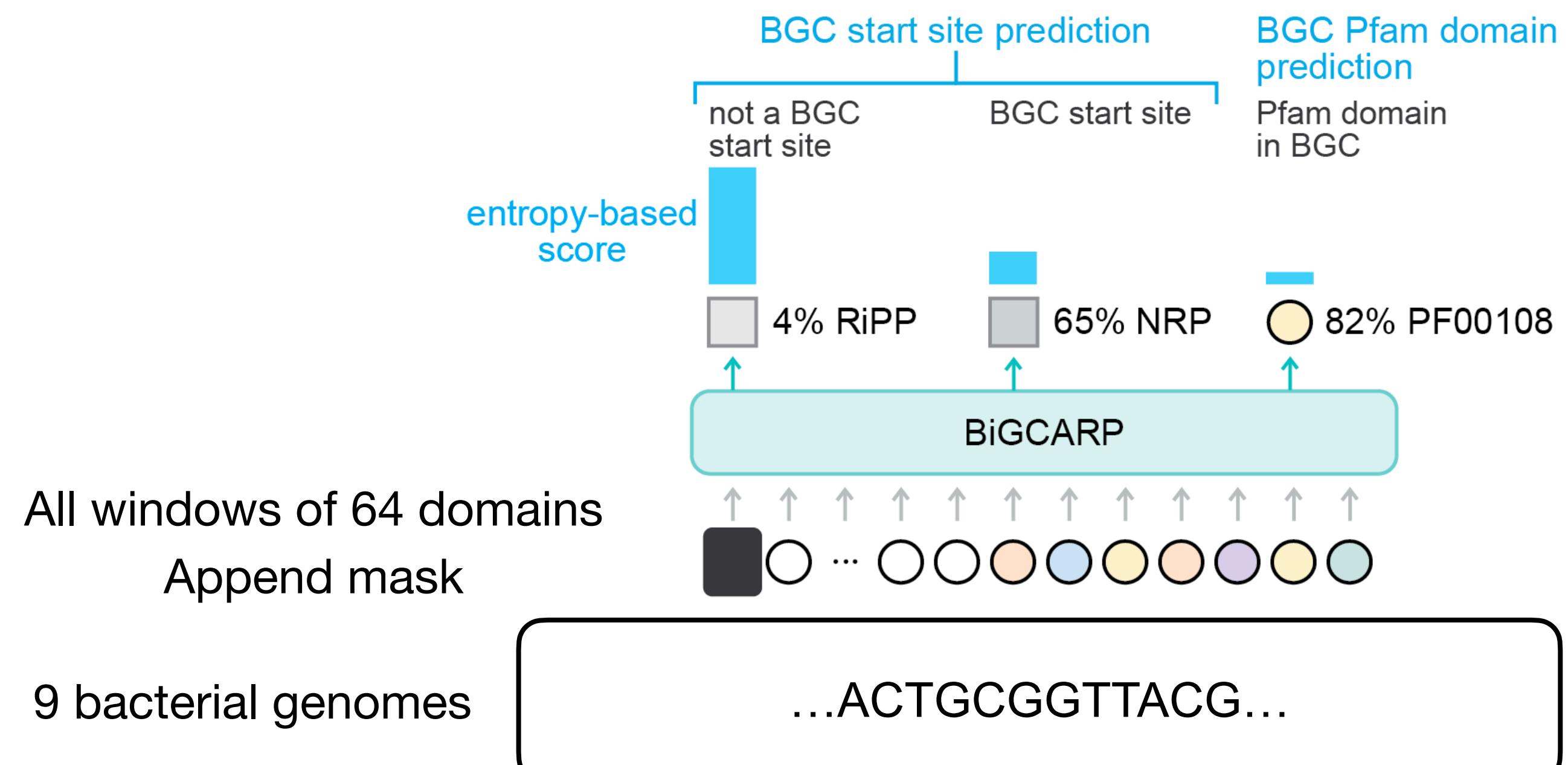
BiGCARP				
	ESM-1b-finetuned	ESM-1b-frozen	random	
9-genomes (AUROC)	start	0.720	0.701	0.723



BiGCARP is an unsupervised BGC detector



BiGCARP			
	ESM-1b-finetuned	ESM-1b-frozen	random
9-genomes (AUROC)	start domain	0.720 0.876	0.701 0.611 0.856



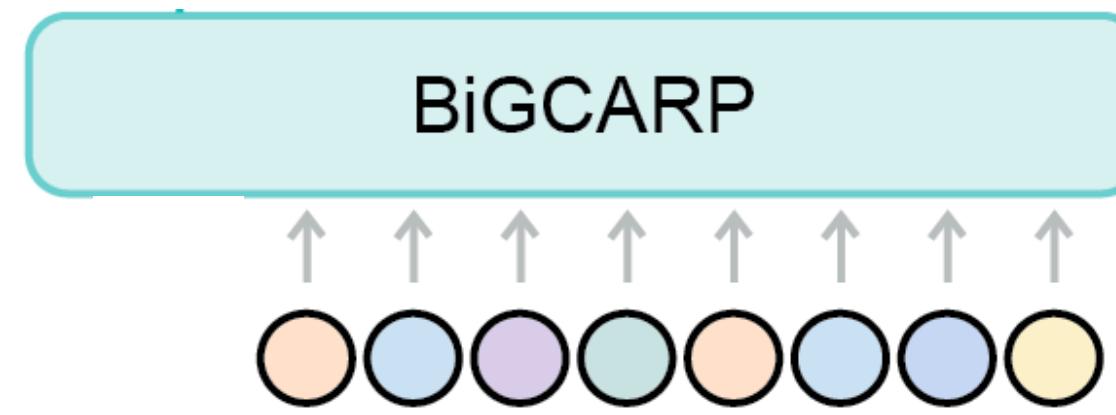
Adding negative examples improves detection

Adding negative examples improves detection

1406 BGCs from MIBiG
+ 10128 negative clusters

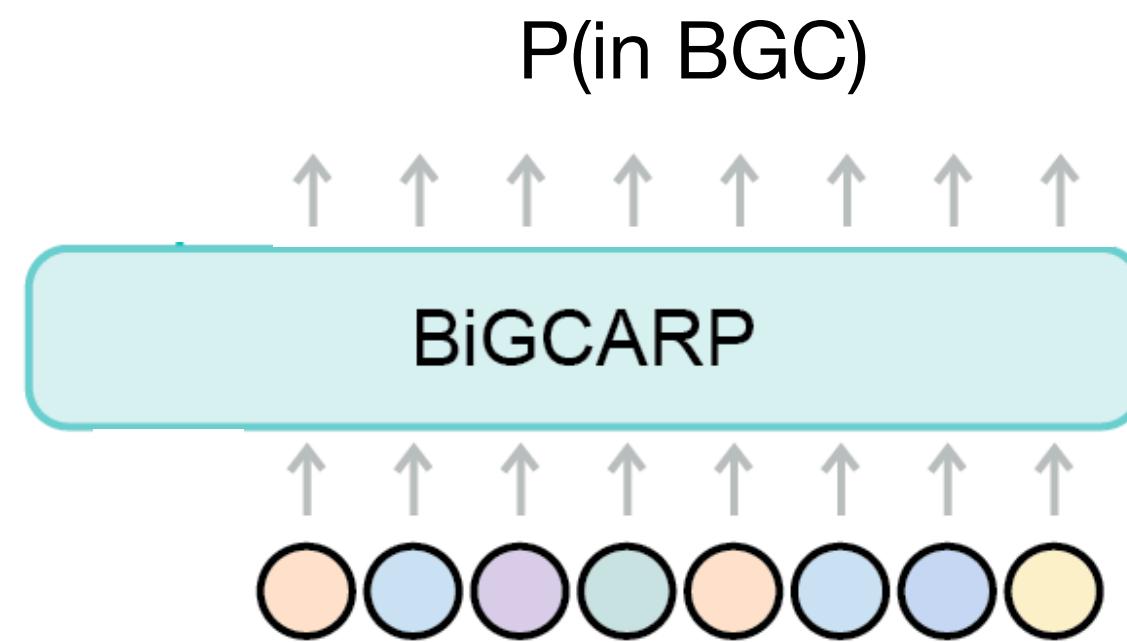
Adding negative examples improves detection

1406 BGCs from MIBiG
+ 10128 negative clusters

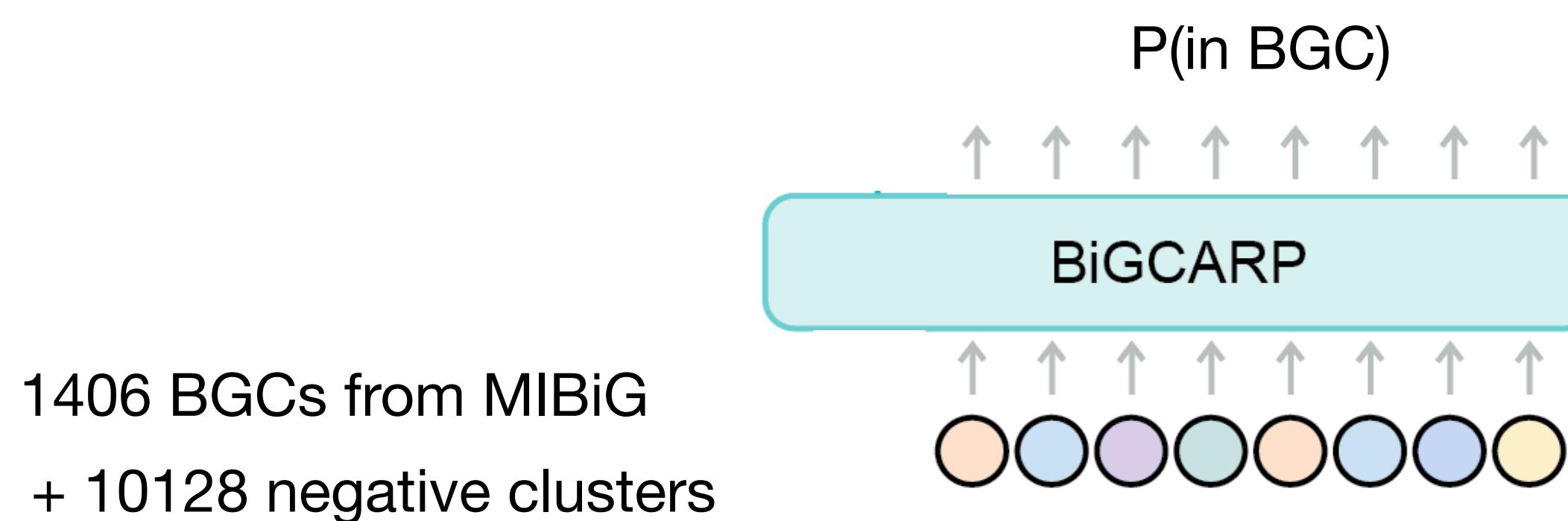


Adding negative examples improves detection

1406 BGCs from MIBiG
+ 10128 negative clusters



Adding negative examples improves detection



pretraining

BiGCARP-ESM-1b-finetuned

BiGCARP-ESM-1b-frozen

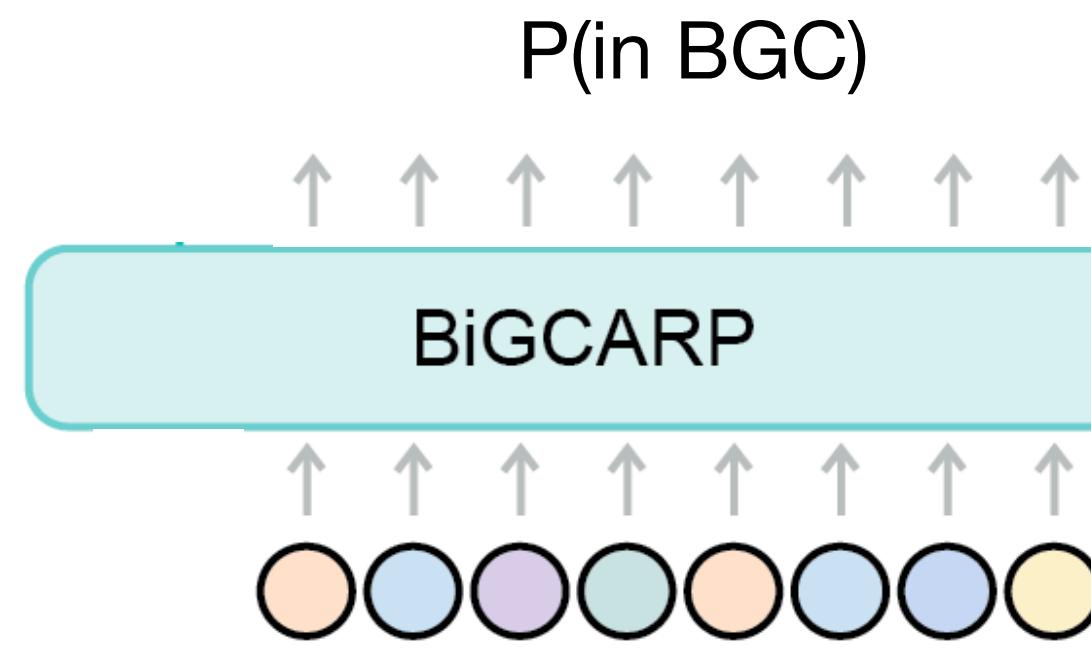
BiGCARP-random

none

DeepBGC

Adding negative examples improves detection

1406 BGCs from MIBiG
+ 10128 negative clusters



pretraining

BiGCARP-ESM-1b-finetuned

BiGCARP-ESM-1b-frozen

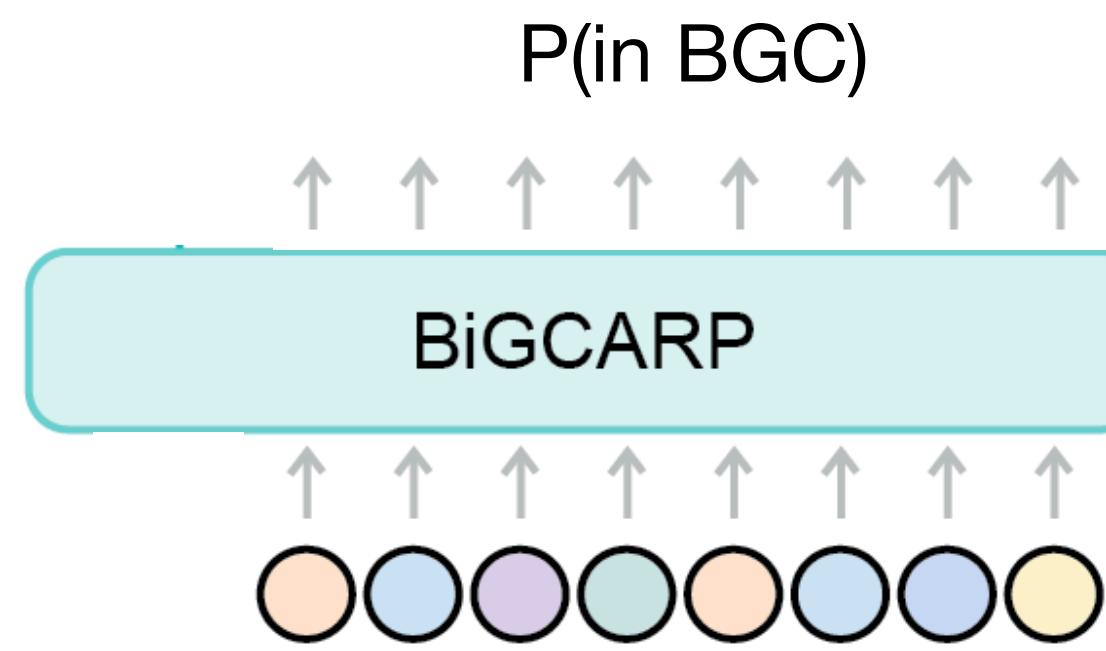
BiGCARP-random

none

DeepBGC **(supervised LSTM)**

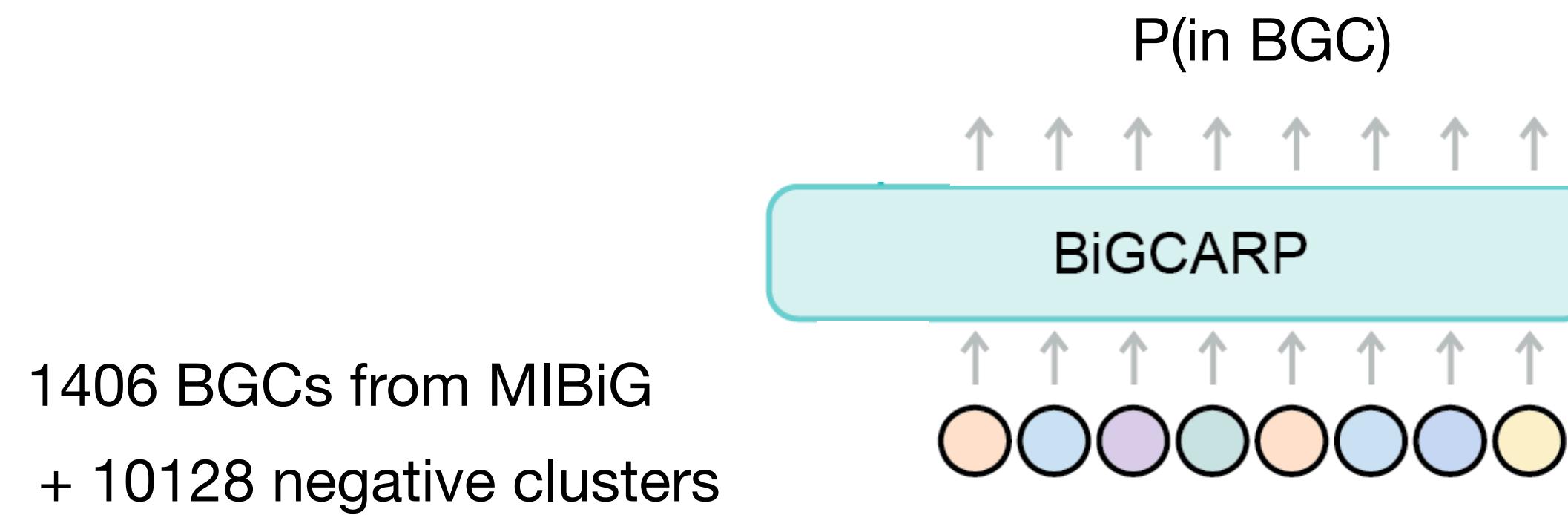
Adding negative examples improves detection

1406 BGCs from MIBiG
+ 10128 negative clusters



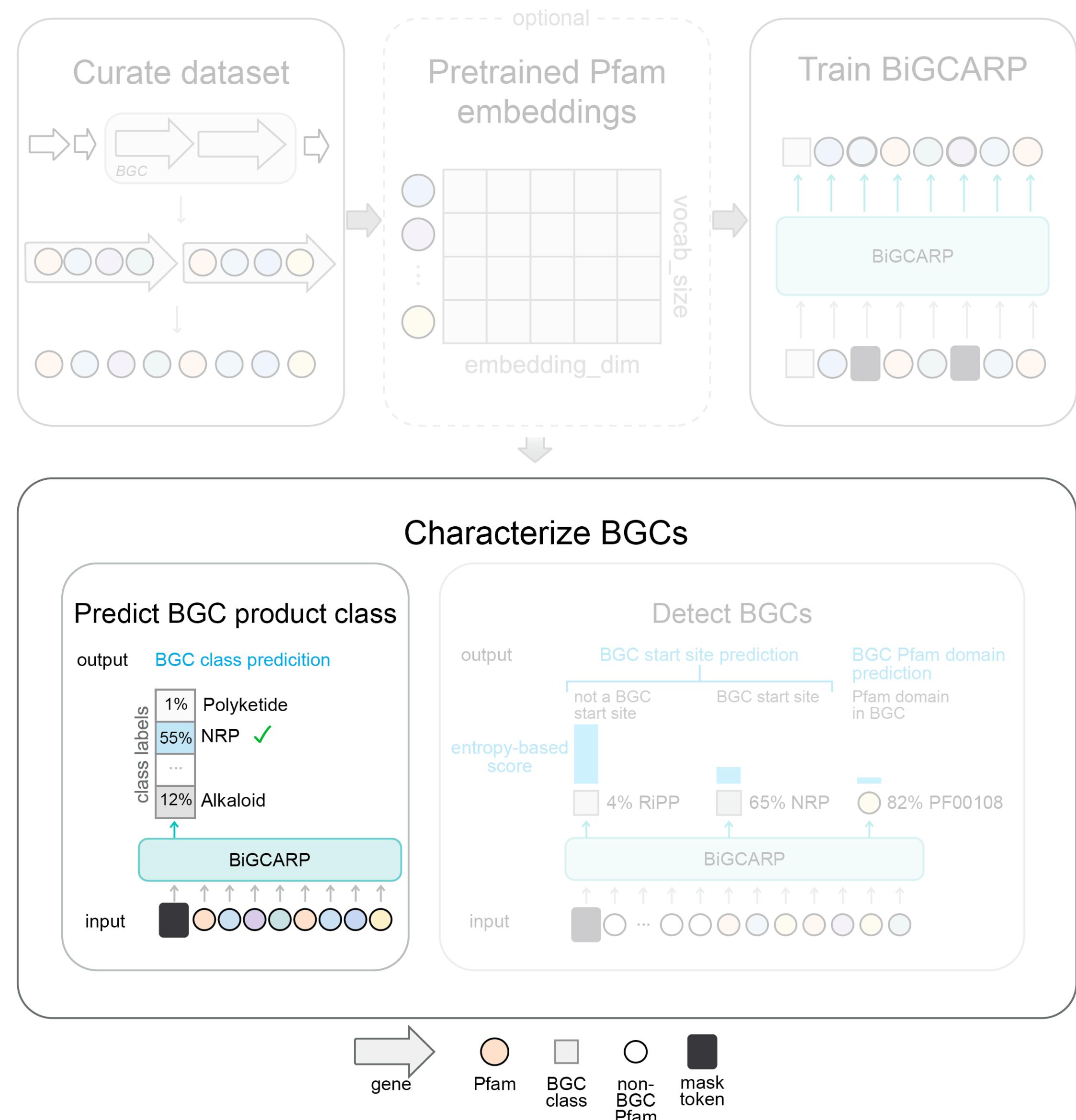
pretraining	Validation
BiGCARP-ESM-1b-finetuned	0.941
BiGCARP-ESM-1b-frozen	0.940
BiGCARP-random	0.936
none	0.937
DeepBGC (supervised LSTM)	0.934

Adding negative examples improves detection

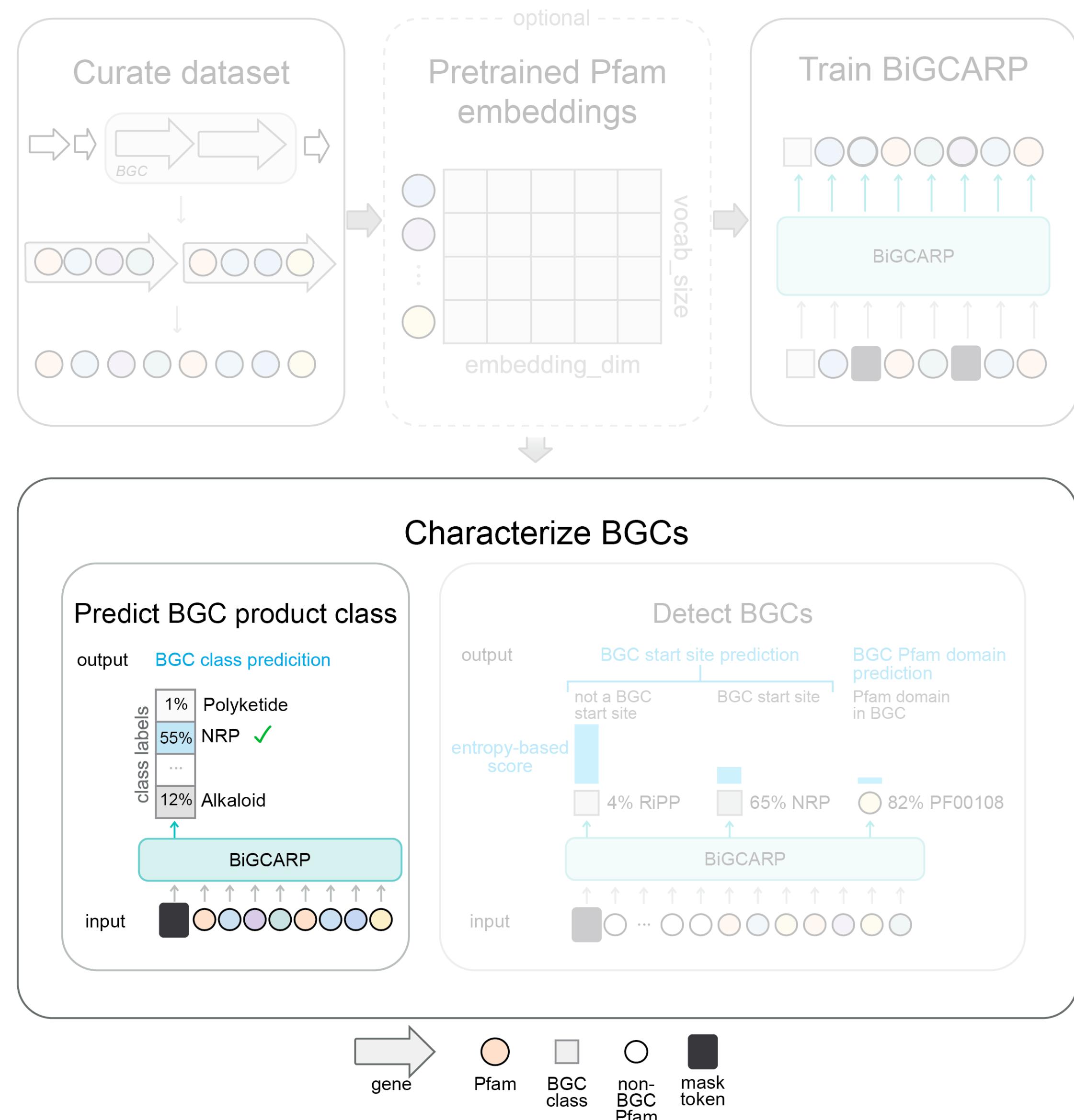


pretraining	Test	Validation
BiGCARP-ESM-1b-finetuned	0.950	0.941
BiGCARP-ESM-1b-frozen	0.946	0.940
BiGCARP-random	0.943	0.936
none	0.950	0.937
DeepBGC (supervised LSTM)	0.921	0.934

BiGCARP predicts BGC product classes



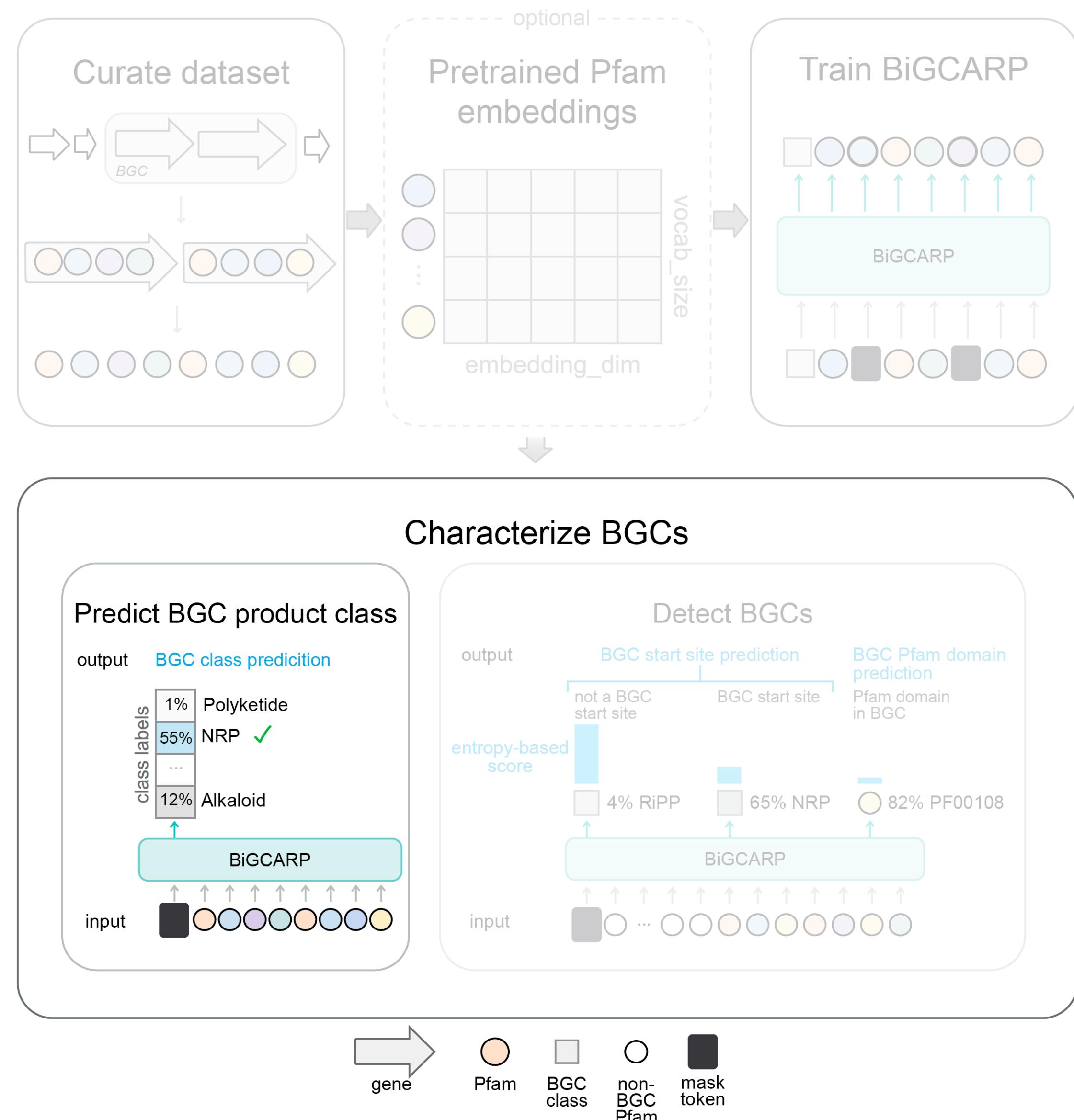
BiGCARP predicts BGC product classes



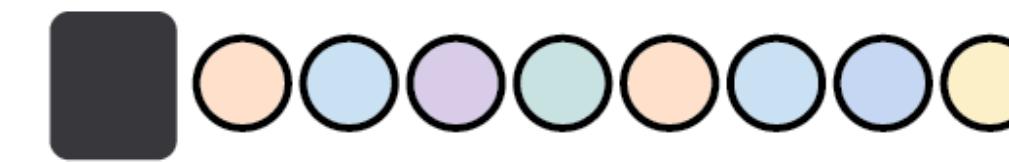
1406 BGCs from MIBiG



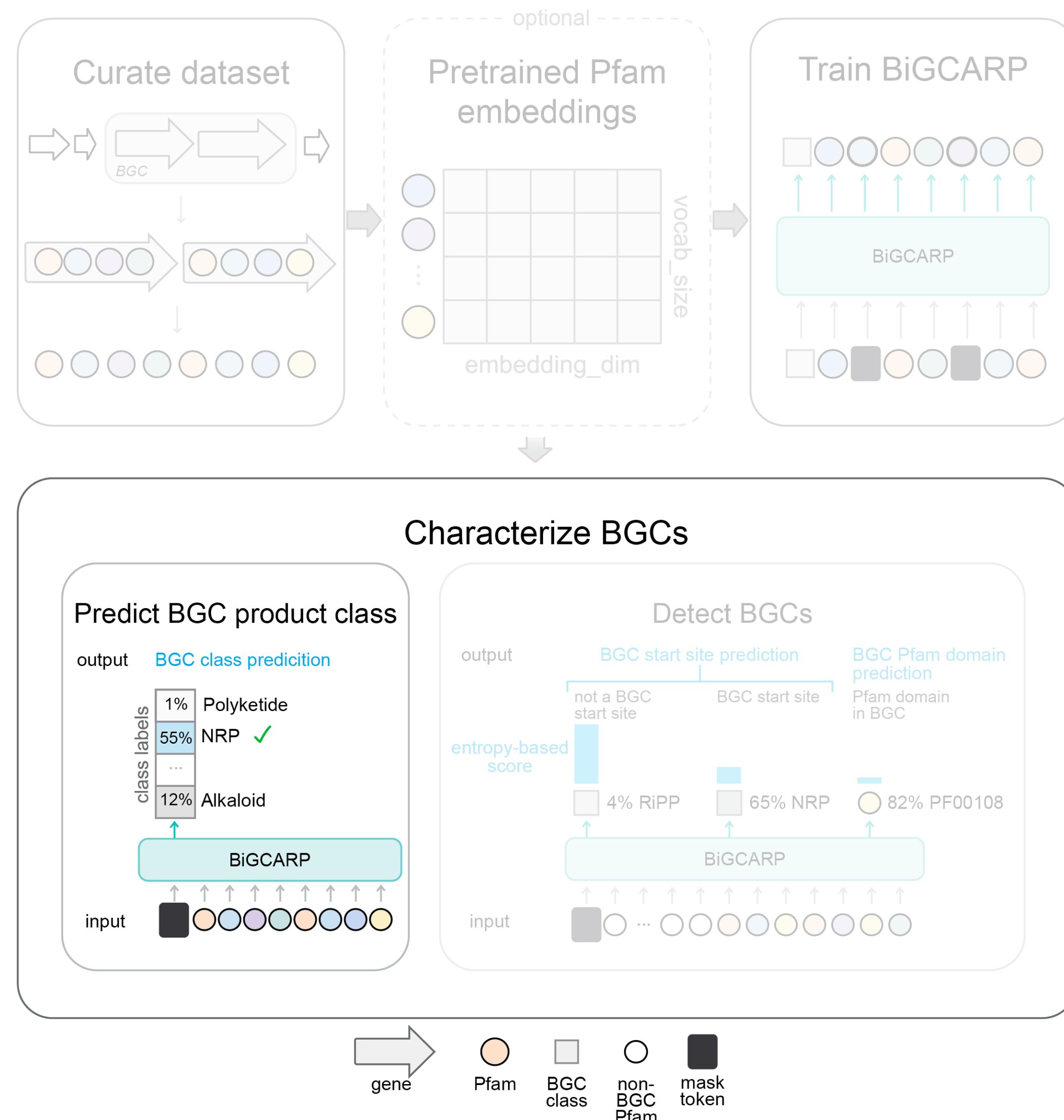
BiGCARP predicts BGC product classes



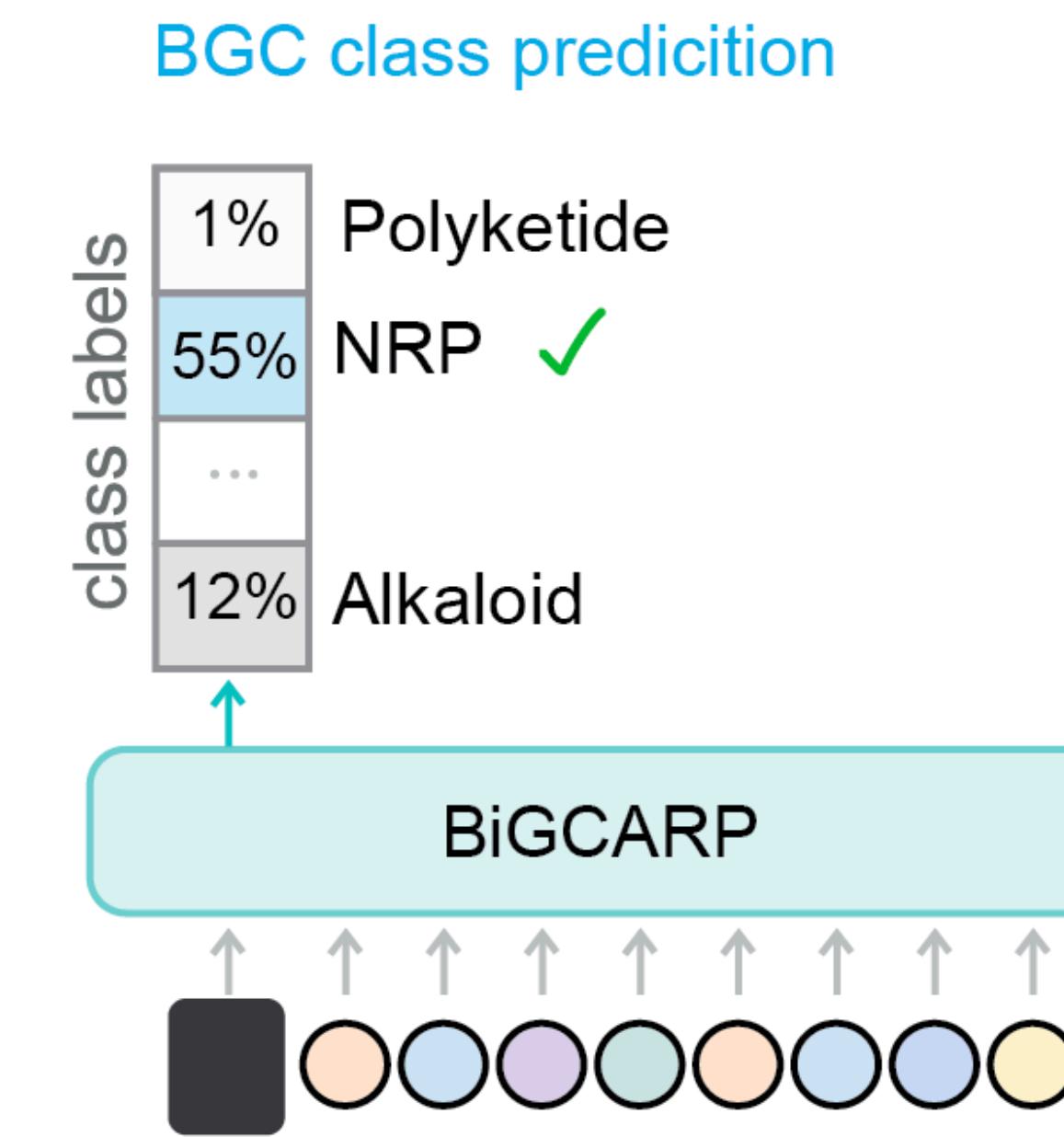
1406 BGCs from MIBiG
Append mask



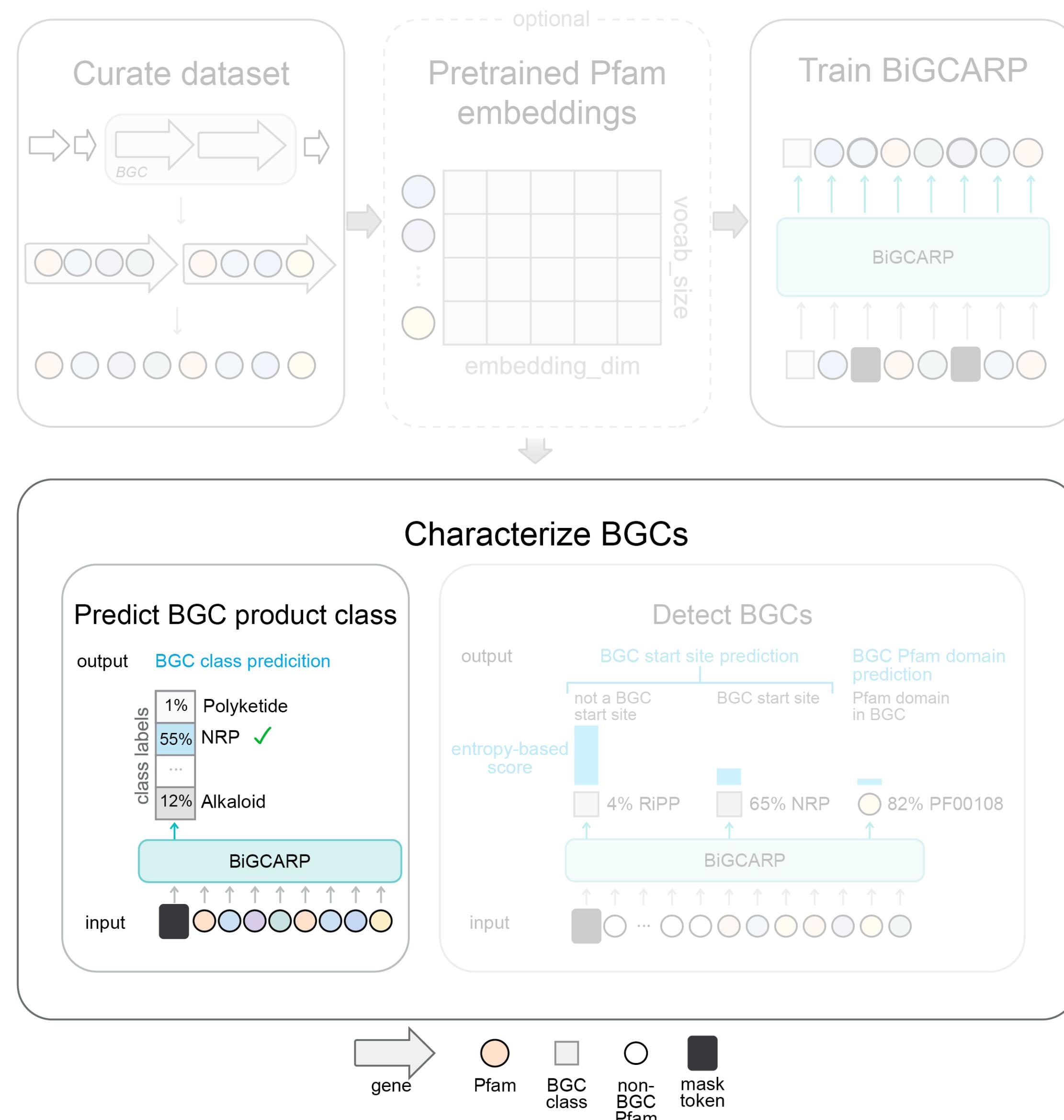
BiGCARP predicts BGC product classes



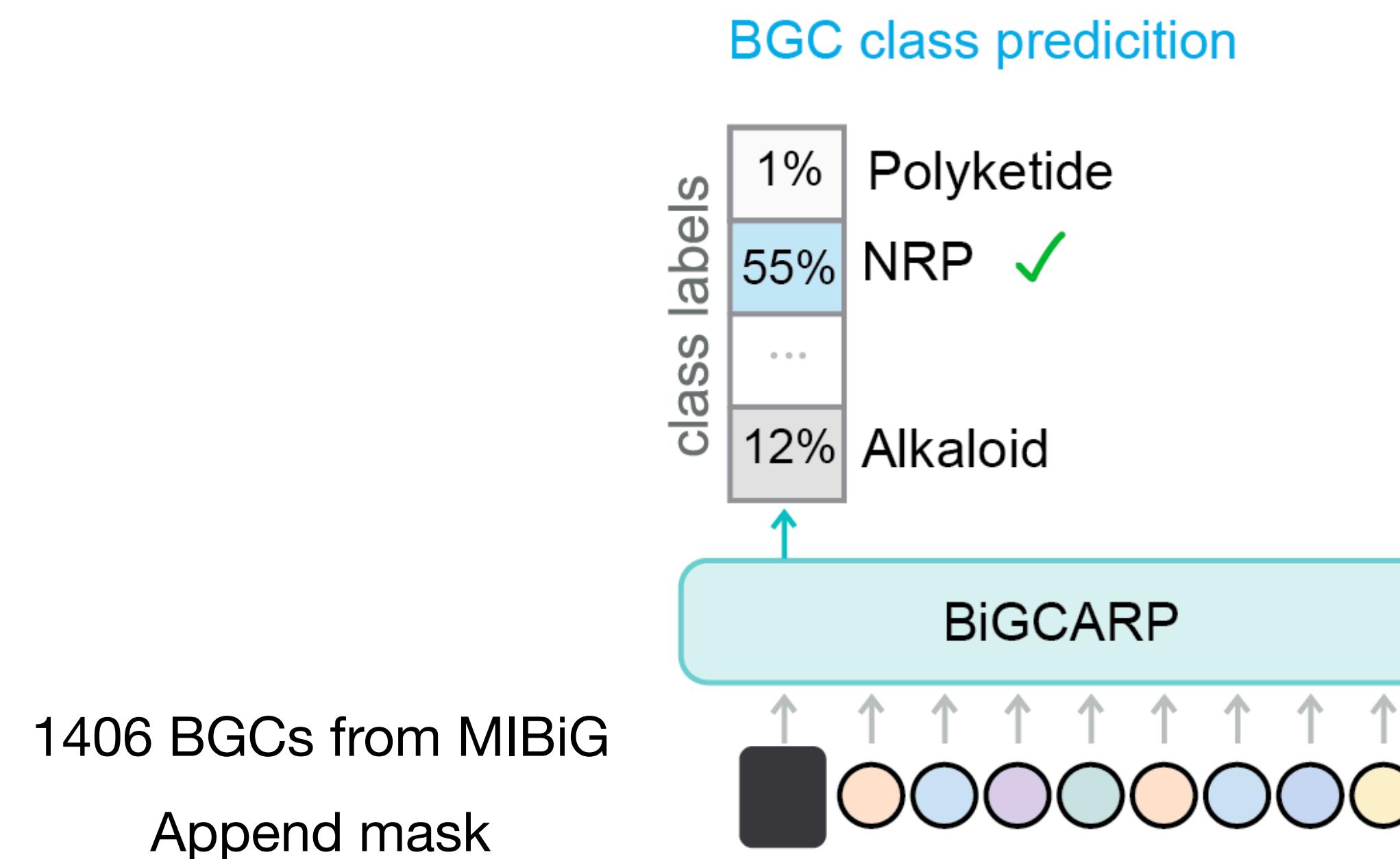
1406 BGCs from MIBiG
Append mask



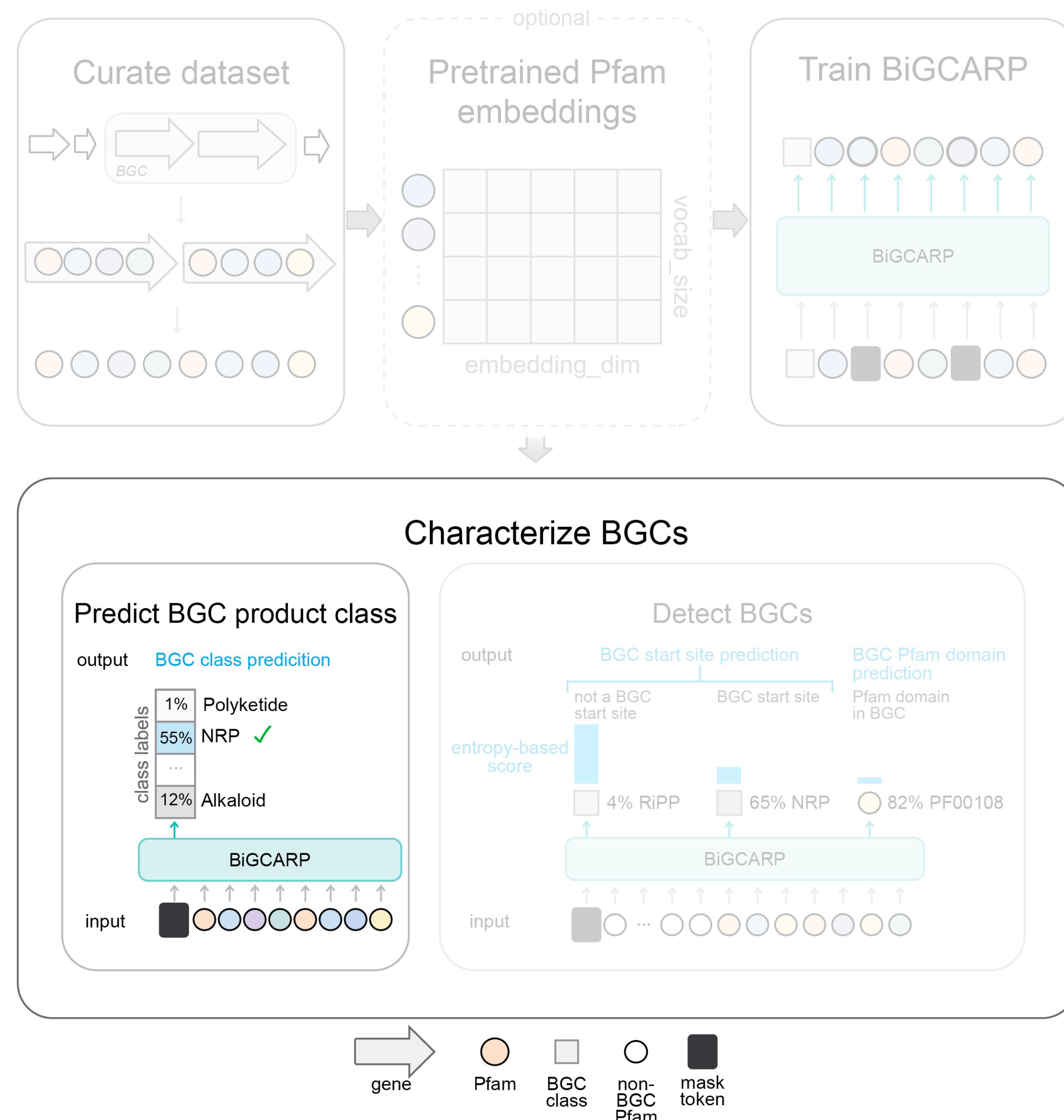
BiGCARP predicts BGC product classes



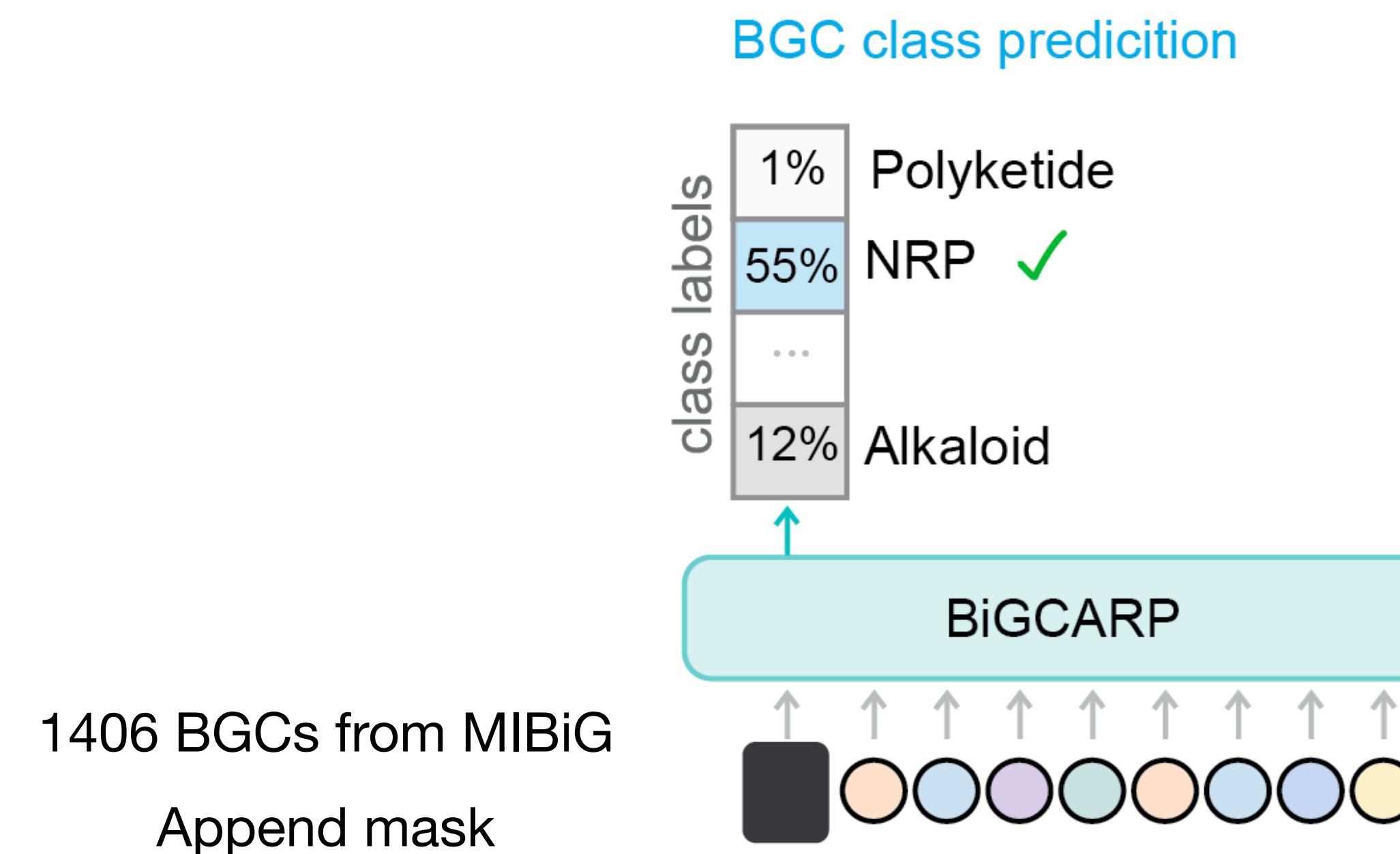
	BiGCARP-ensemble	DeepBGC
polyketide		
NRP		
RiPP		
saccharide		
other		
terpene		
alkaloid		
average		



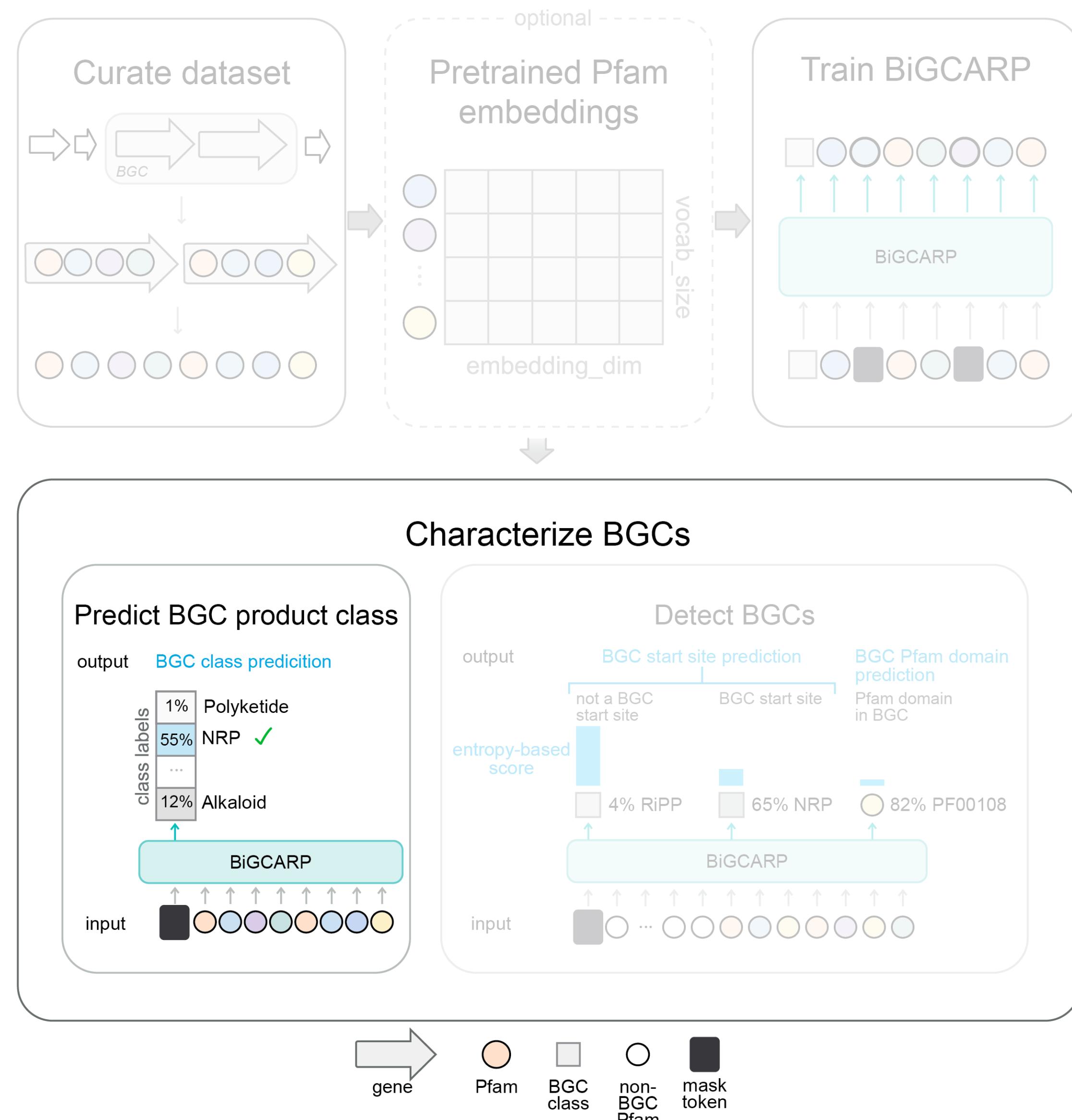
BiGCARP predicts BGC product classes



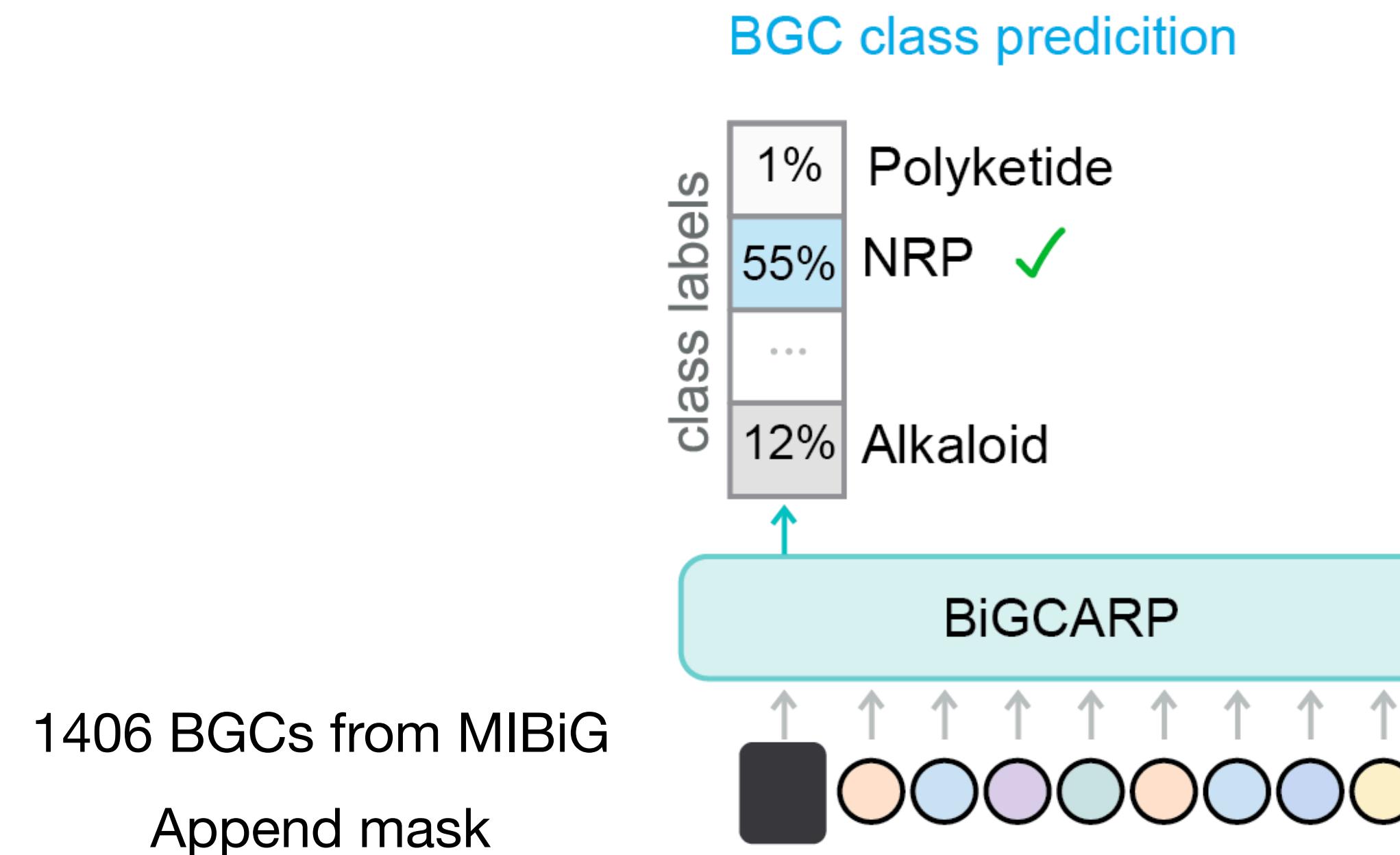
	BiGCARP-ensemble	DeepBGC
polyketide		
NRP		
RiPP		
saccharide		
other		
terpene		
alkaloid		
average	0.855	0.792



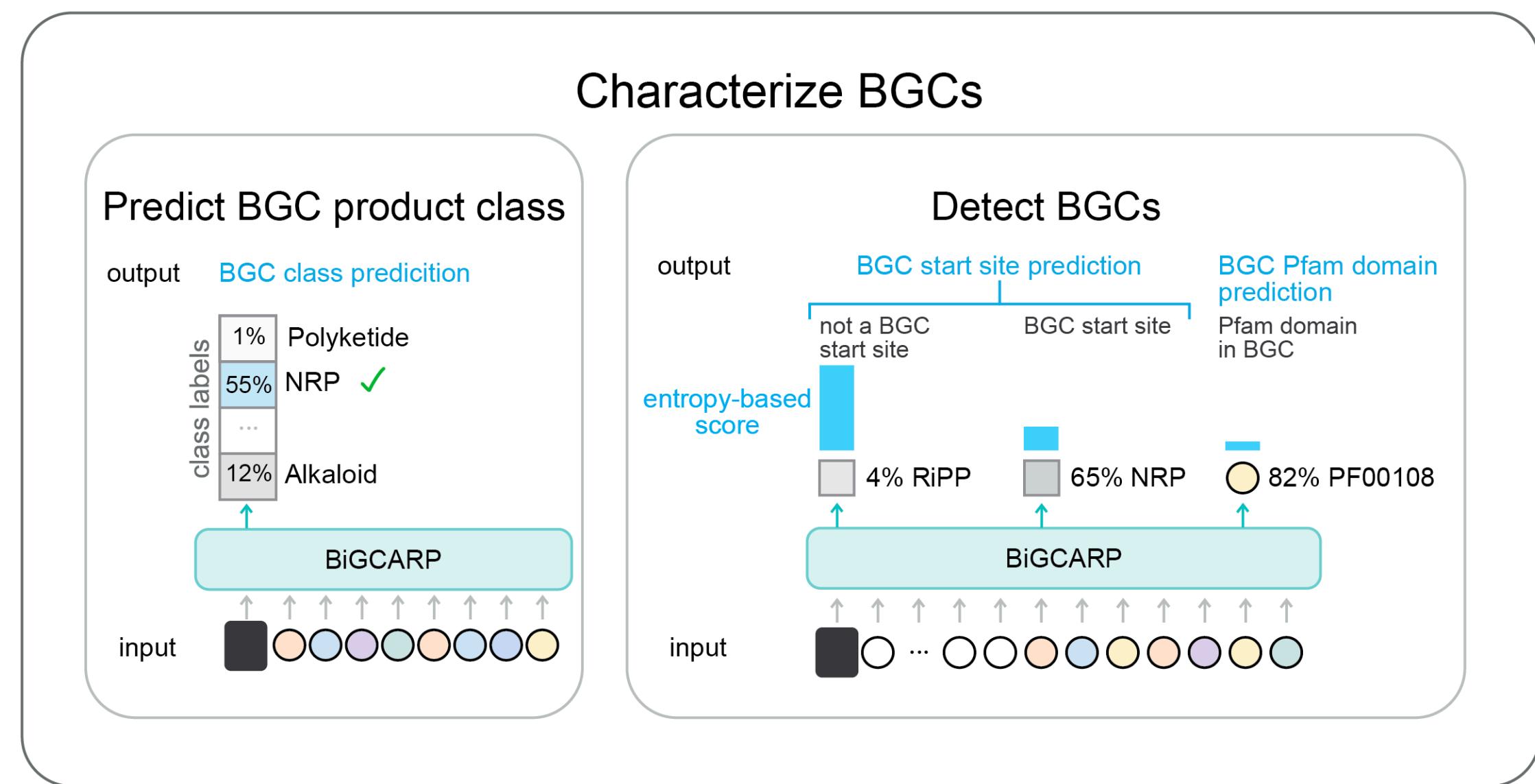
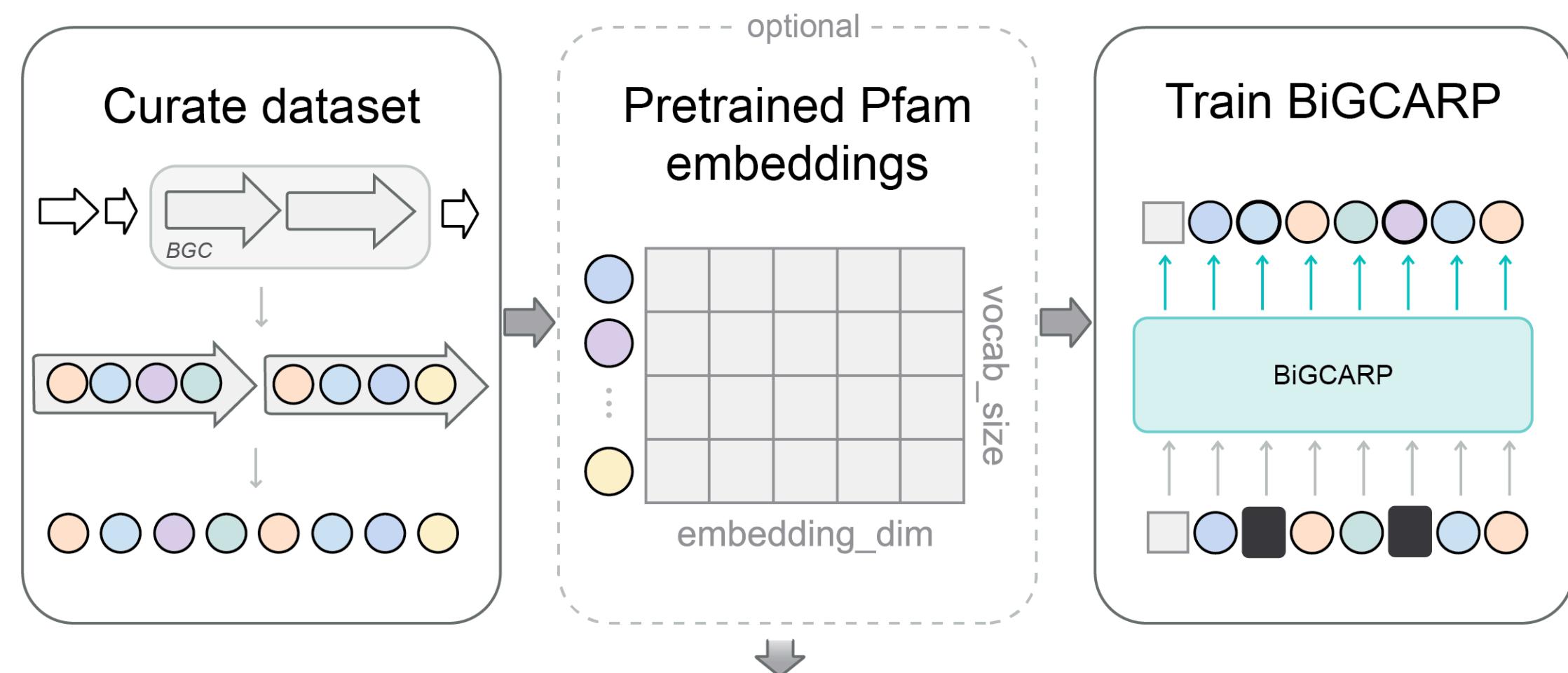
BiGCARP predicts BGC product classes



	BiGCARP-ensemble	DeepBGC
polyketide	0.898	0.903
NRP	0.898	0.907
RiPP	0.963	0.907
saccharide	0.773	0.811
other	0.763	0.583
terpene	0.869	0.824
alkaloid	0.820	0.607
average	0.855	0.792



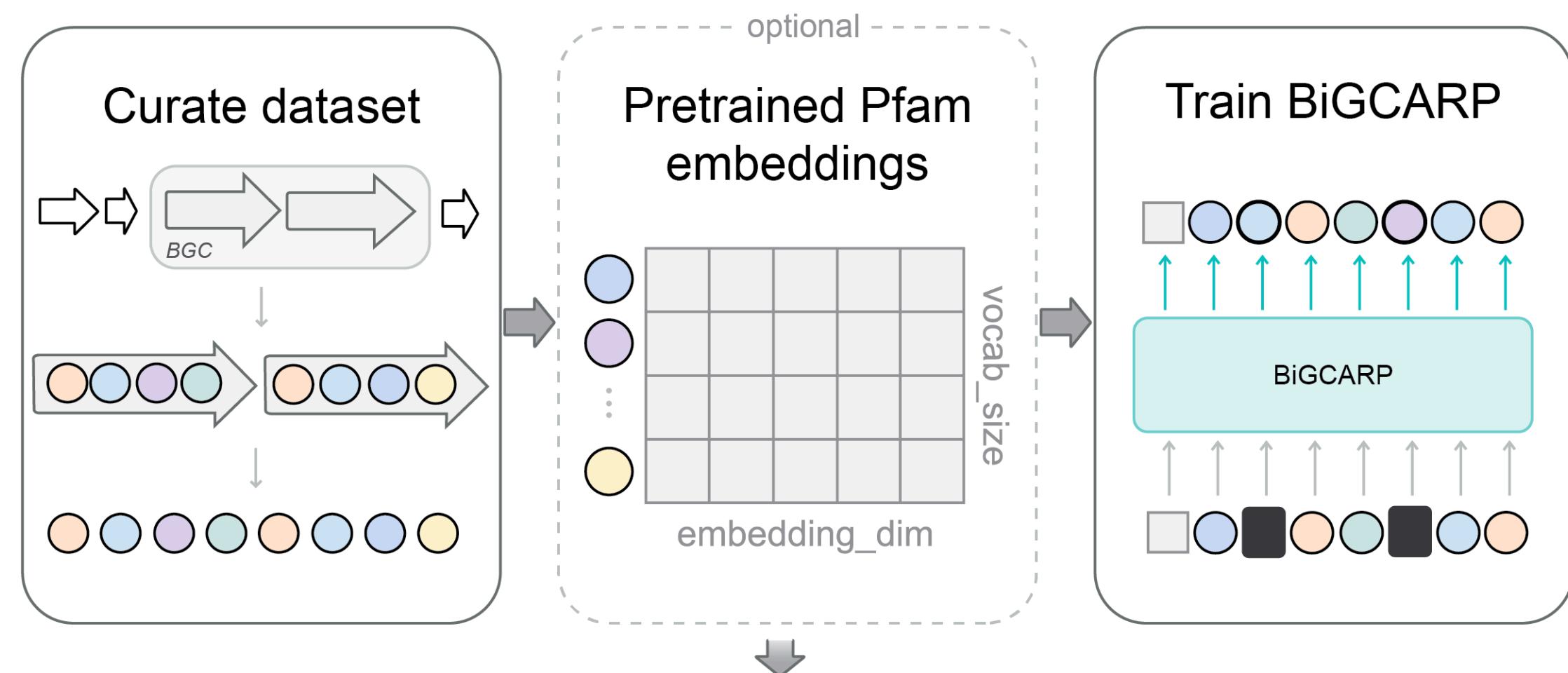
BiGCARP detects and classifies BGCs



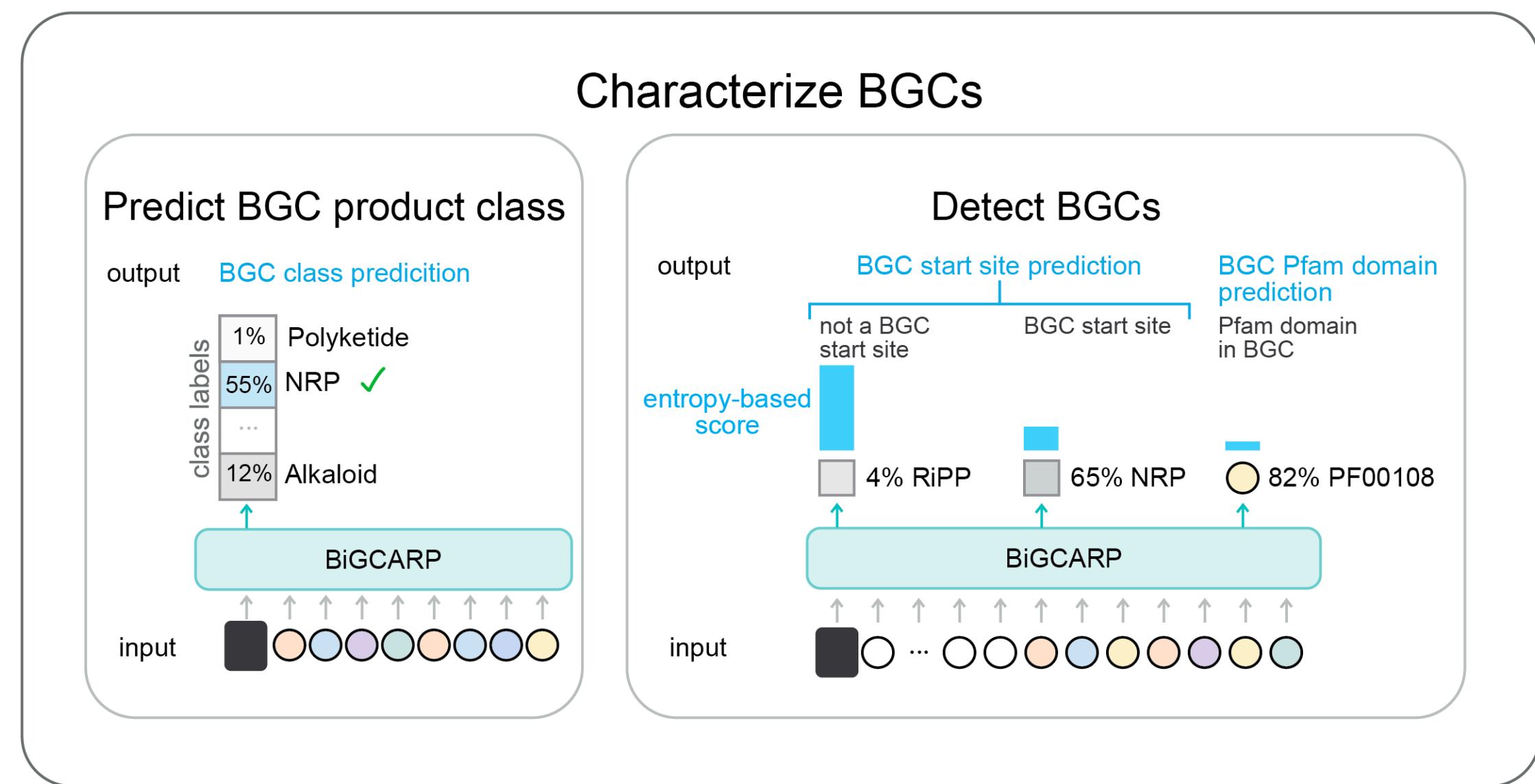
gene →

Pfam	BGC class	non-BGC Pfam	mask token
------	-----------	--------------	------------

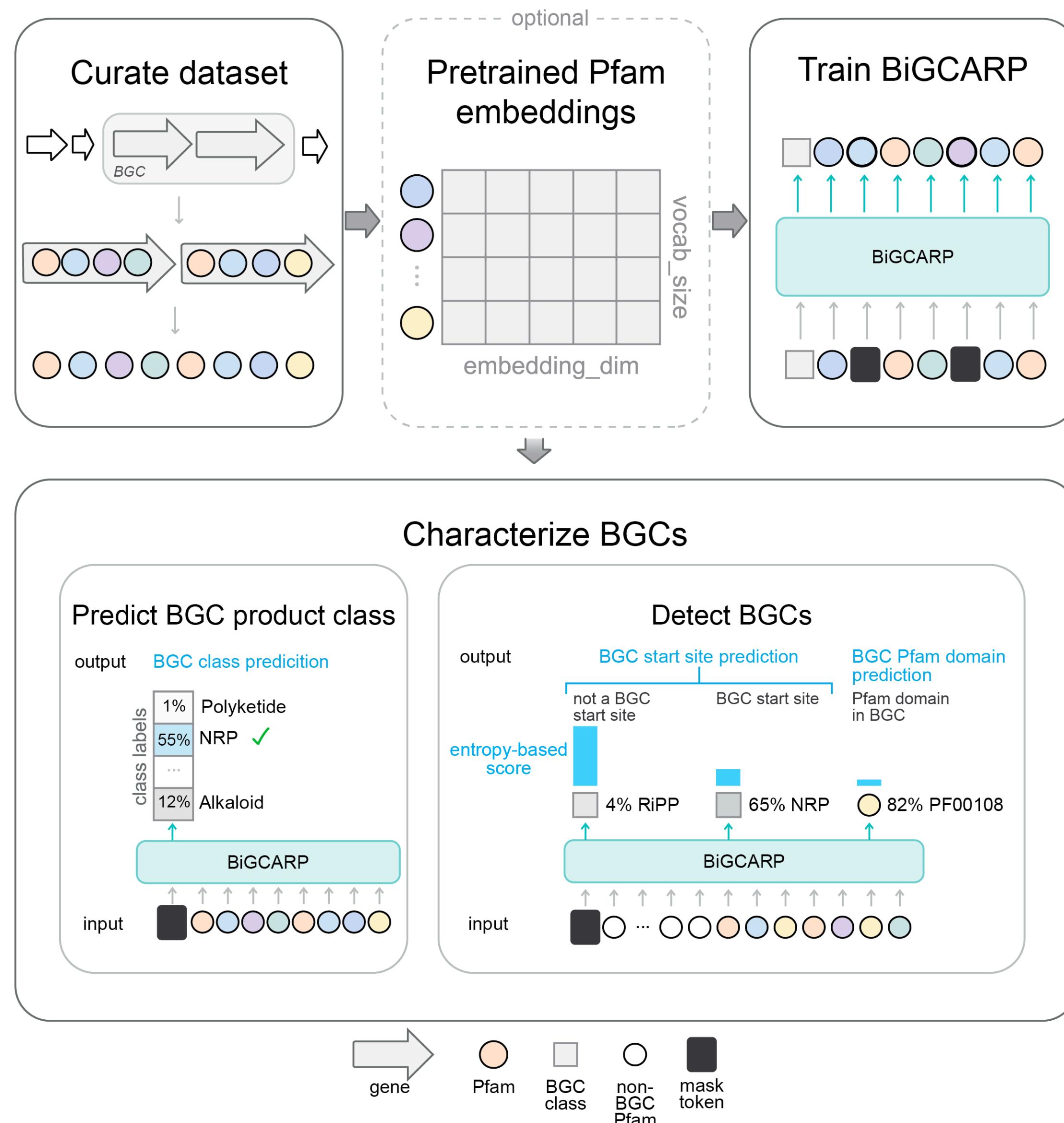
BiGCARP detects and classifies BGCs



- First masked language model on domains and for BGCs

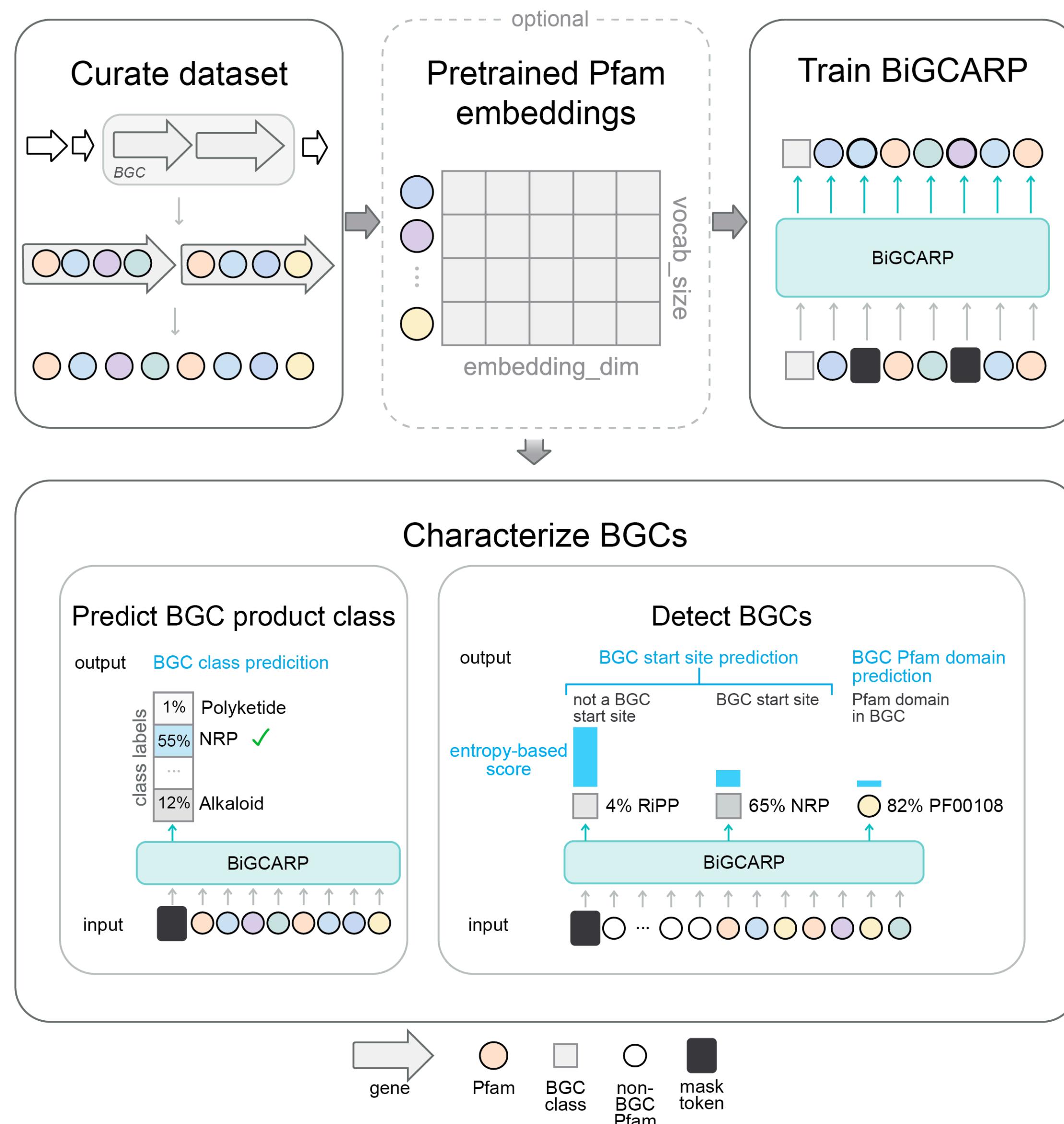


BiGCARP detects and classifies BGCs



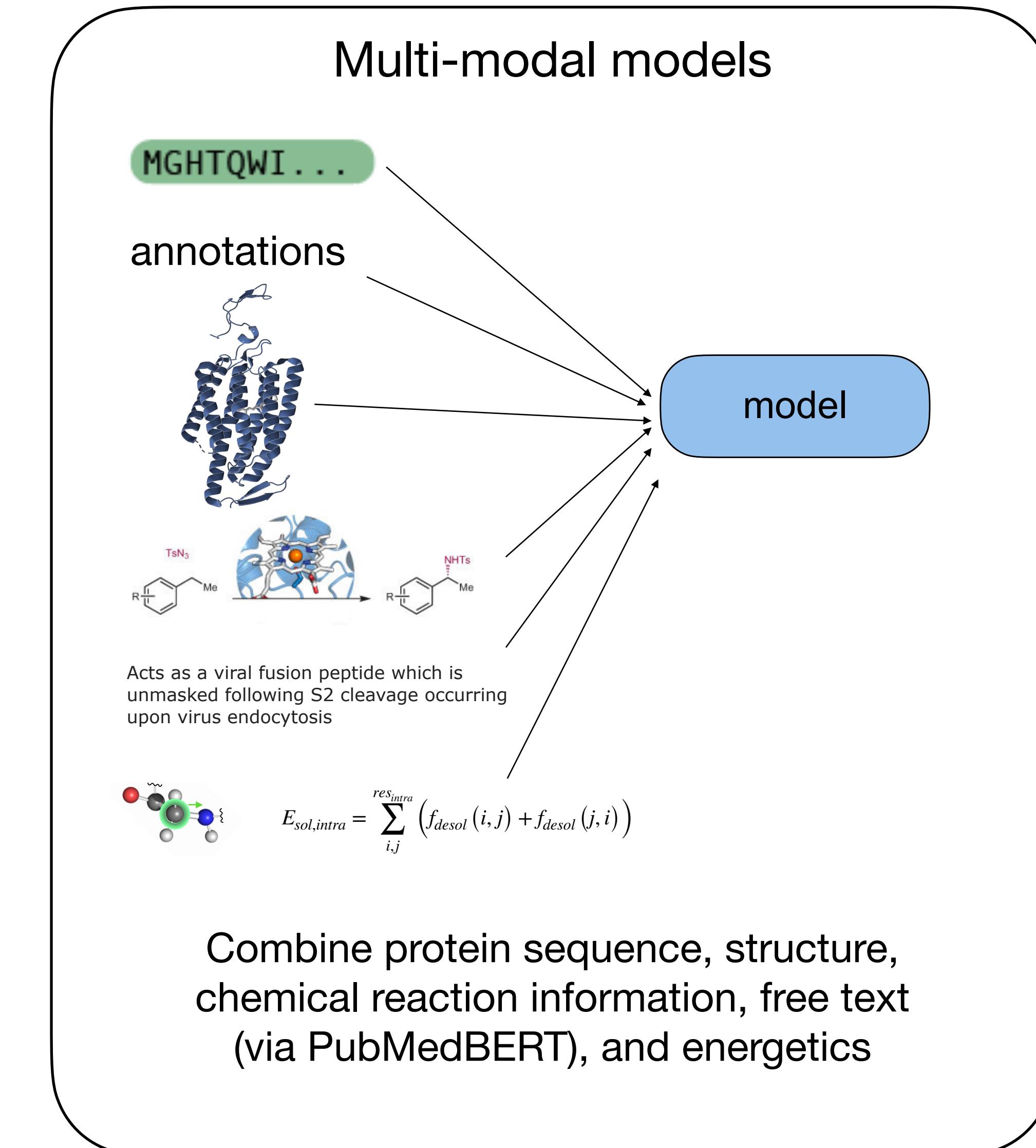
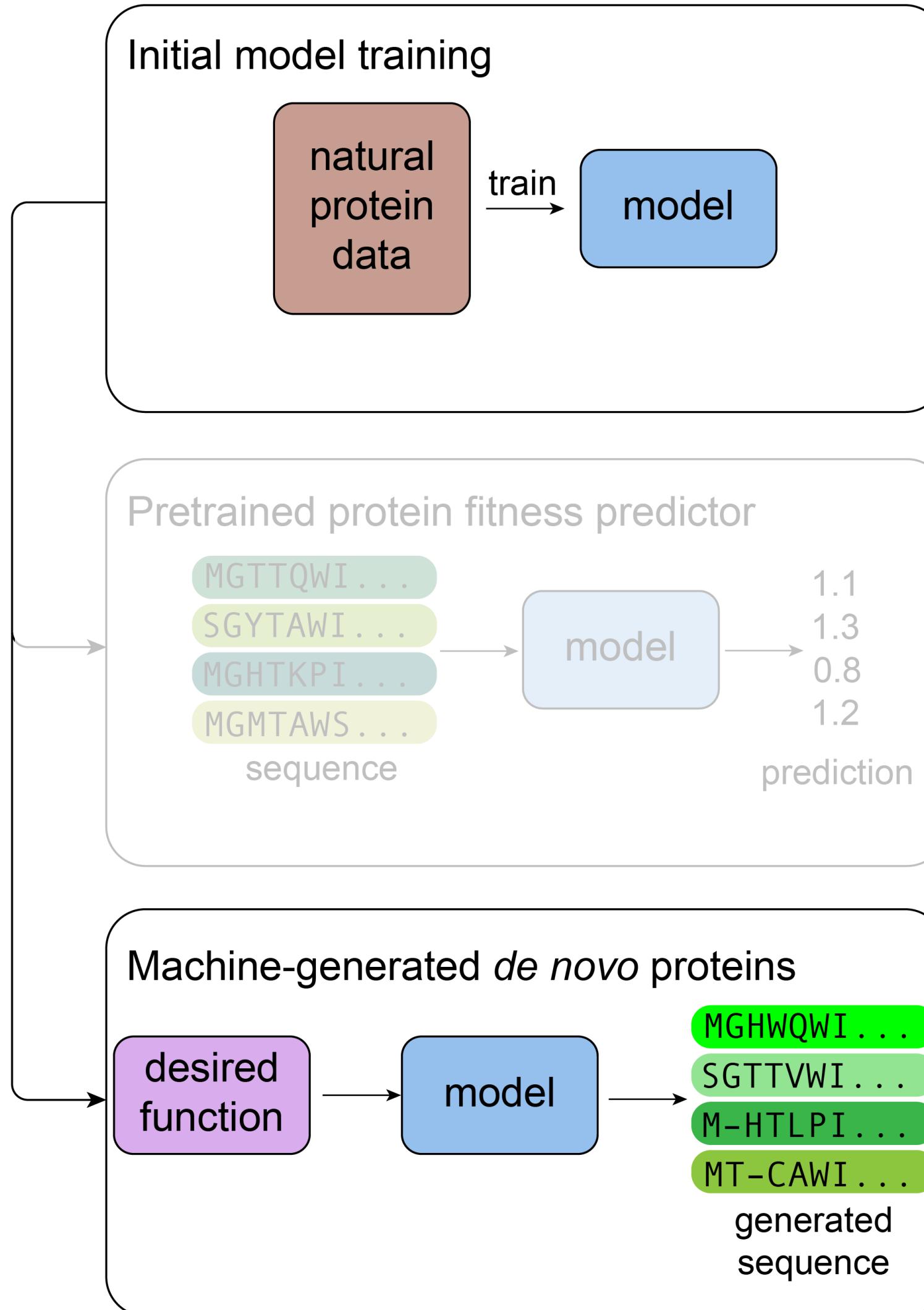
- First masked language model on domains and for BGCs
- Unsupervised BGC detection

BiGCARP detects and classifies BGCs

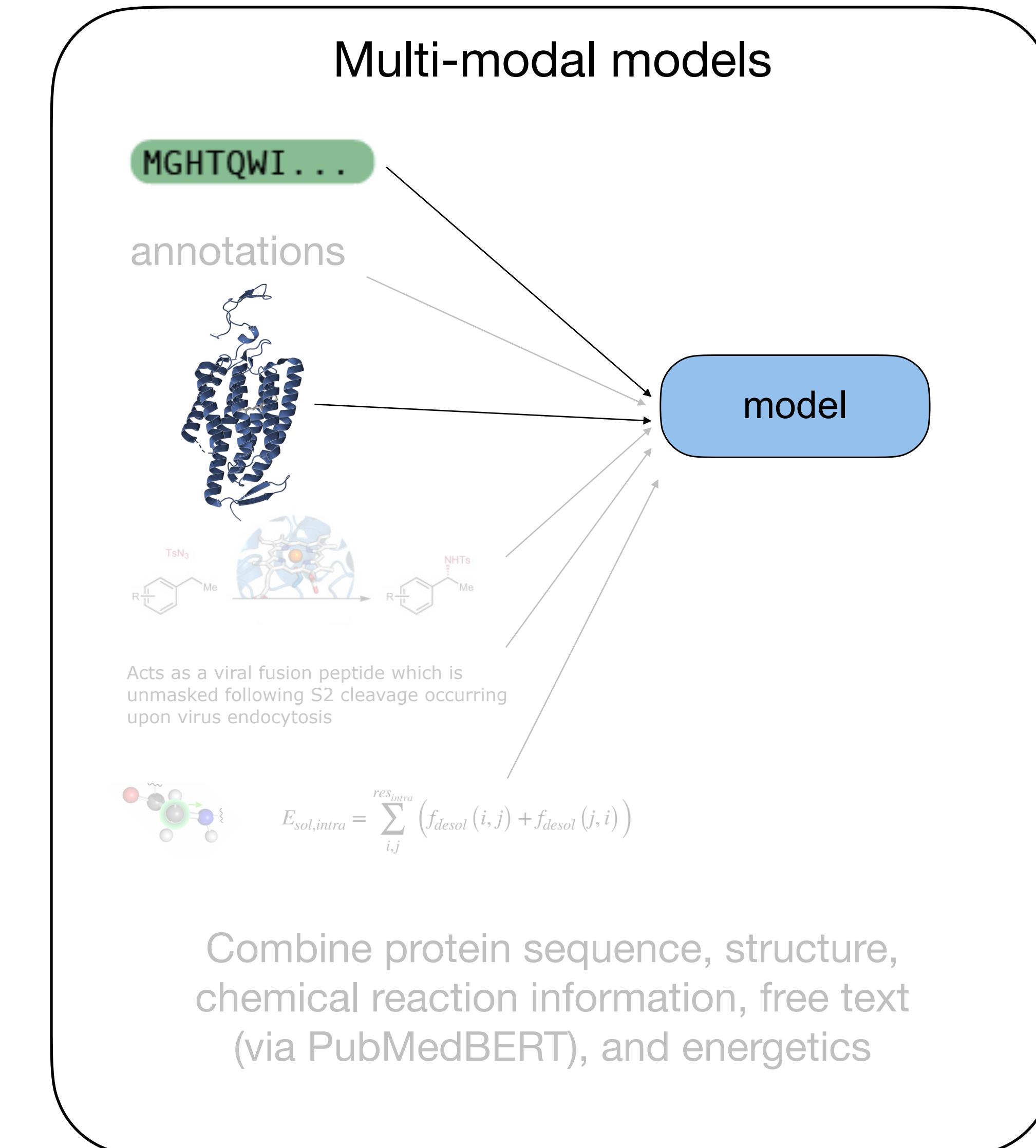
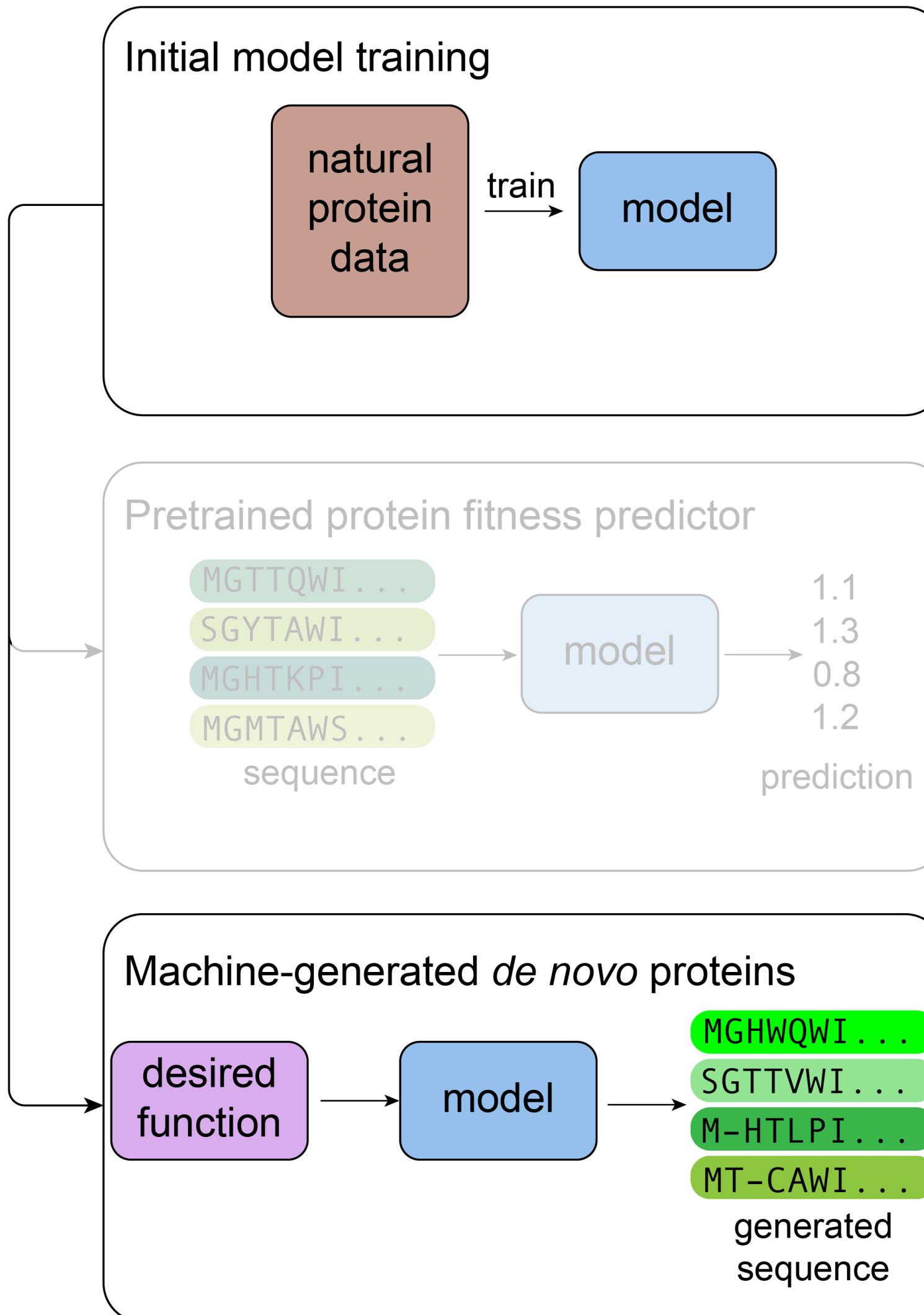


- First masked language model on domains and for BGCs
- Unsupervised BGC detection
- Learn product classes and BGC structure simultaneously

Use multiple data modalities to discover and design proteins



Use multiple data modalities to discover and design proteins

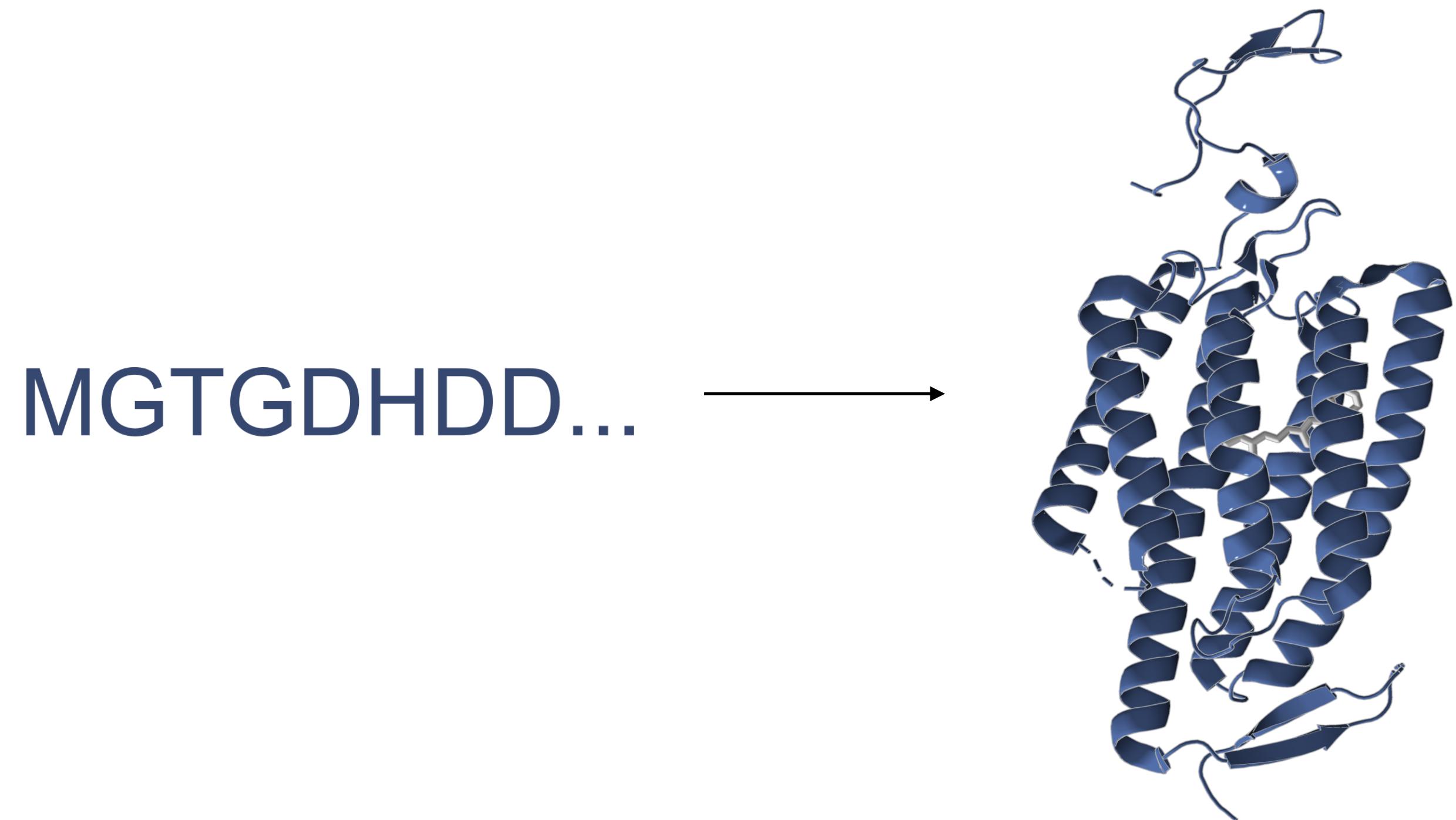


Sequence determines structure determines function

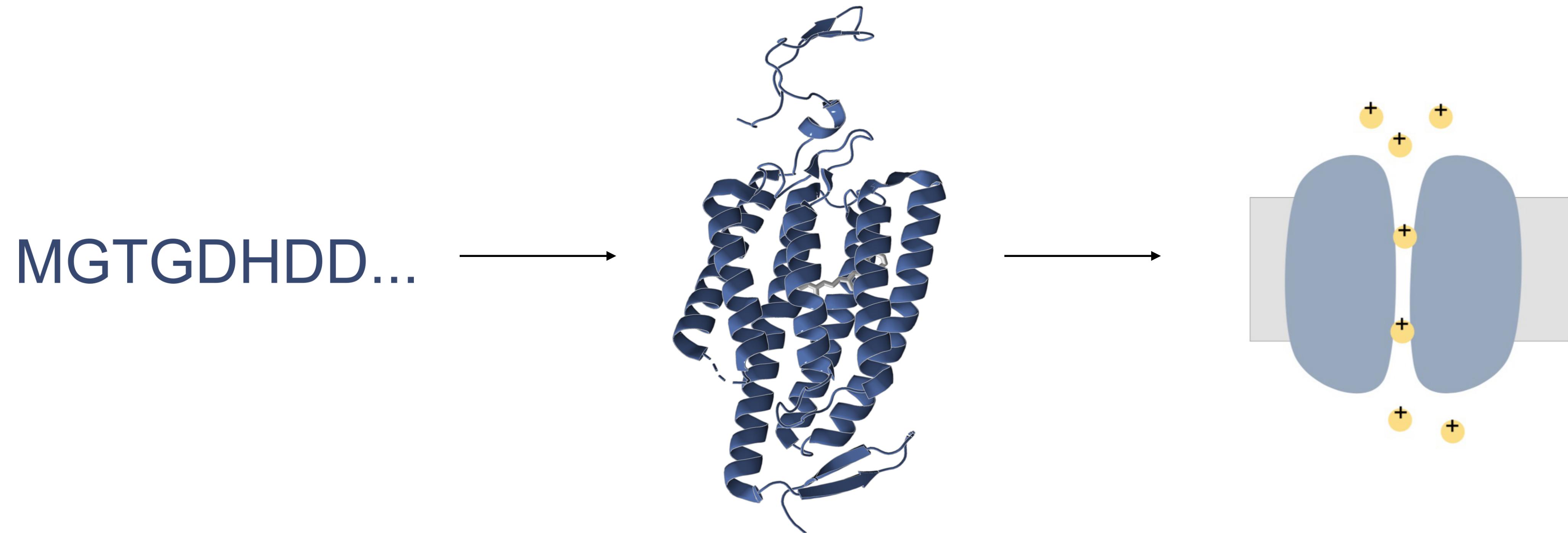
Sequence determines structure determines function

MGTGDHDD...

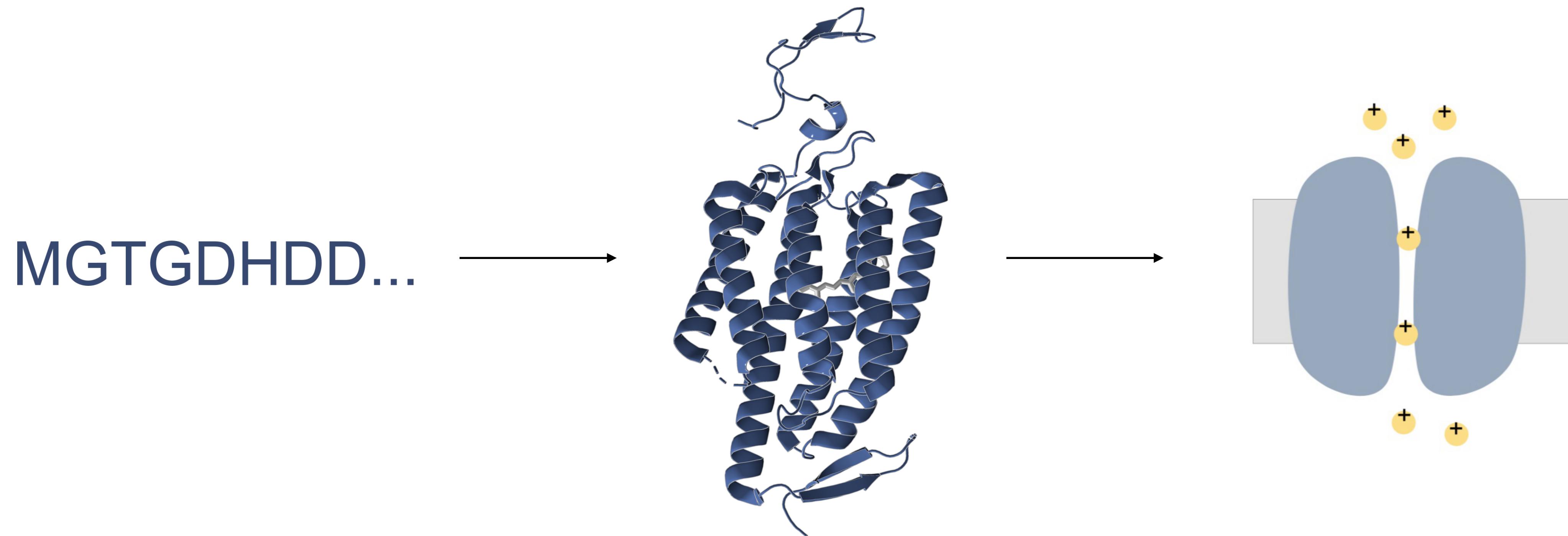
Sequence determines structure determines function



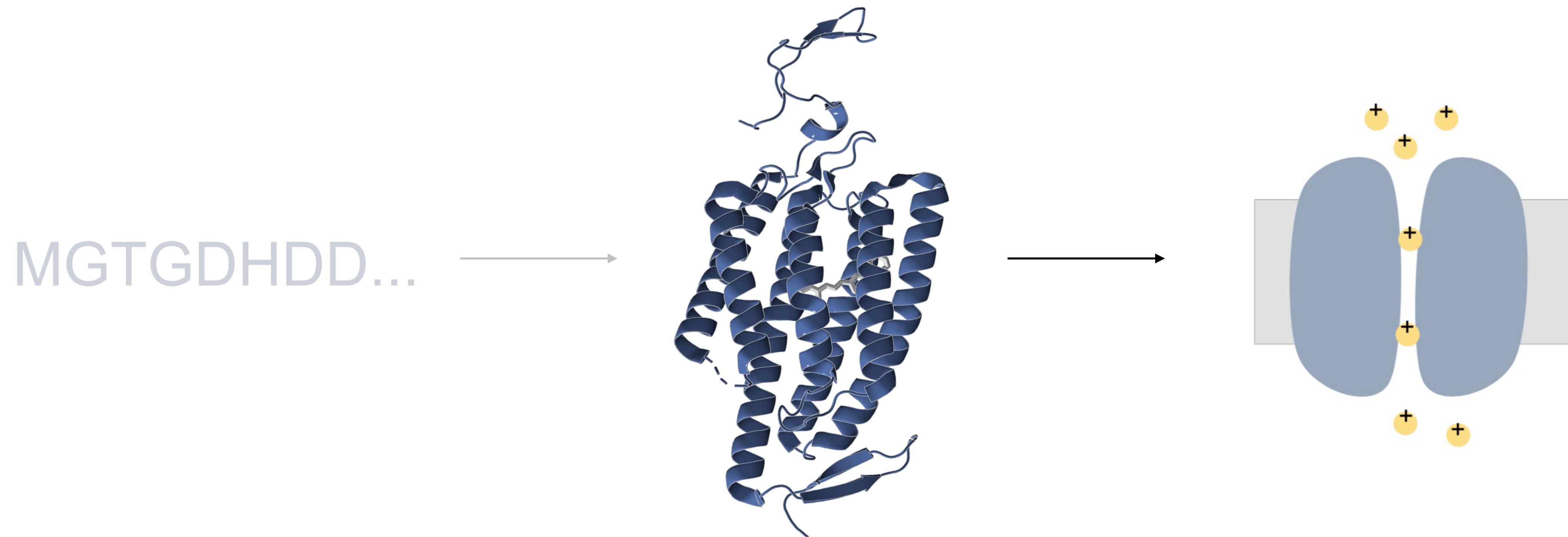
Sequence determines structure determines function



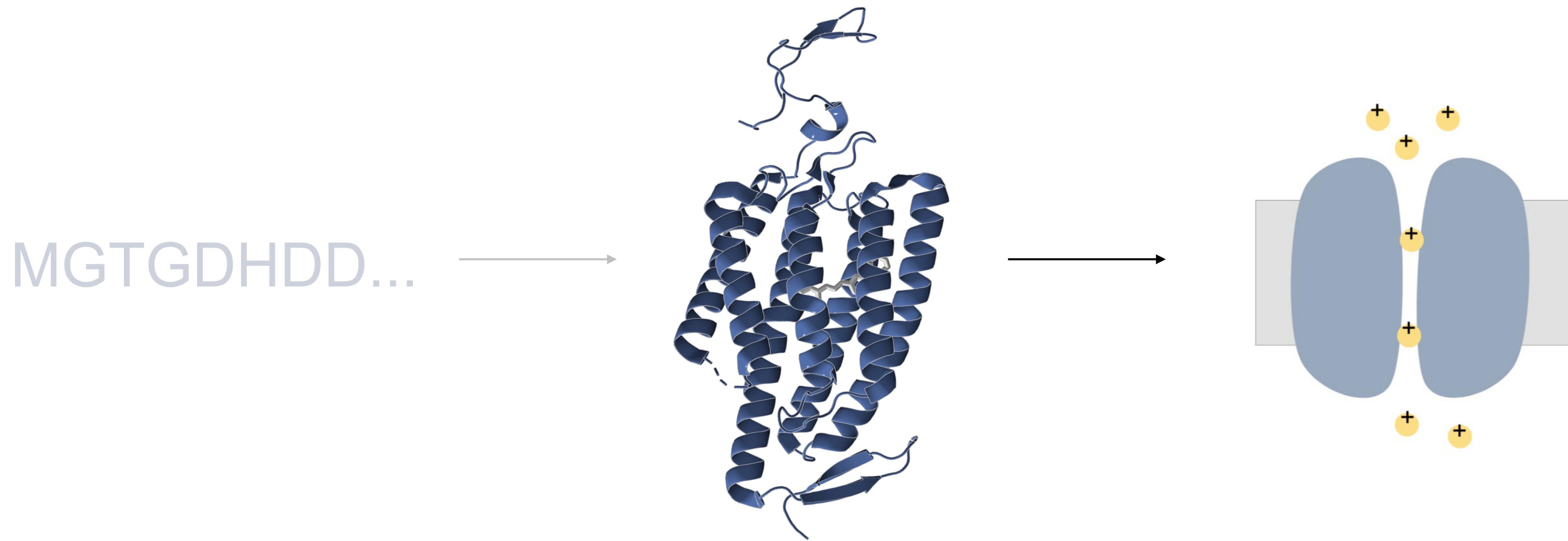
Generating new, designable structures expands functional space



Generating new, designable structures expands functional space

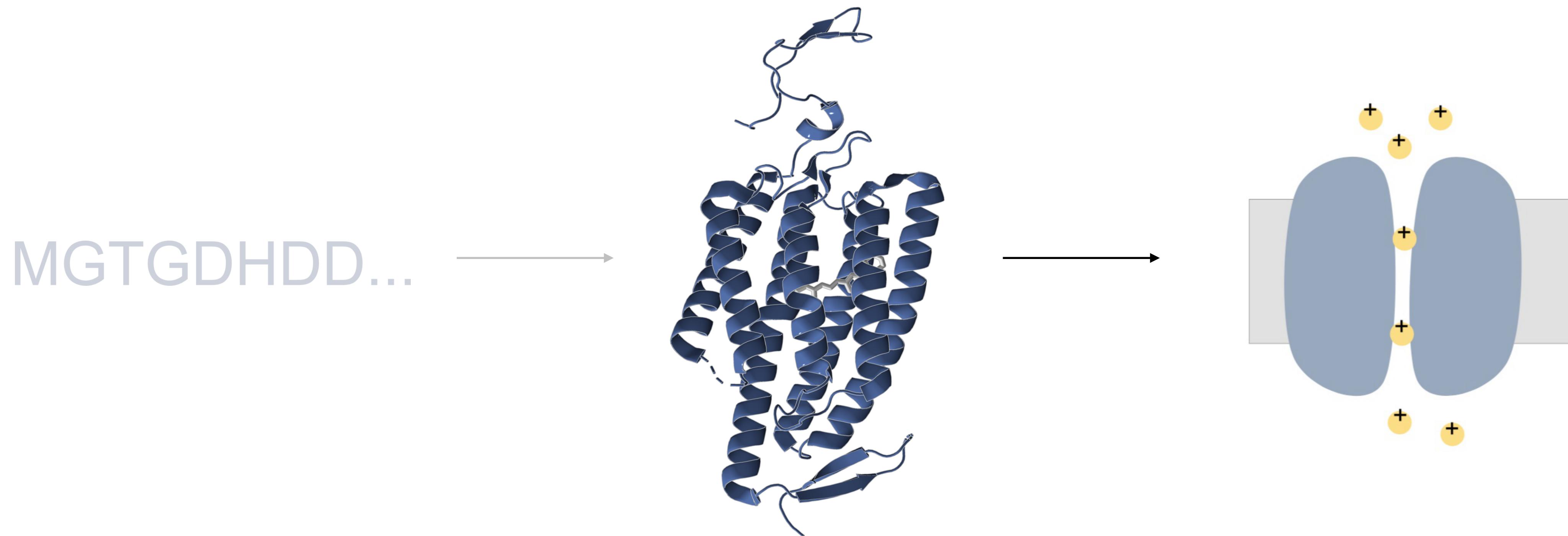


Generating new, designable structures expands functional space



Challenge: Generate diverse and designable structures

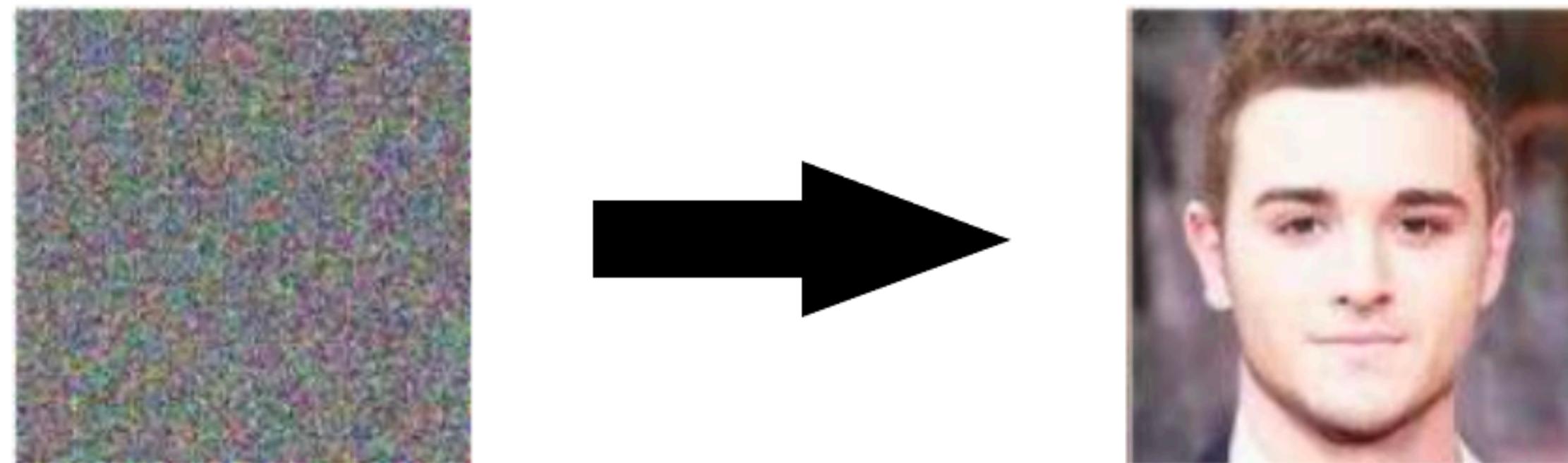
Generating new, designable structures expands functional space



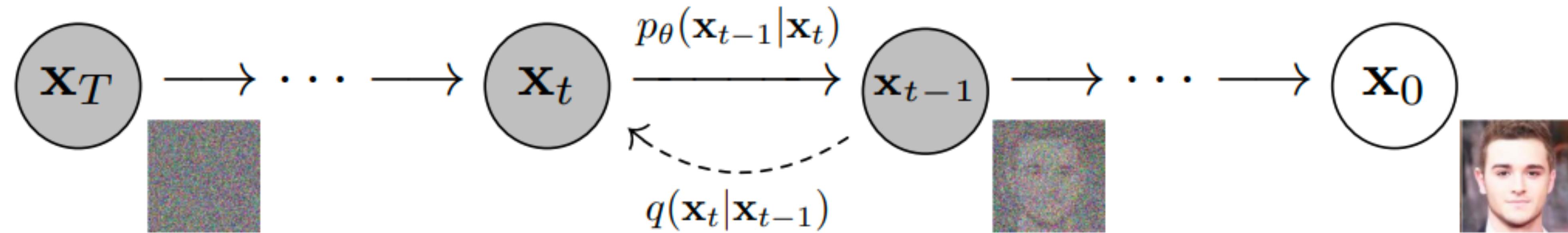
Challenge: Generate diverse and designable structures



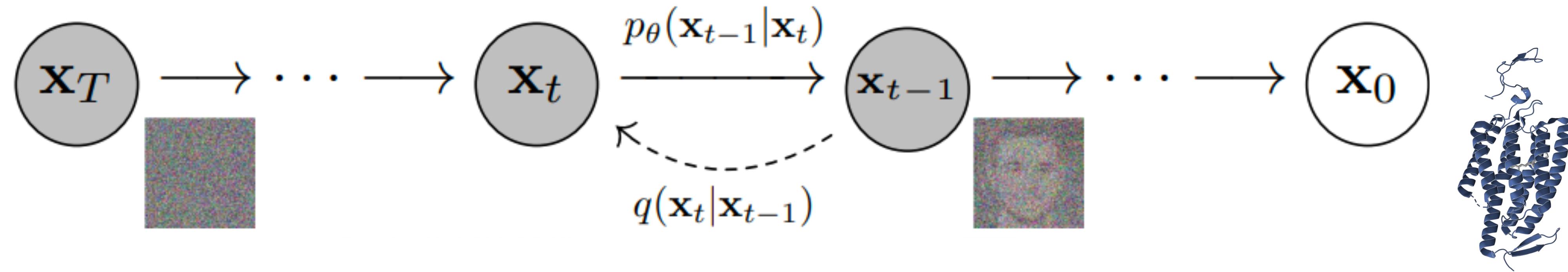
Diffusion models learn to generate objects from noise



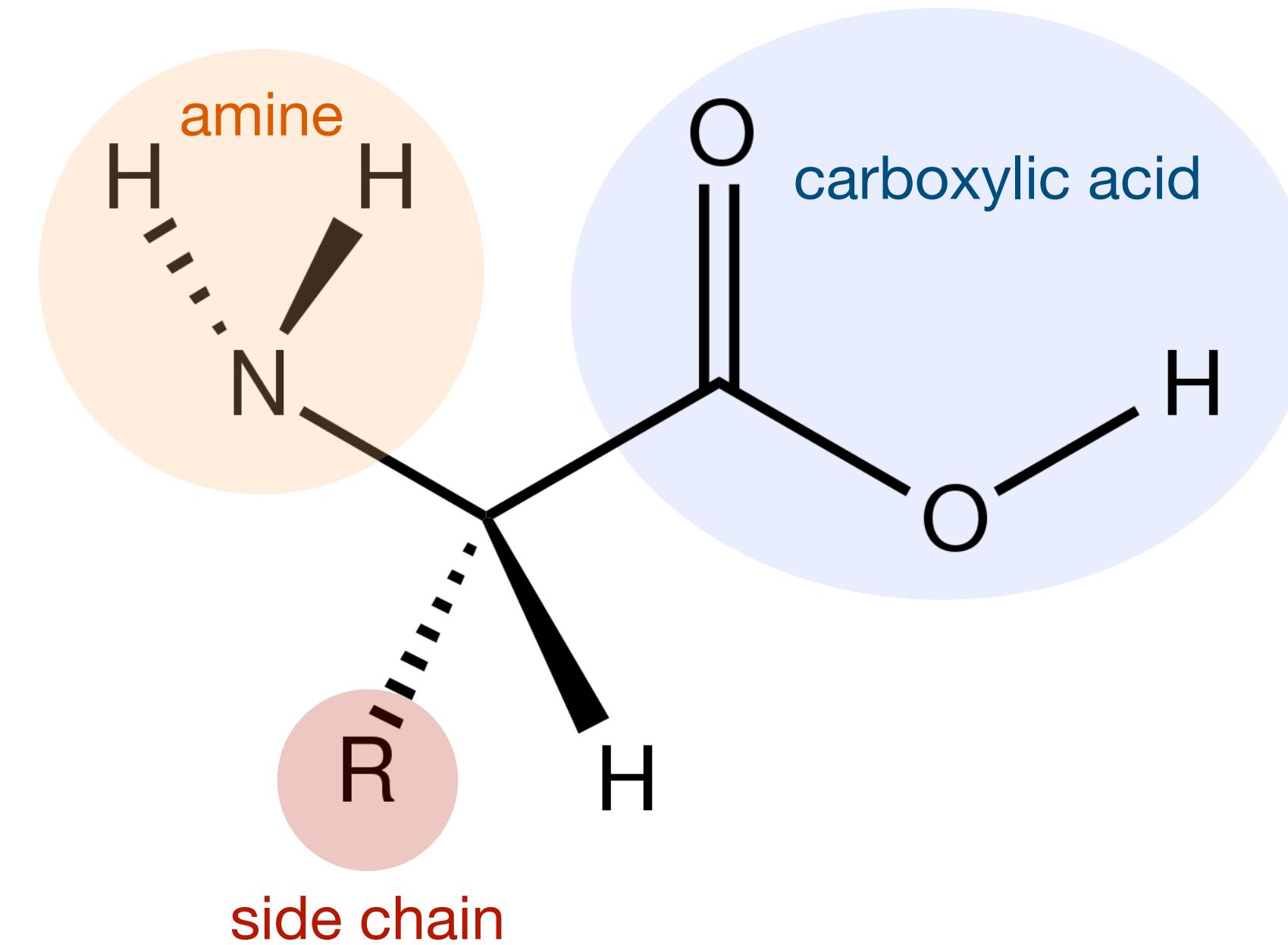
Diffusion models learn to generate objects from noise



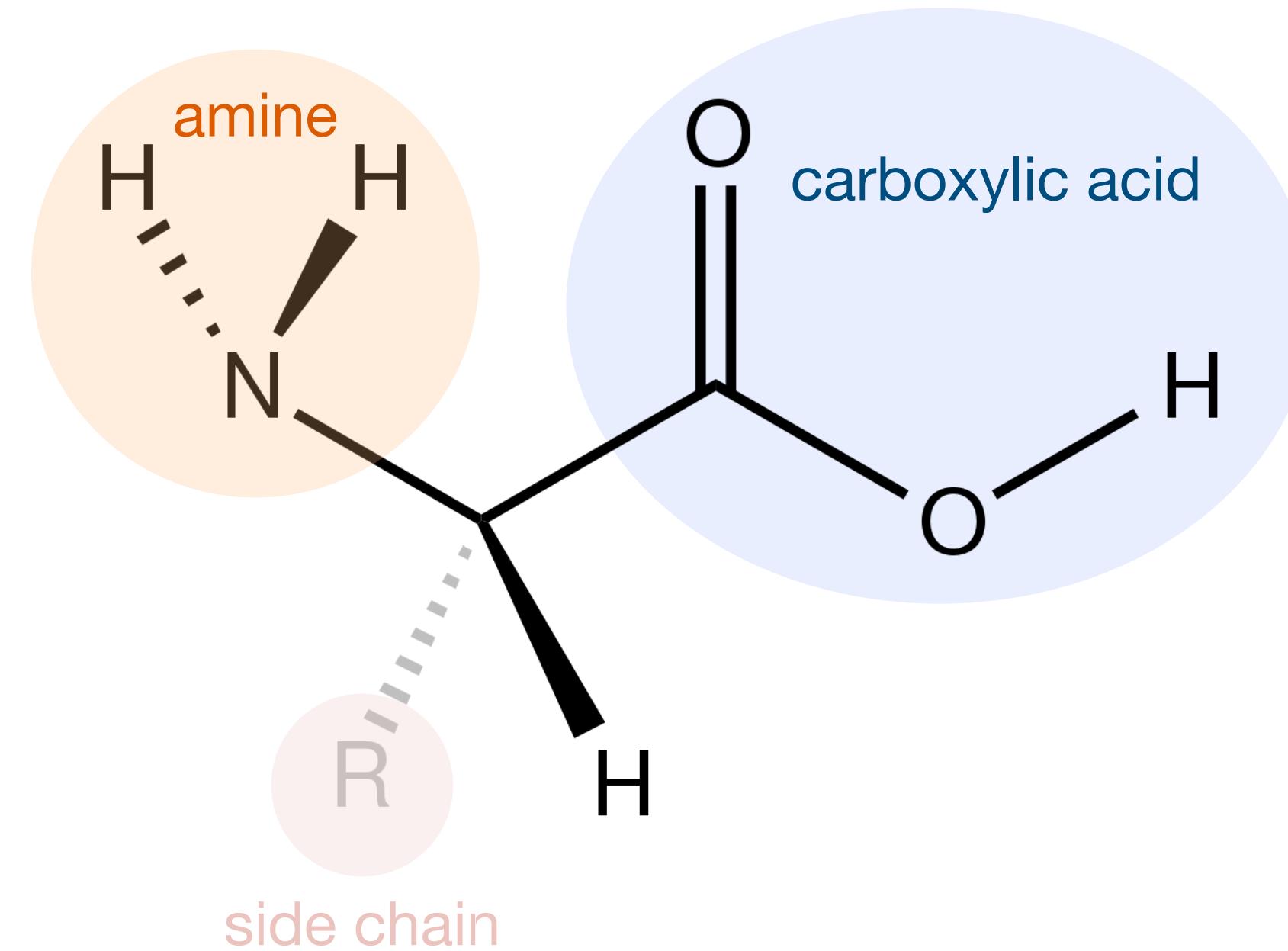
Diffusion models learn to generate objects from noise



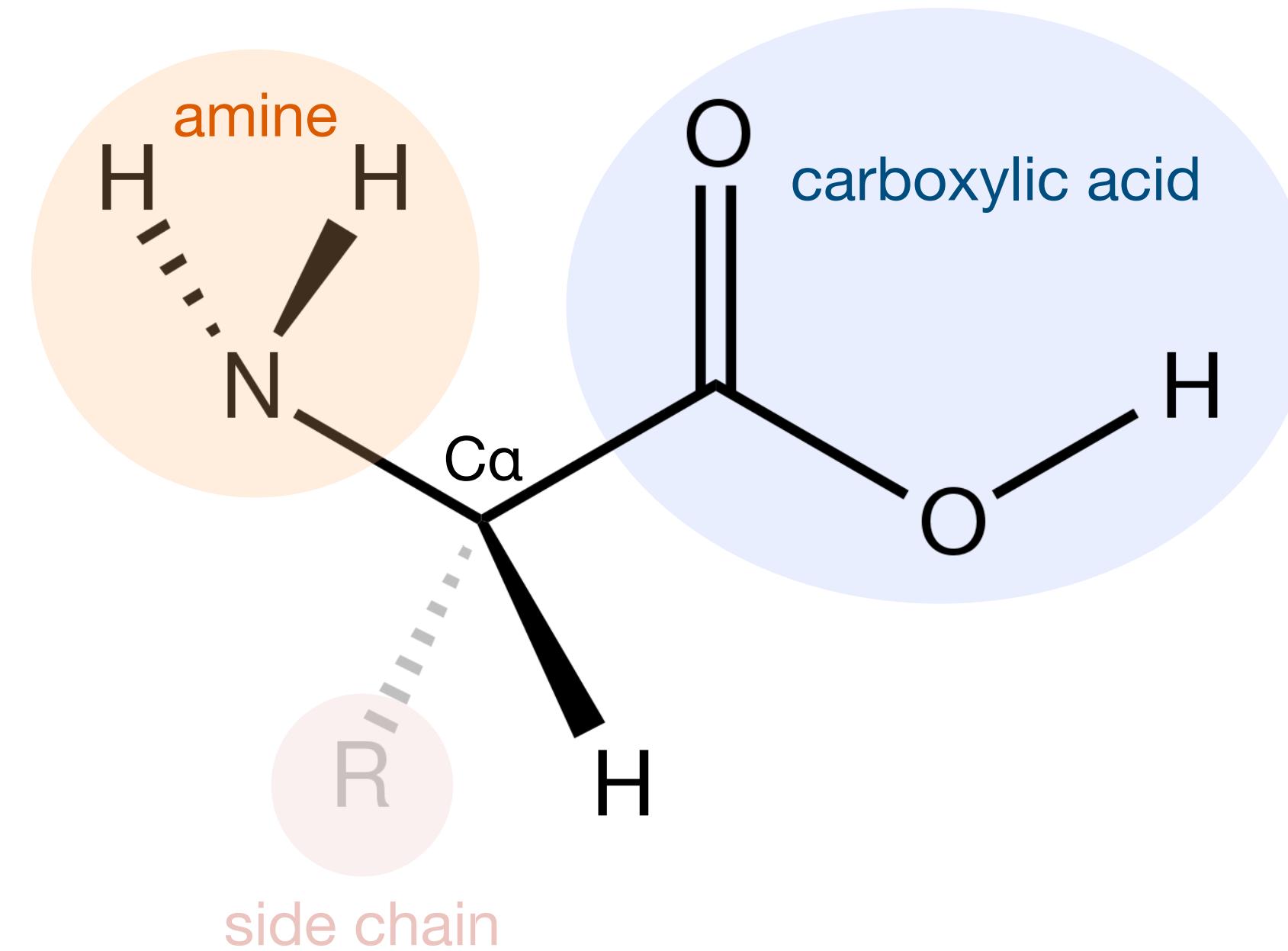
Generate backbone structures



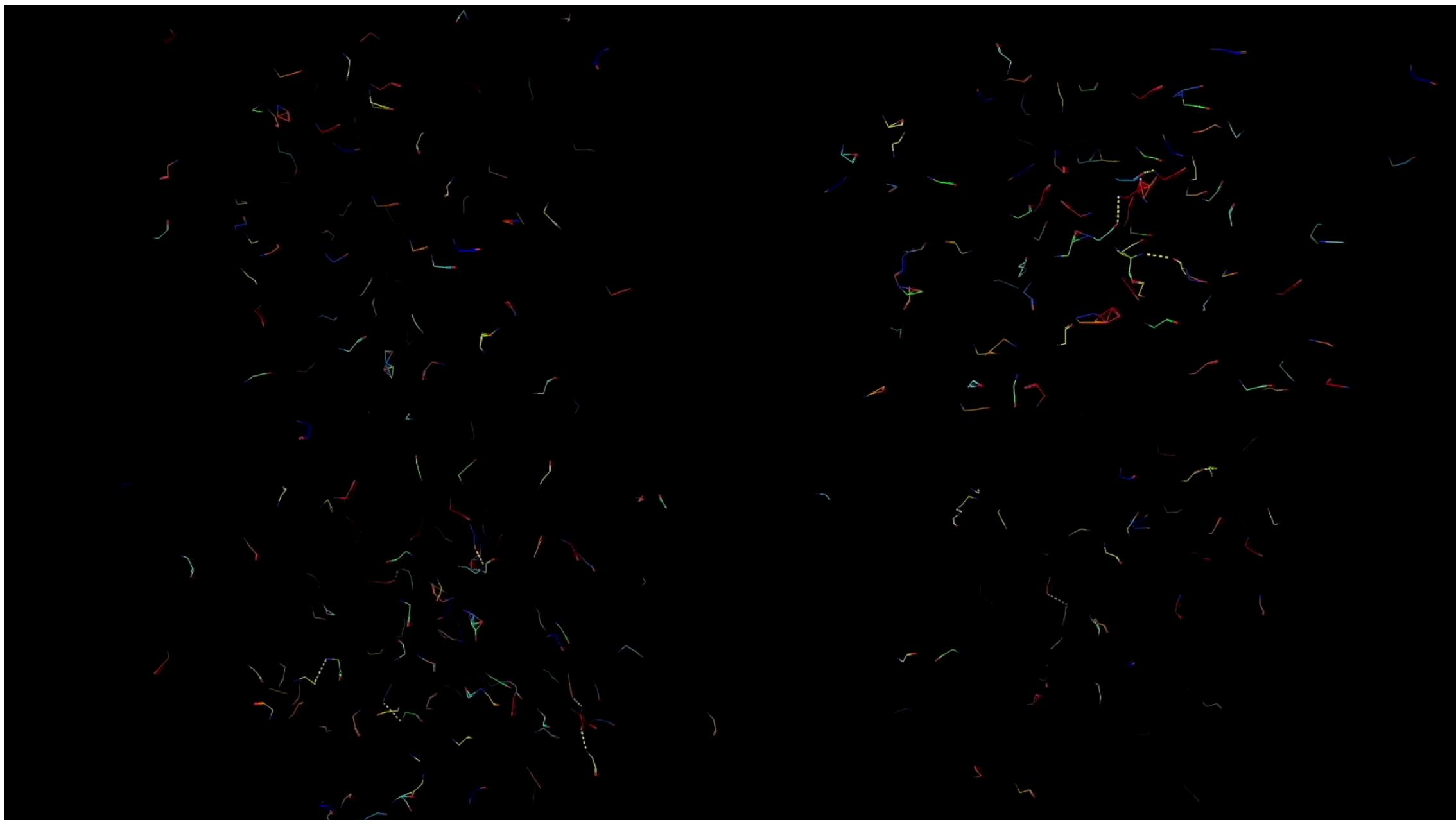
Generate backbone structures



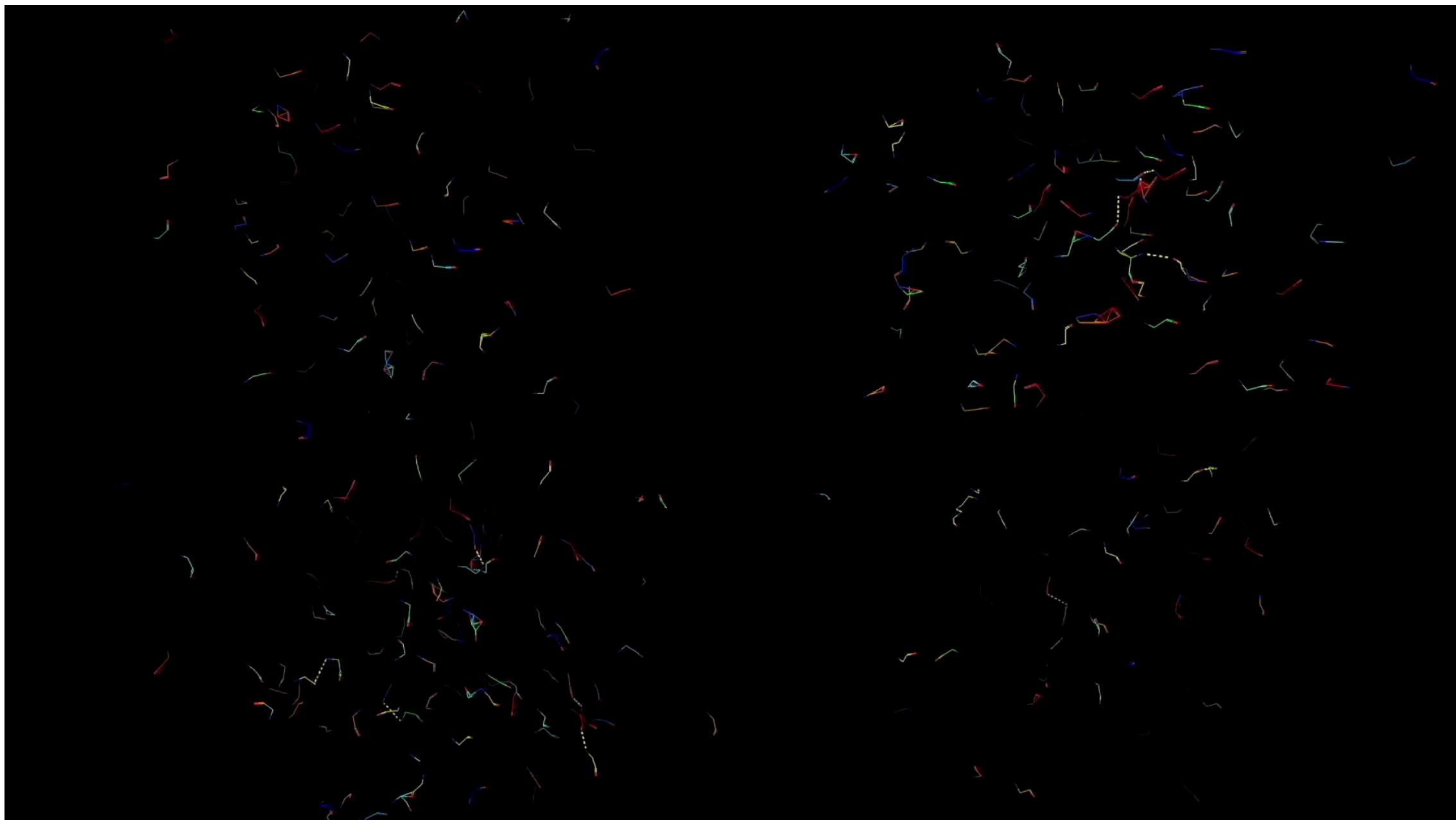
Generate backbone structures



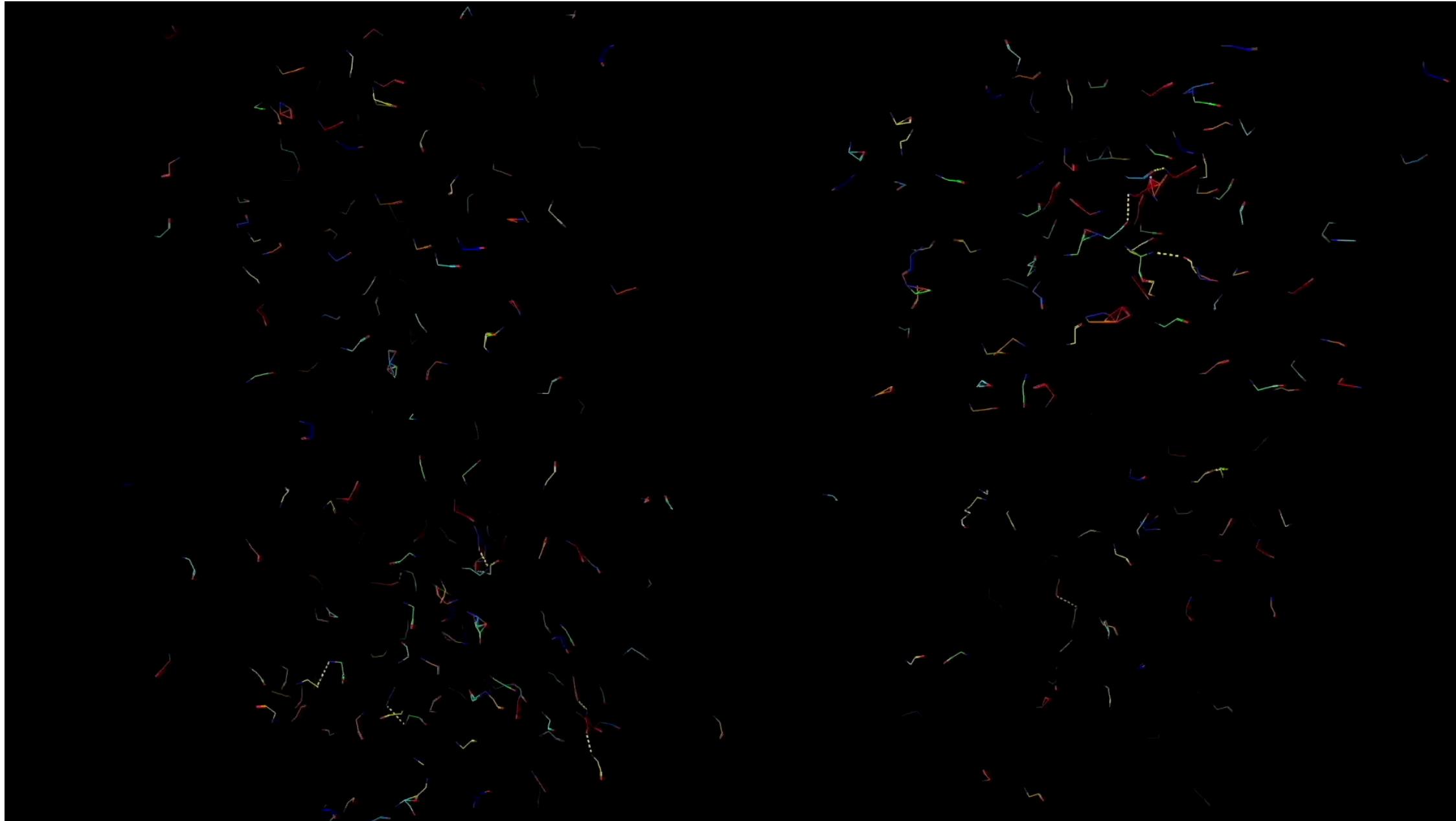
Diffusion on 3D coordinates requires equivariances



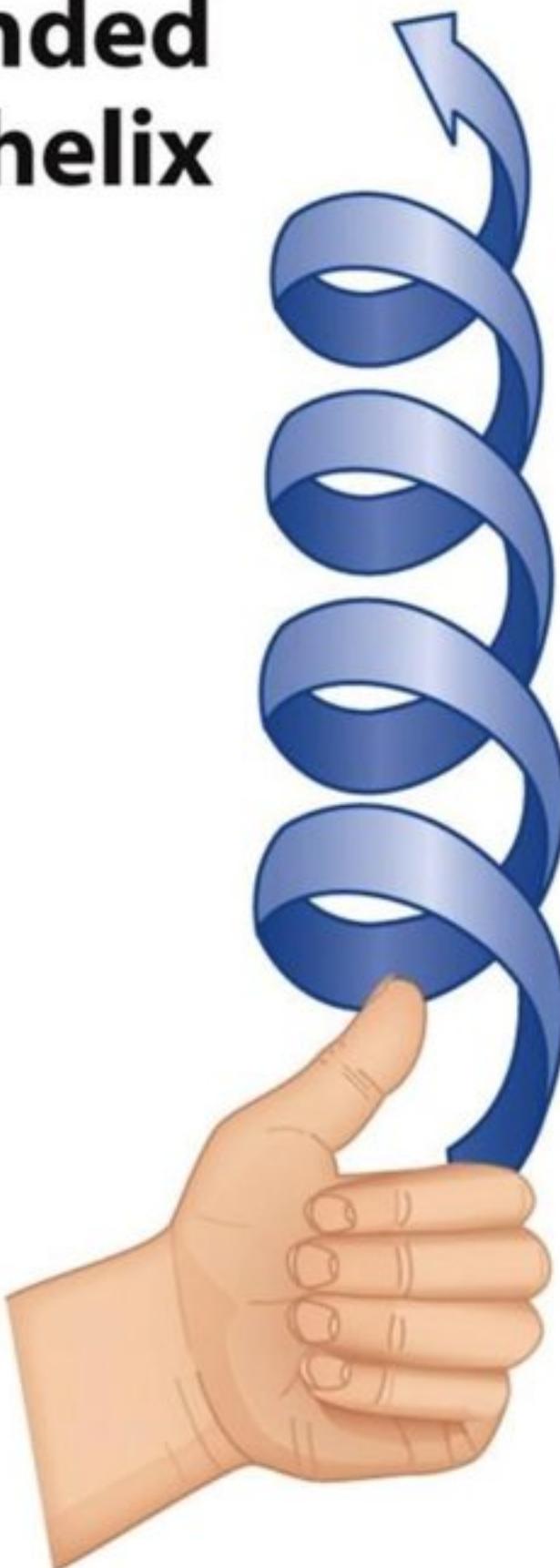
Diffusion on 3D coordinates requires equivariances



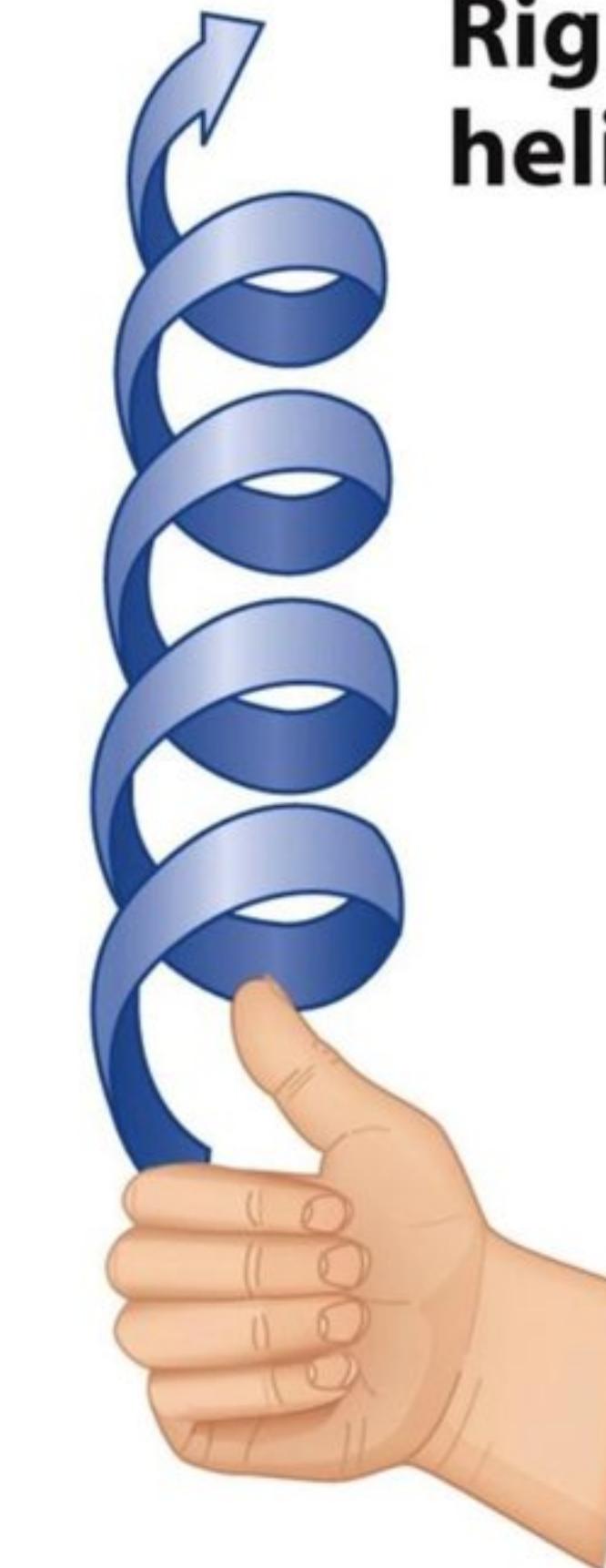
Diffusion on 3D coordinates requires equivariances



**Left-handed
helix**



**Right-handed
helix**



Box 4-1
Lehninger Principles of Biochemistry, Fifth Edition
© 2008 W.H. Freeman and Company

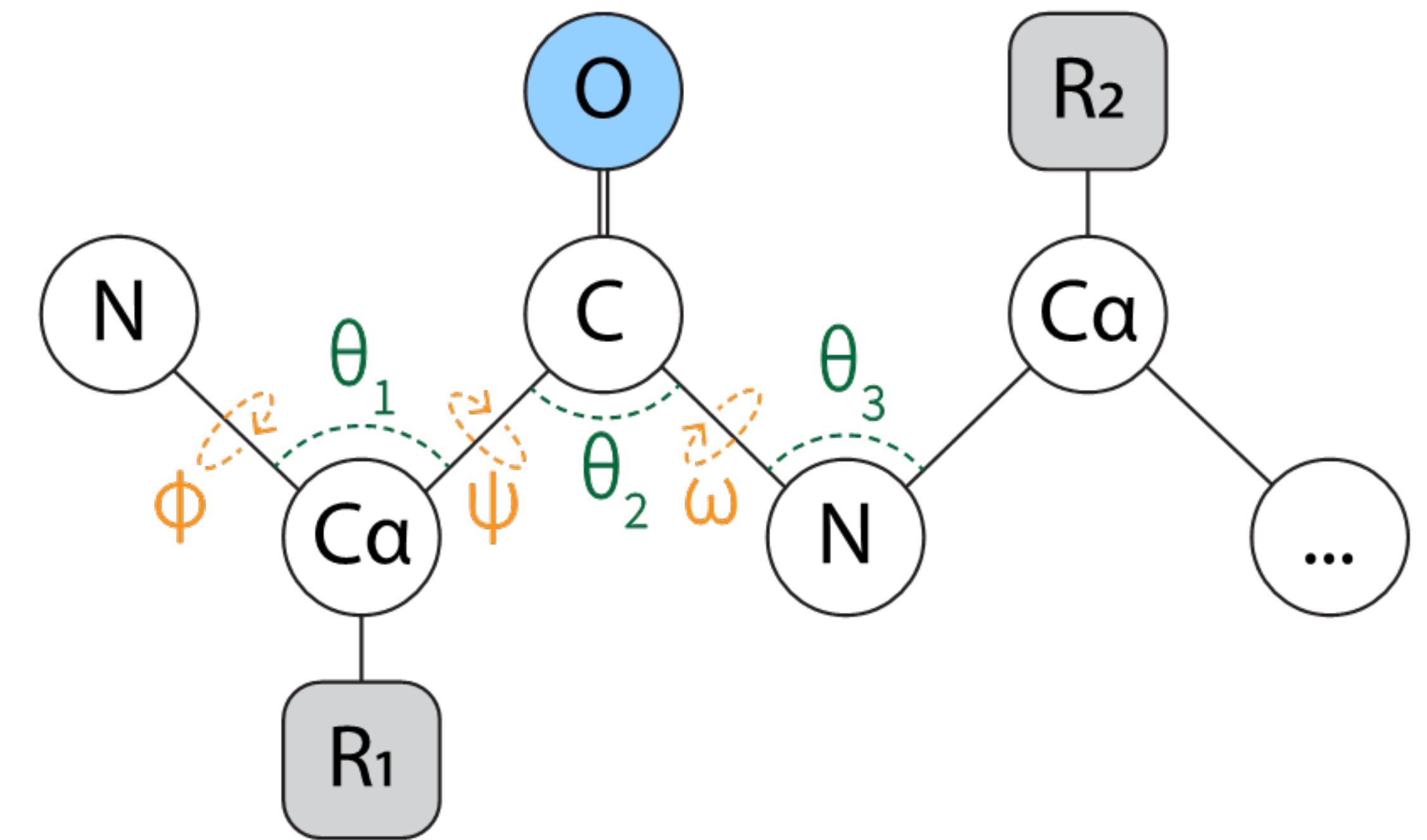
Folding Diffusion uses coordinates inspired by protein folding



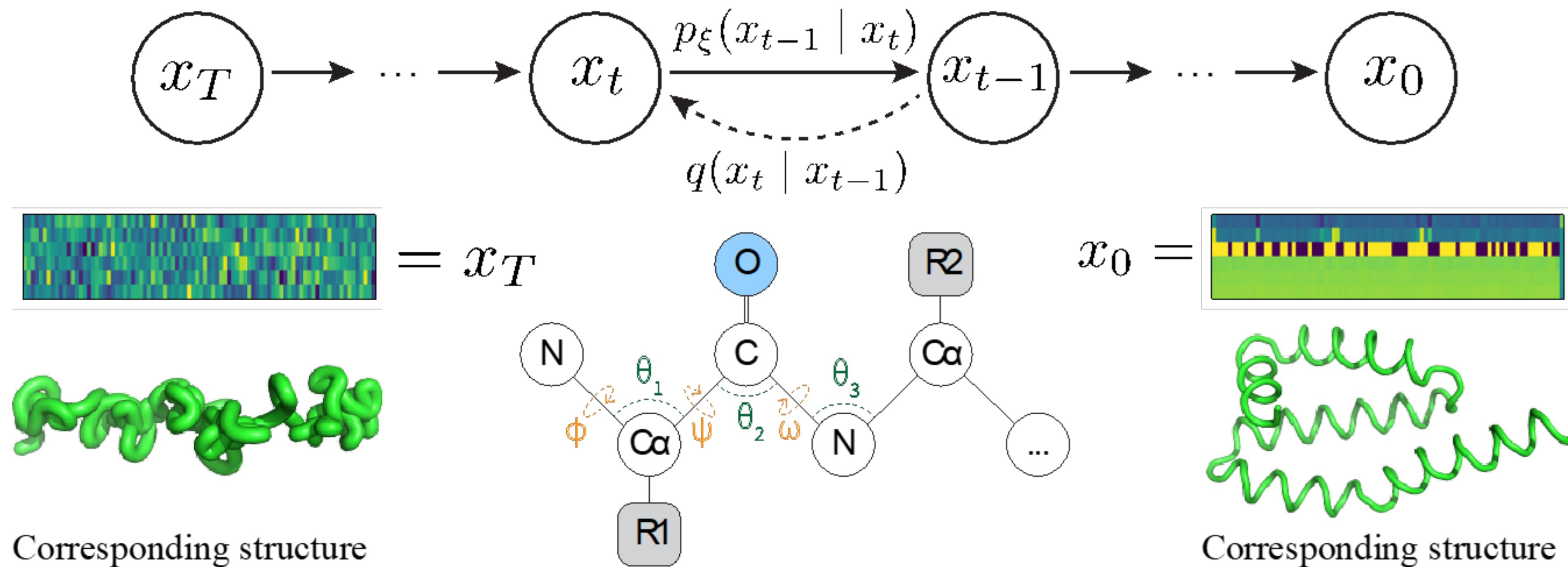
Folding Diffusion uses coordinates inspired by protein folding



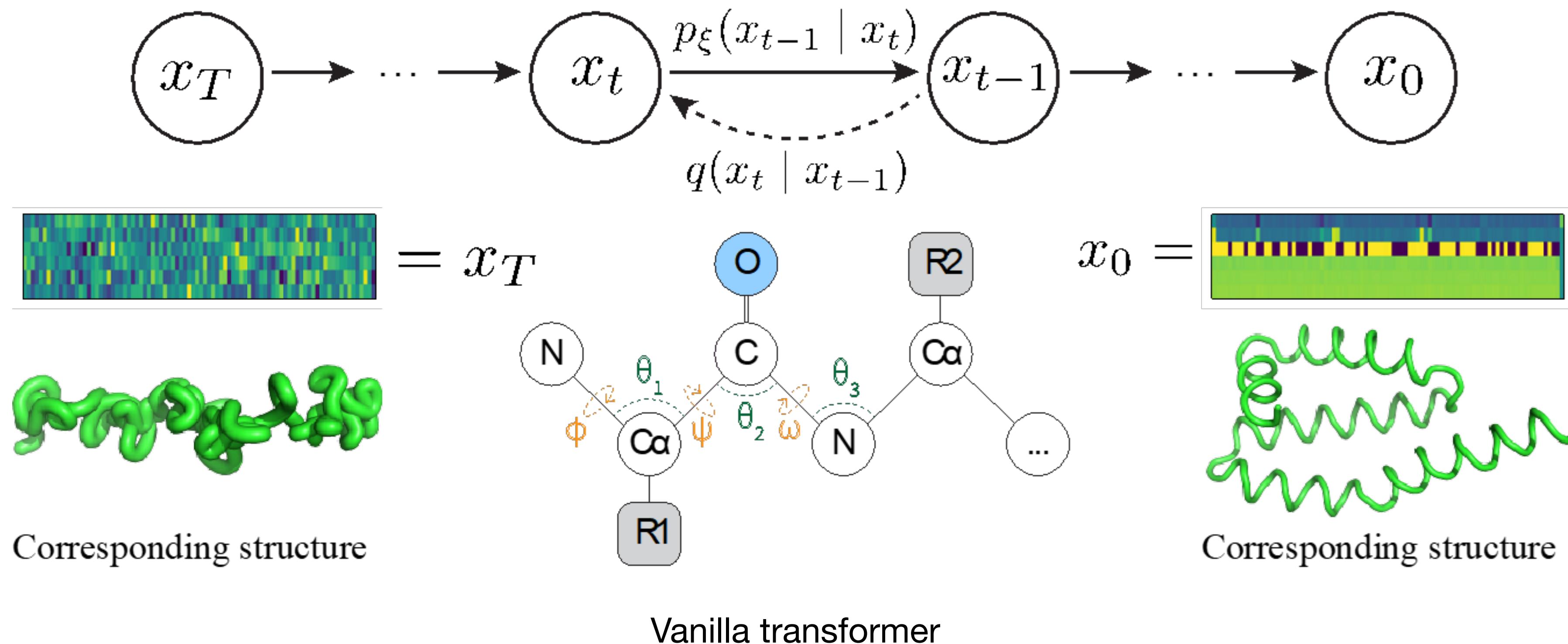
Folding Diffusion uses coordinates inspired by protein folding



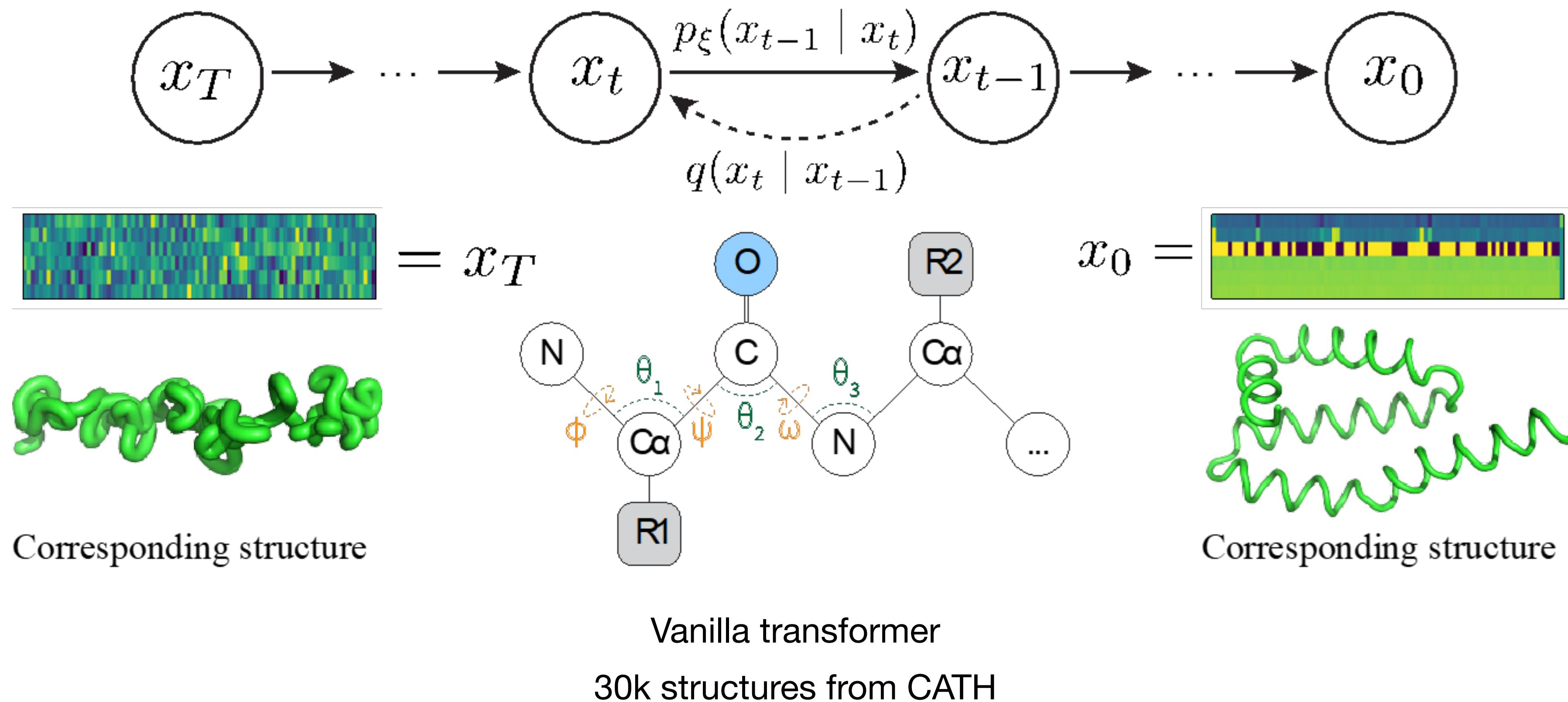
Folding Diffusion uses coordinates inspired by protein folding



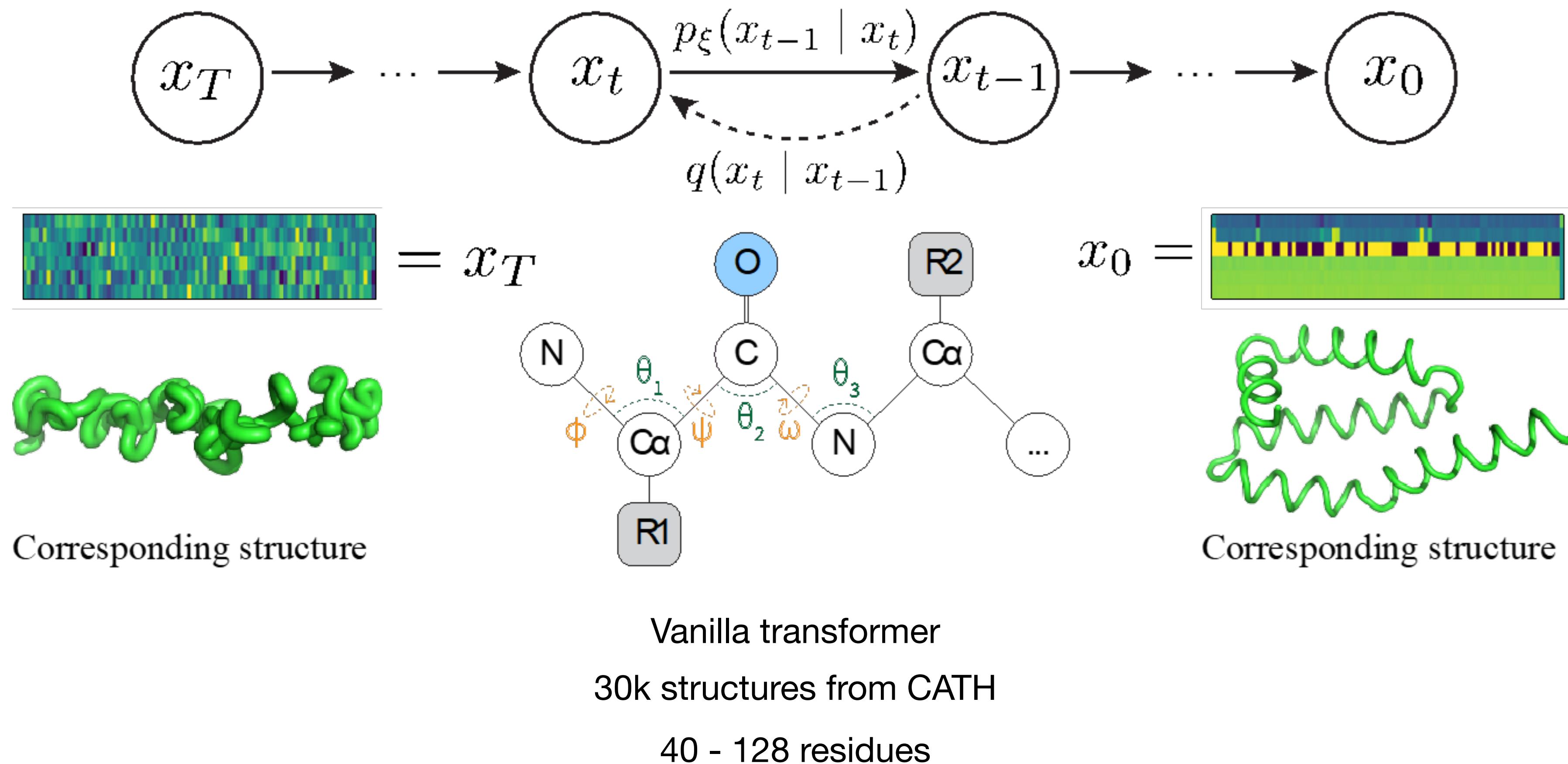
Folding Diffusion uses coordinates inspired by protein folding



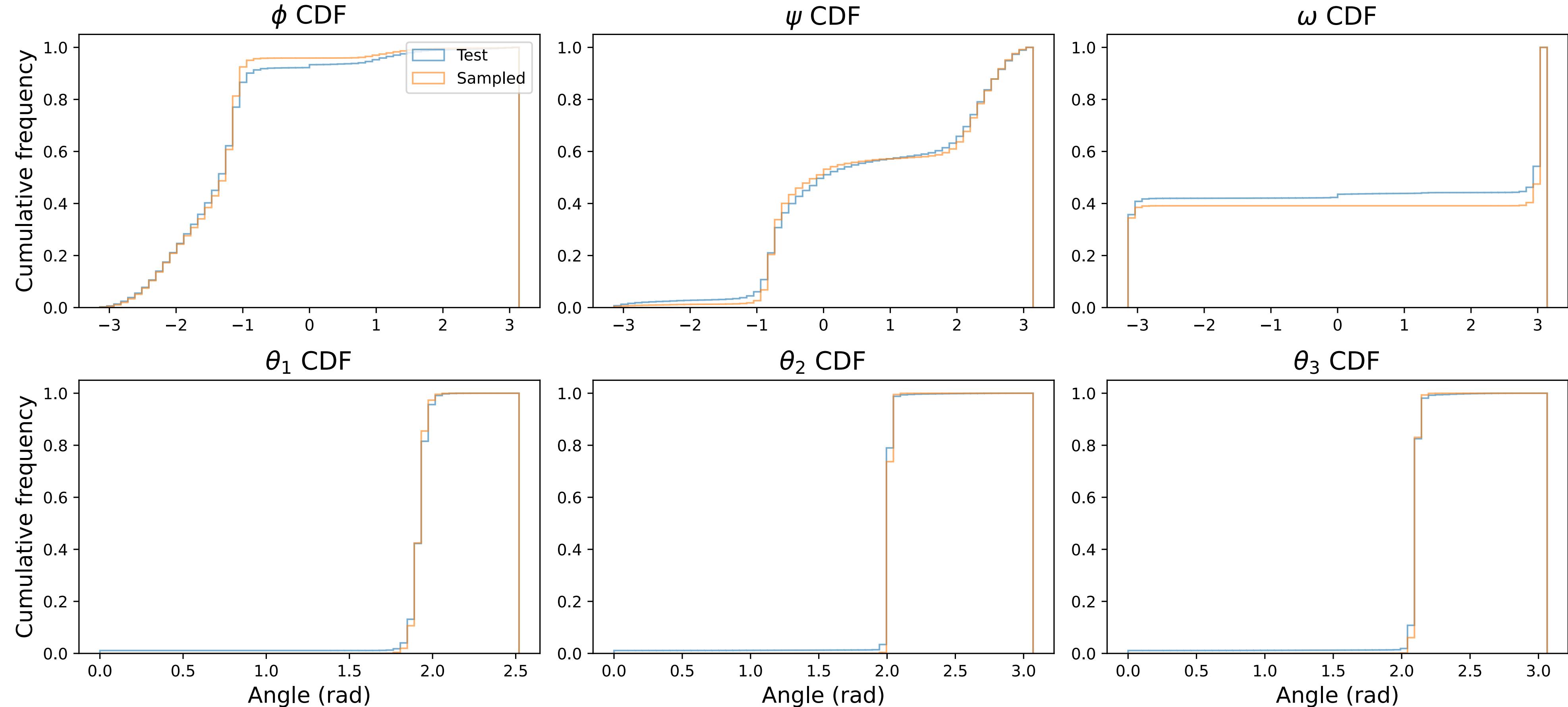
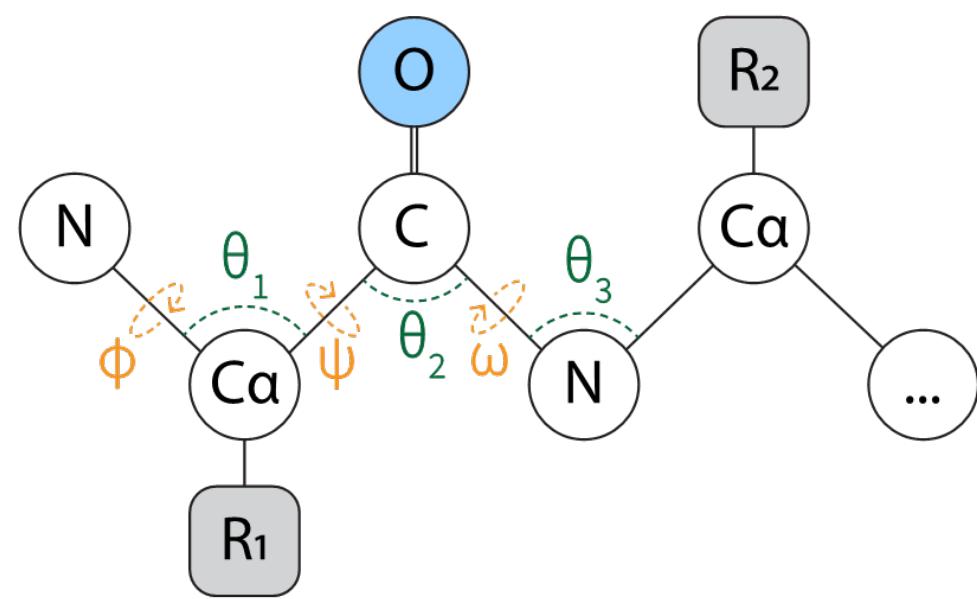
Folding Diffusion uses coordinates inspired by protein folding



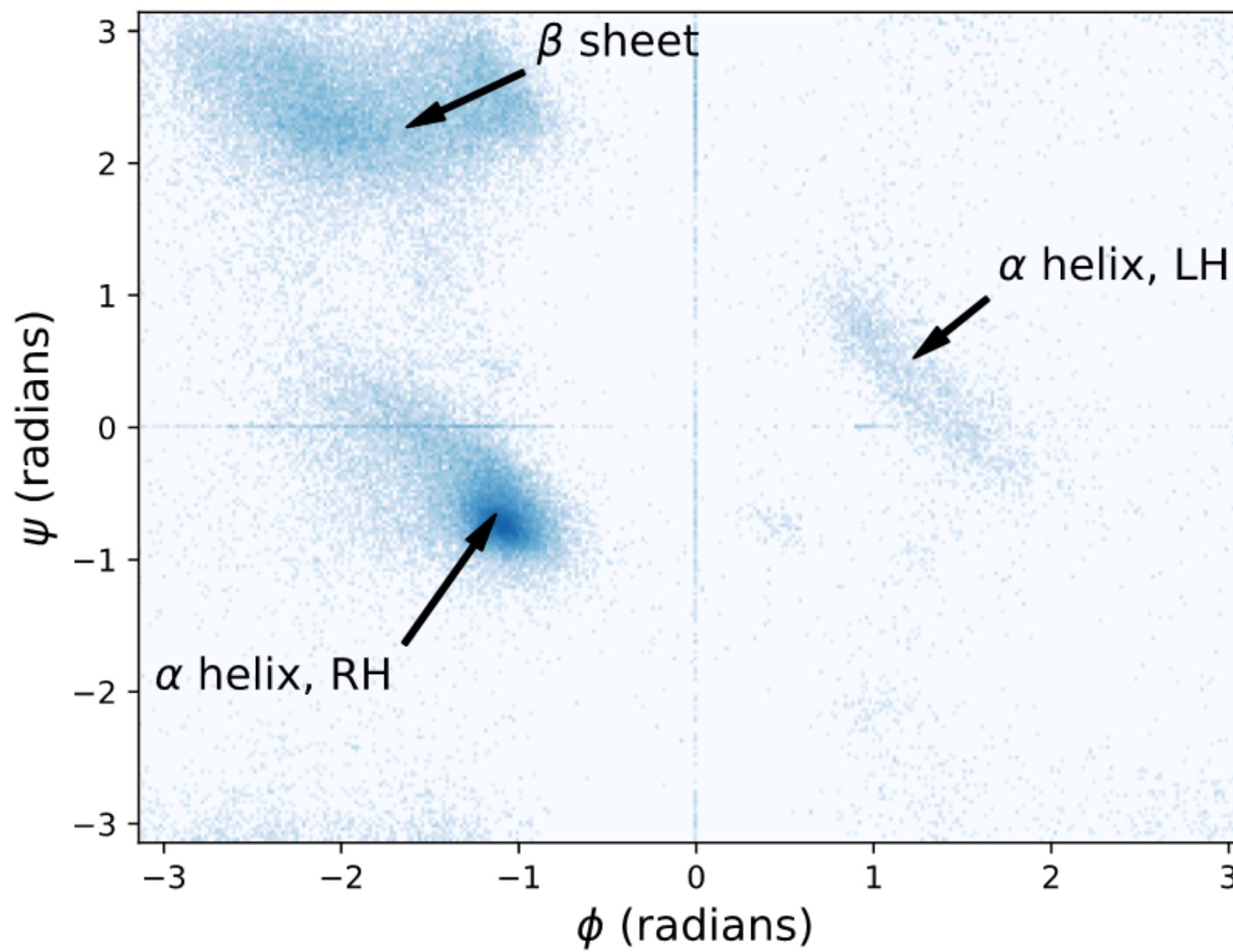
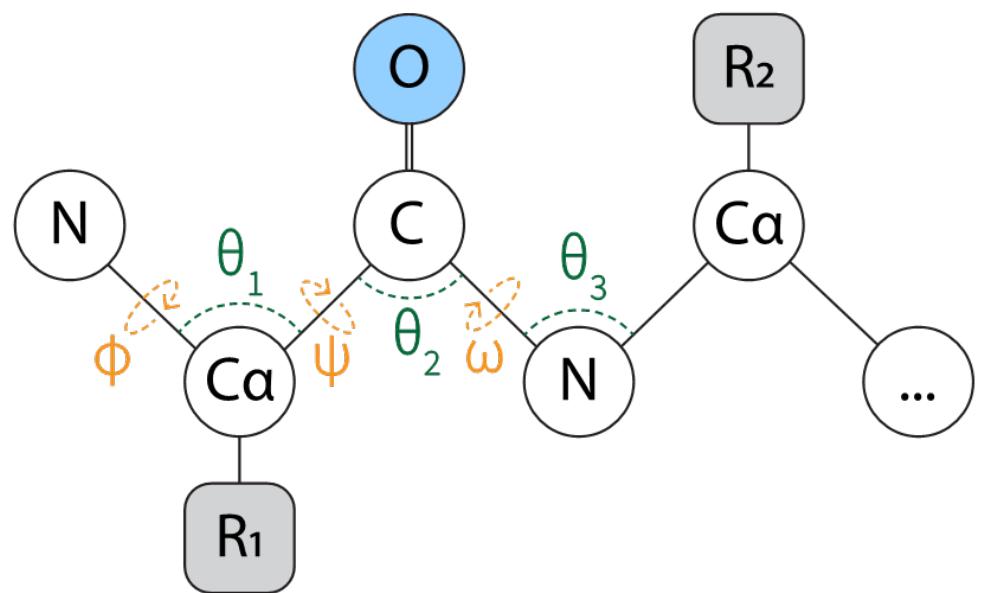
Folding Diffusion uses coordinates inspired by protein folding



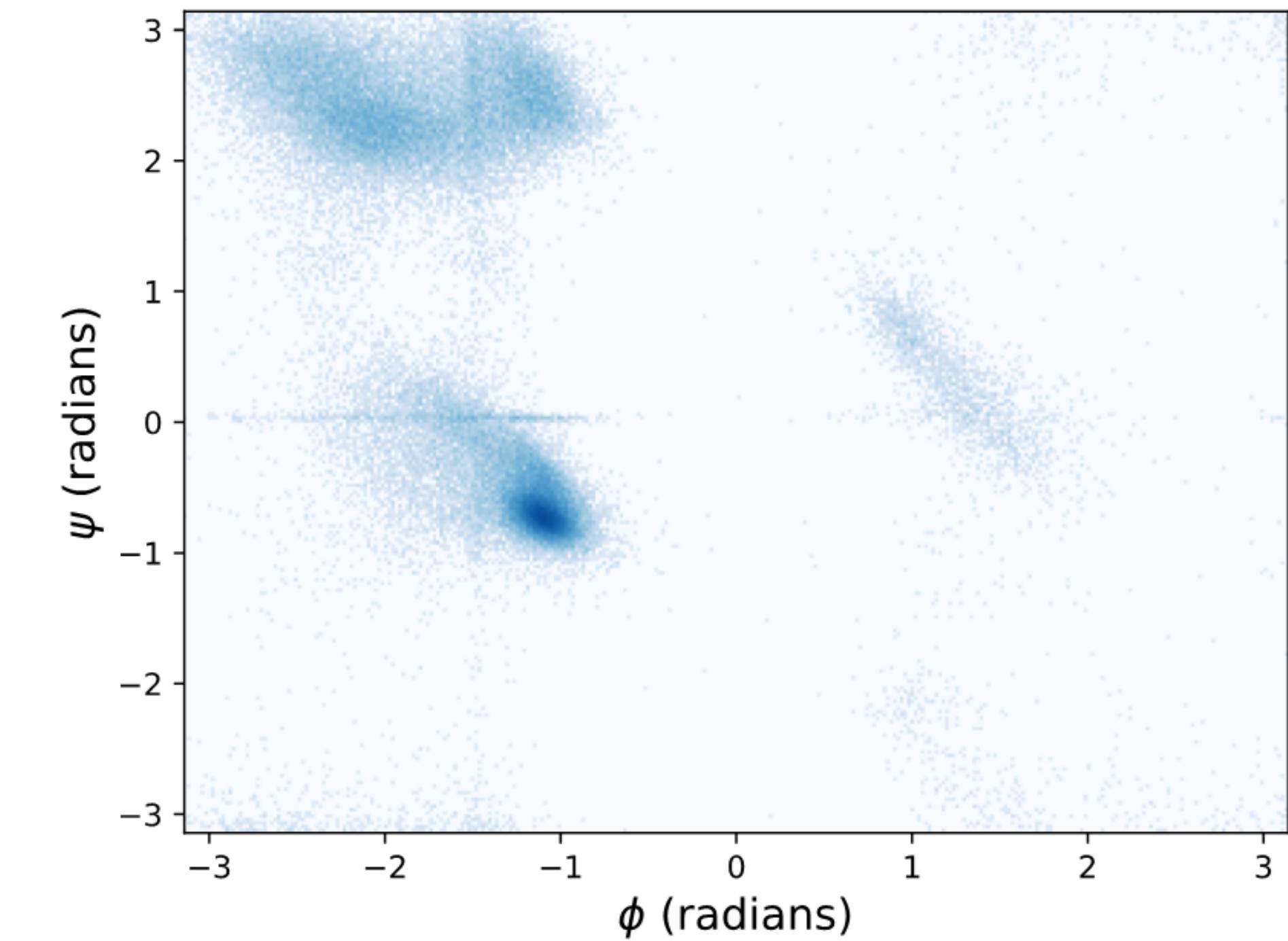
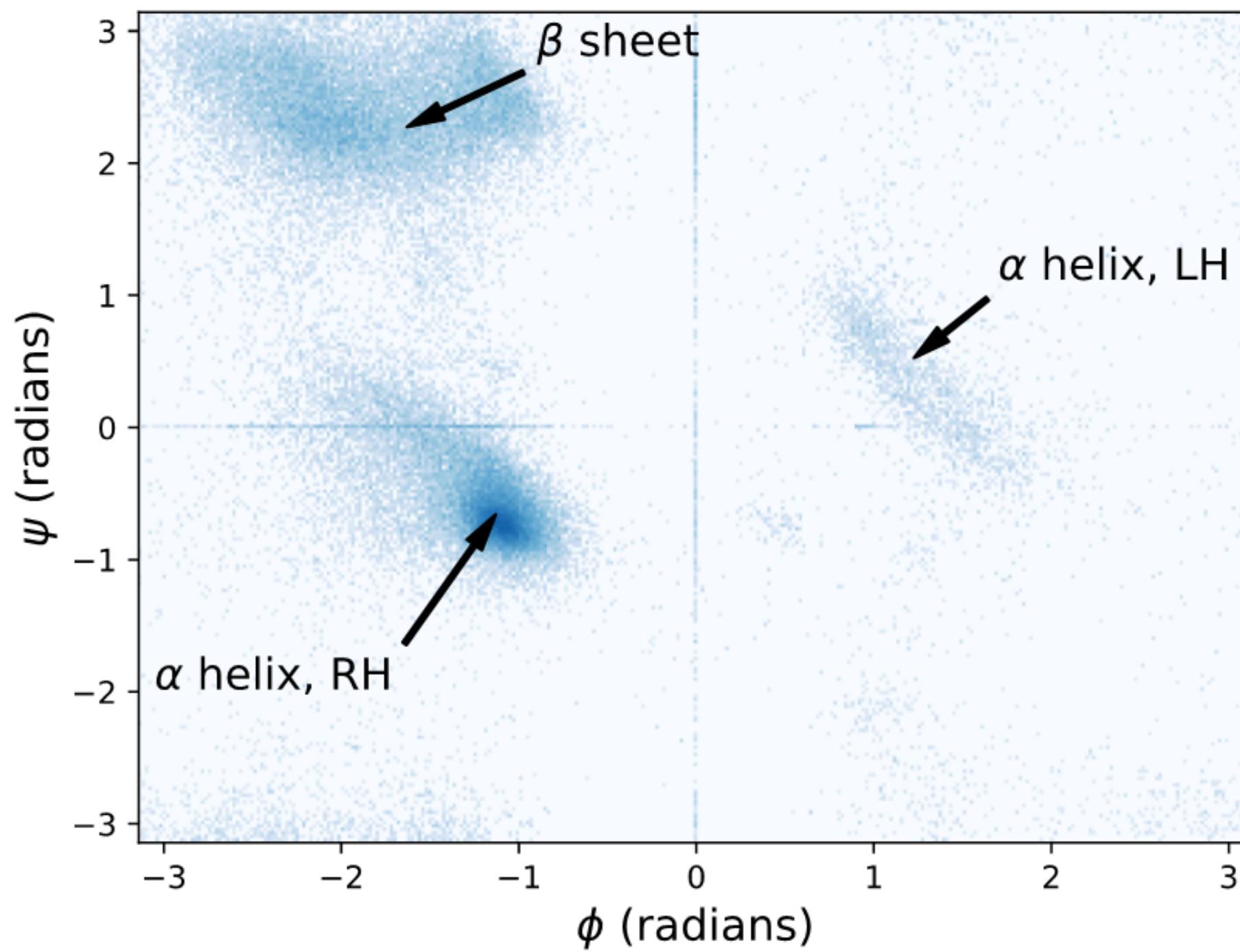
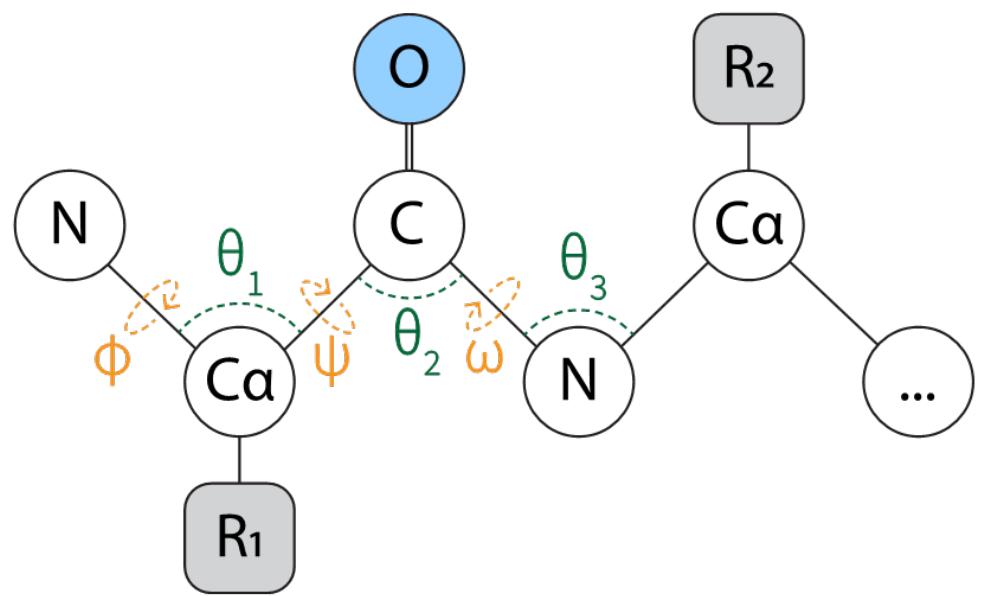
Generated angles match test distribution



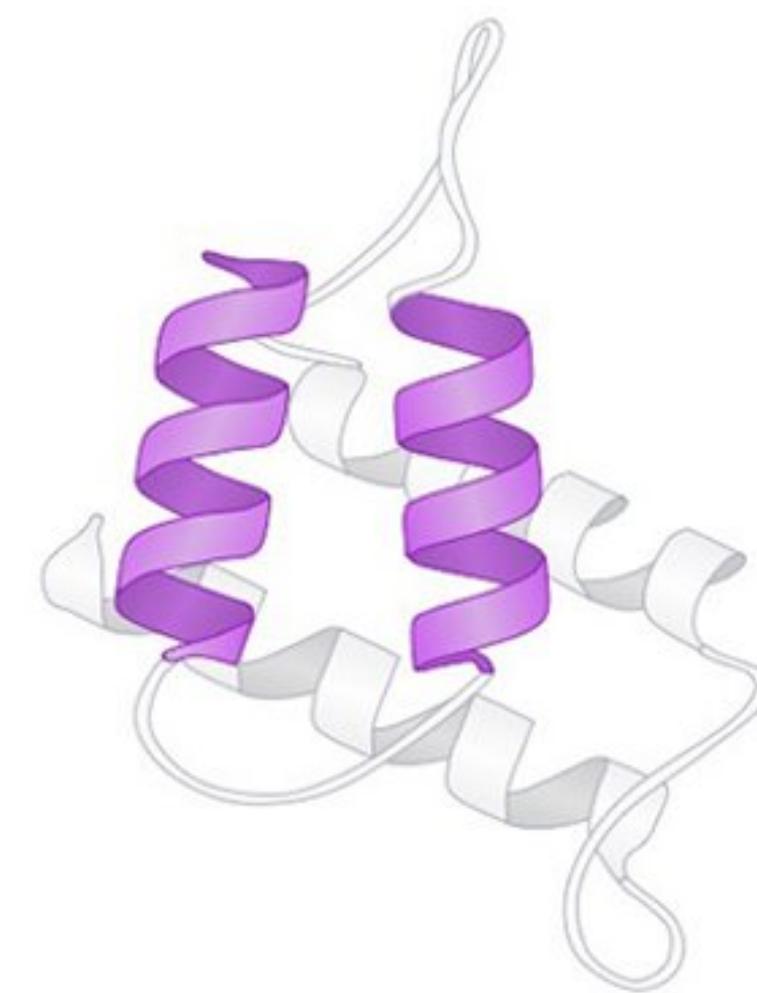
Generated Ramachandran plots match test structures



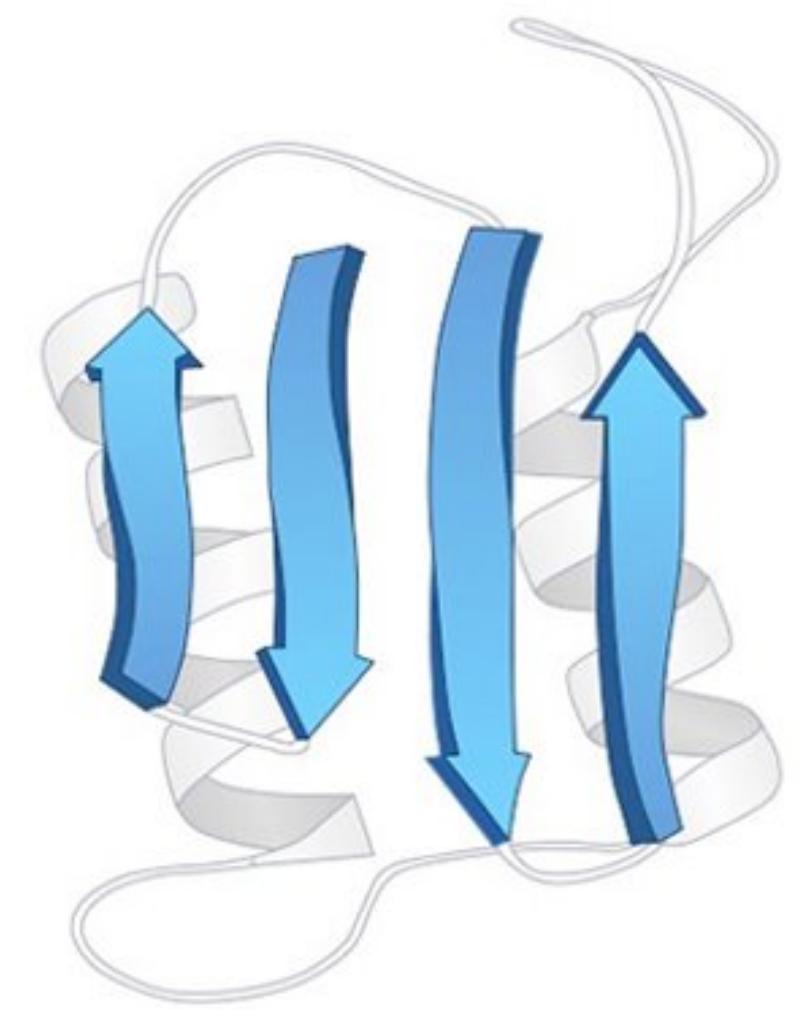
Generated Ramachandran plots match test structures



Generated secondary structures match test structures

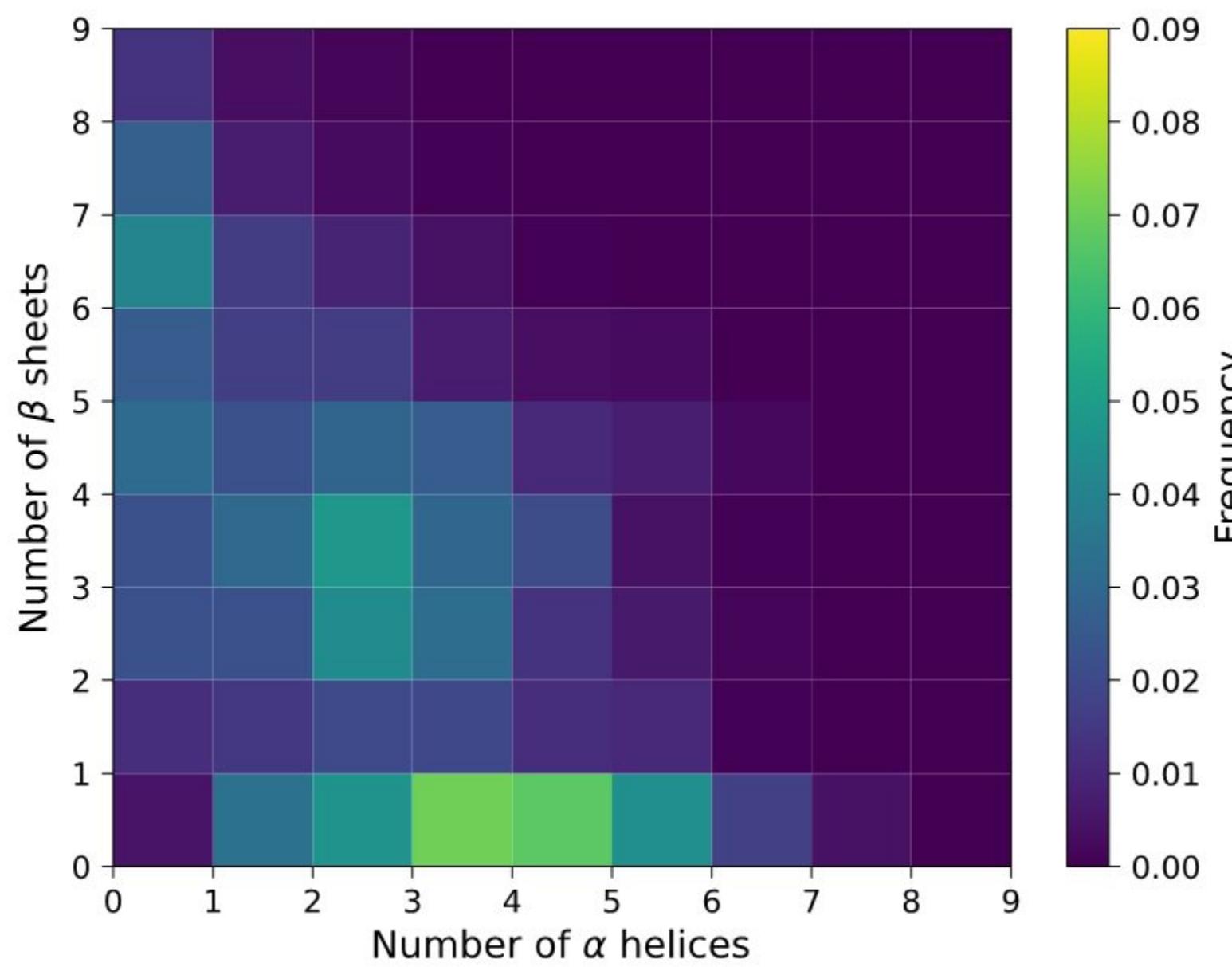
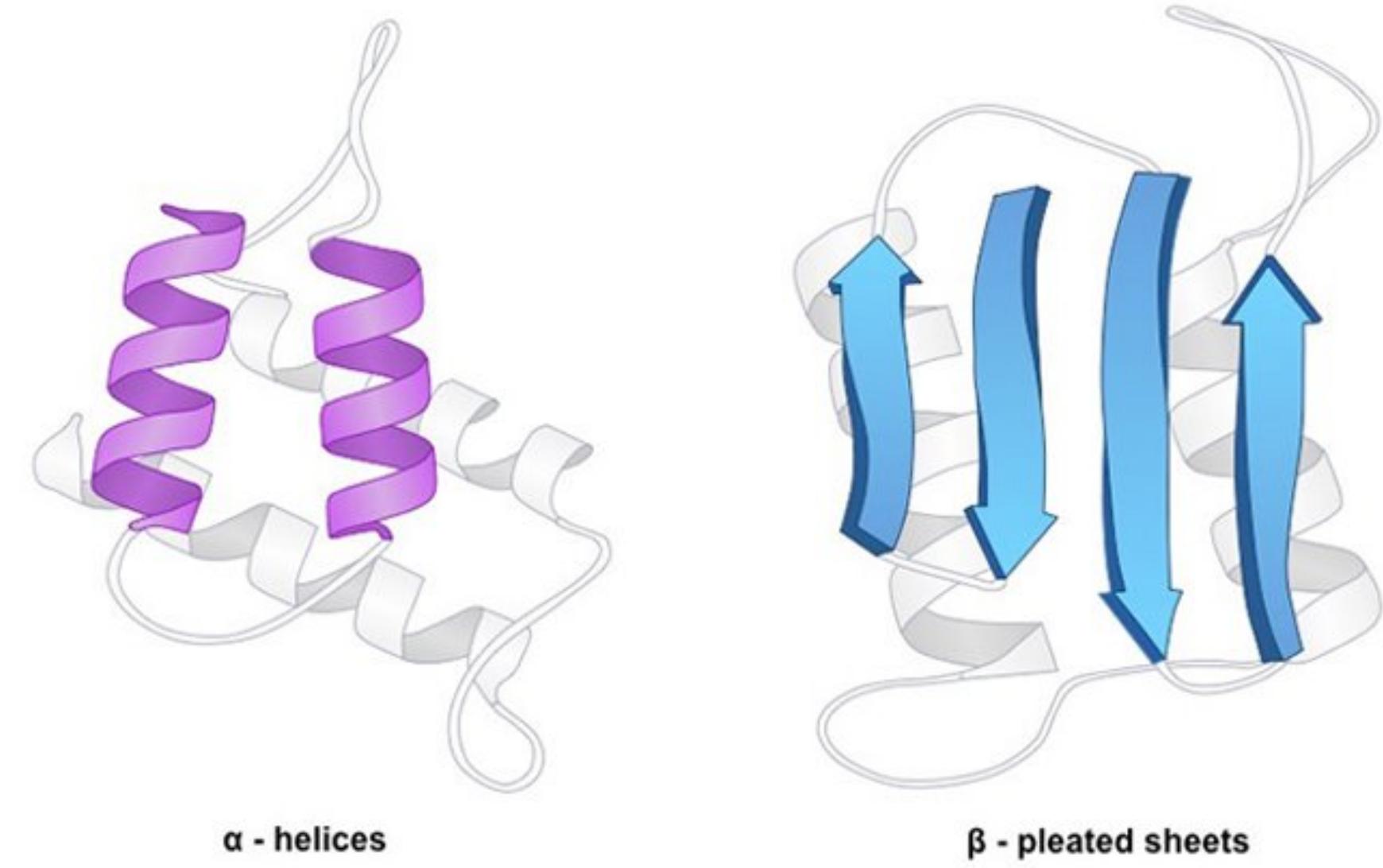


α - helices

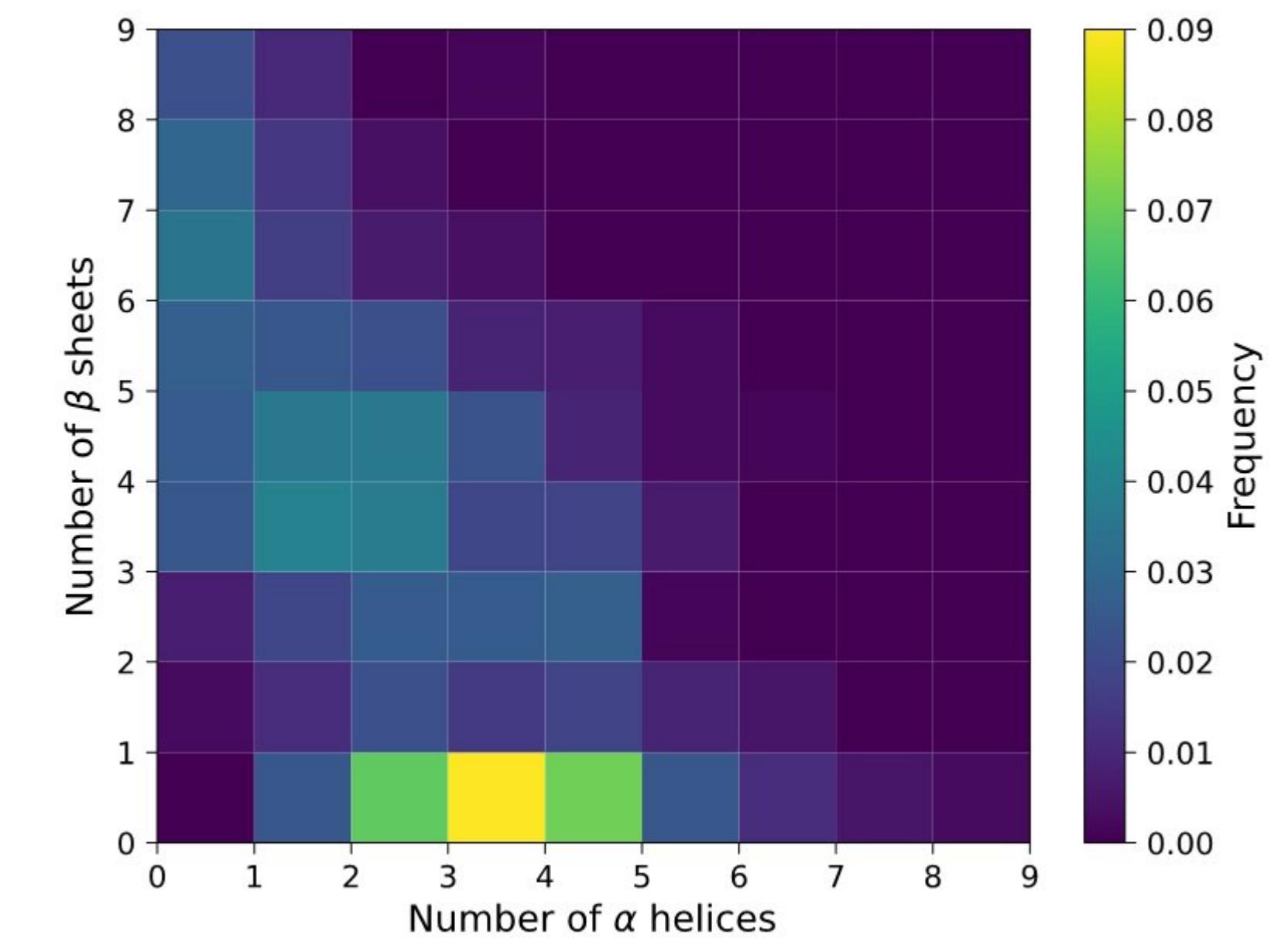
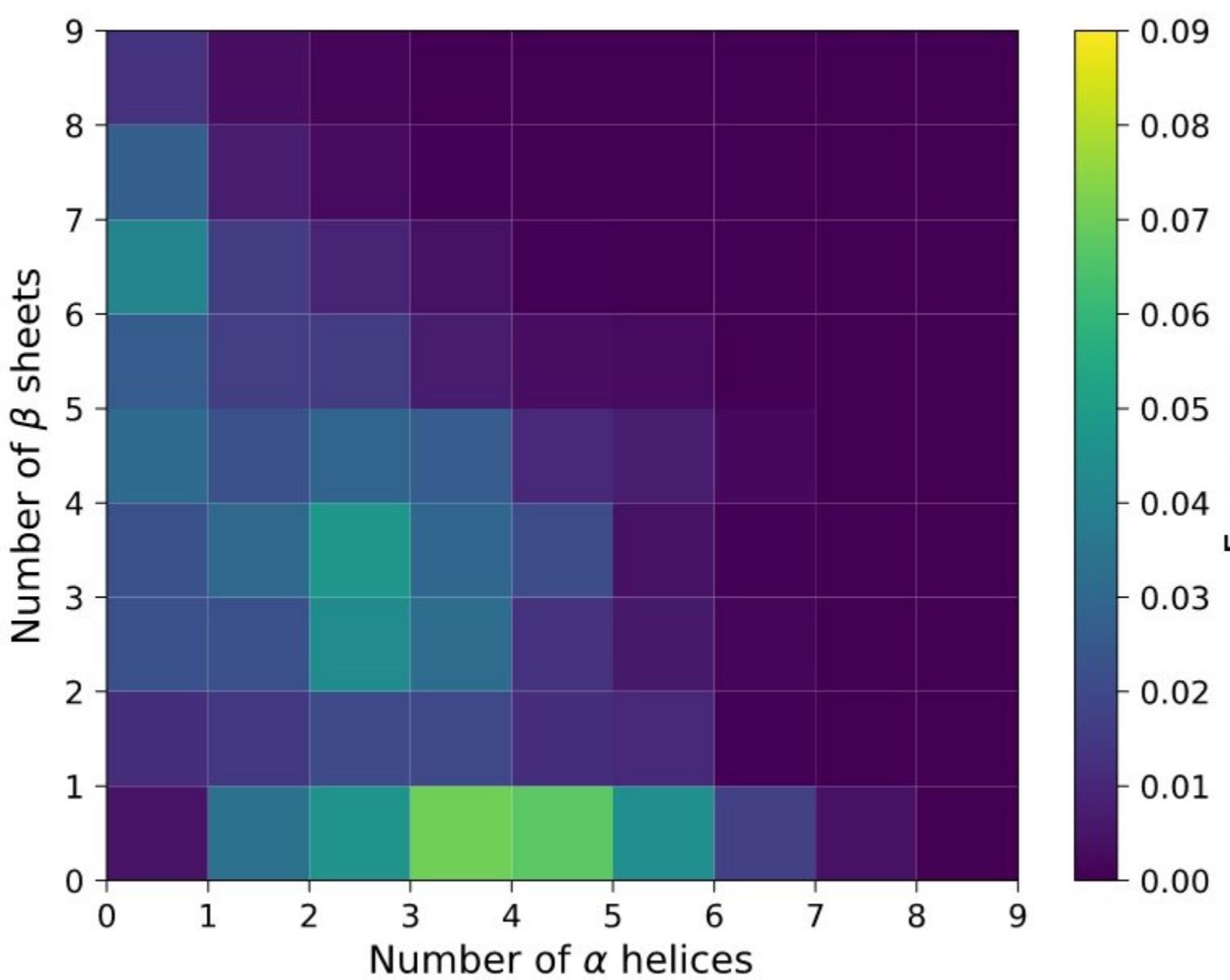
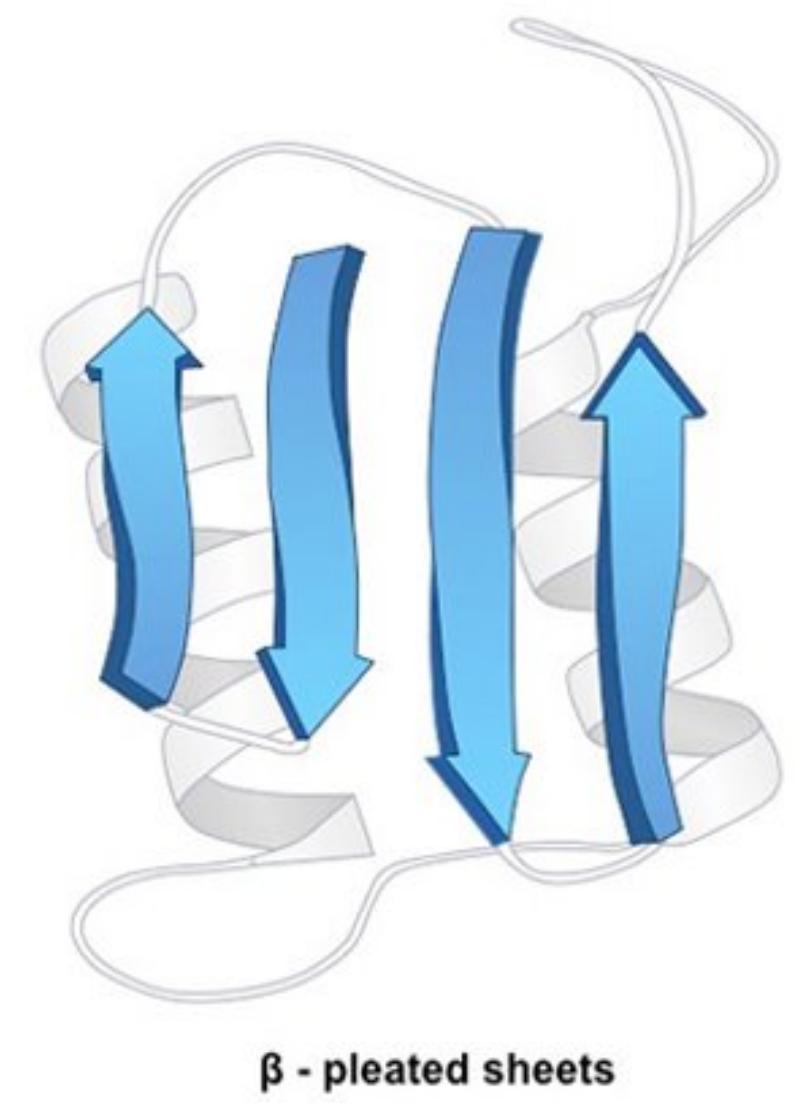
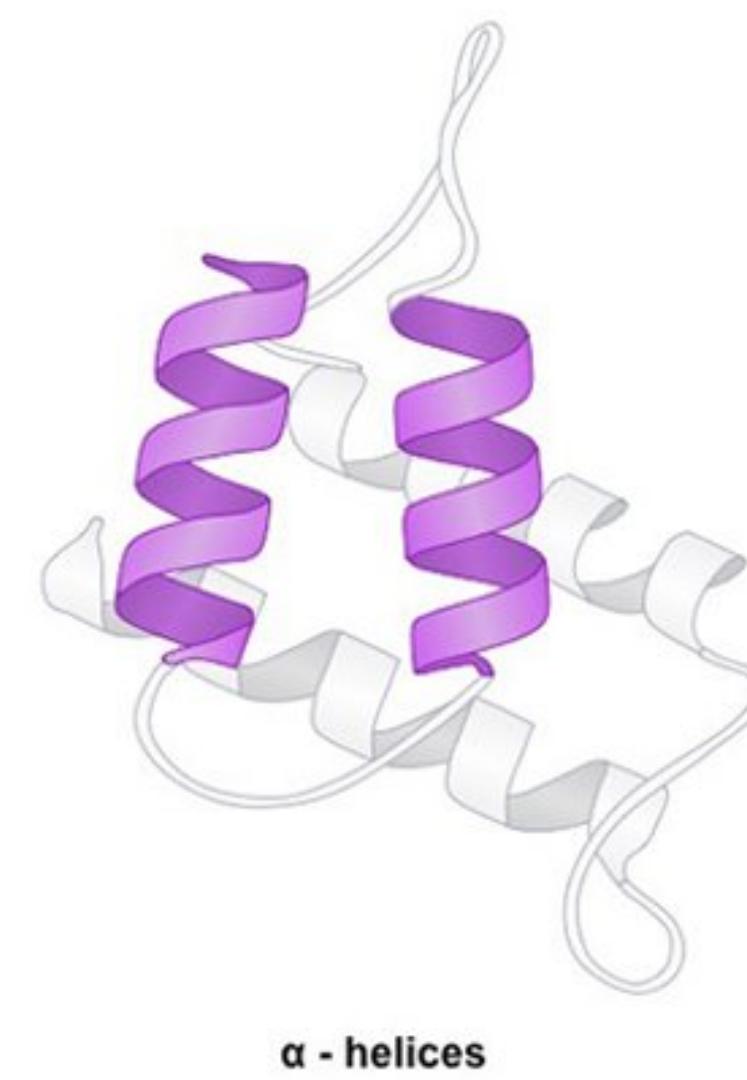


β - pleated sheets

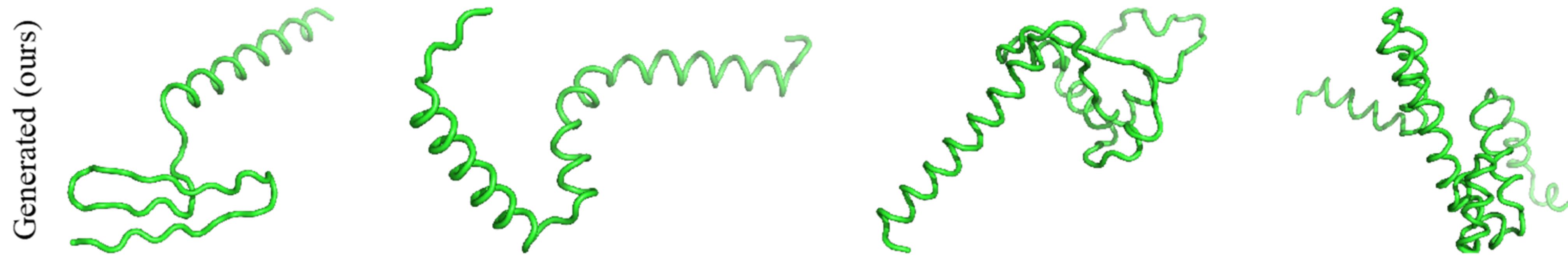
Generated secondary structures match test structures



Generated secondary structures match test structures

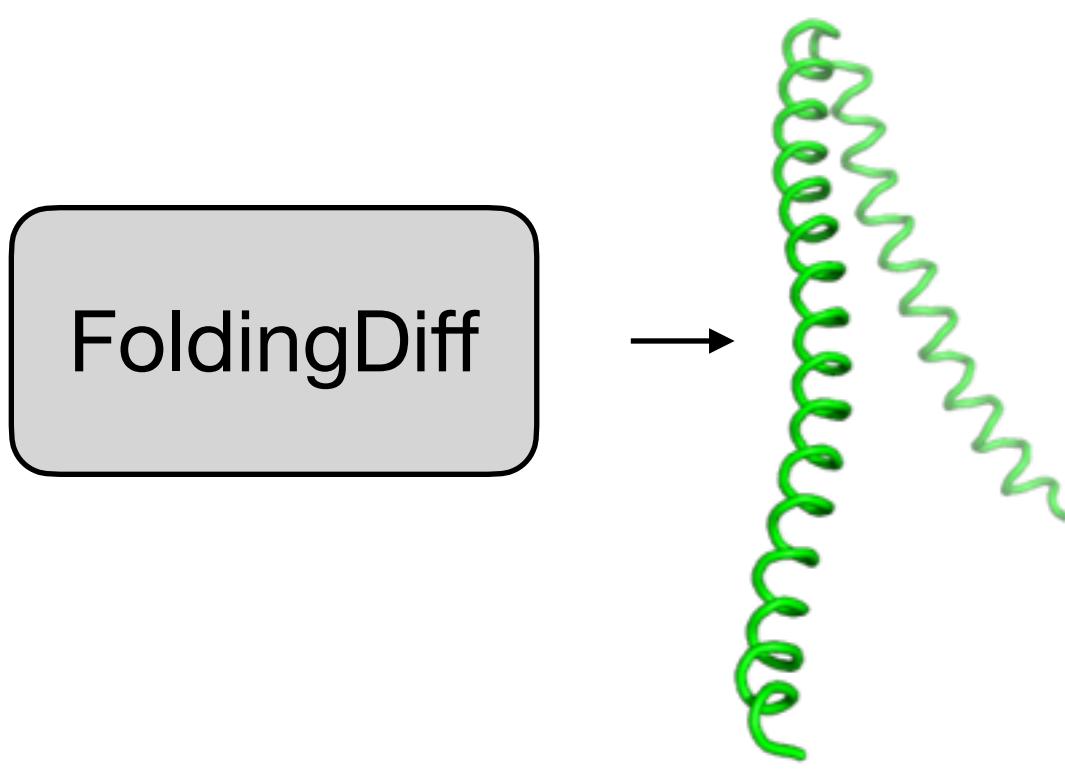


Generated structures look reasonable

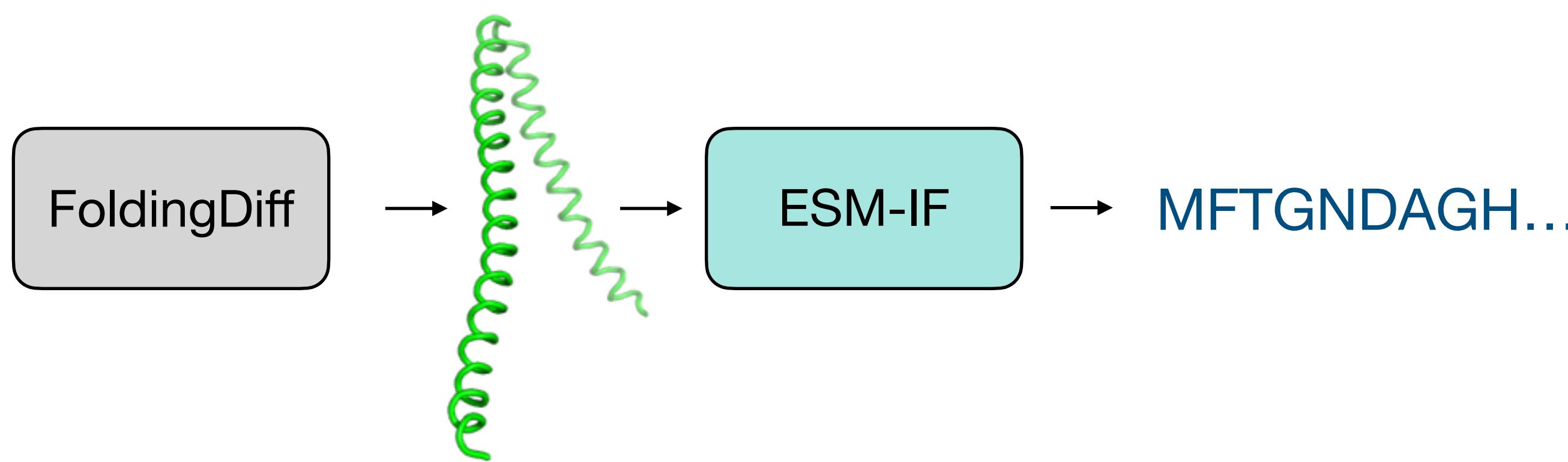


Measure design ability of structures with self-consistency TMscore

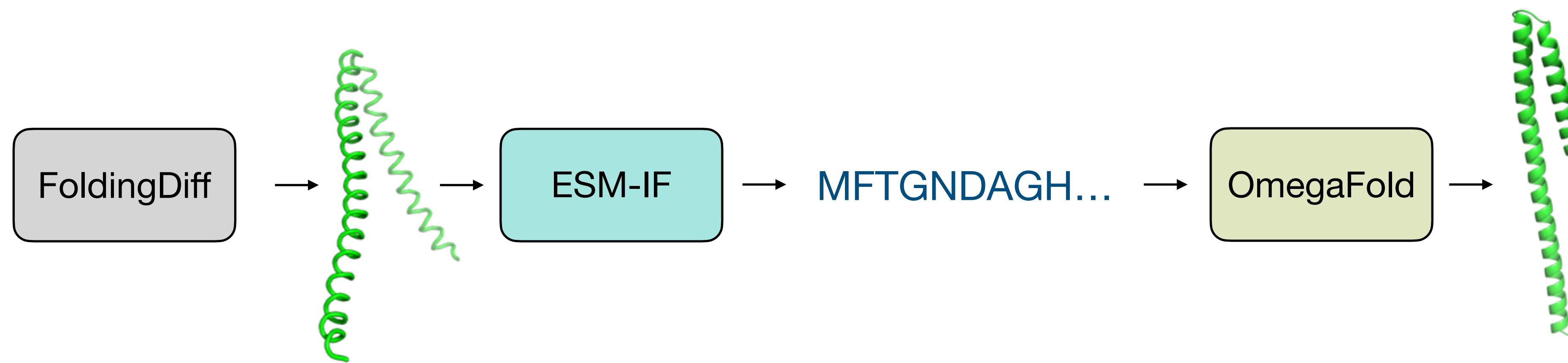
Measure design ability of structures with self-consistency TMscore



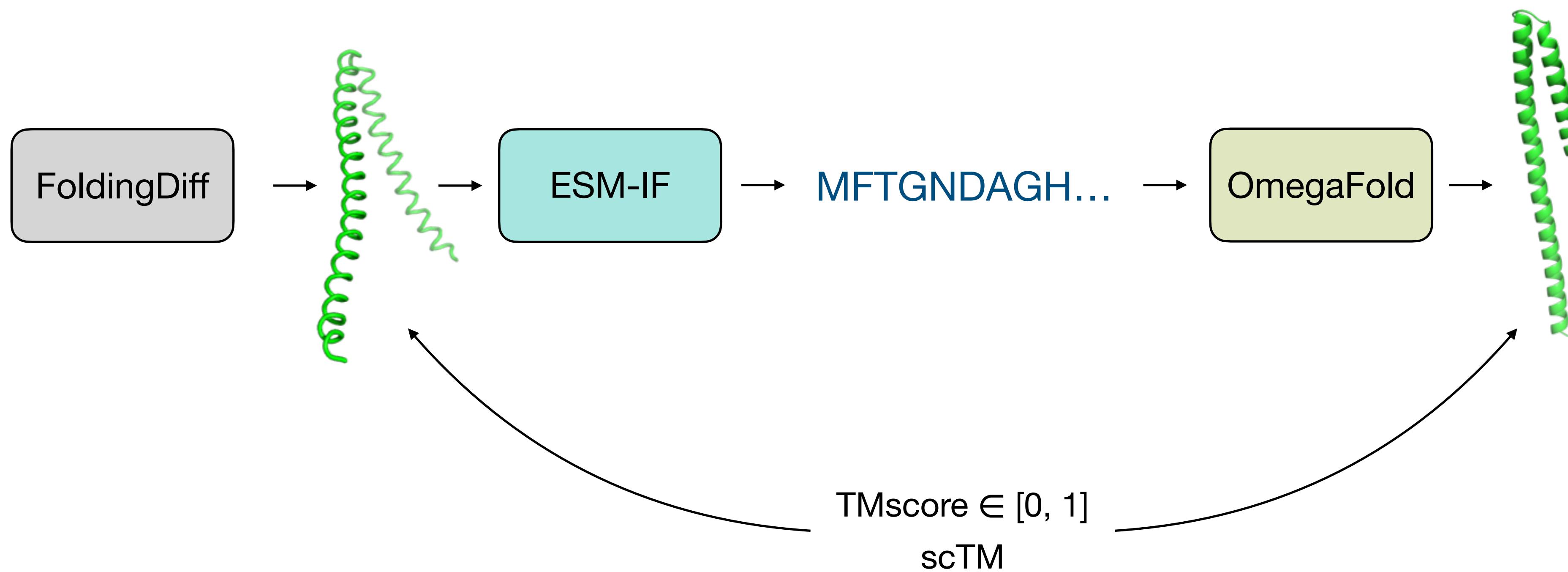
Measure design ability of structures with self-consistency TMscore



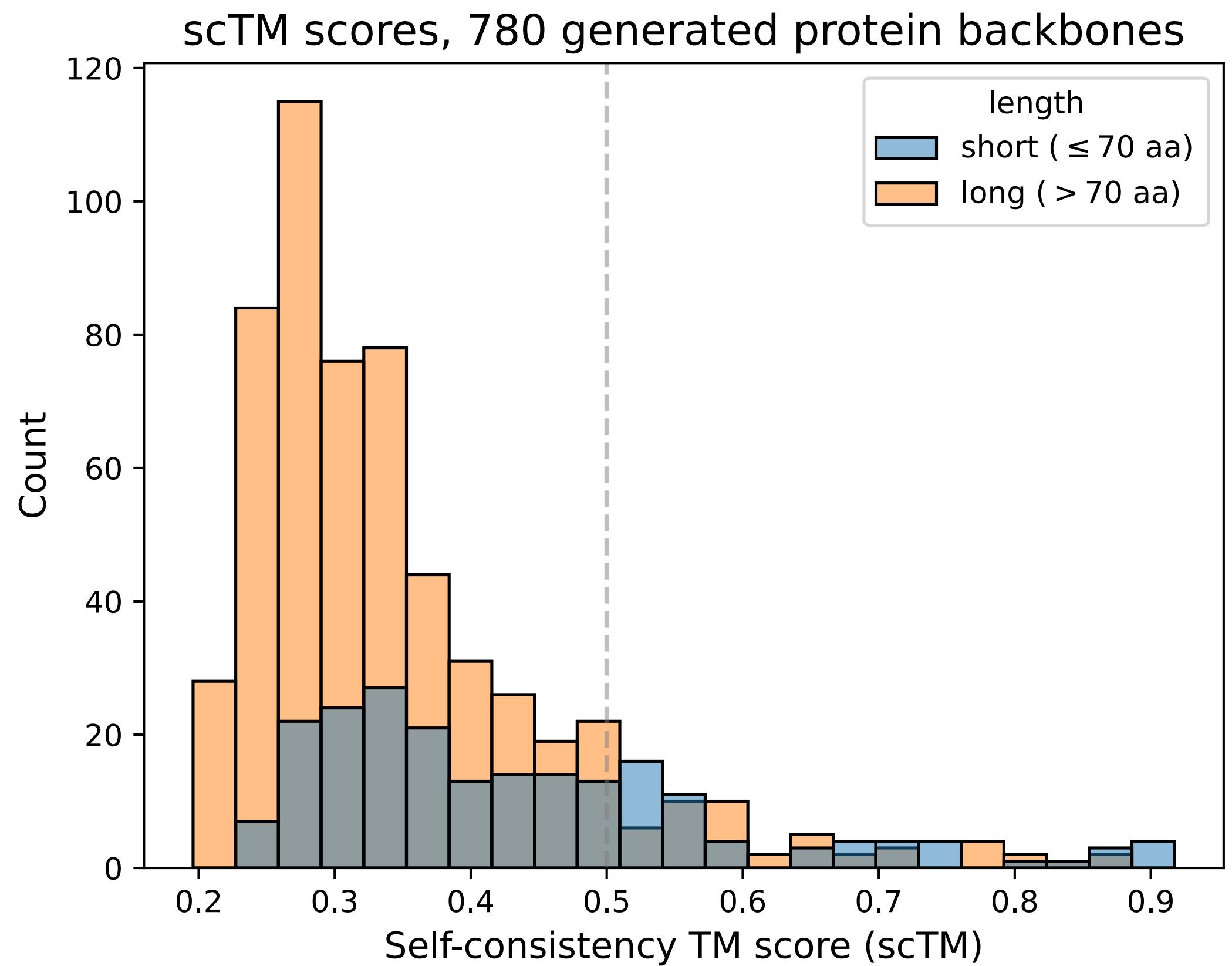
Measure design ability of structures with self-consistency TMscore



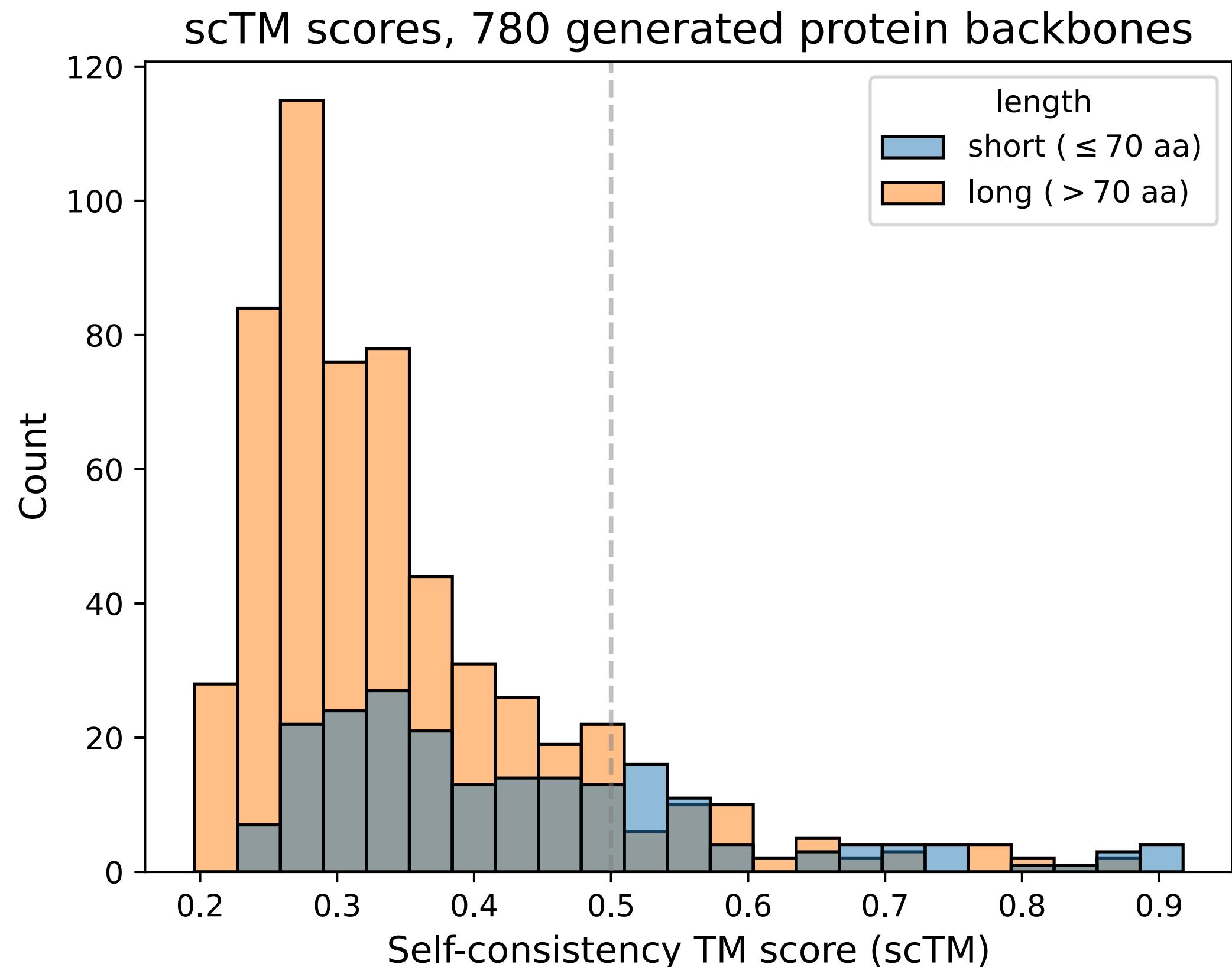
Measure design ability of structures with self-consistency TMscore



Many generated structures are designable

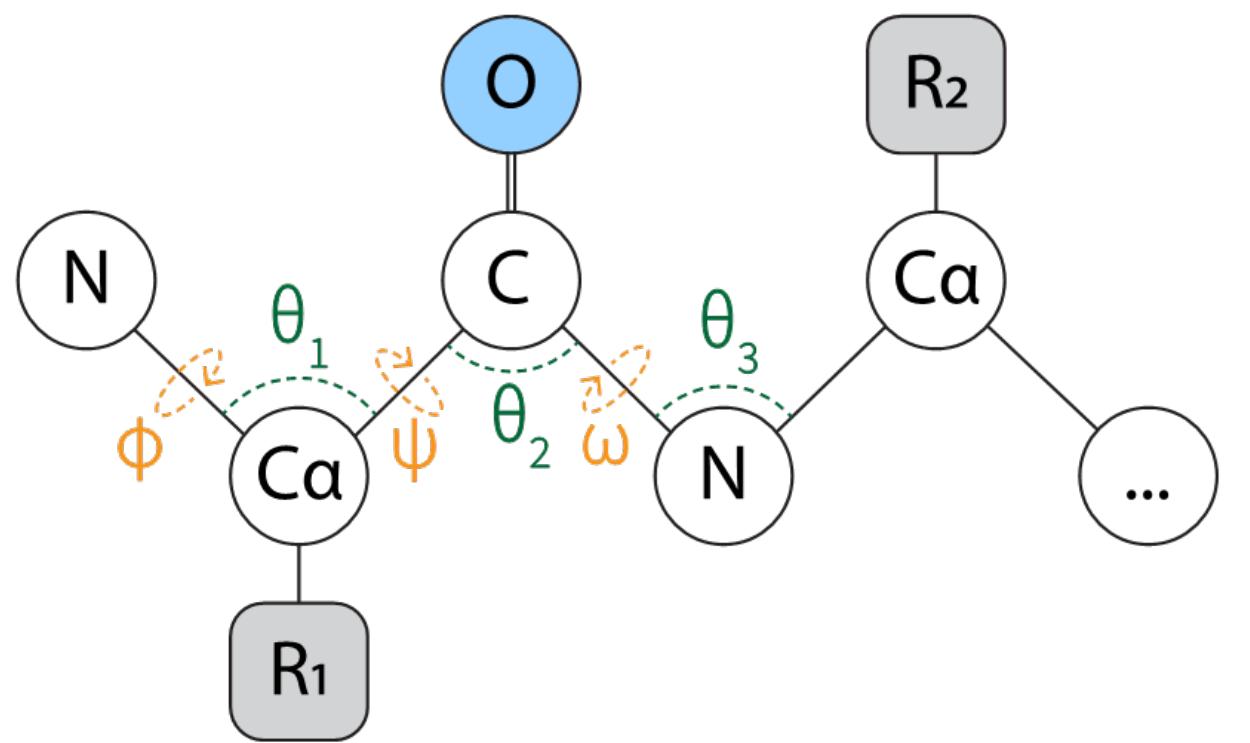


Many generated structures are designable



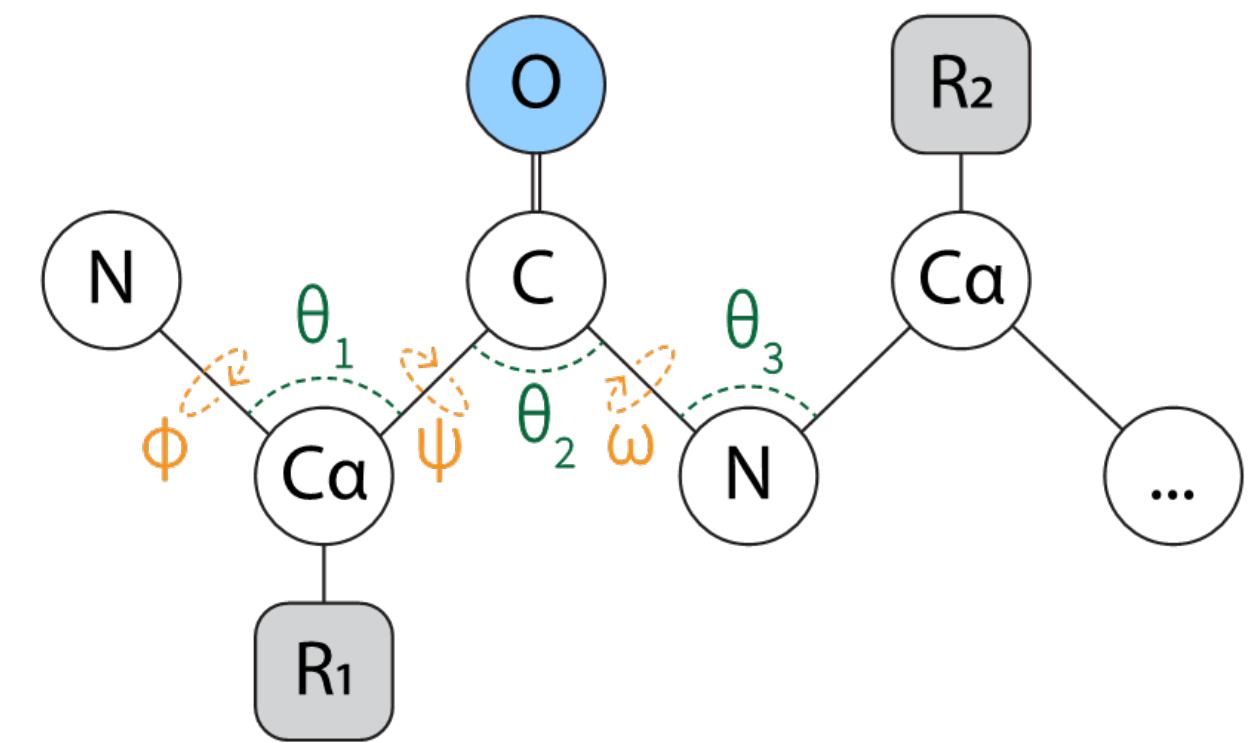
scTM > 0.5	≤ 70 aa	> 70 aa
FoldingDiff	57/210	54/570
ProtDiff (Trippe et al.)	36/210	51/570

FoldingDiff structures are better than random baseline

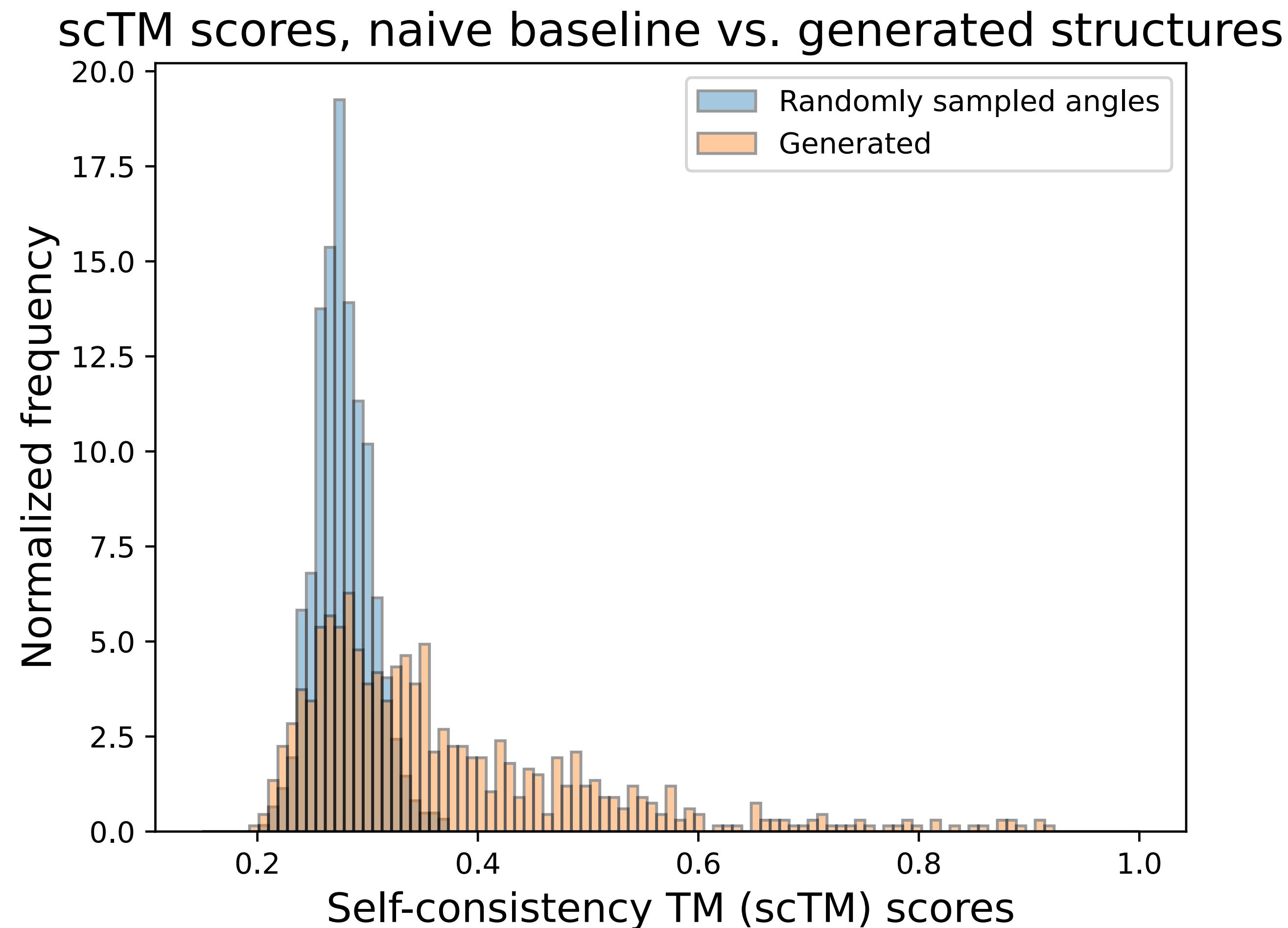


Sample sets of angles
Preserves Ramachandran plot

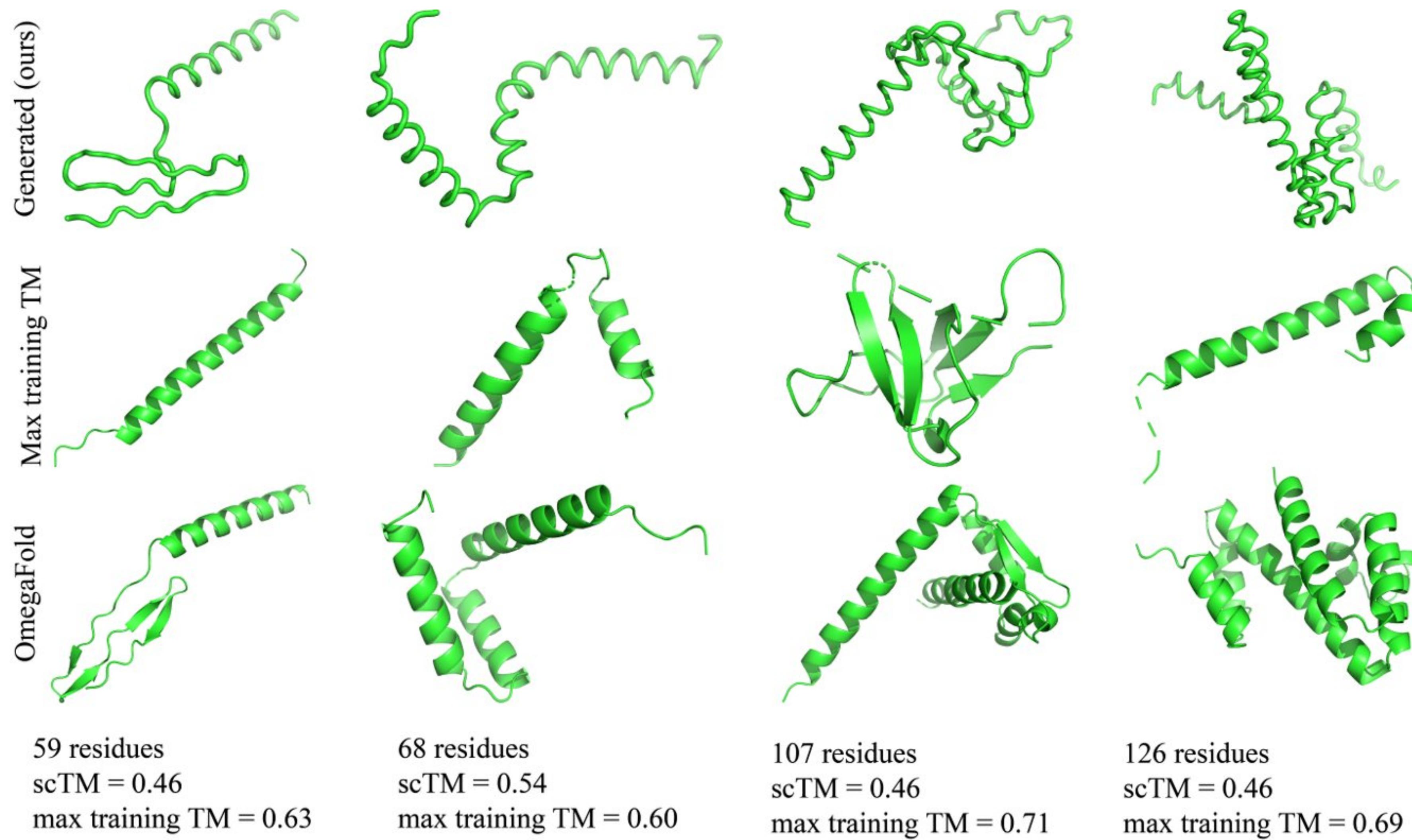
FoldingDiff structures are better than random baseline



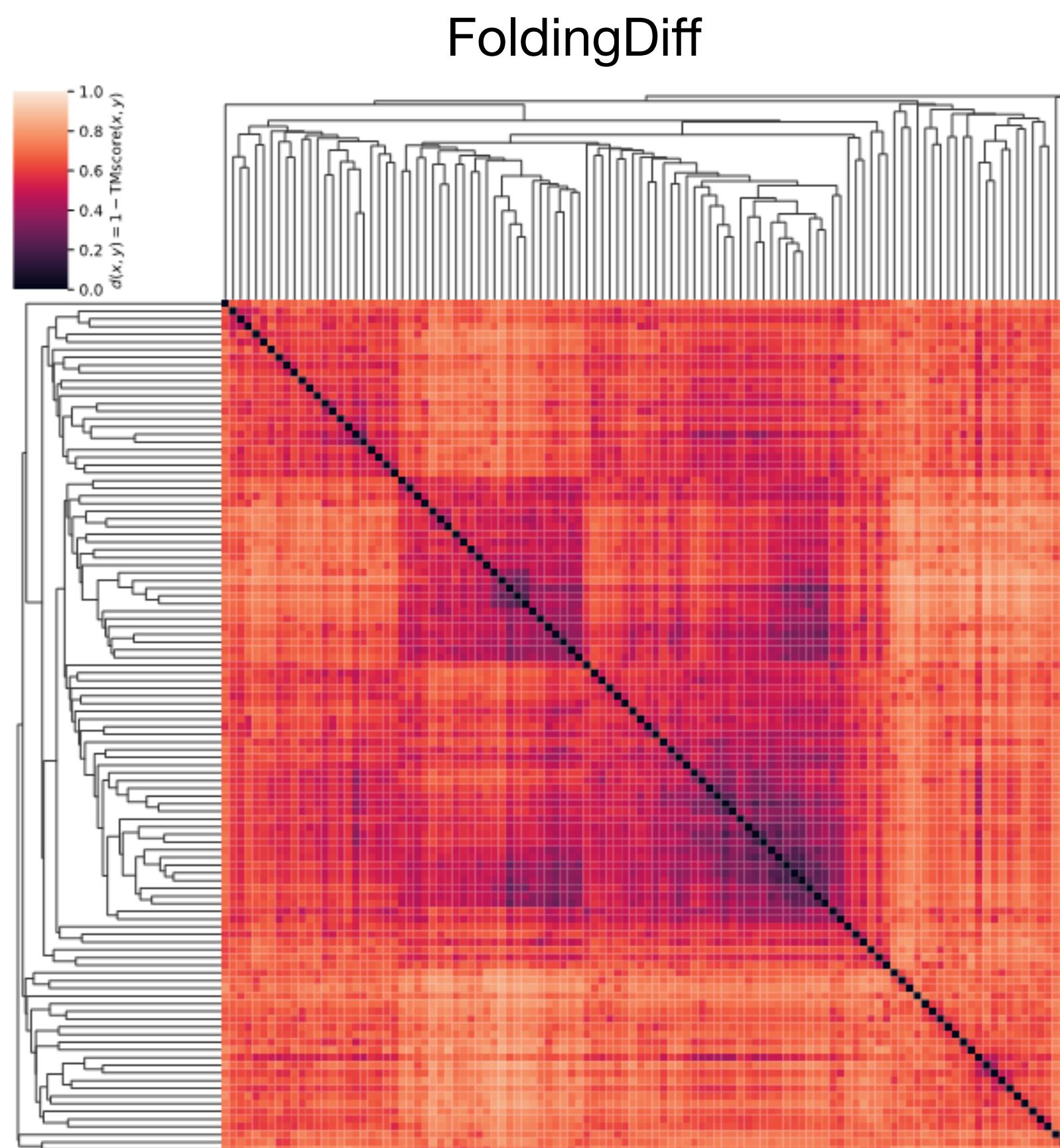
Sample sets of angles
Preserves Ramachandran plot



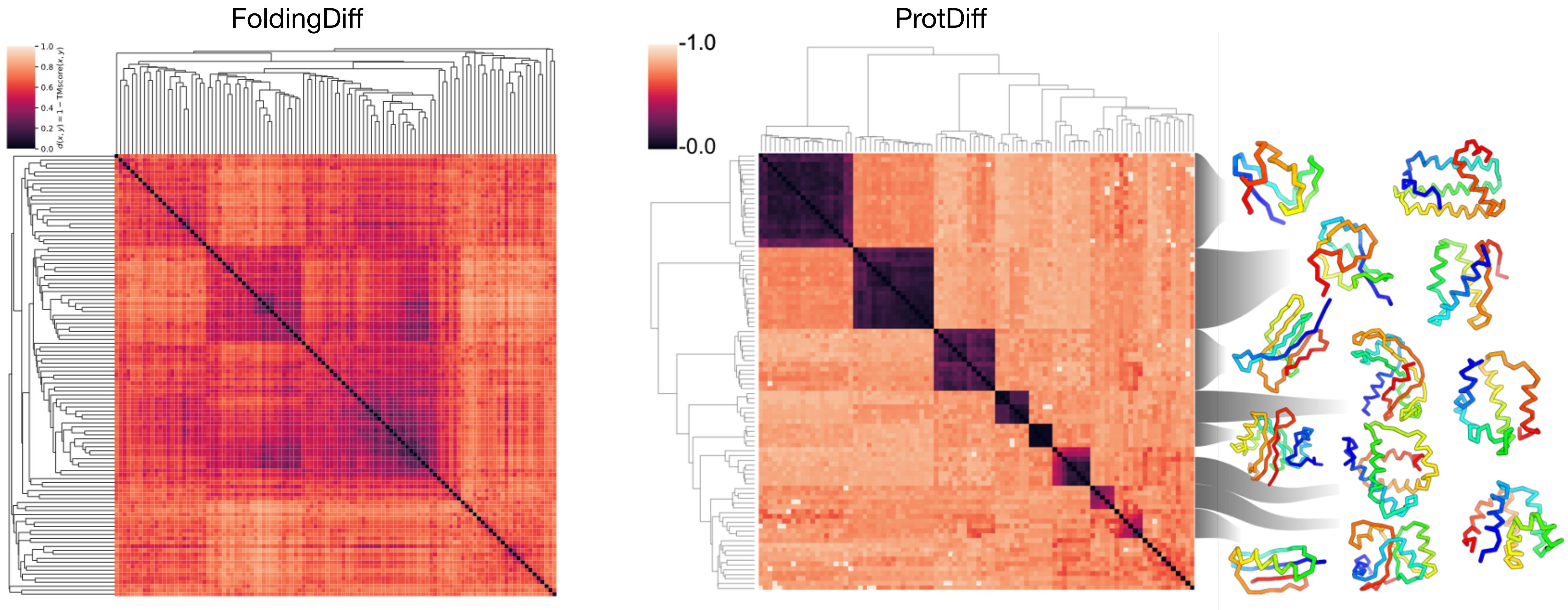
Generated structures are diverse



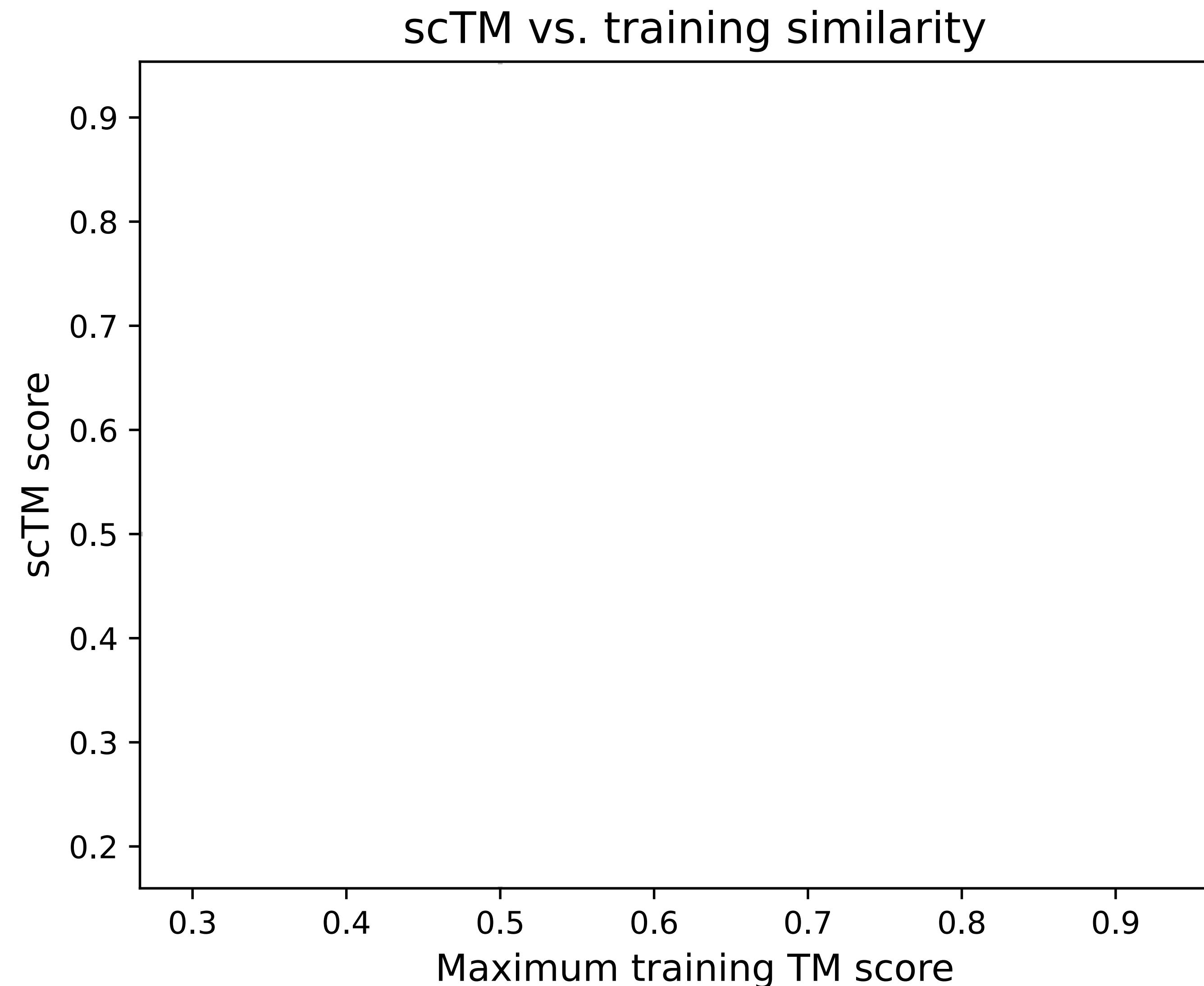
Generated structures are diverse



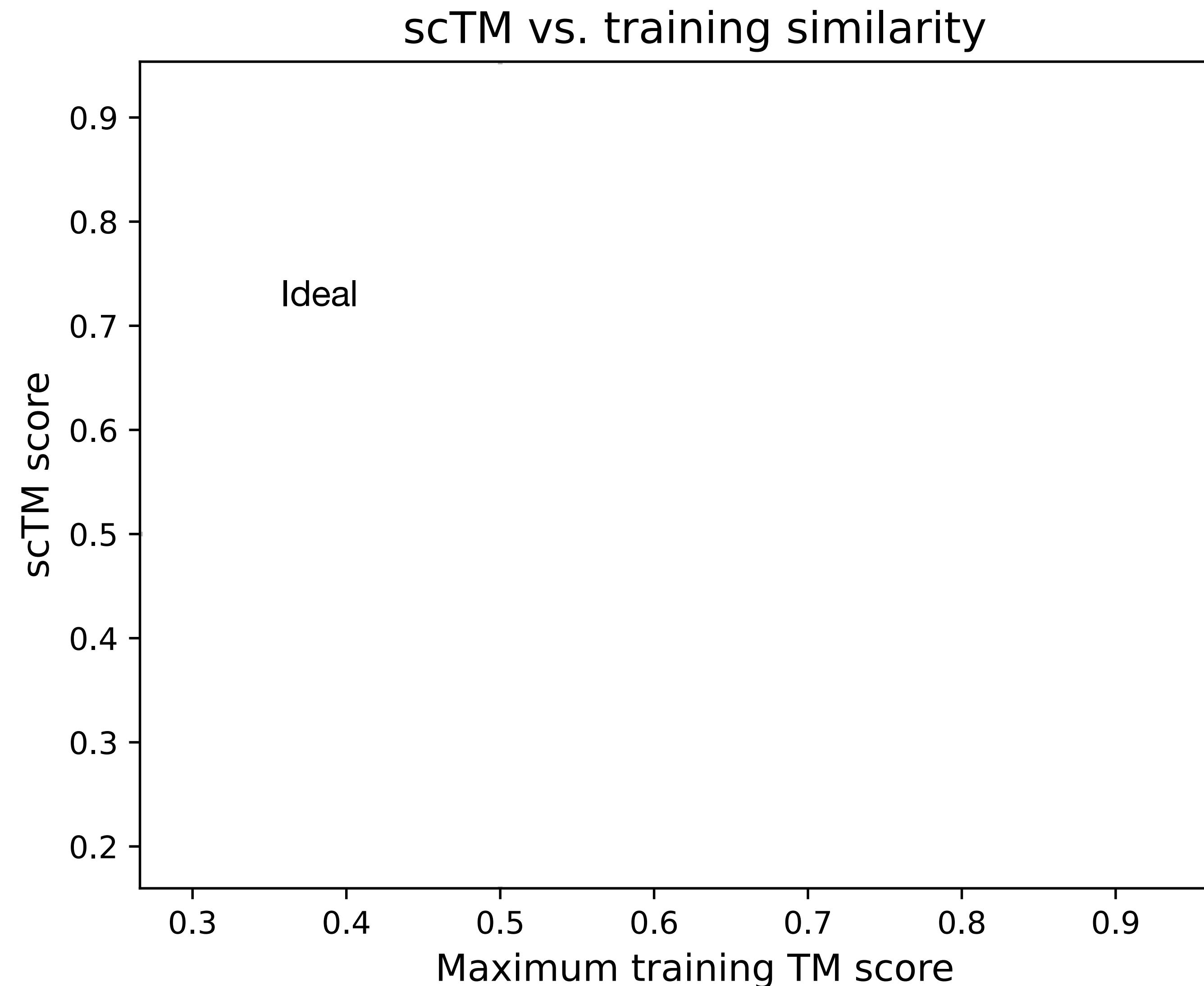
Generated structures are diverse



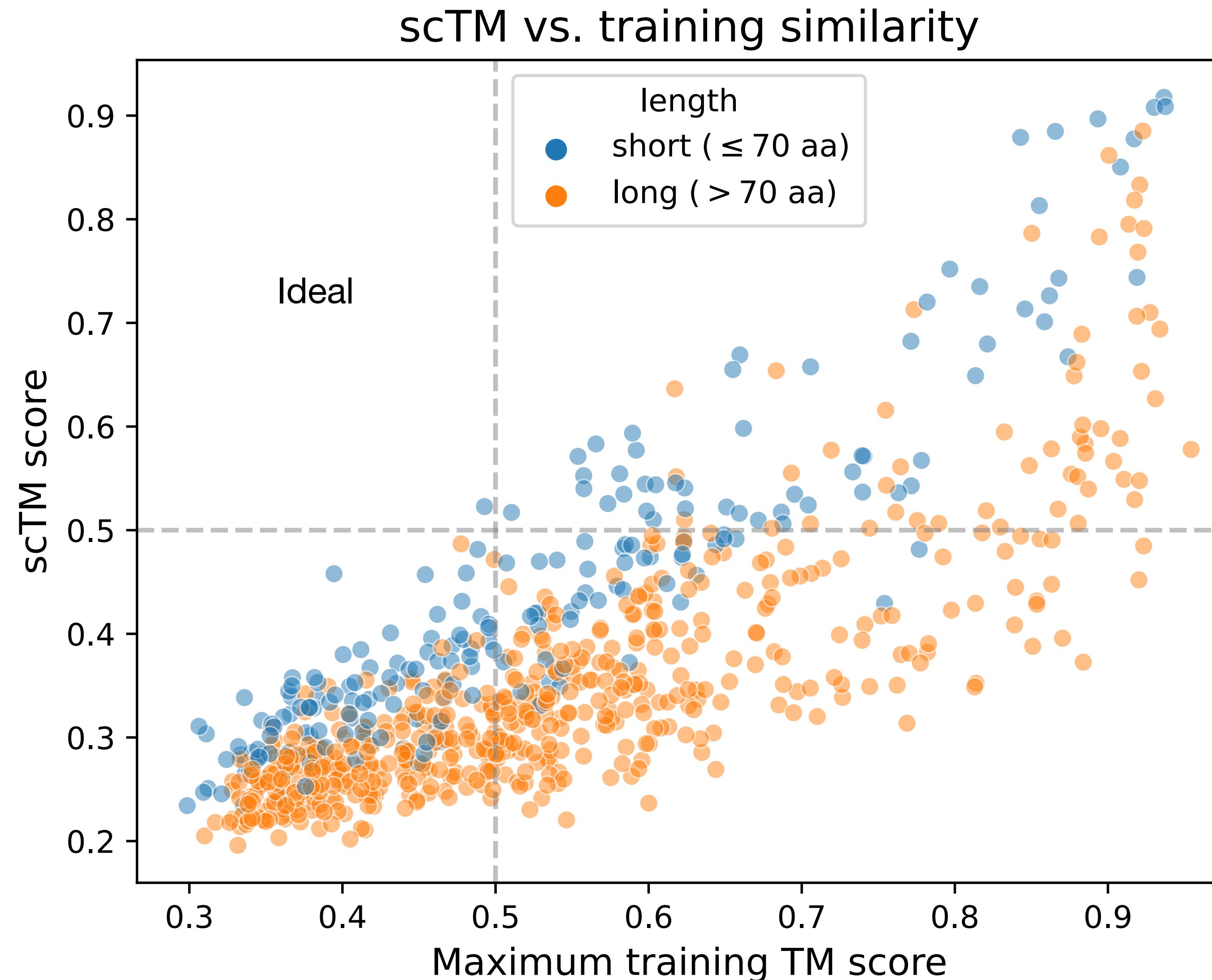
Generated structures are diverse



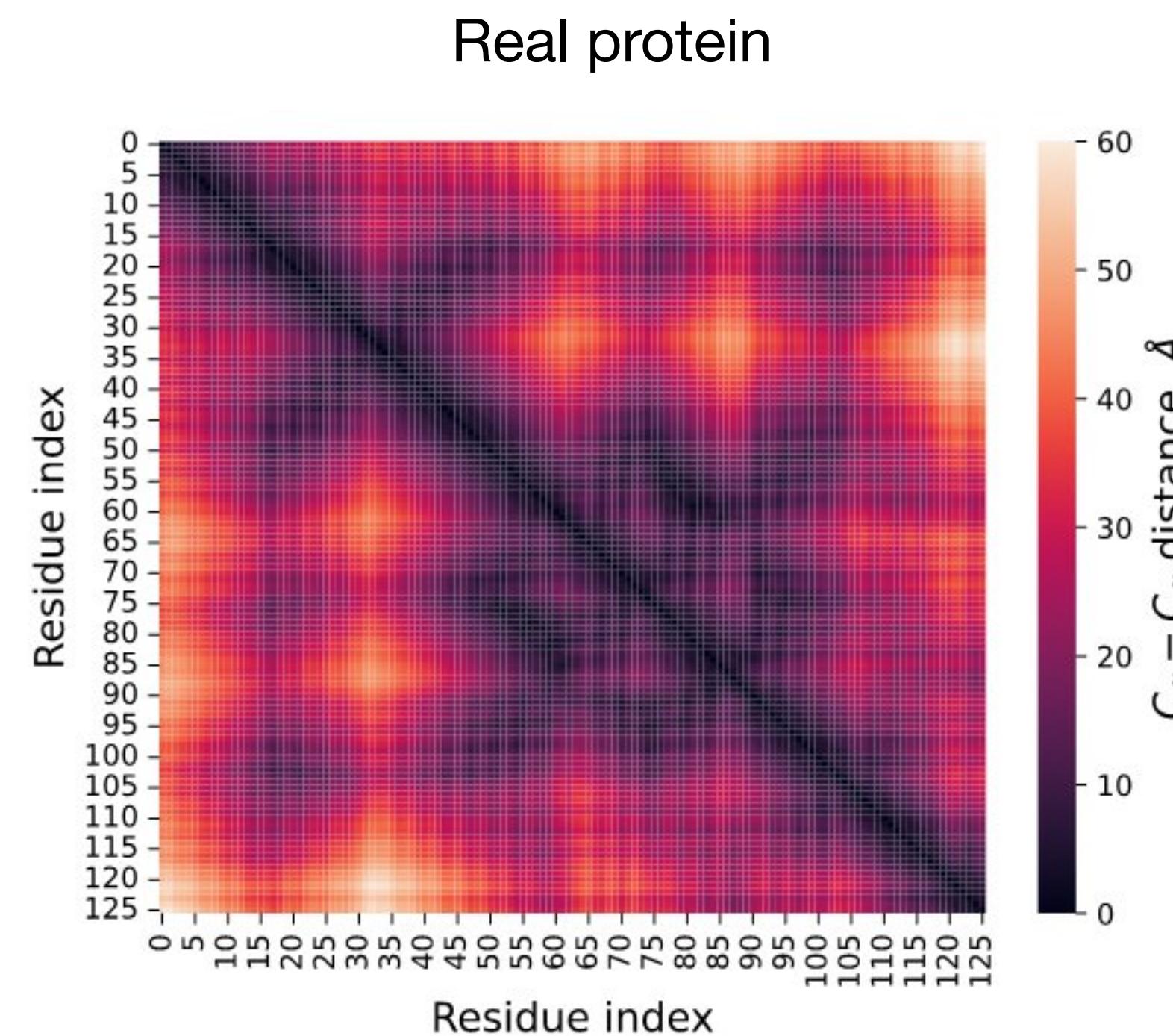
Generated structures are diverse



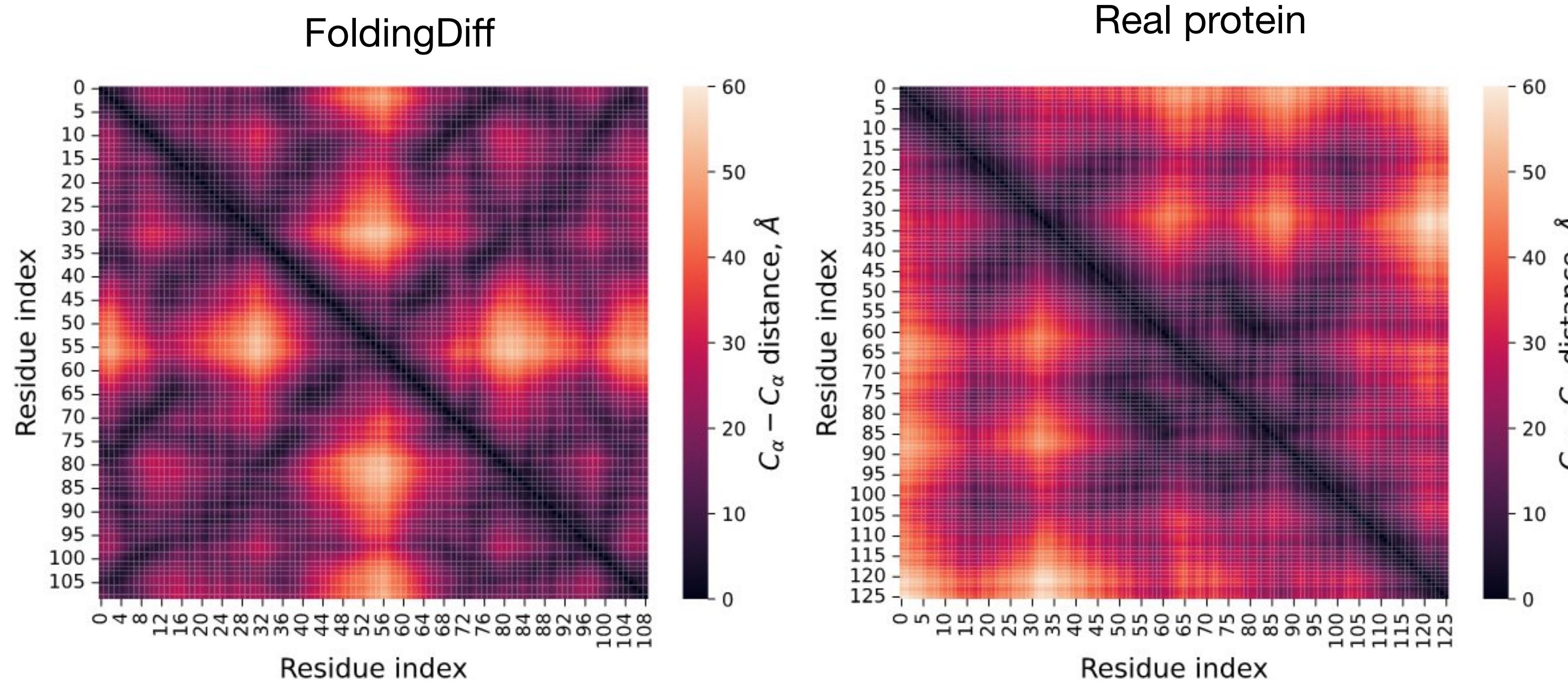
Generated structures are diverse



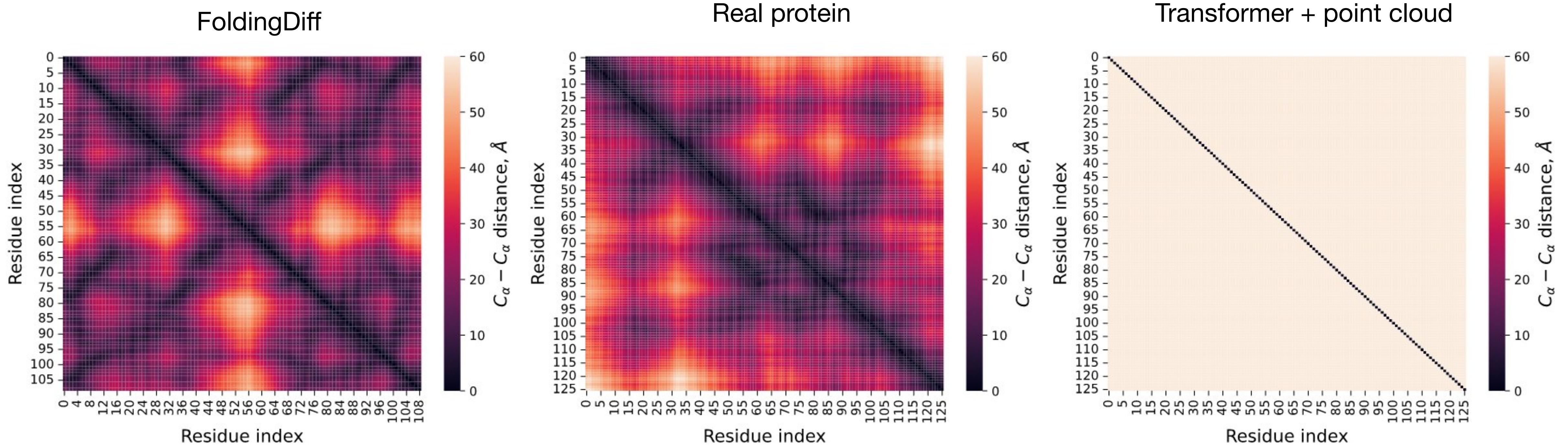
Angle formulation is enables generation from vanilla transformers



Angle formulation is enables generation from vanilla transformers

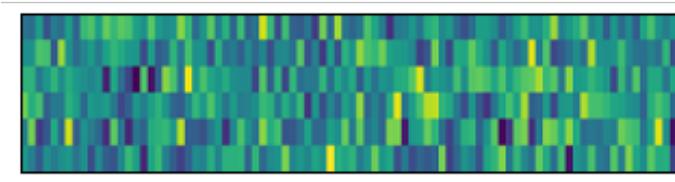
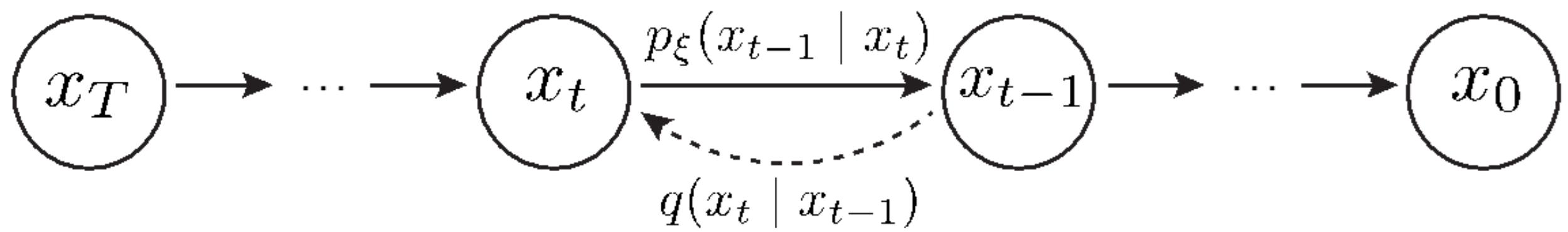


Angle formulation is enables generation from vanilla transformers

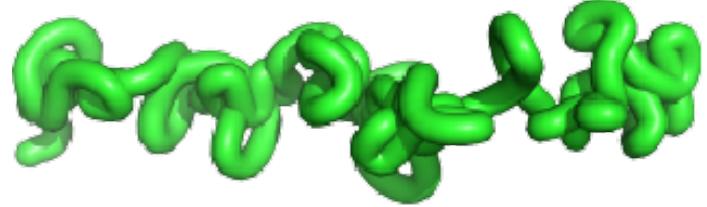


FoldingDiff is first step towards generating new functions

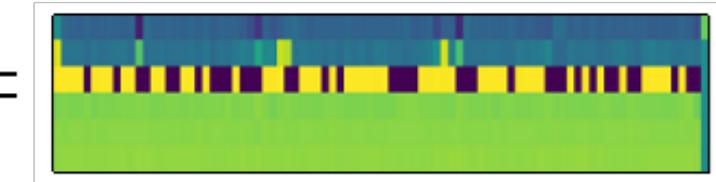
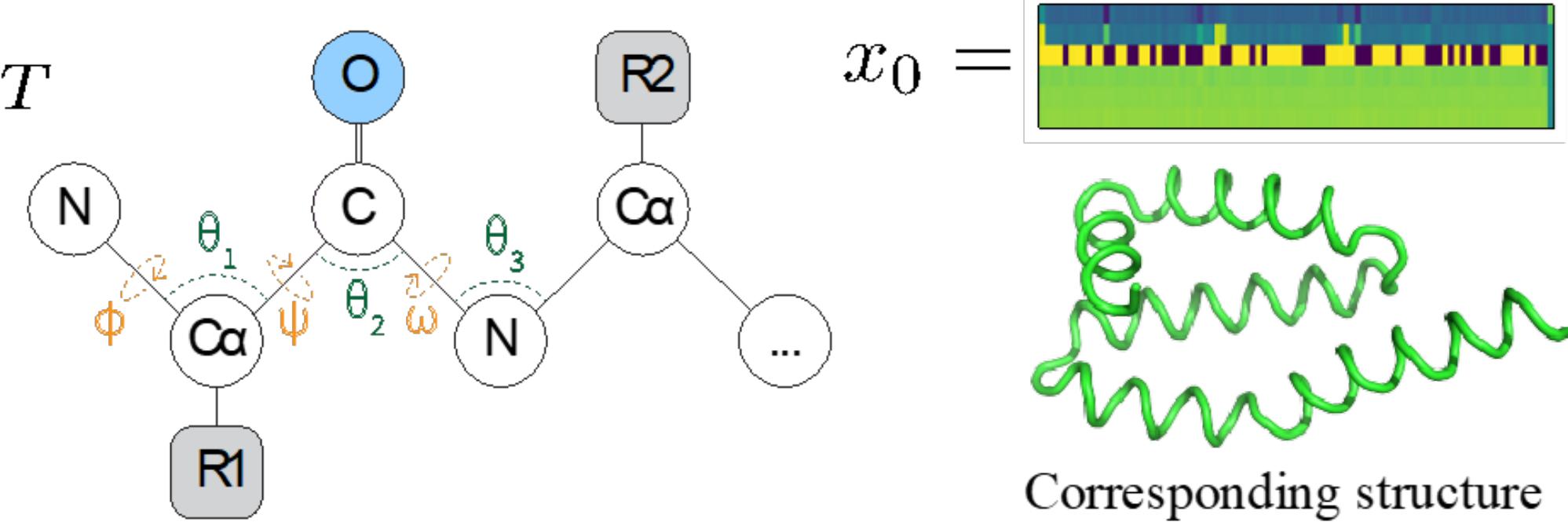
FoldingDiff is first step towards generating new functions



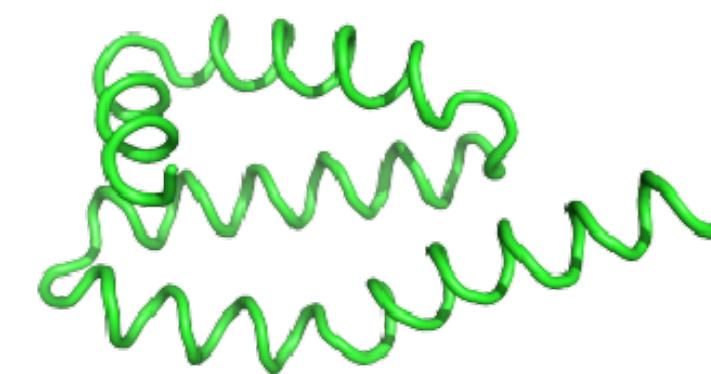
$$= x_T$$



Corresponding structure

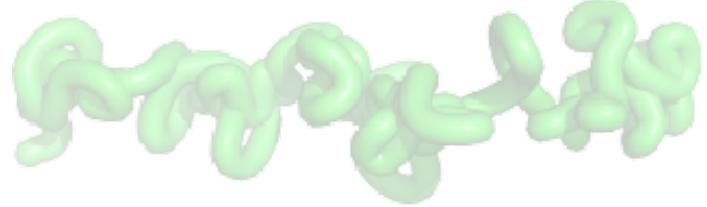
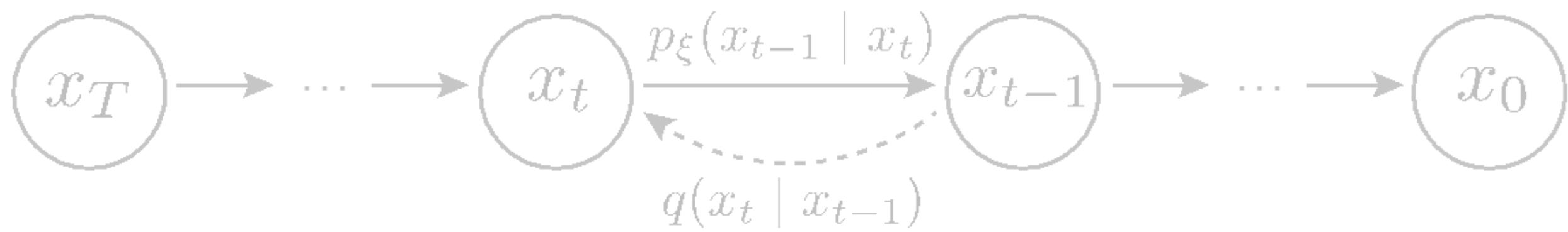


$$x_0 =$$

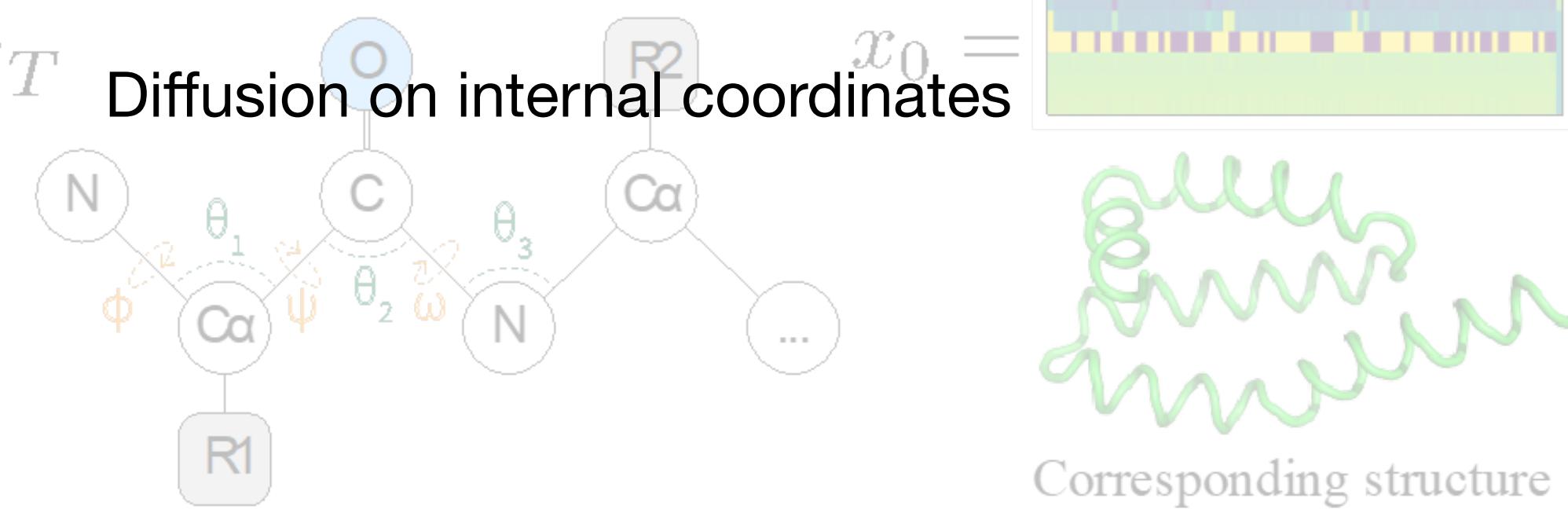


Corresponding structure

FoldingDiff is first step towards generating new functions

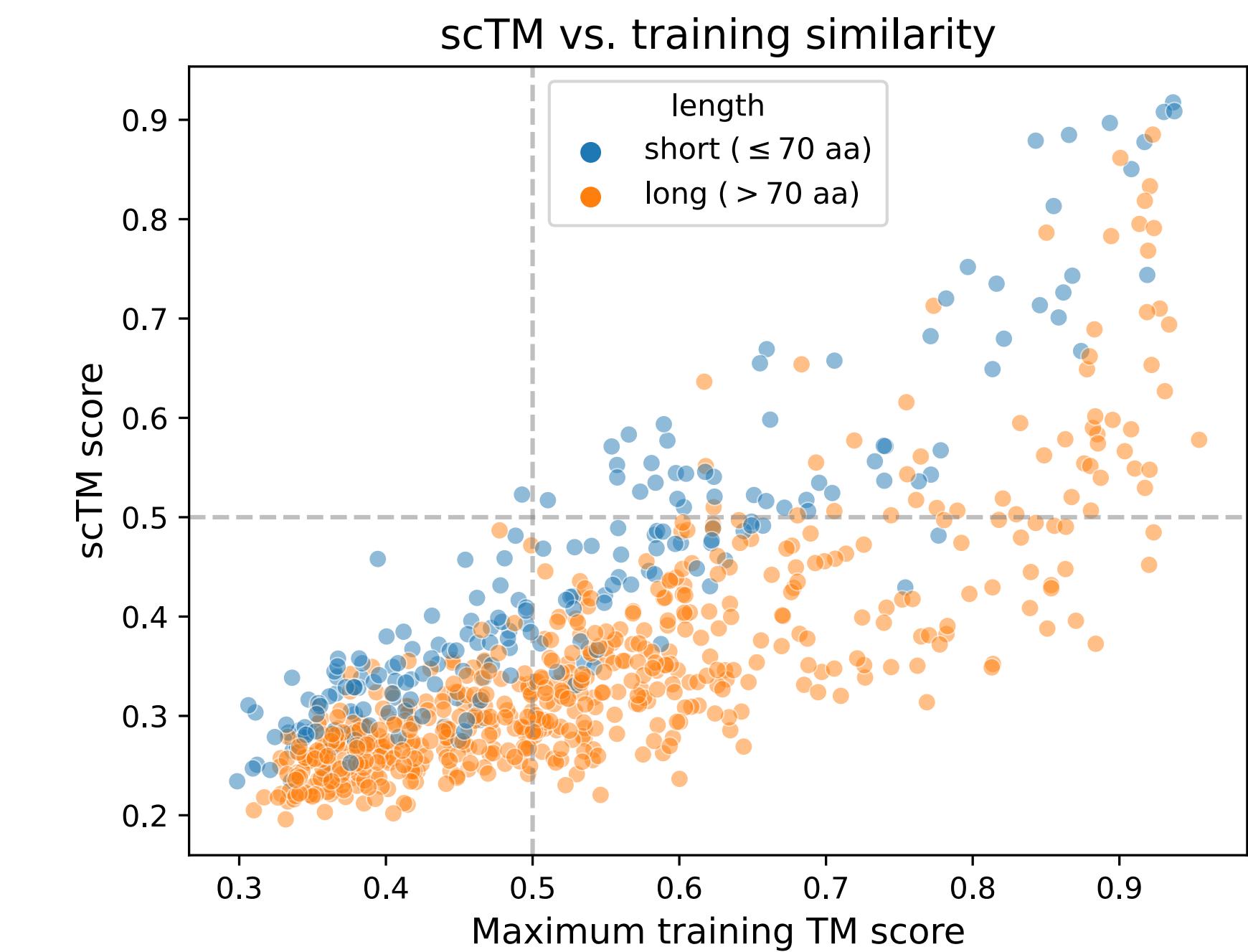
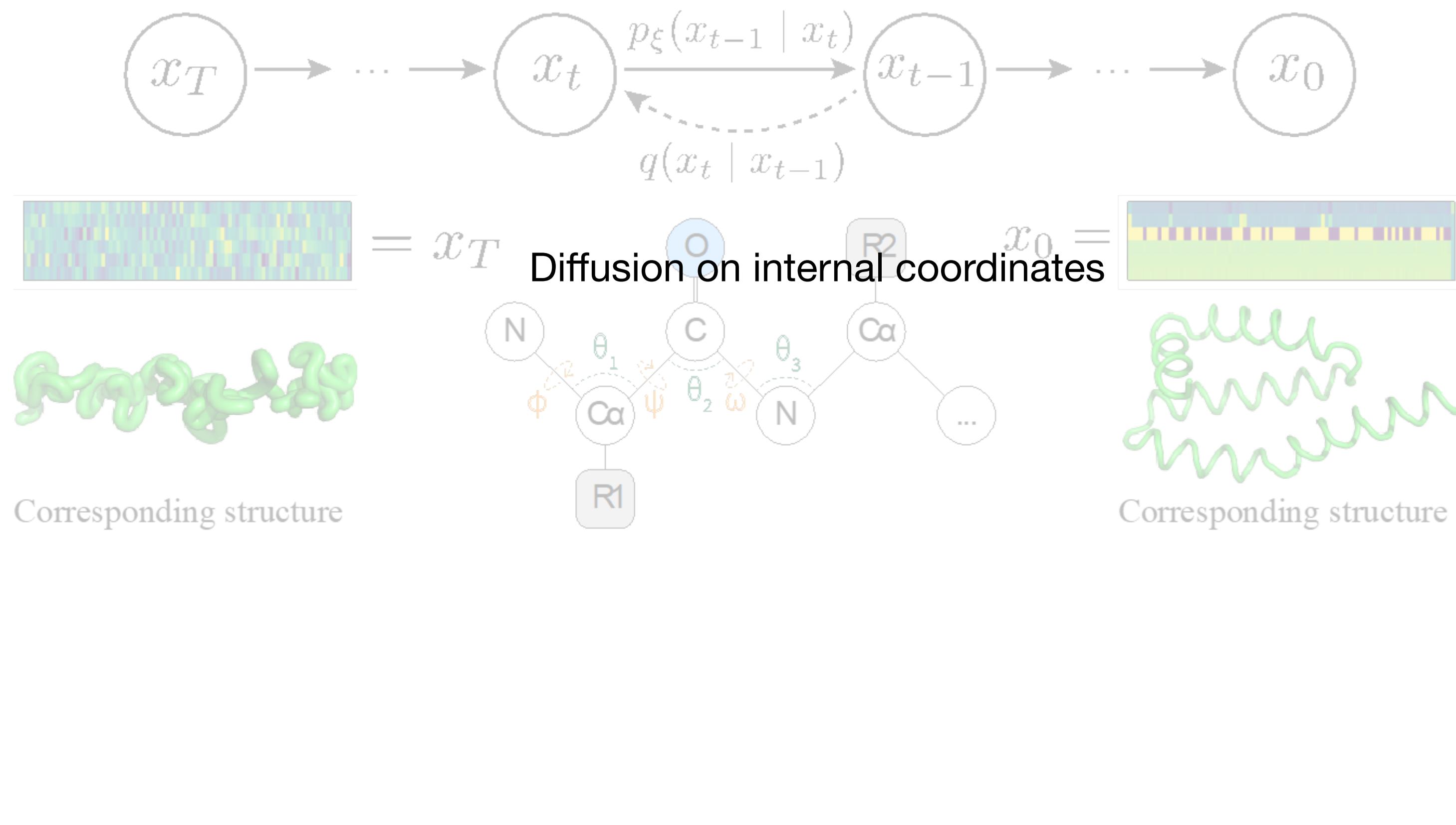


Corresponding structure

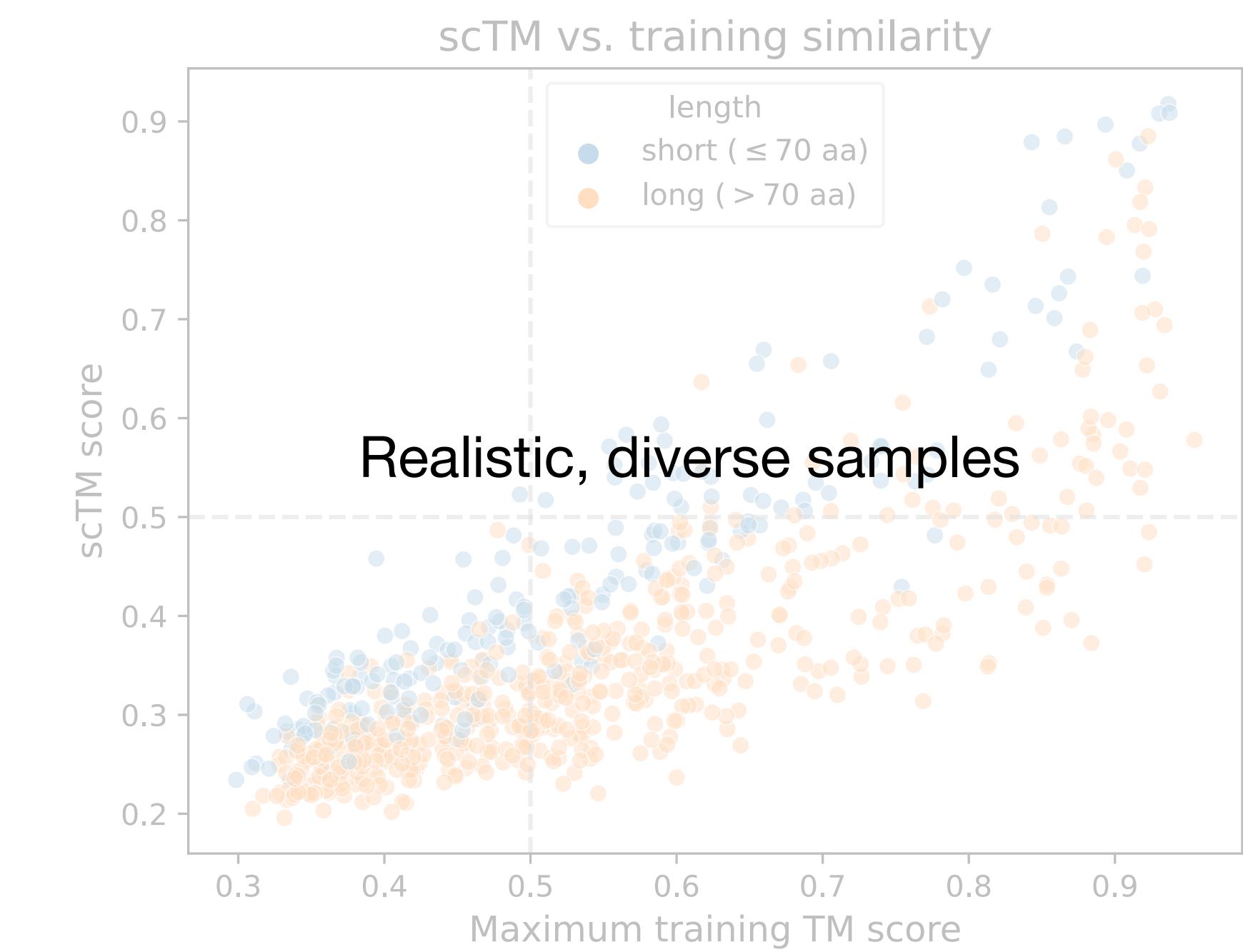
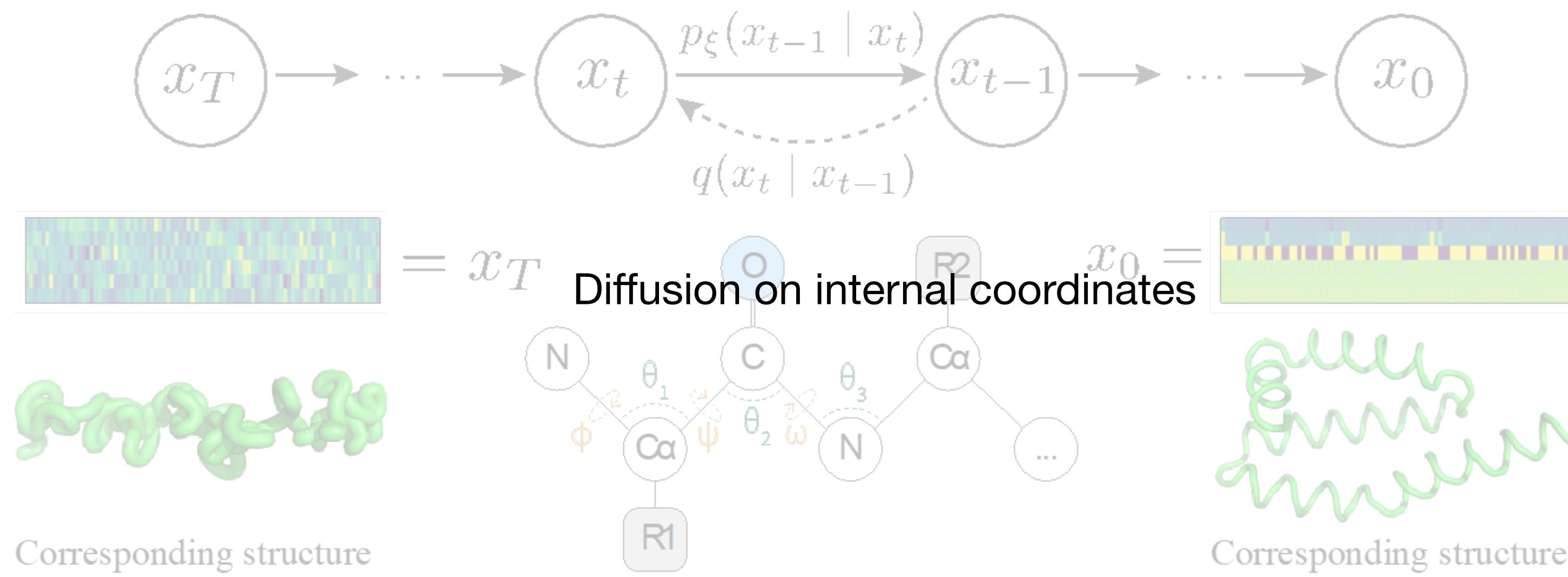


Corresponding structure

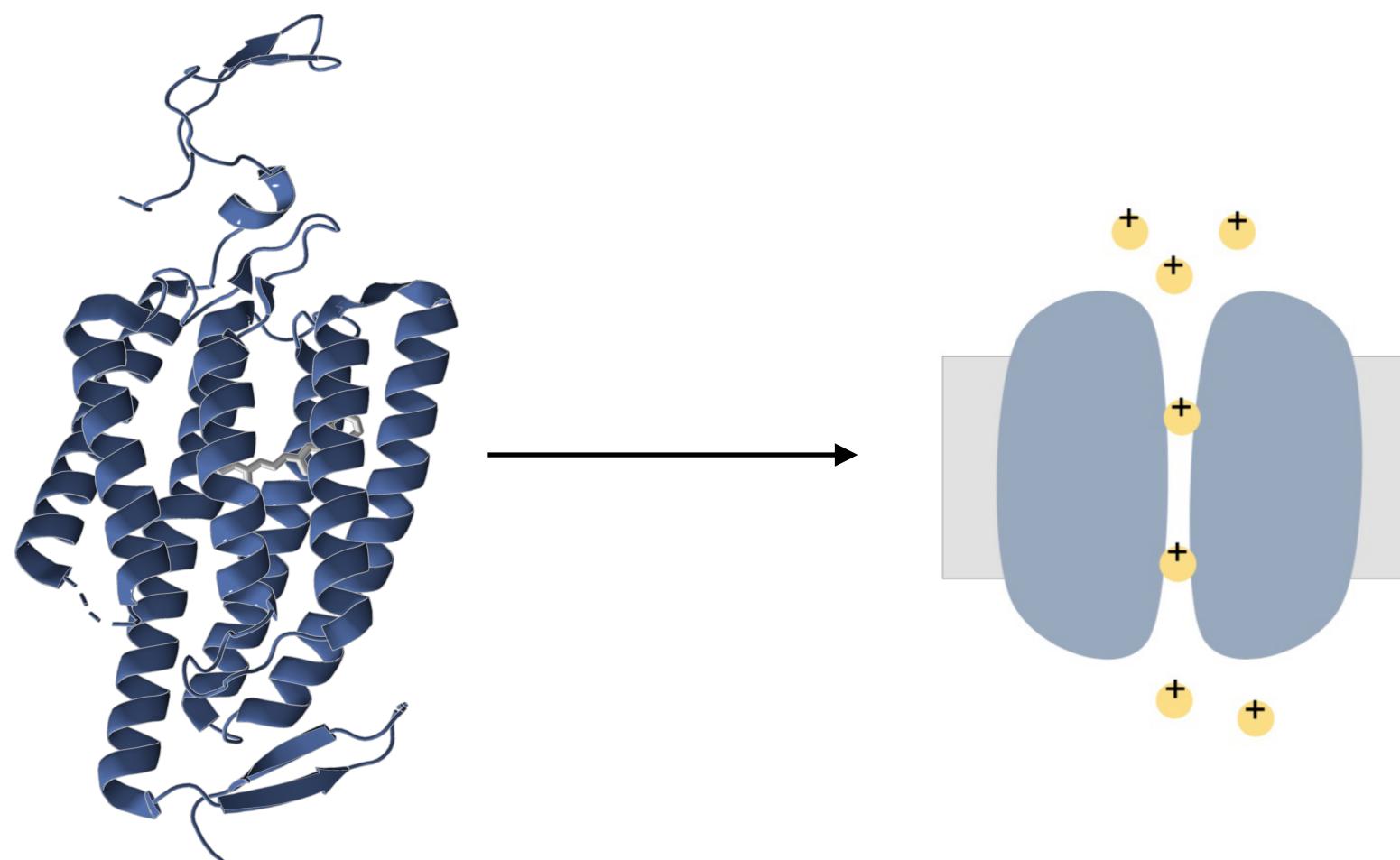
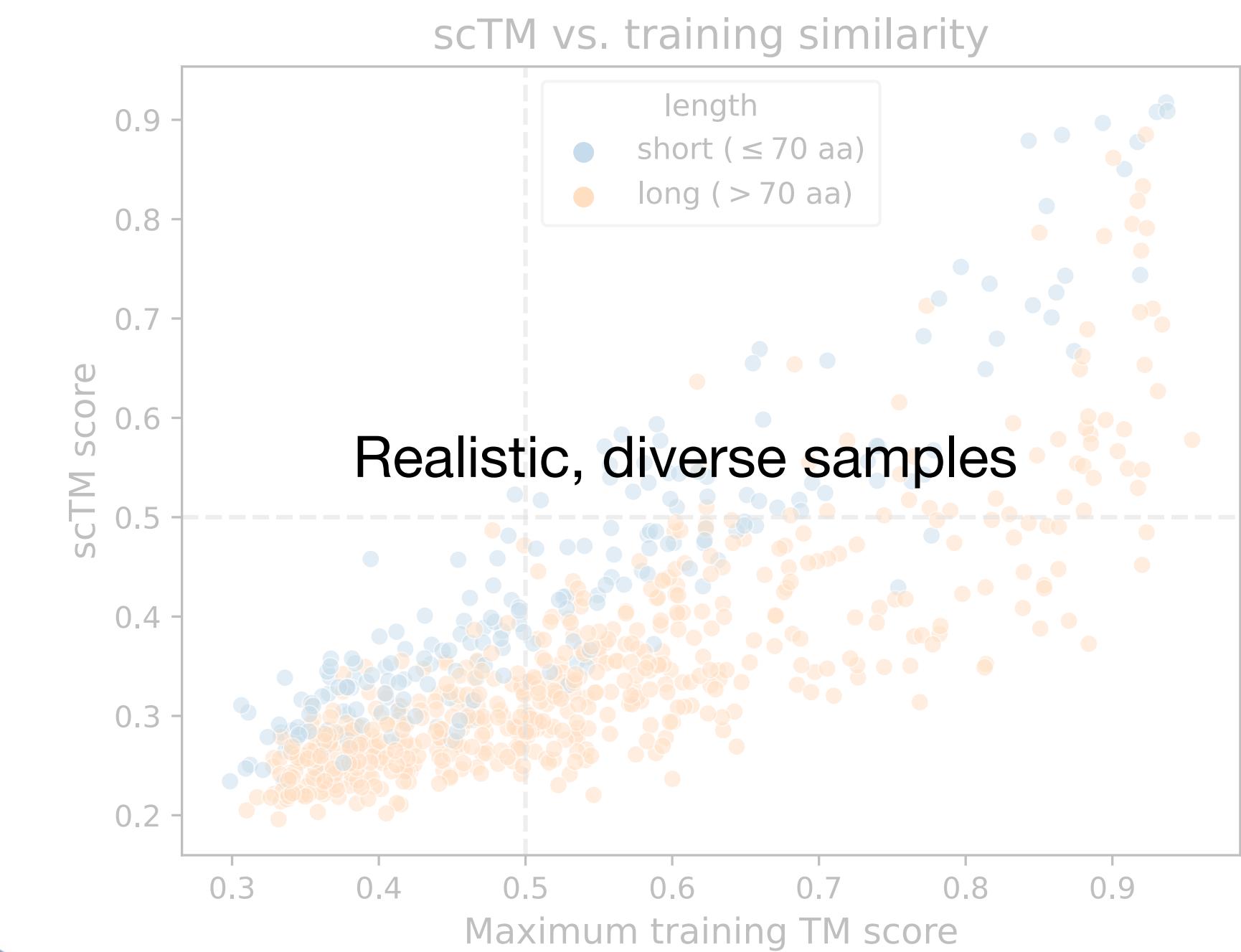
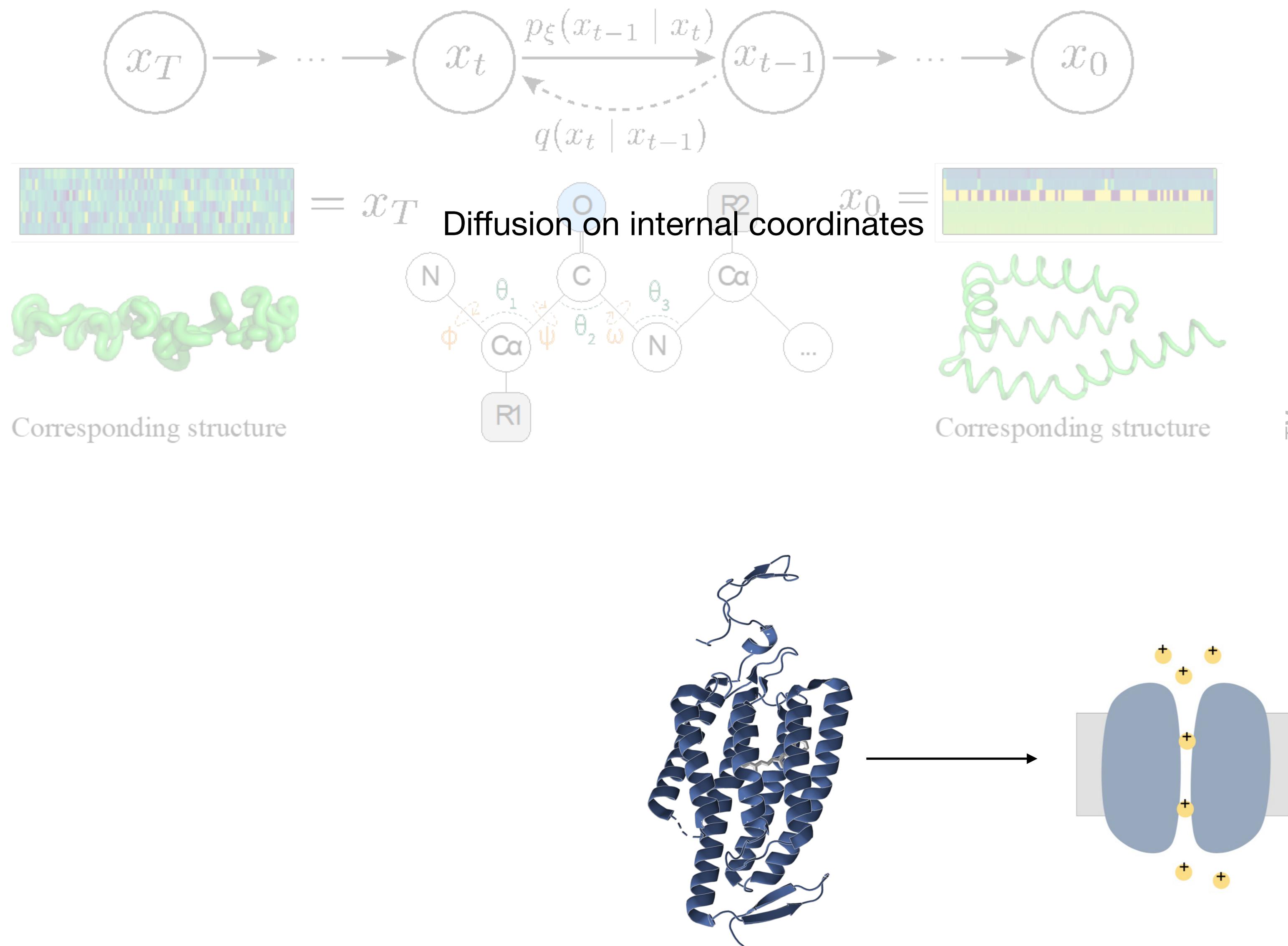
FoldingDiff is first step towards generating new functions



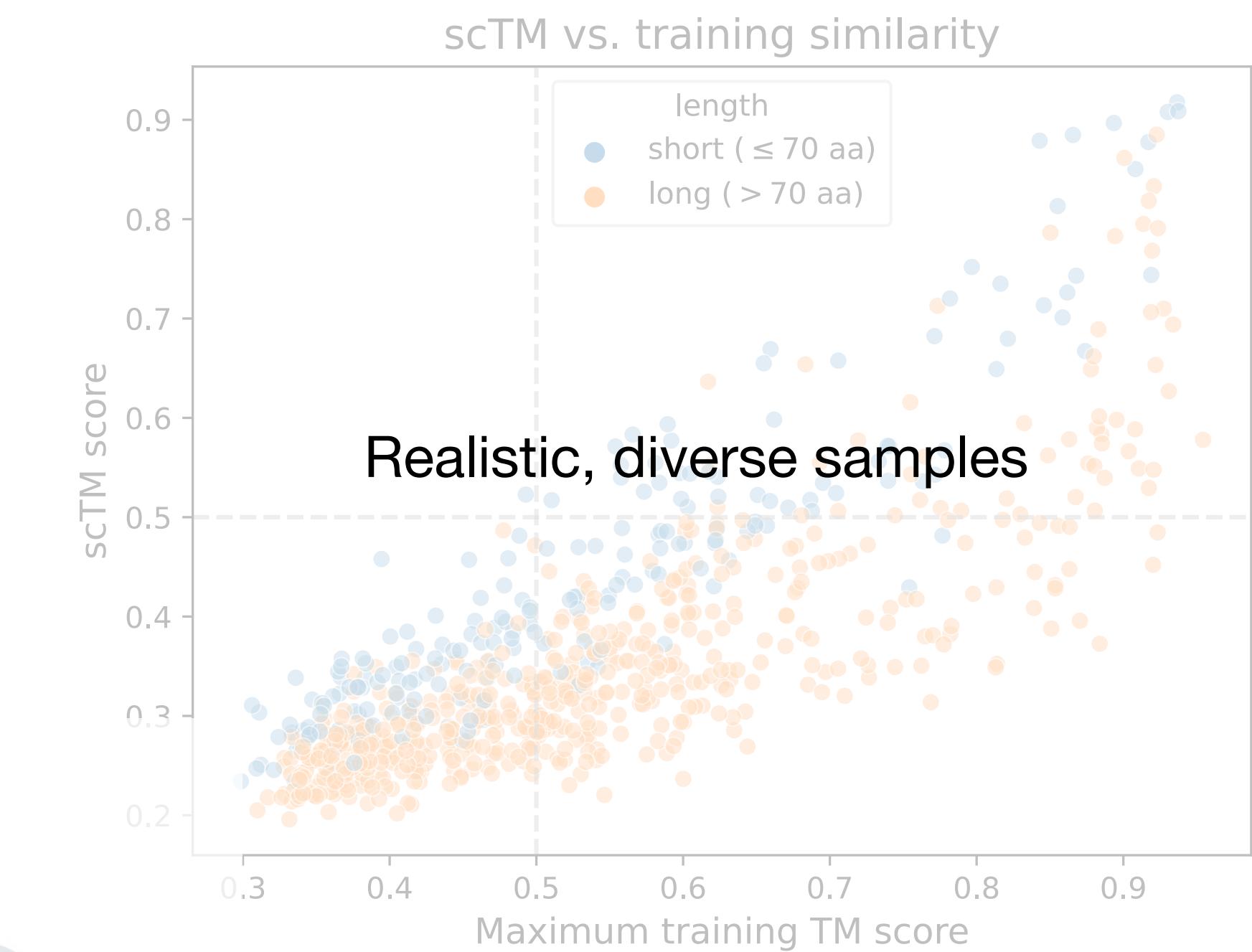
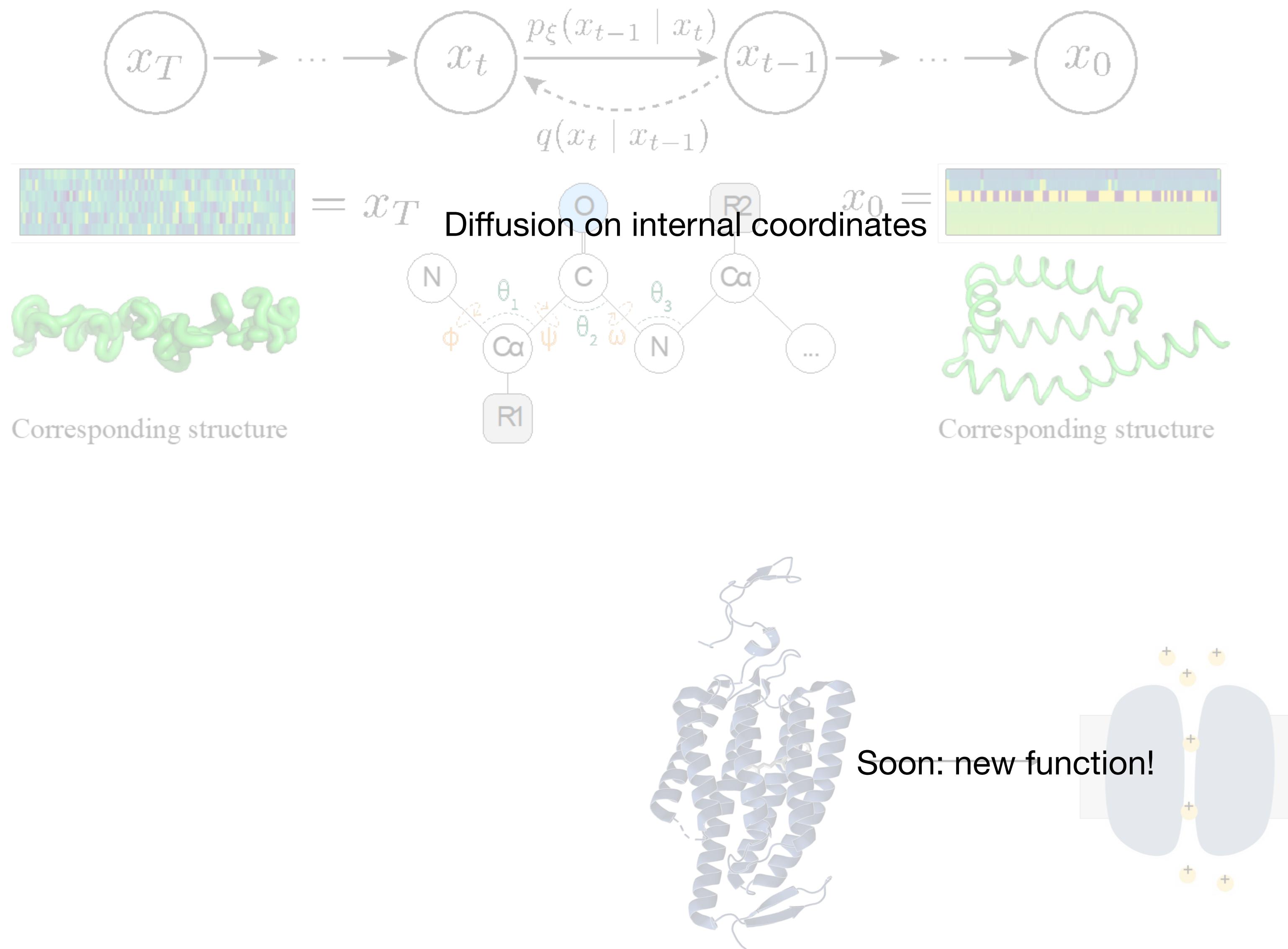
FoldingDiff is first step towards generating new functions



FoldingDiff is first step towards generating new functions



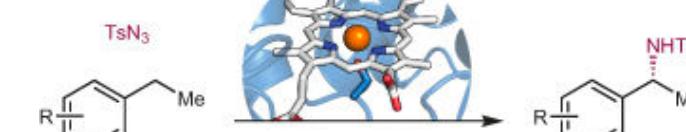
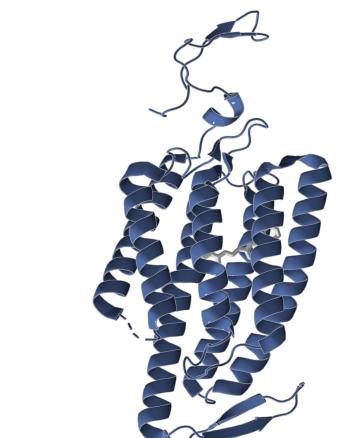
FoldingDiff is first step towards generating new functions



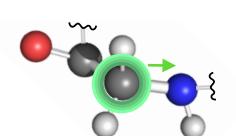
Multi-modal models

MGHTQWI . . .

annotations



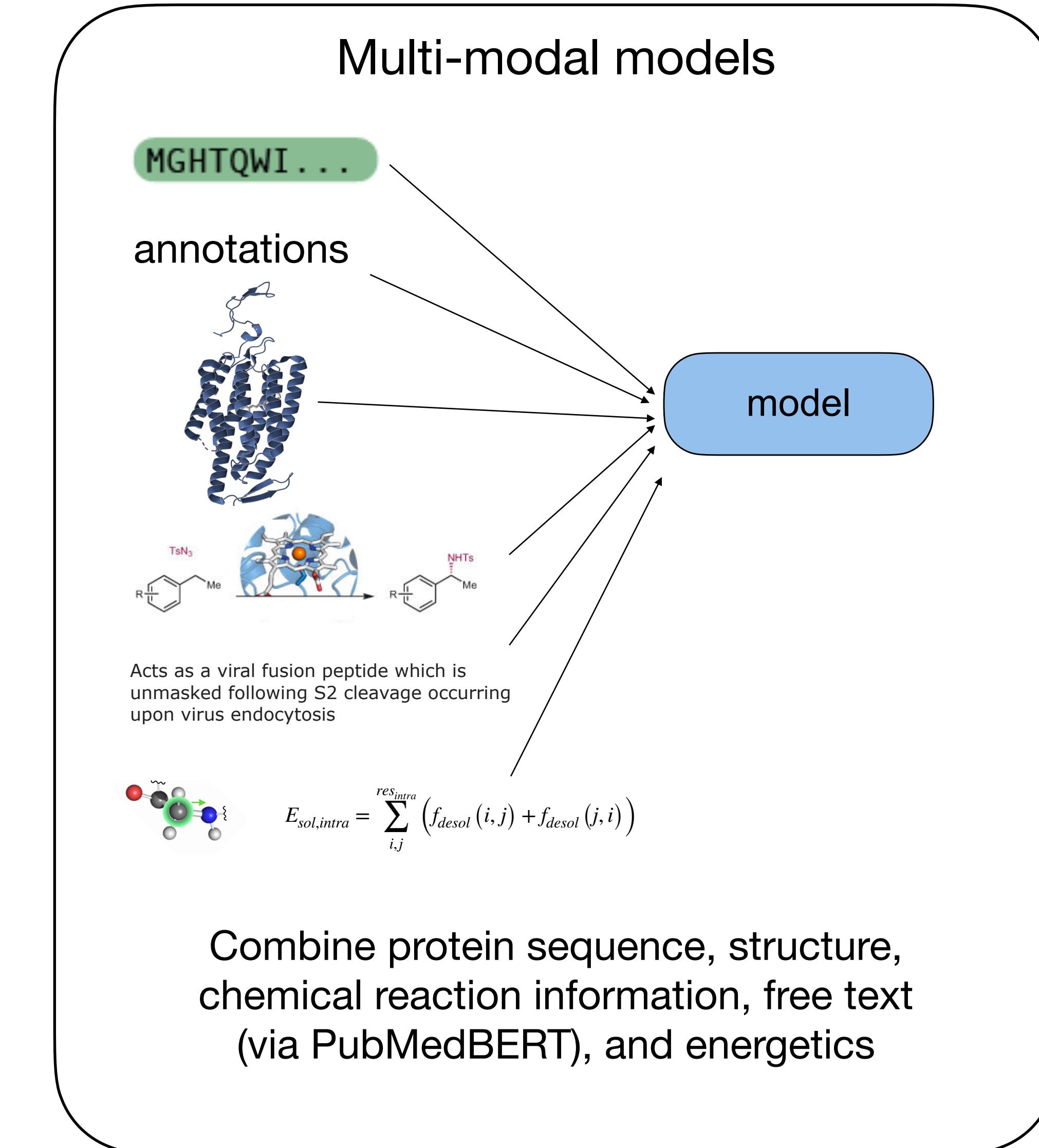
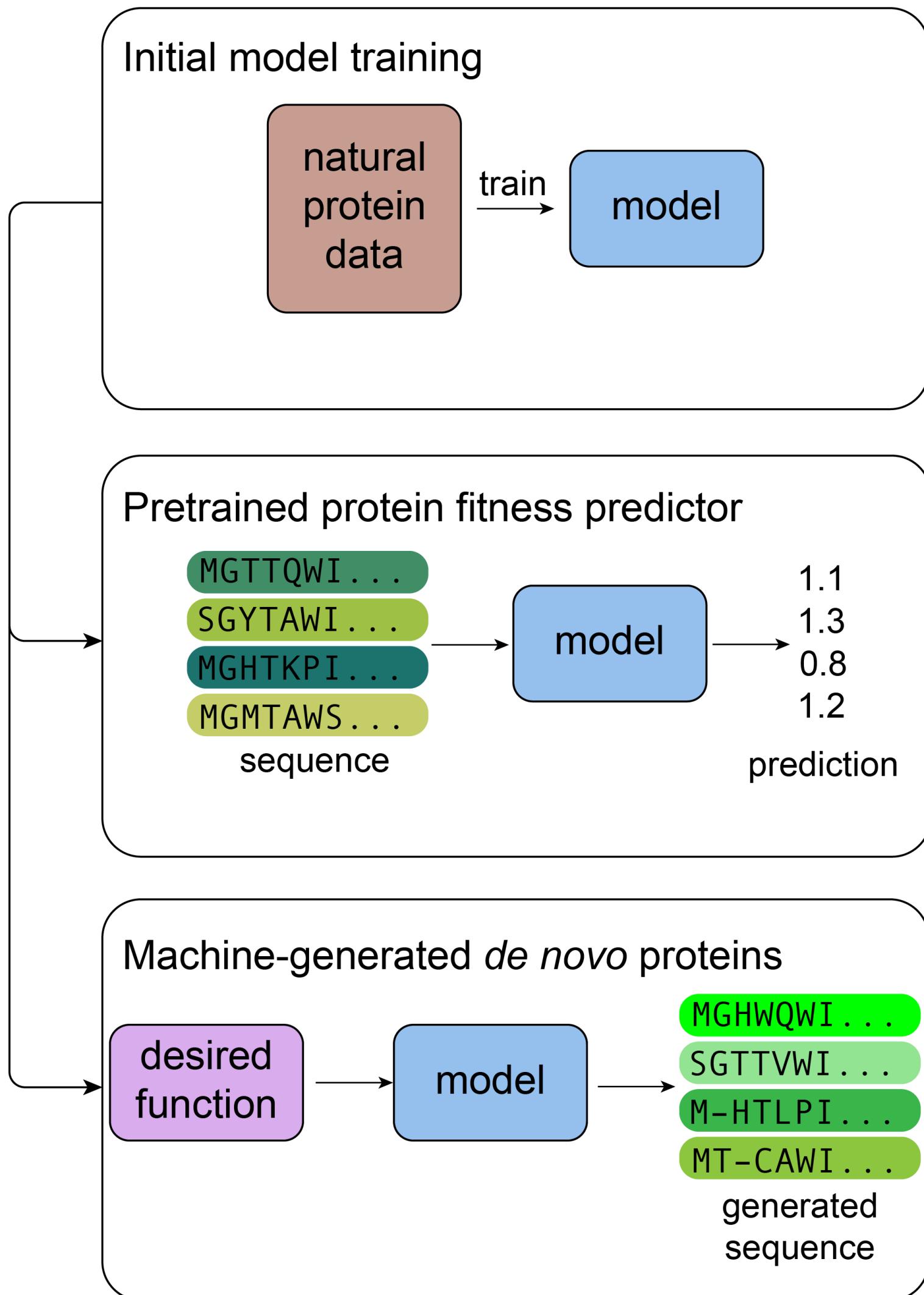
Acts as a viral fusion peptide which is unmasked following S2 cleavage occurring upon virus endocytosis



$$E_{sol,intra} = \sum_{i,j}^{res_{intra}} (f_{desol}(i,j) + f_{desol}(j,i))$$

Combine protein sequence, structure, chemical reaction information, free text (via PubMedBERT), and energetics

Can we predict and generate functional proteins?





BioML at MSR New England

Acknowledgments



BioML at MSR New England