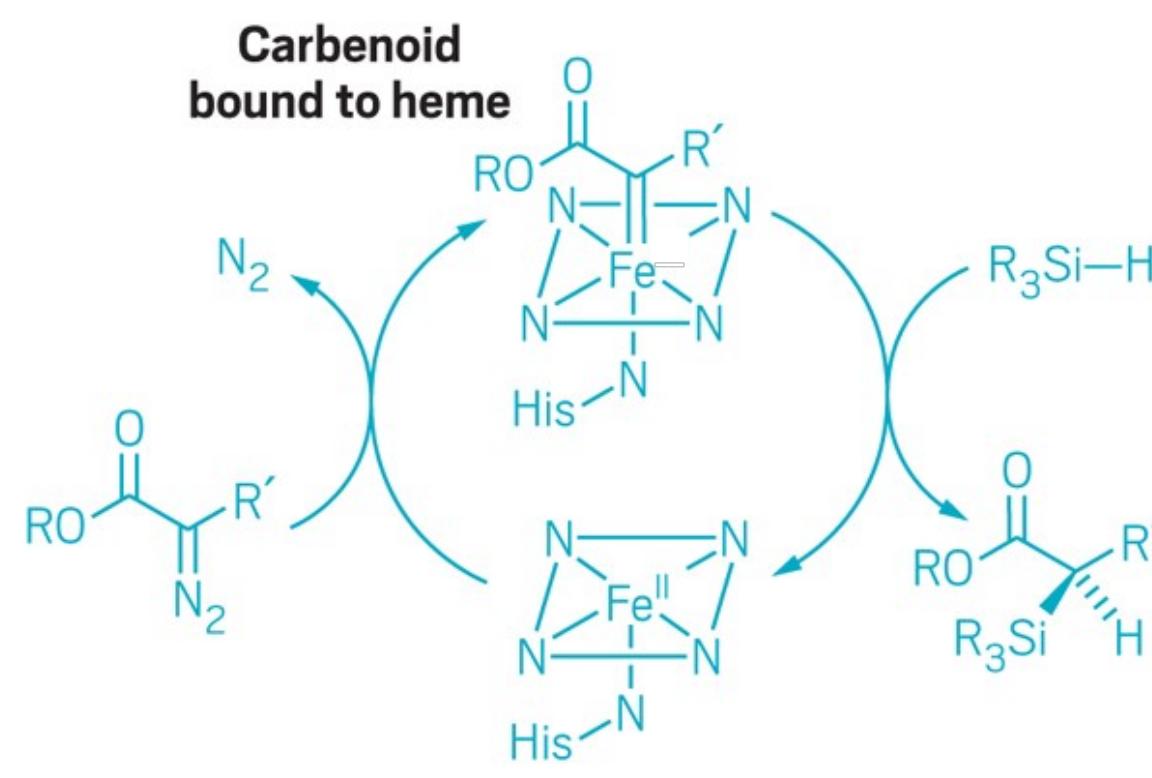


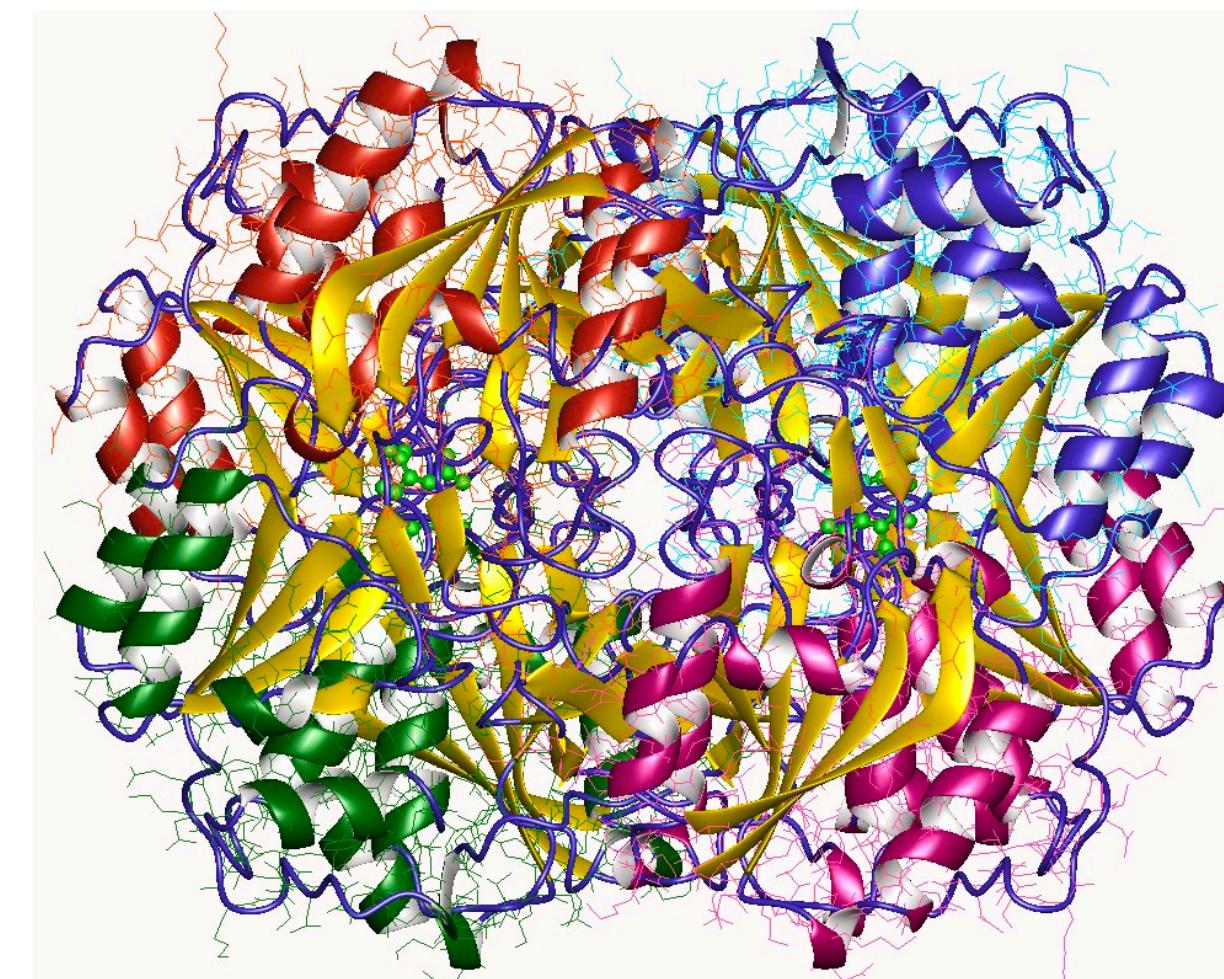
Protein sequence generation with evolutionary diffusion

Kevin Kaichuang Yang
Microsoft Research New England
 @KevinKaichuang

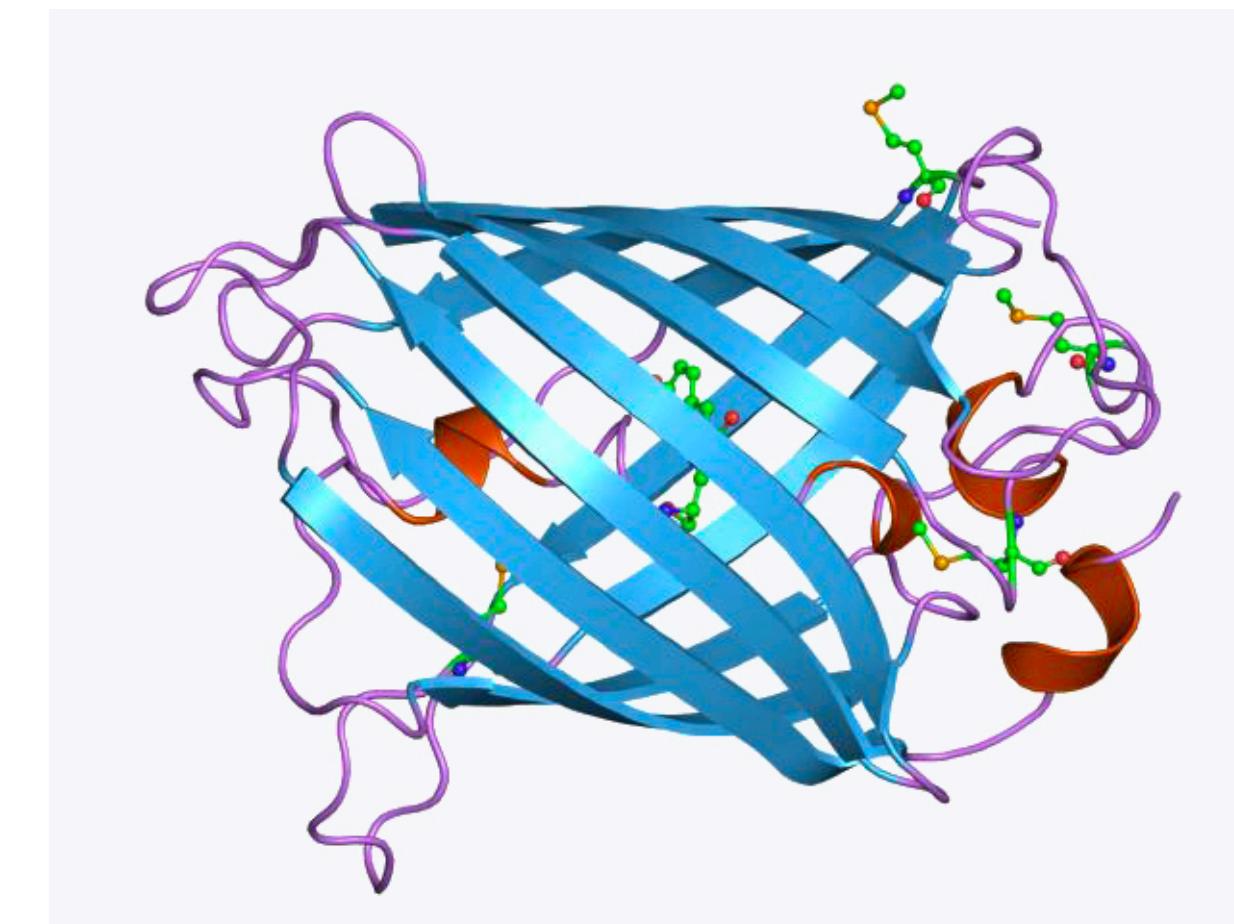
We need proteins with new functions



new chemistry

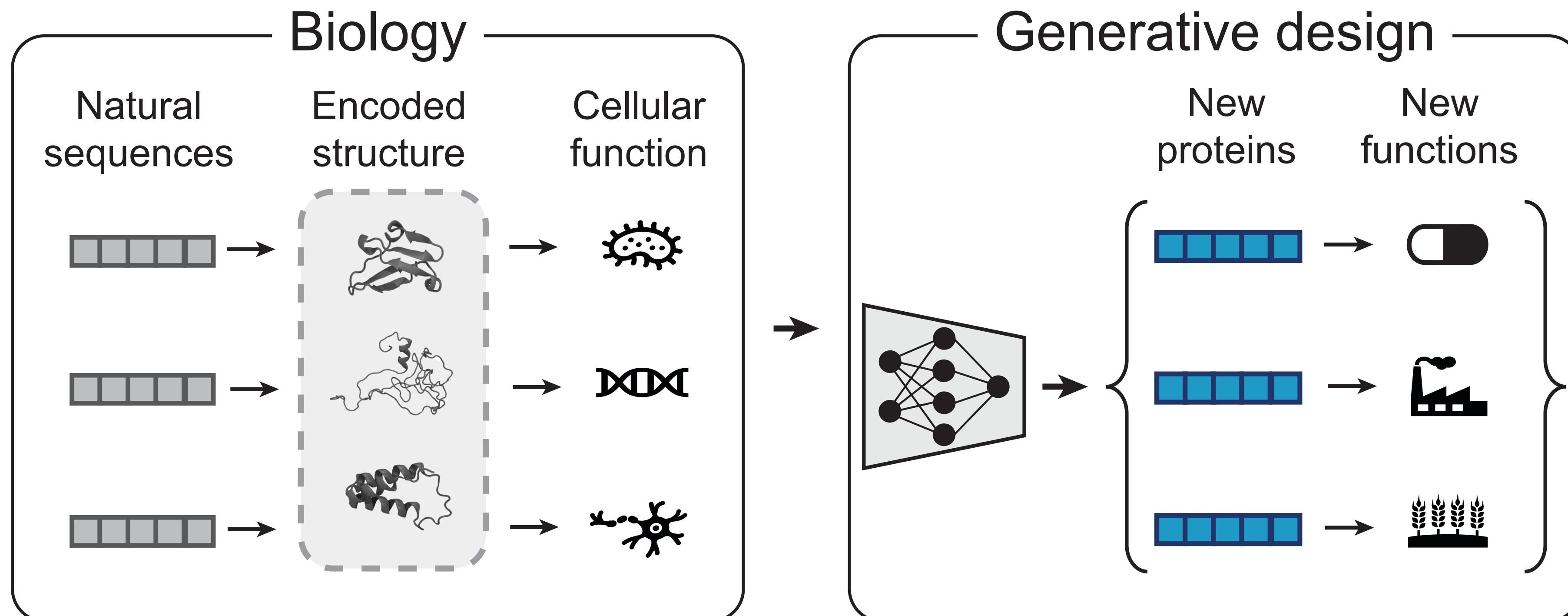


therapeutics

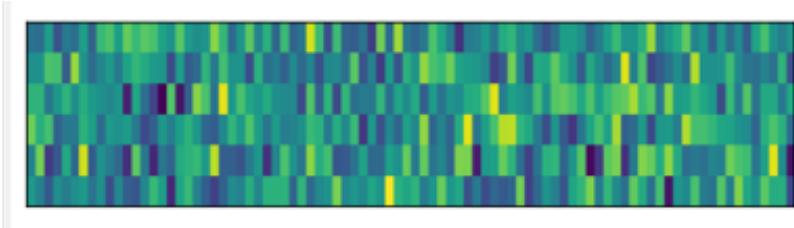
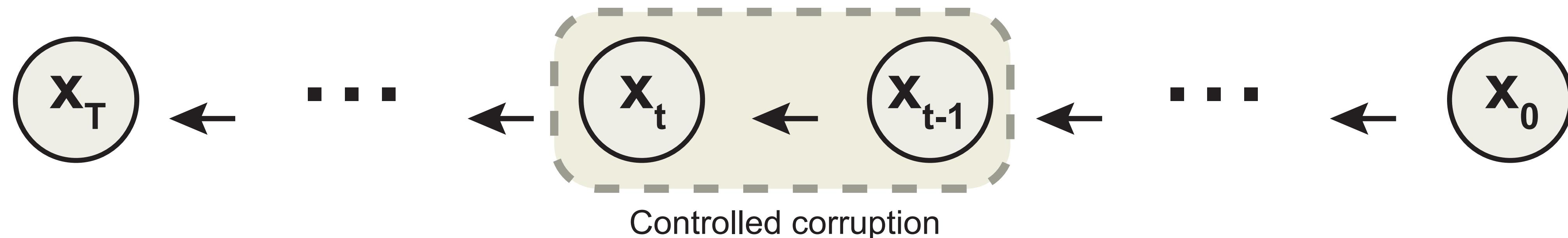


molecular tools

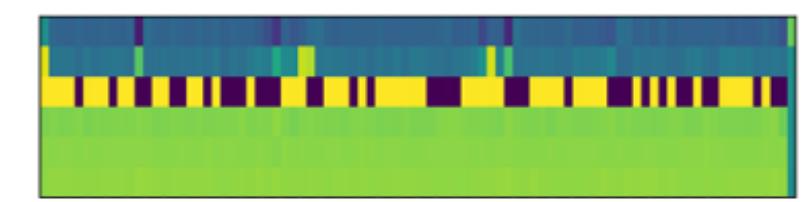
Generate new proteins to expand functional space



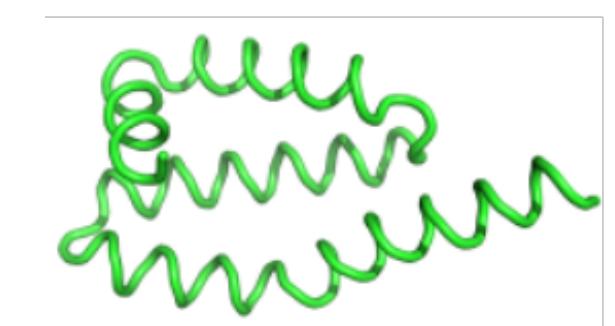
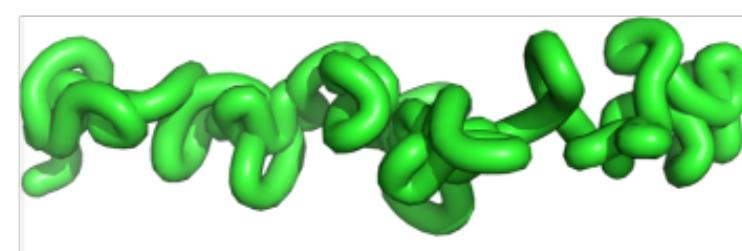
Diffusion models generate realistic structures



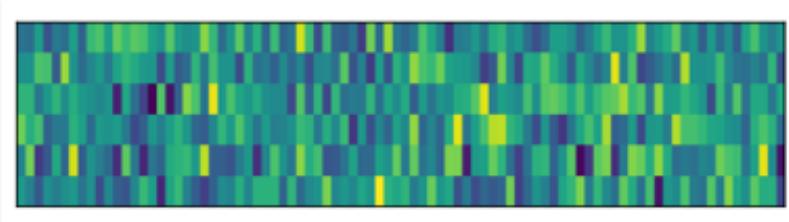
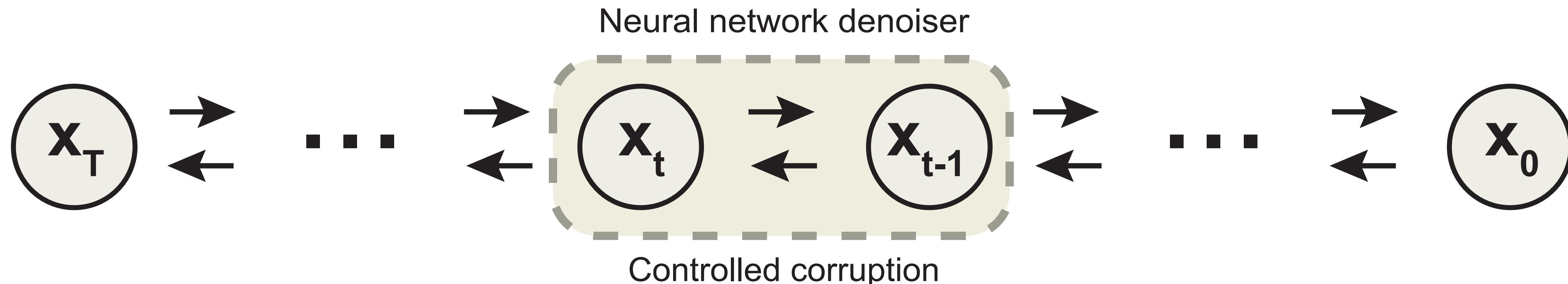
random structures



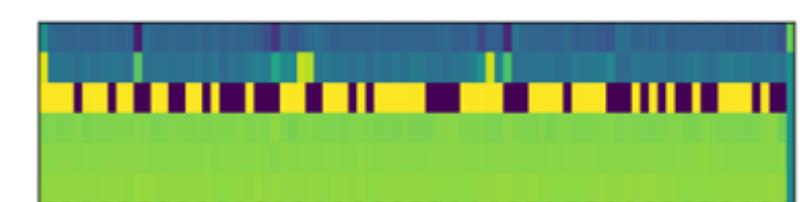
real structures



Diffusion models generate realistic structures



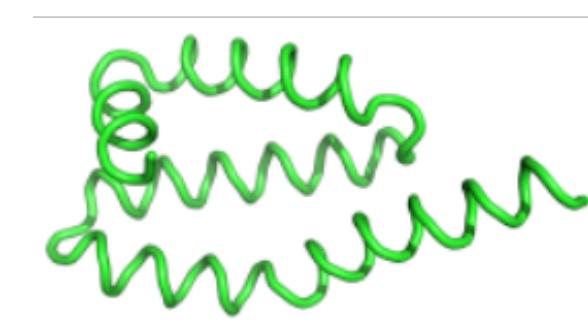
random structures



generated structures

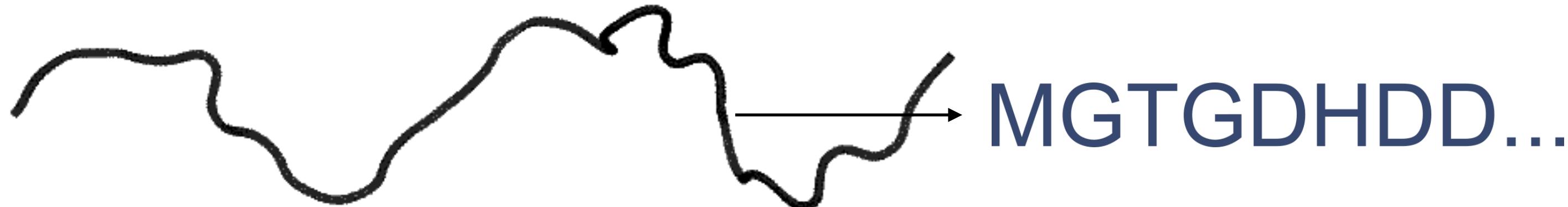


Wu et al., preprint 2022

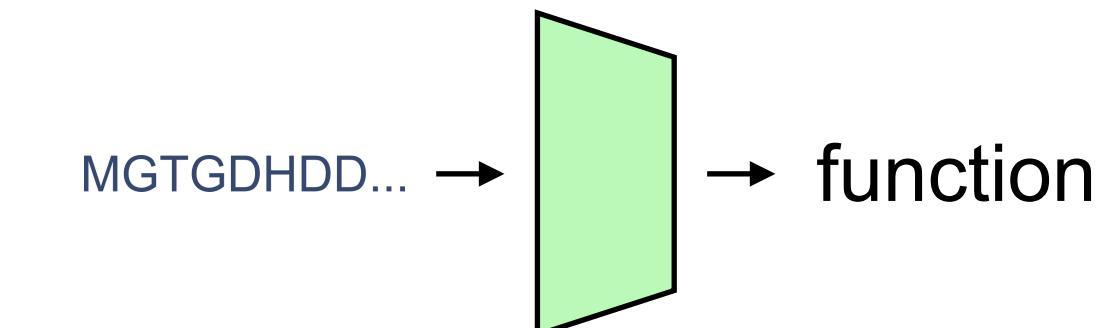


Sequence is the universal protein design space

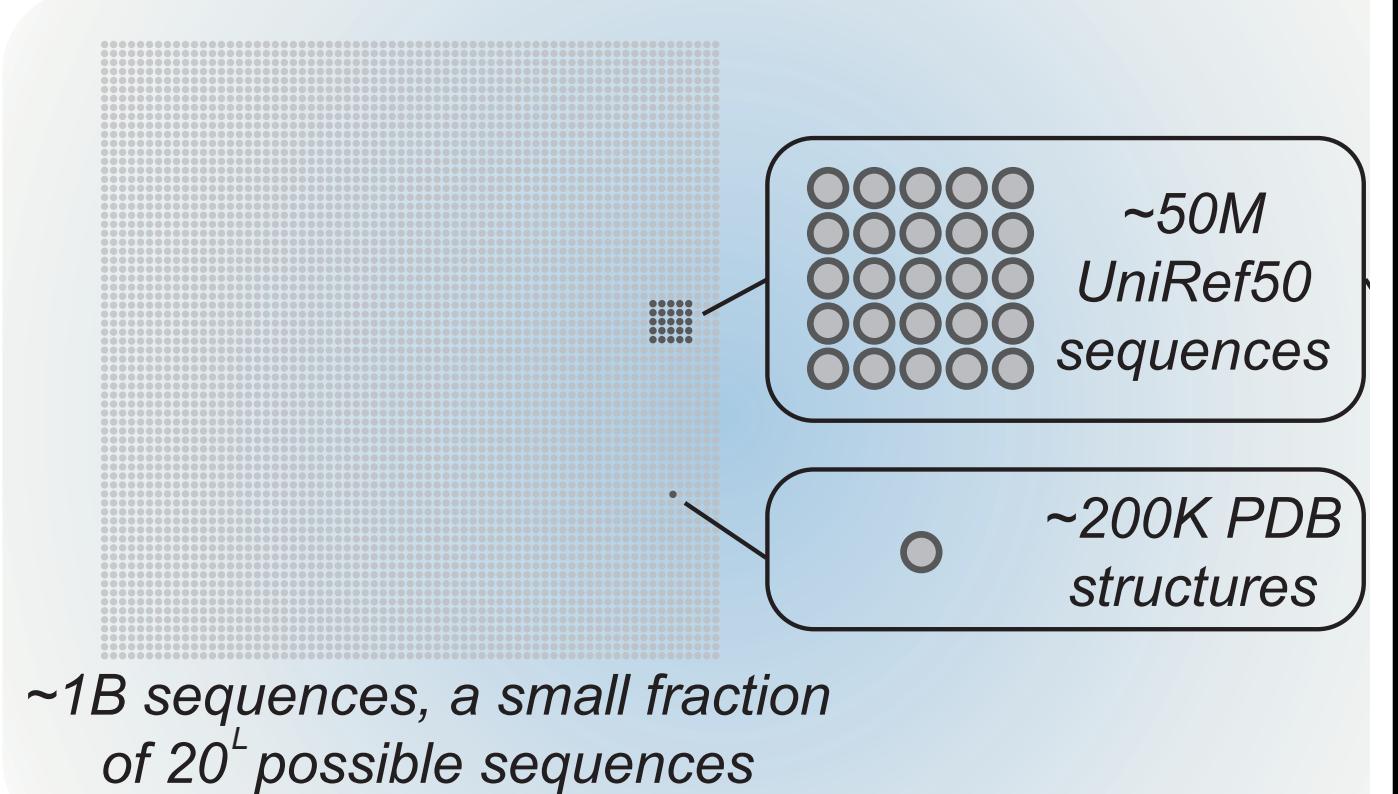
No sequence design



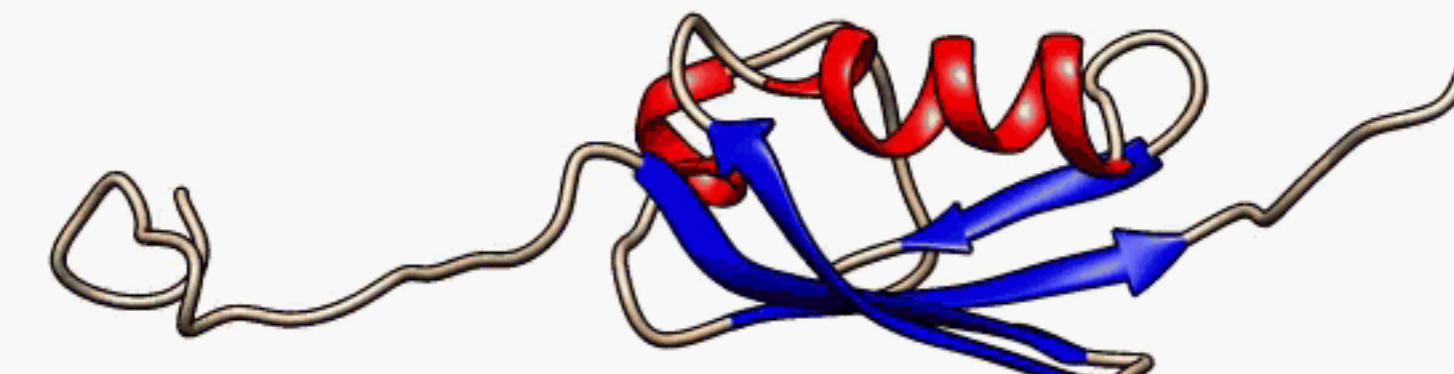
Many function predictors



Many more sequences
than structures

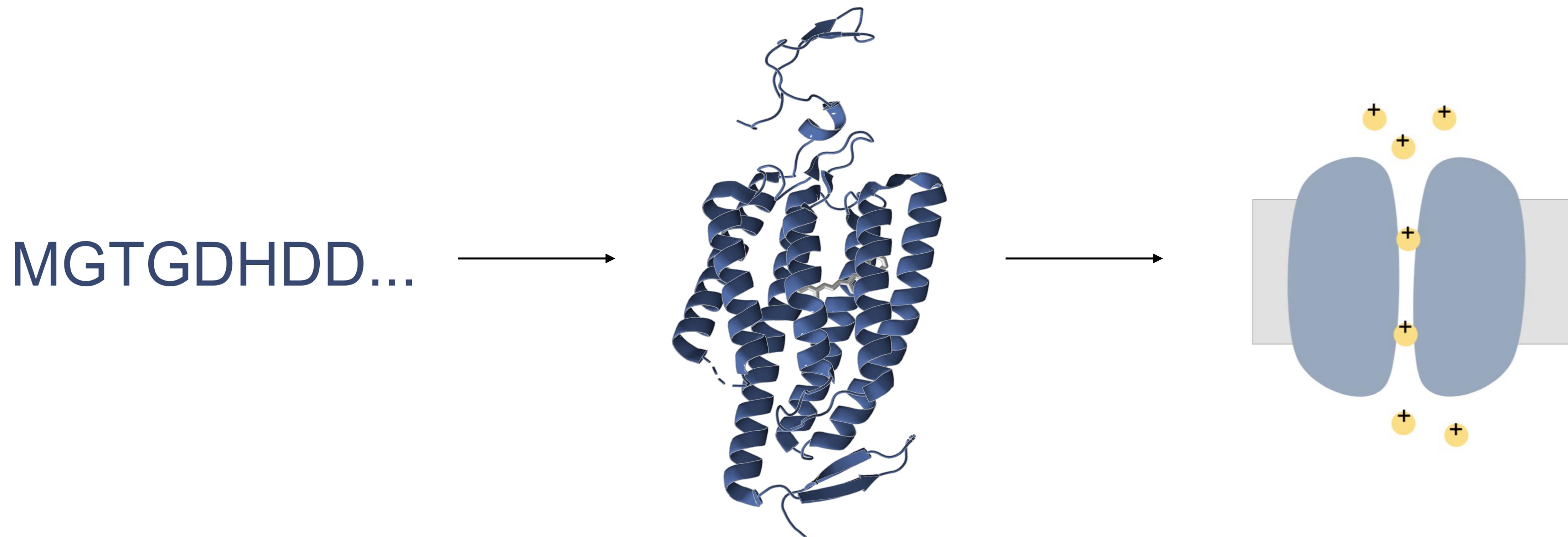


Handle disordered regions

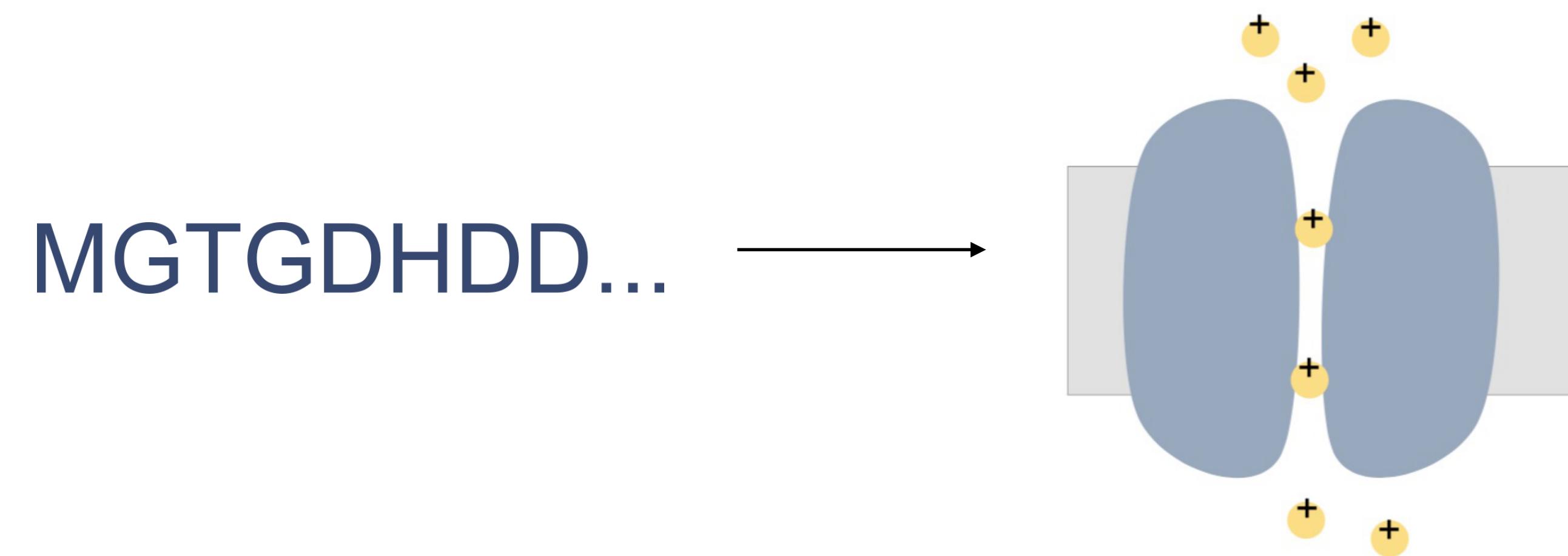


By Lukasz Kozlowski - Own work, CC BY-SA 4.0, <https://commons.wikimedia.org/w/index.php?curid=36298697>

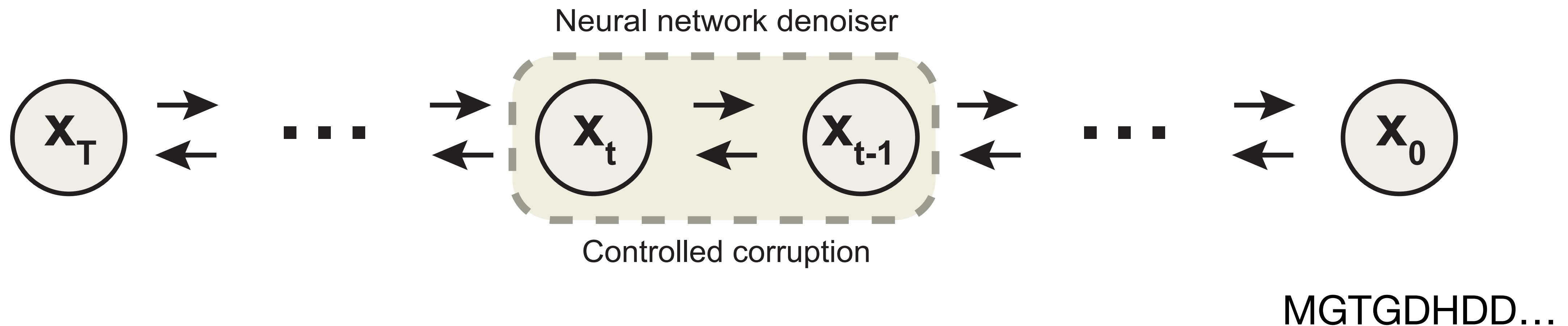
Sequence is the universal protein design space



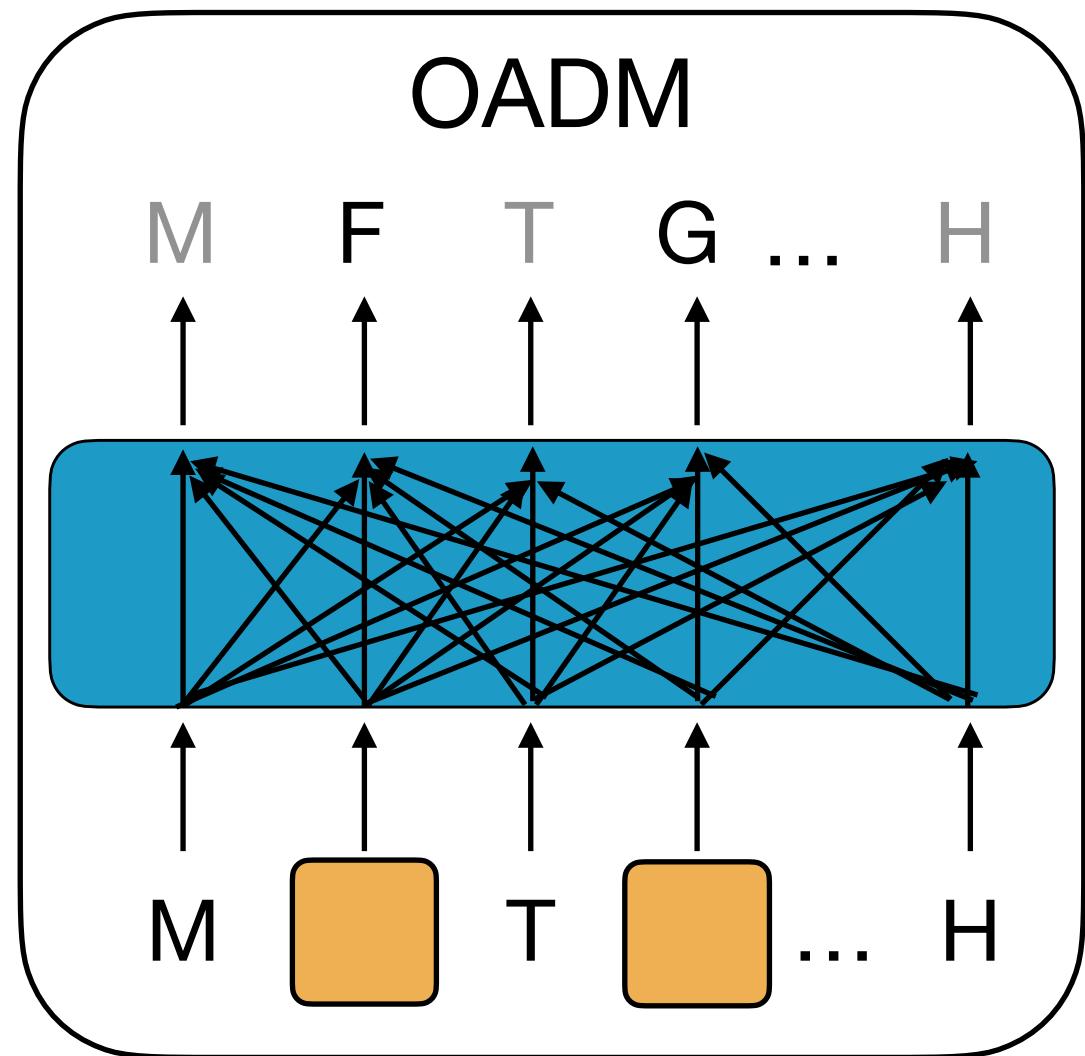
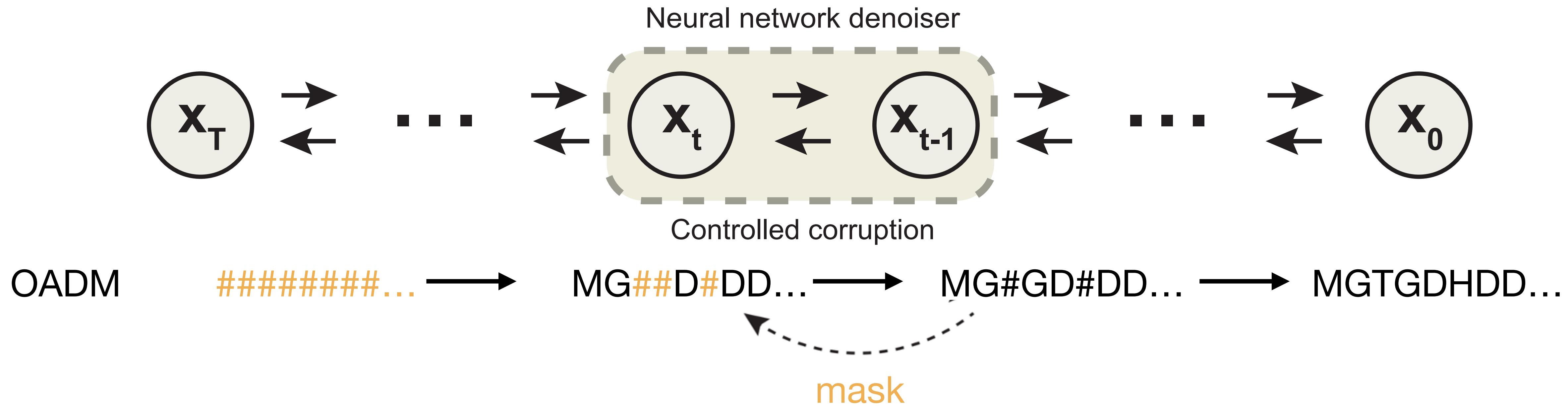
Sequence is the universal protein design space



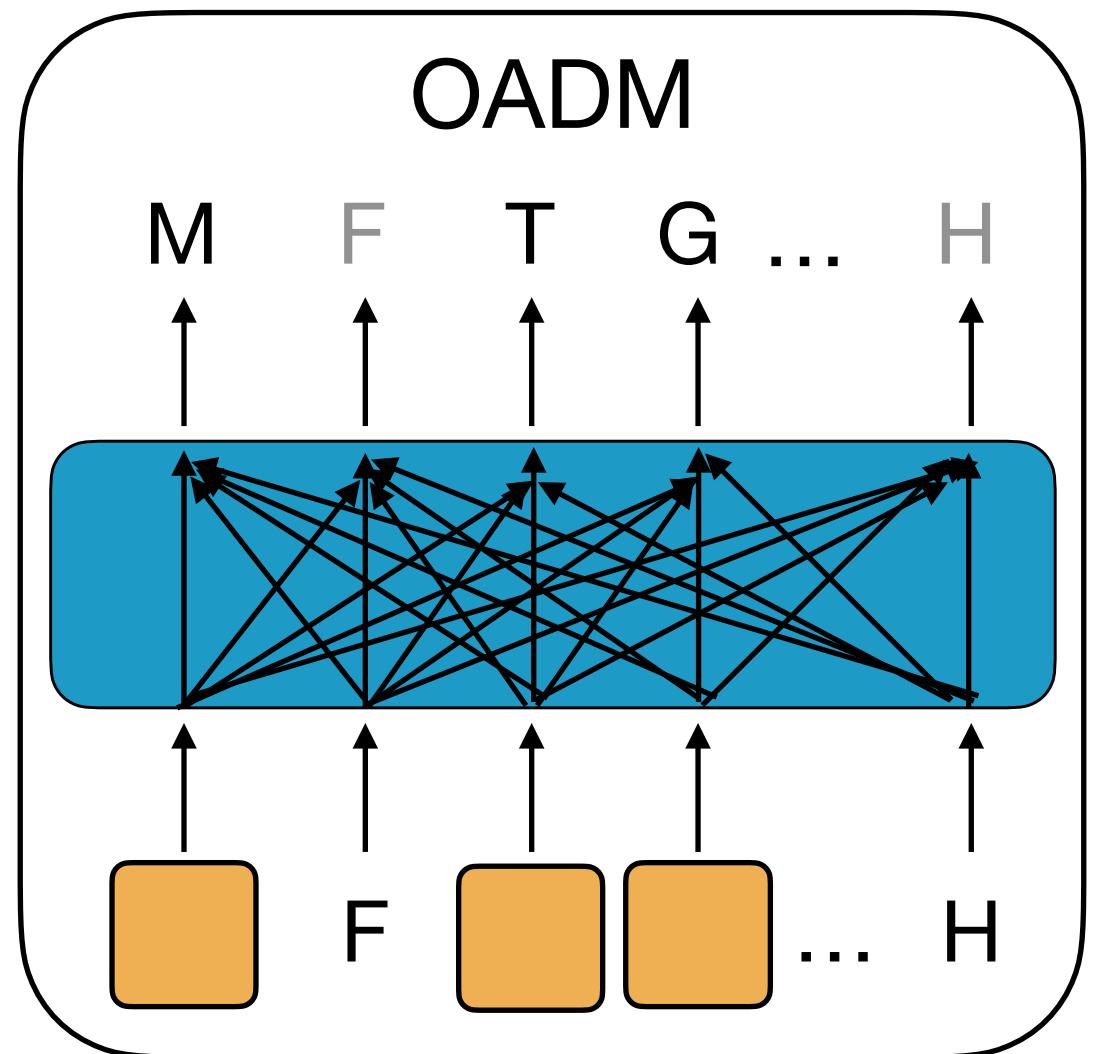
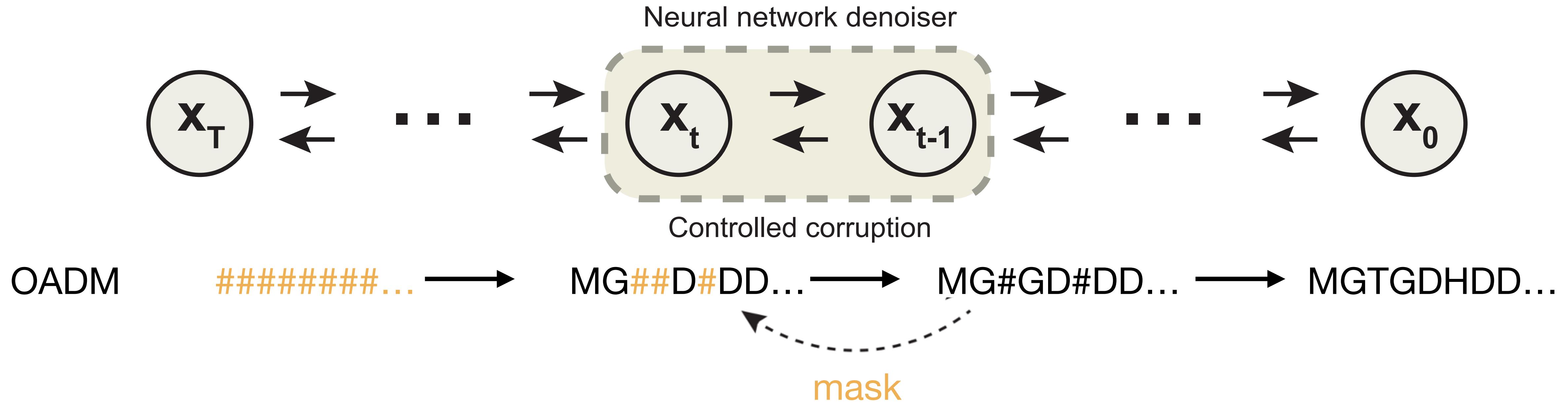
EvoDiff: evolutionary-scale diffusion



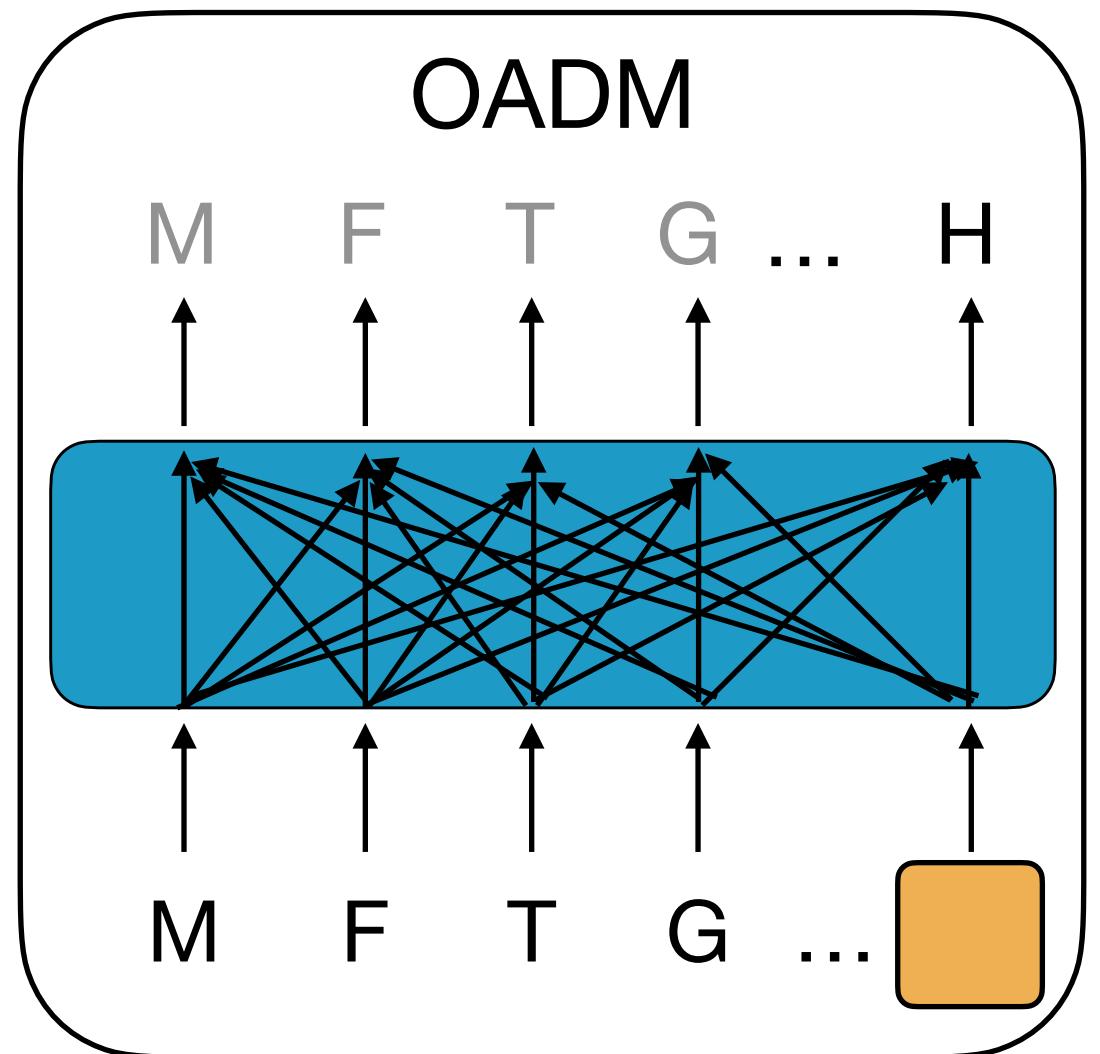
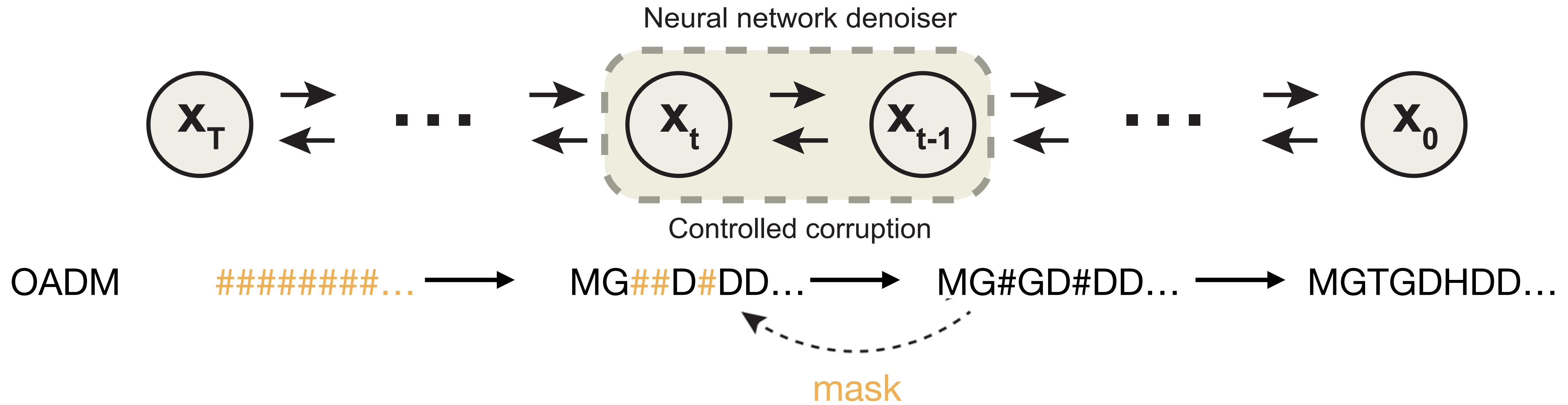
EvoDiff: evolutionary-scale diffusion



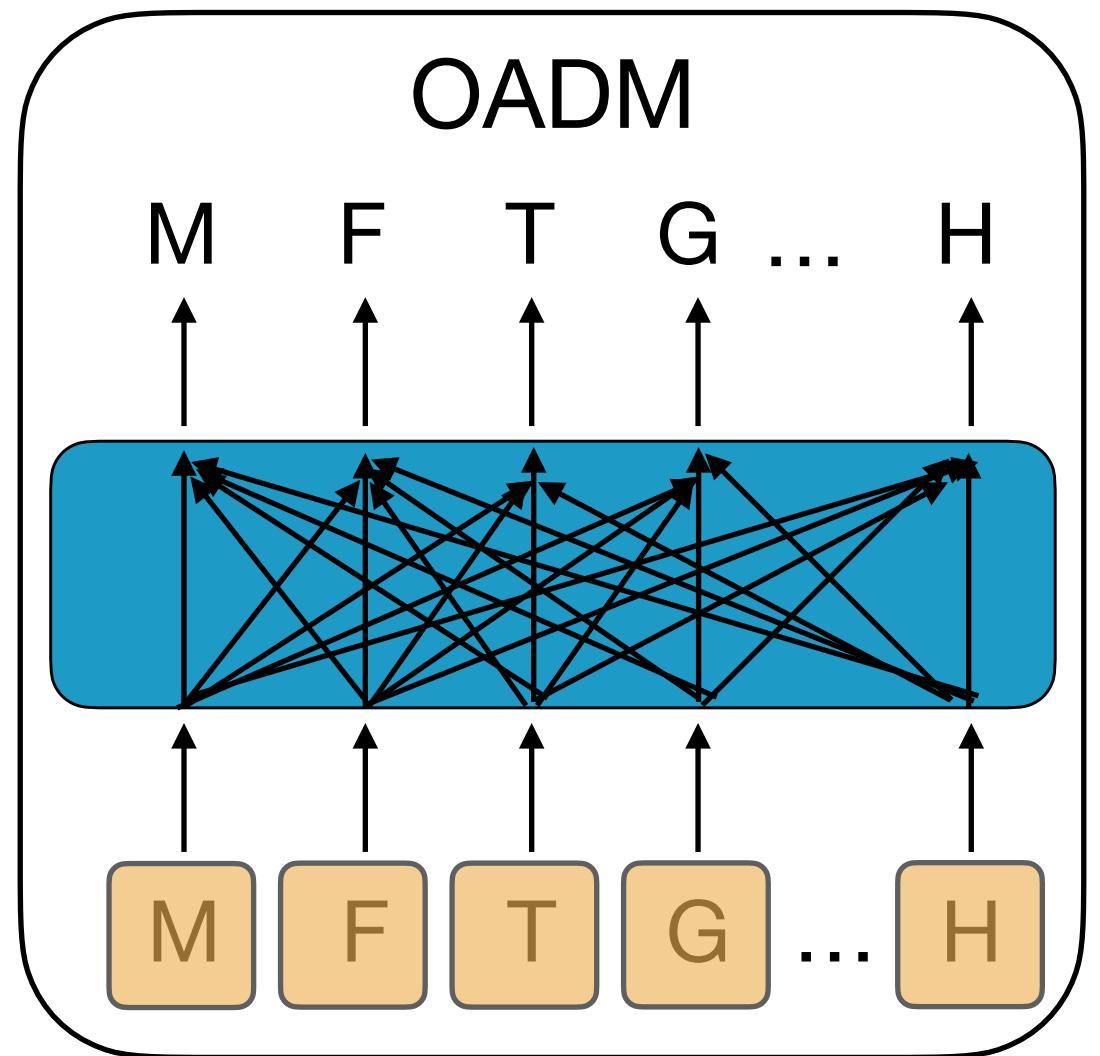
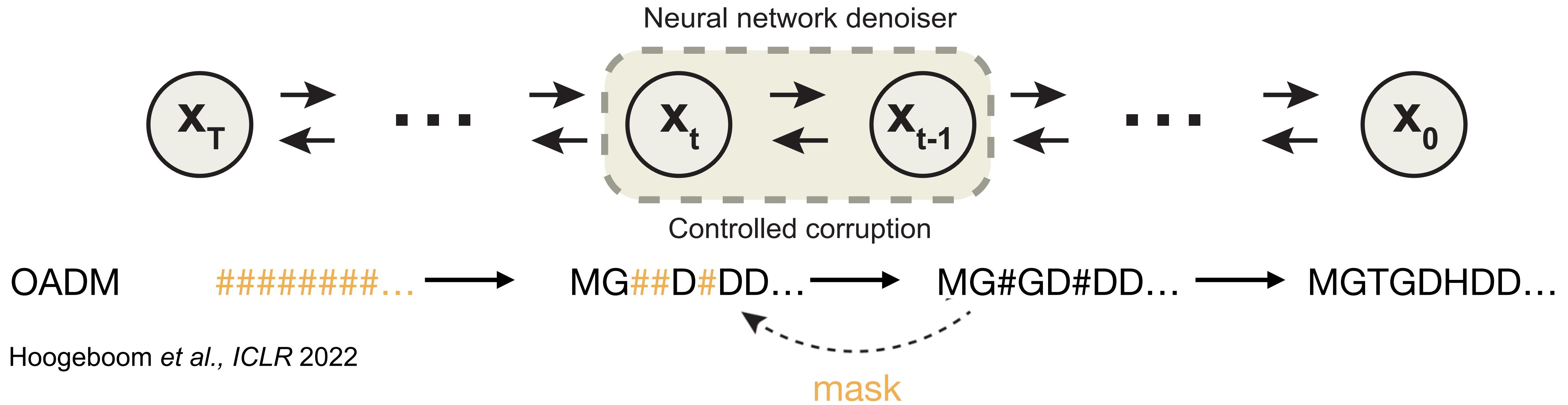
EvoDiff: evolutionary-scale diffusion



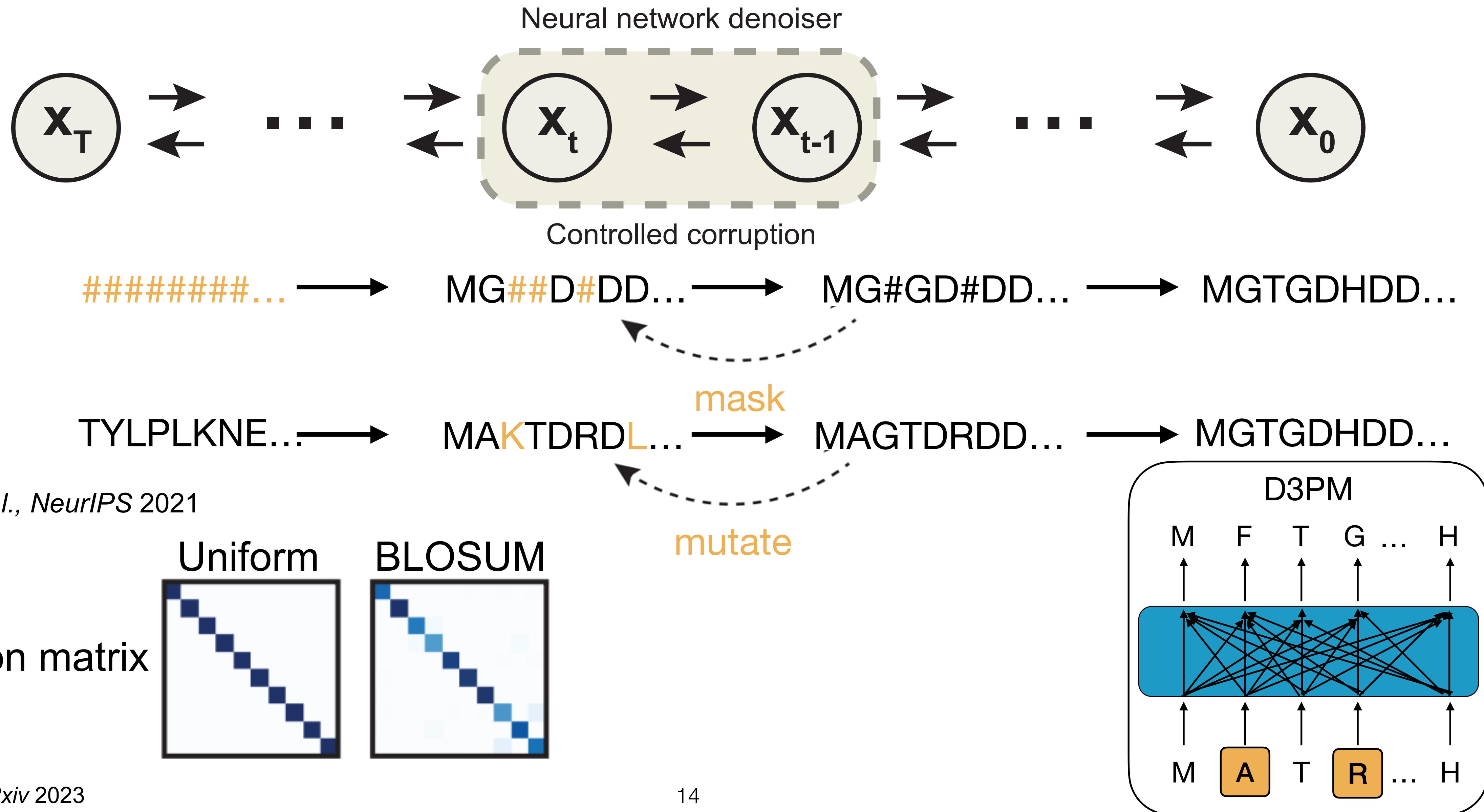
EvoDiff: evolutionary-scale diffusion



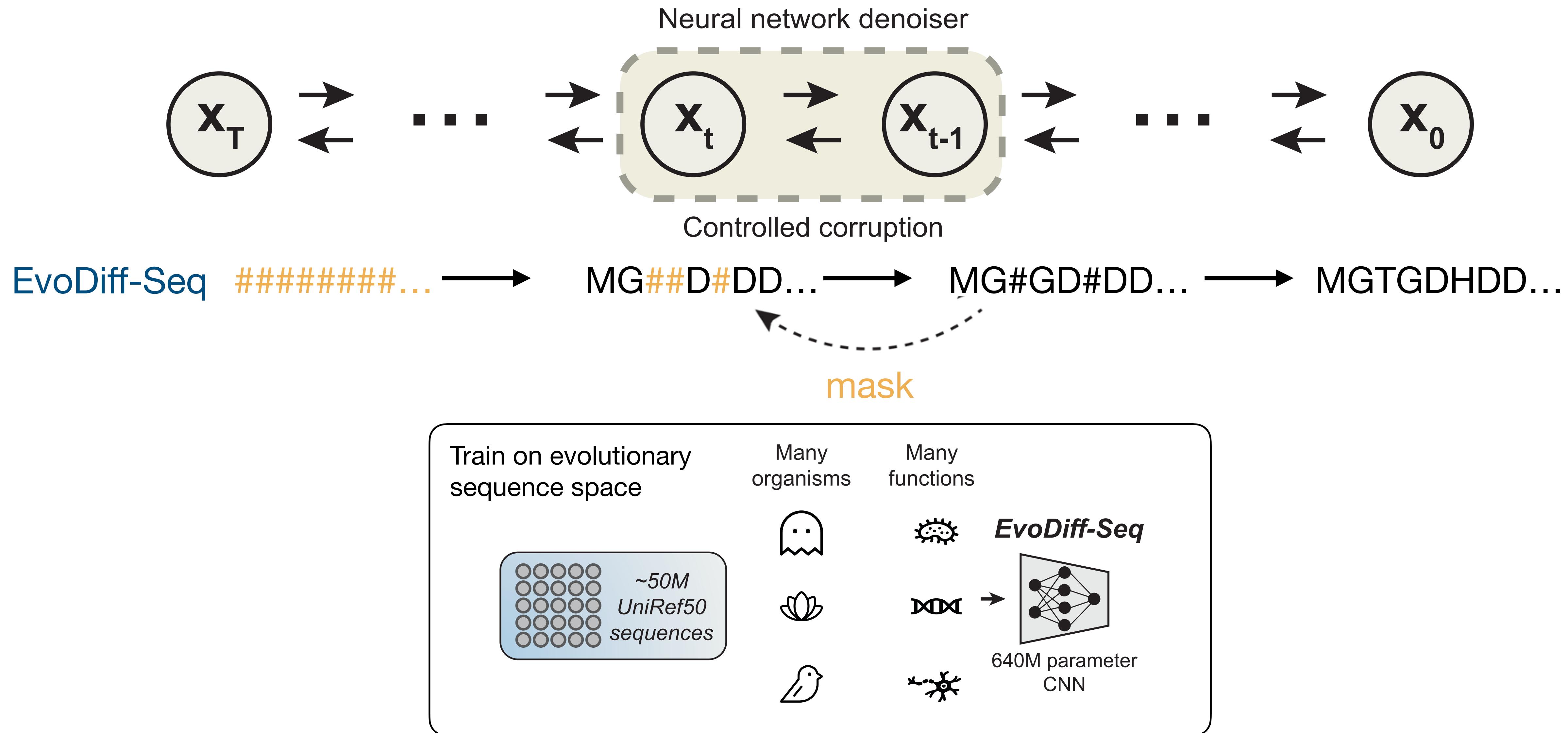
EvoDiff: evolutionary-scale diffusion



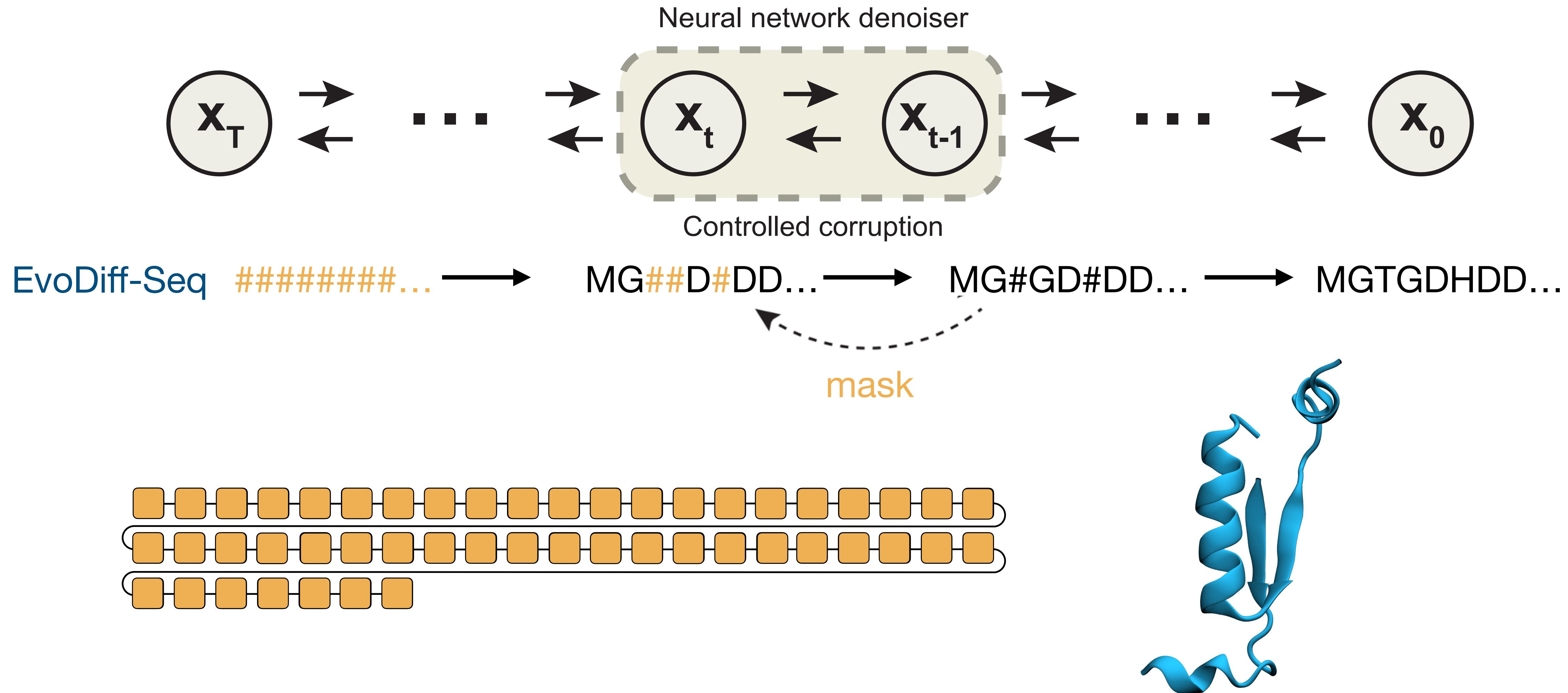
EvoDiff: evolutionary-scale diffusion



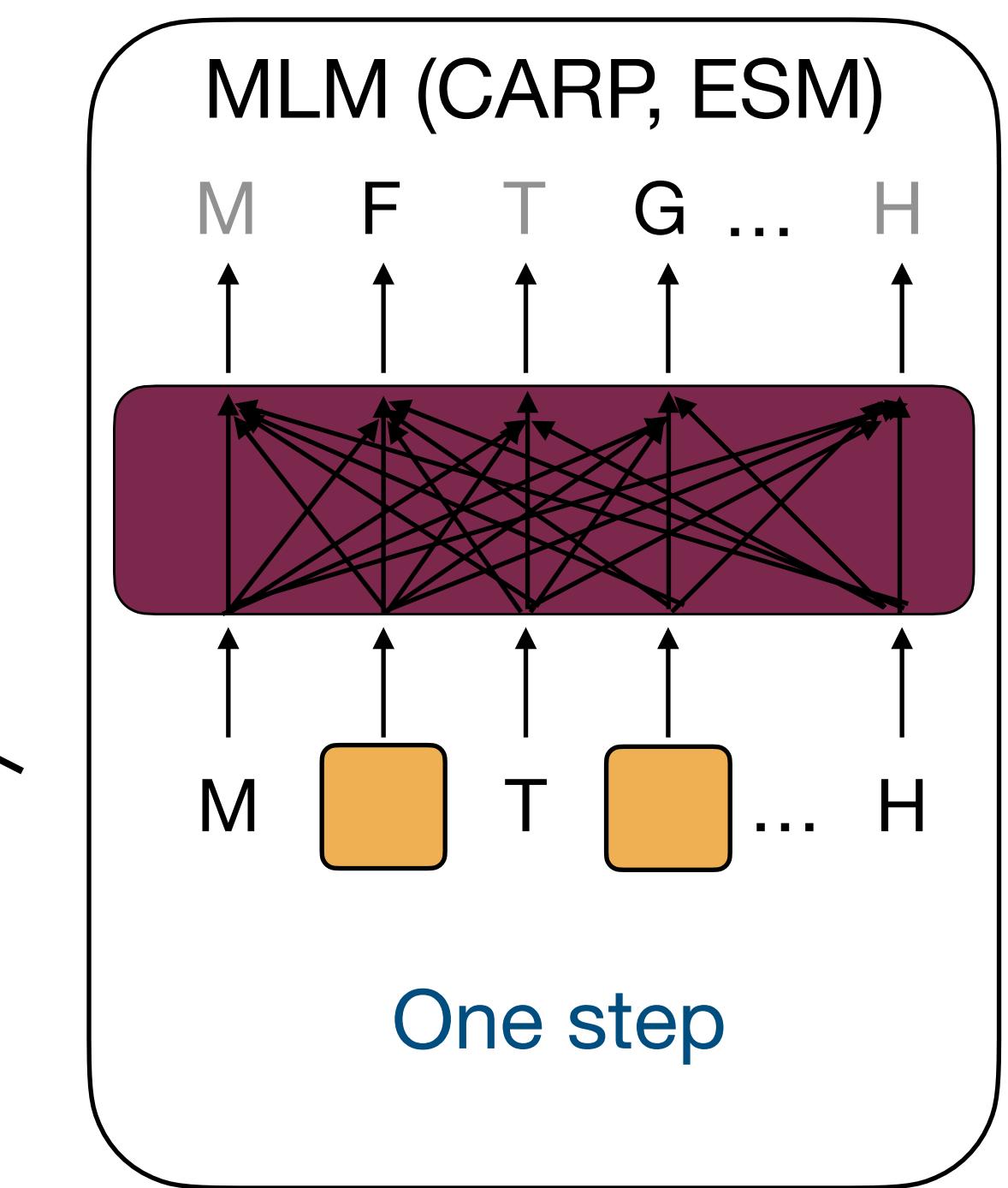
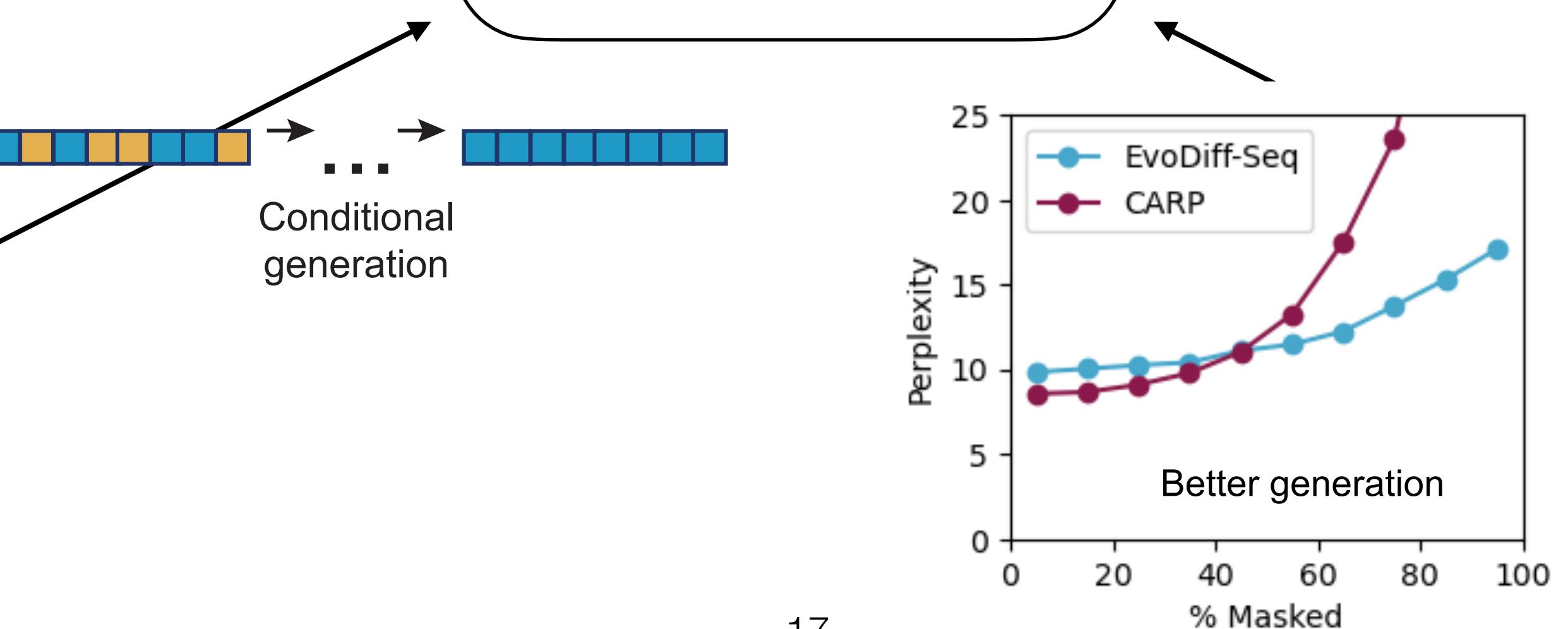
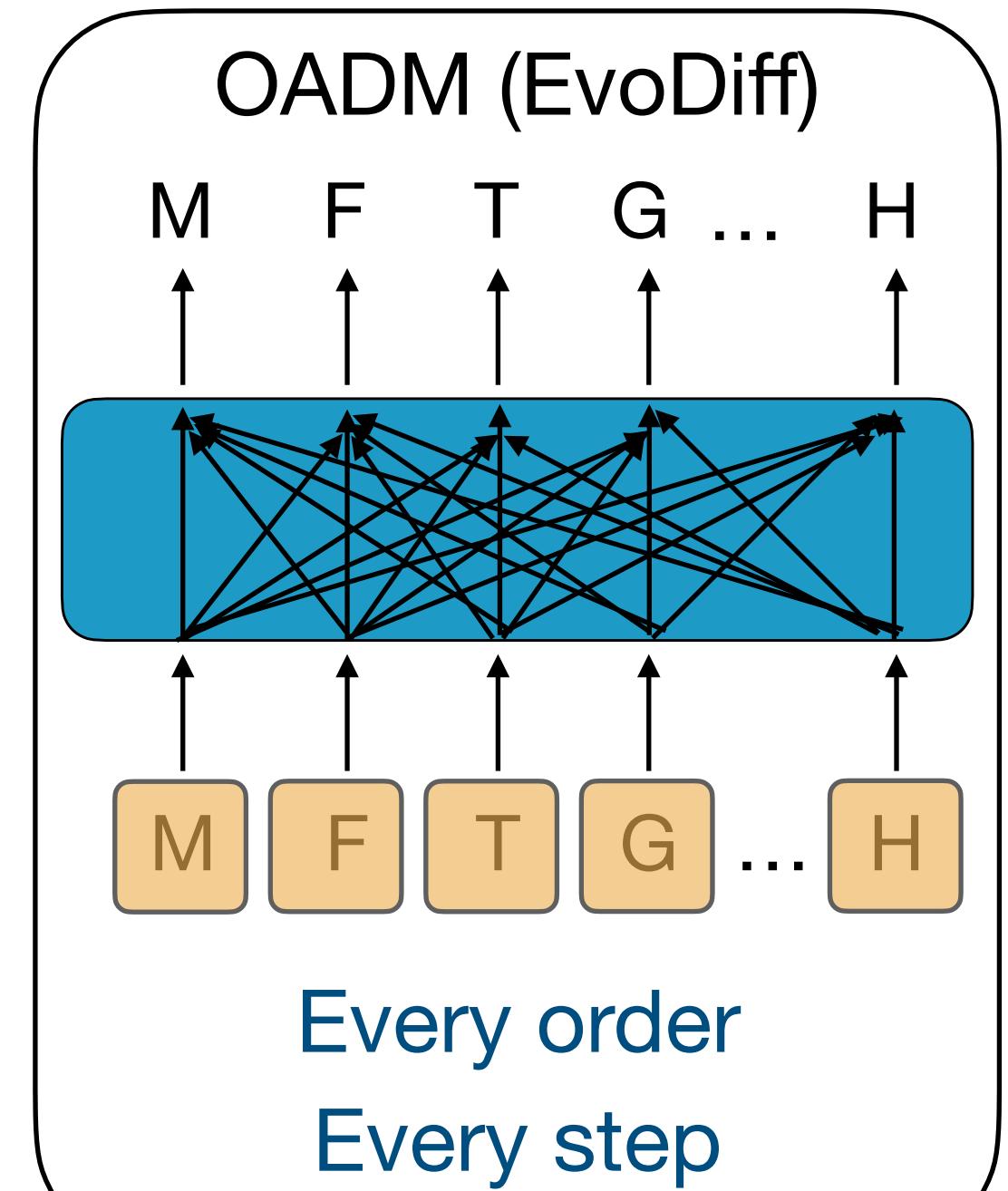
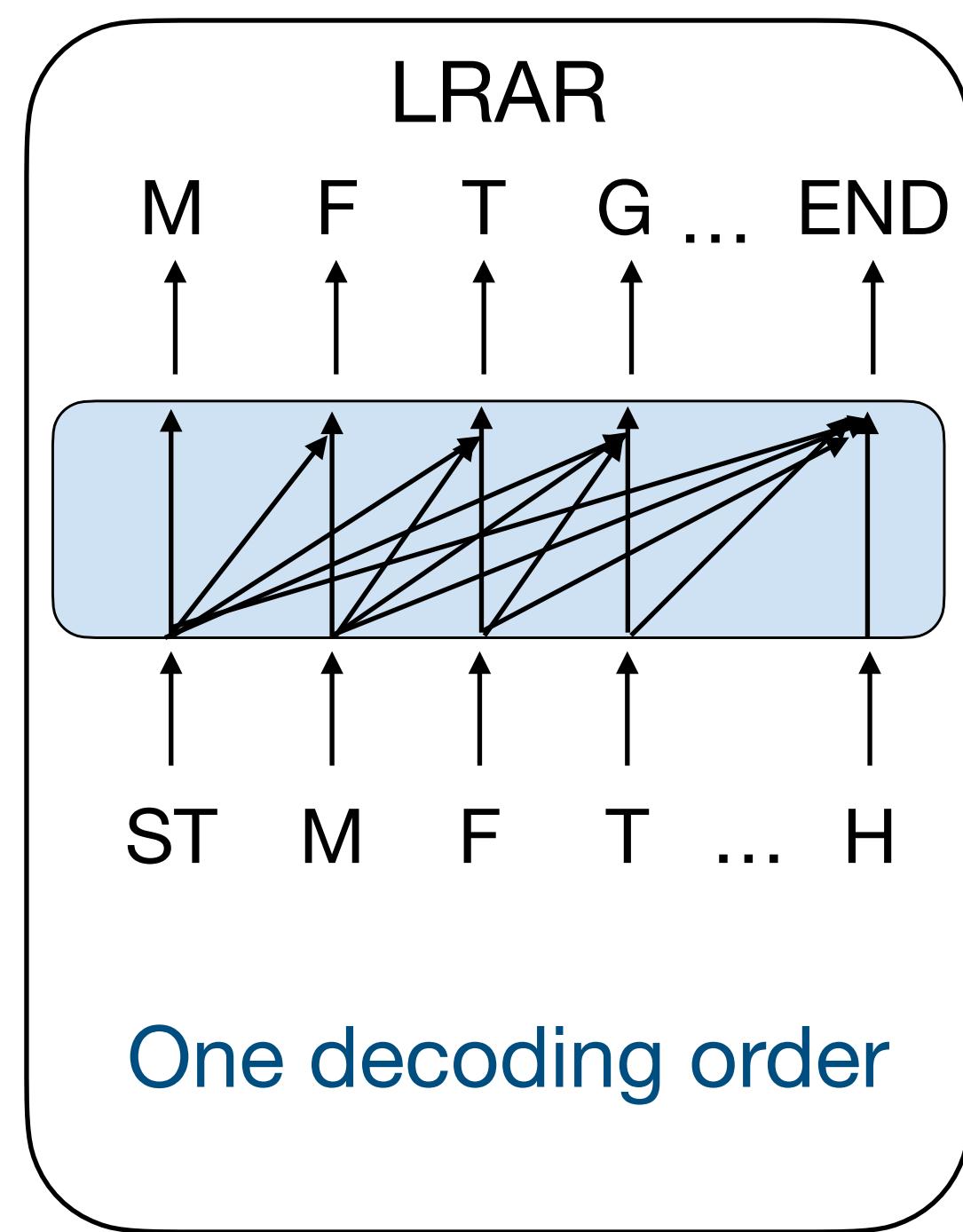
EvoDiff: evolutionary-scale diffusion



EvoDiff: evolutionary-scale diffusion

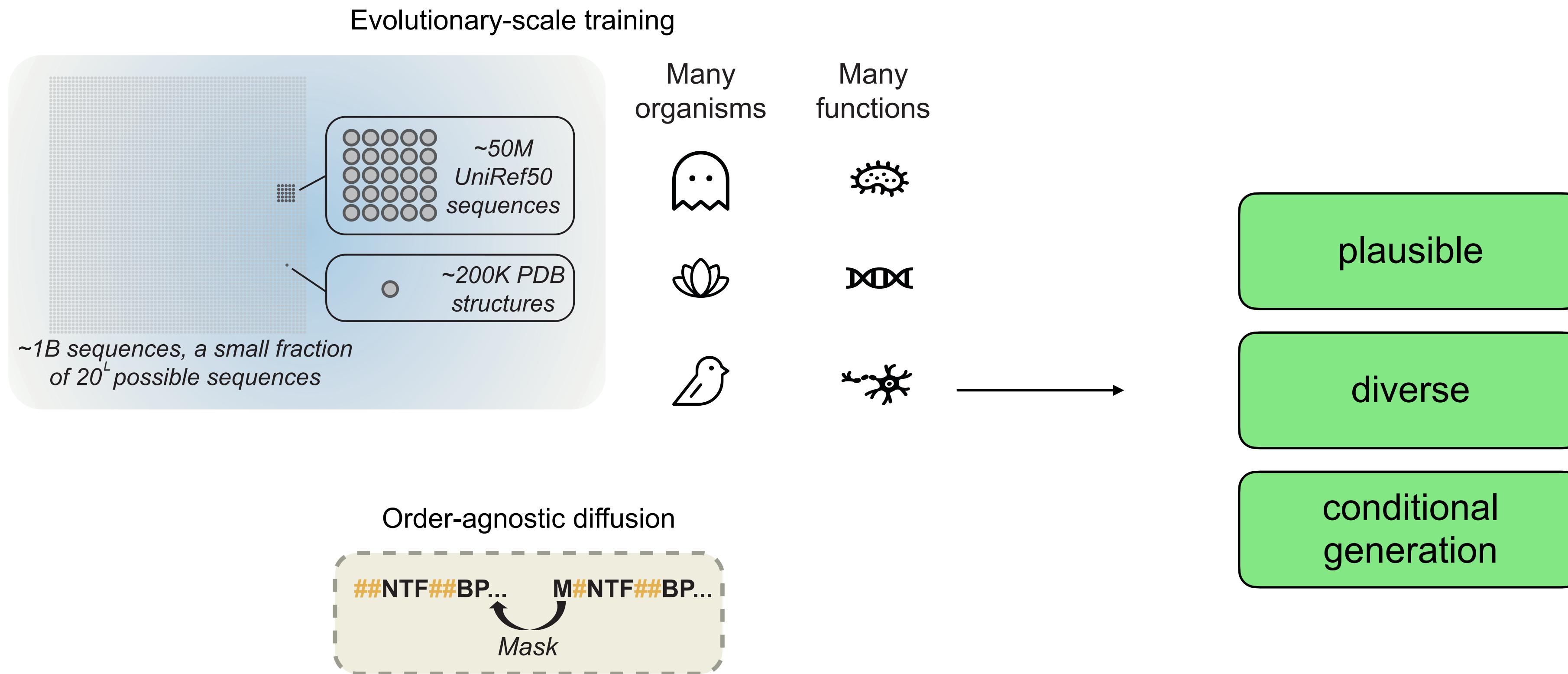


EvoDiff-Seq generalizes masked and autoregressive language models

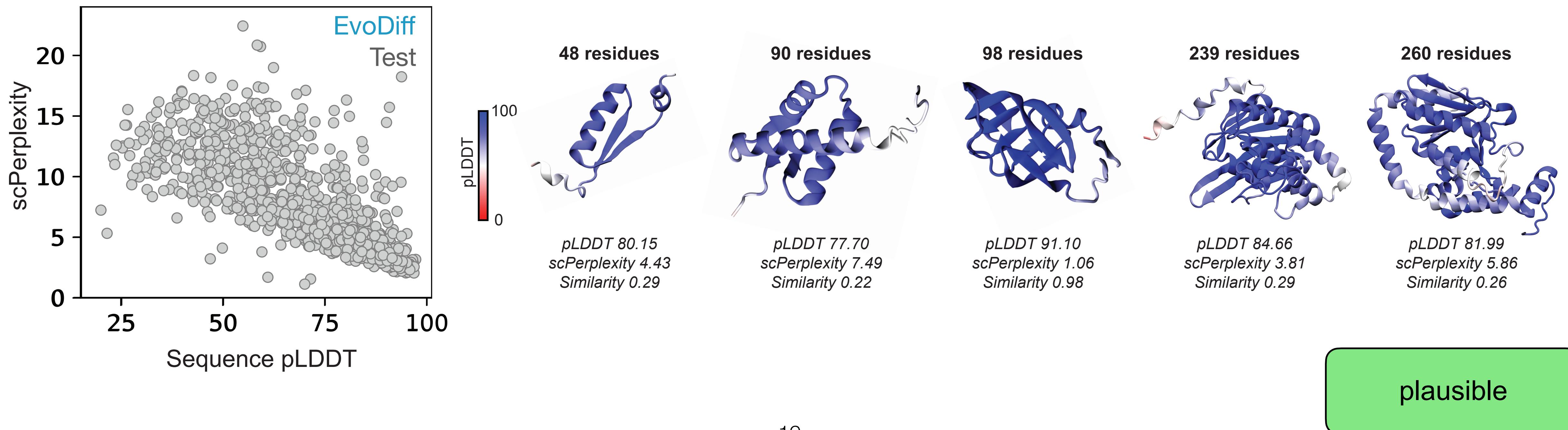
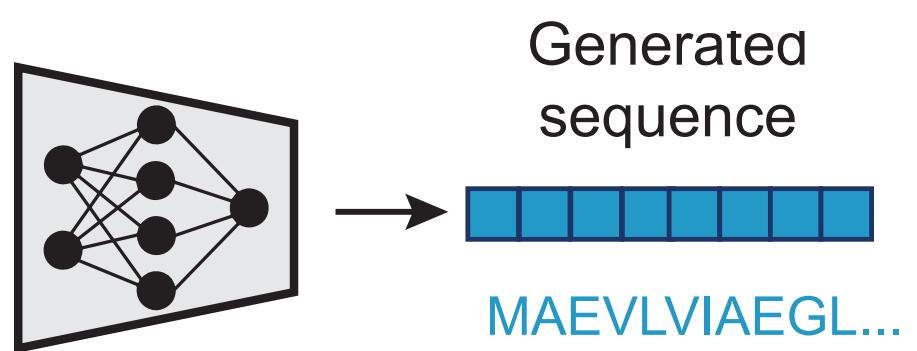


CARP: Yang et al., preprint, 2022
ESM: Rives et al., PNAS, 2021

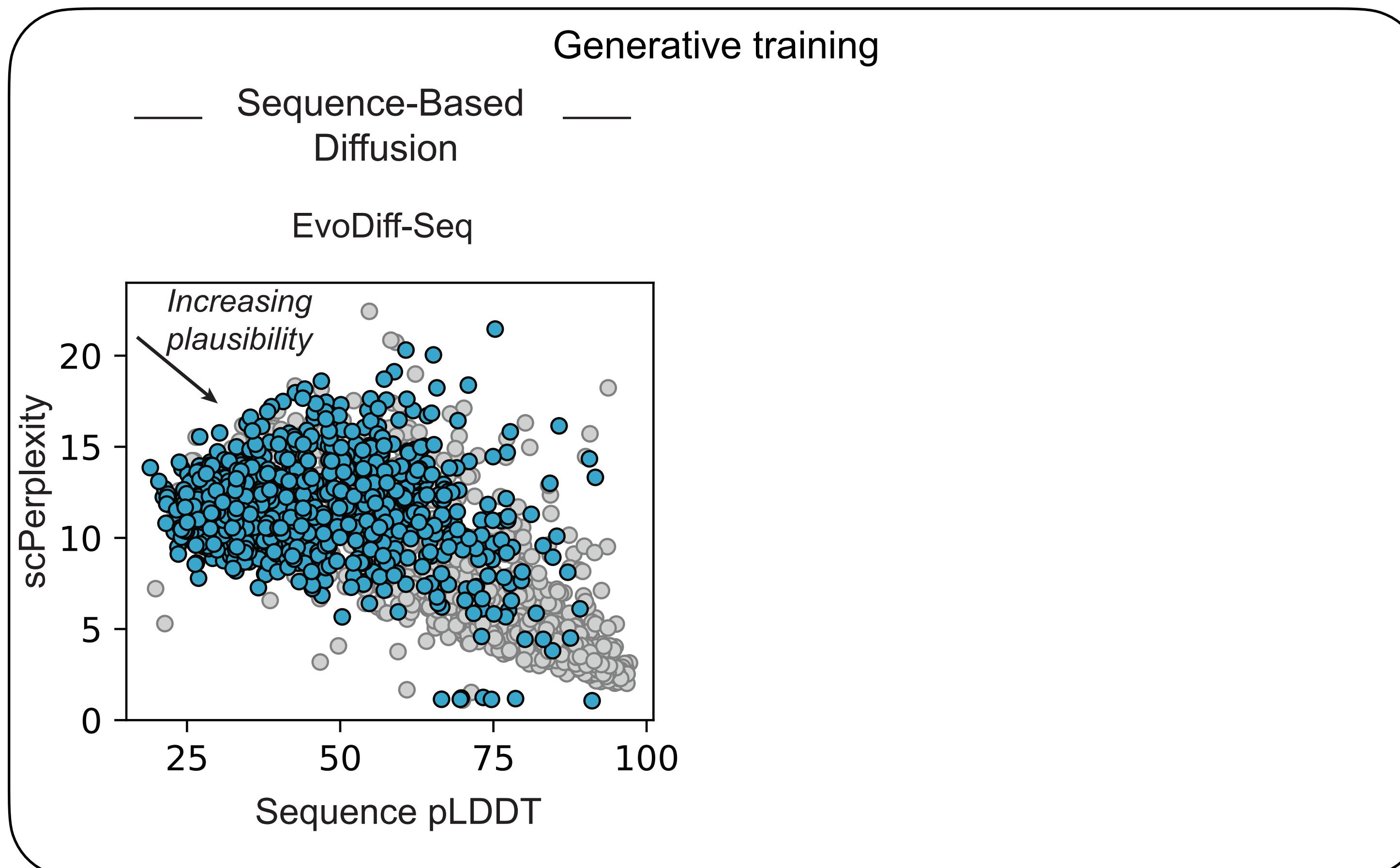
EvoDiff enables controllable generation of plausible, diverse proteins



EvoDiff-Seq generates highly-plausible proteins

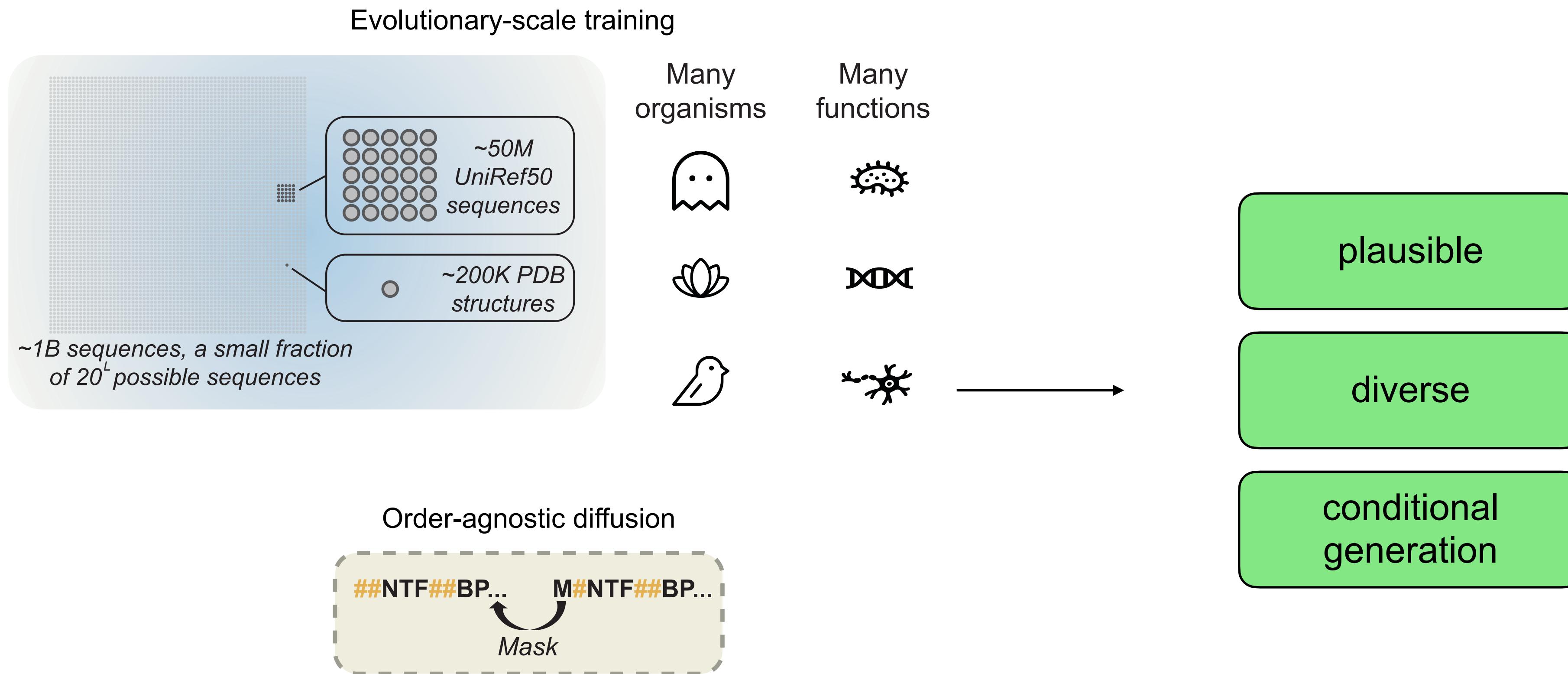


Generative training results in better sequences

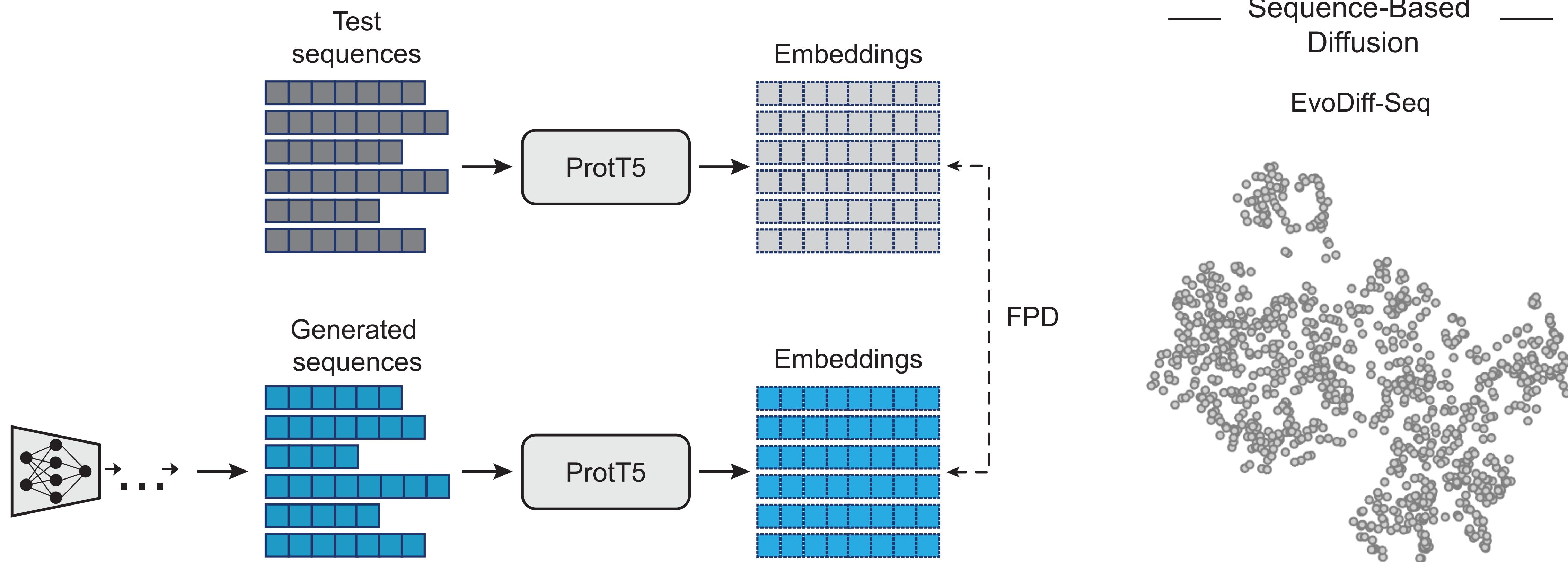


plausible

EvoDiff enables controllable generation of plausible, diverse proteins



EvoDiff-Seq recapitulates natural functional distribution

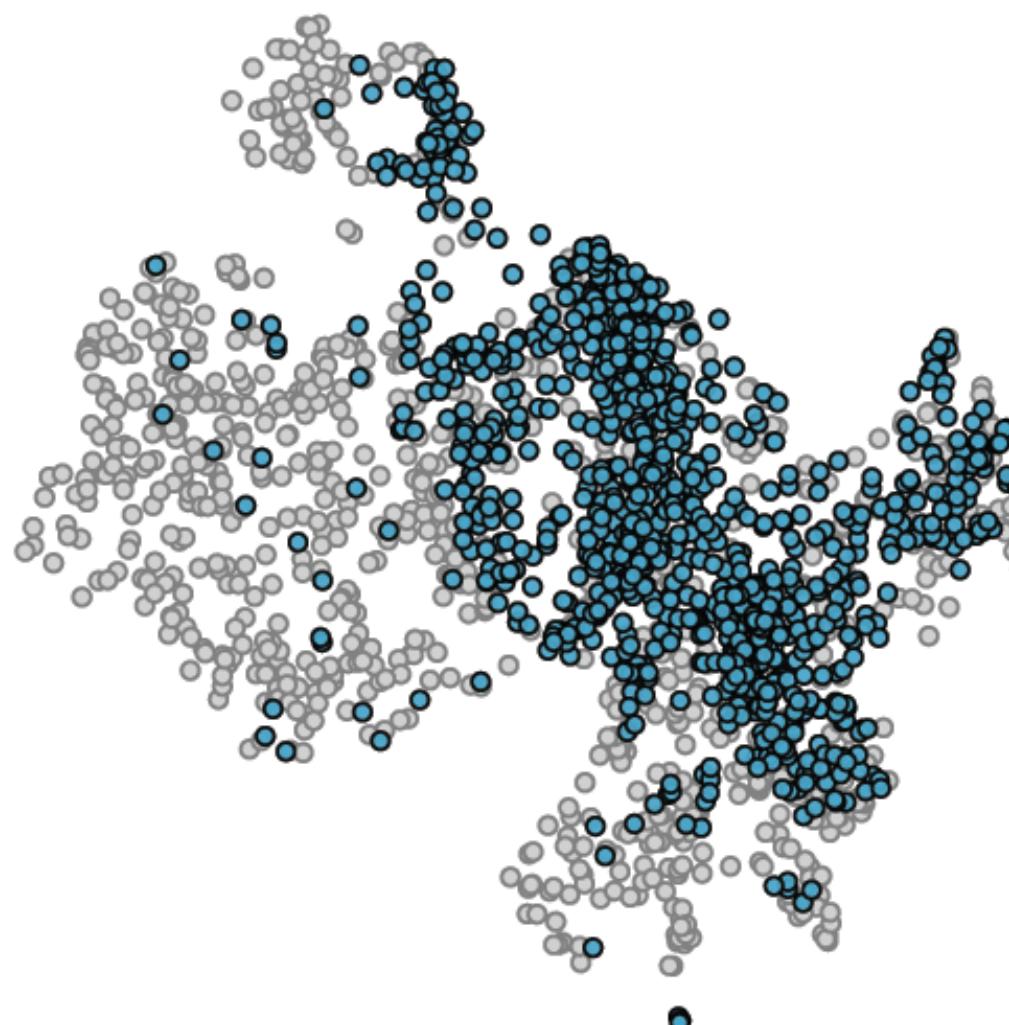


diverse

Evolutionary-scale diffusion improves FPD

Sequence-Based Diffusion

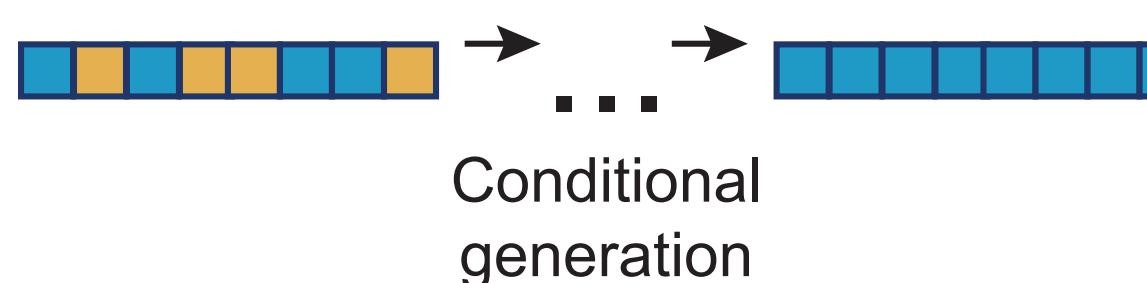
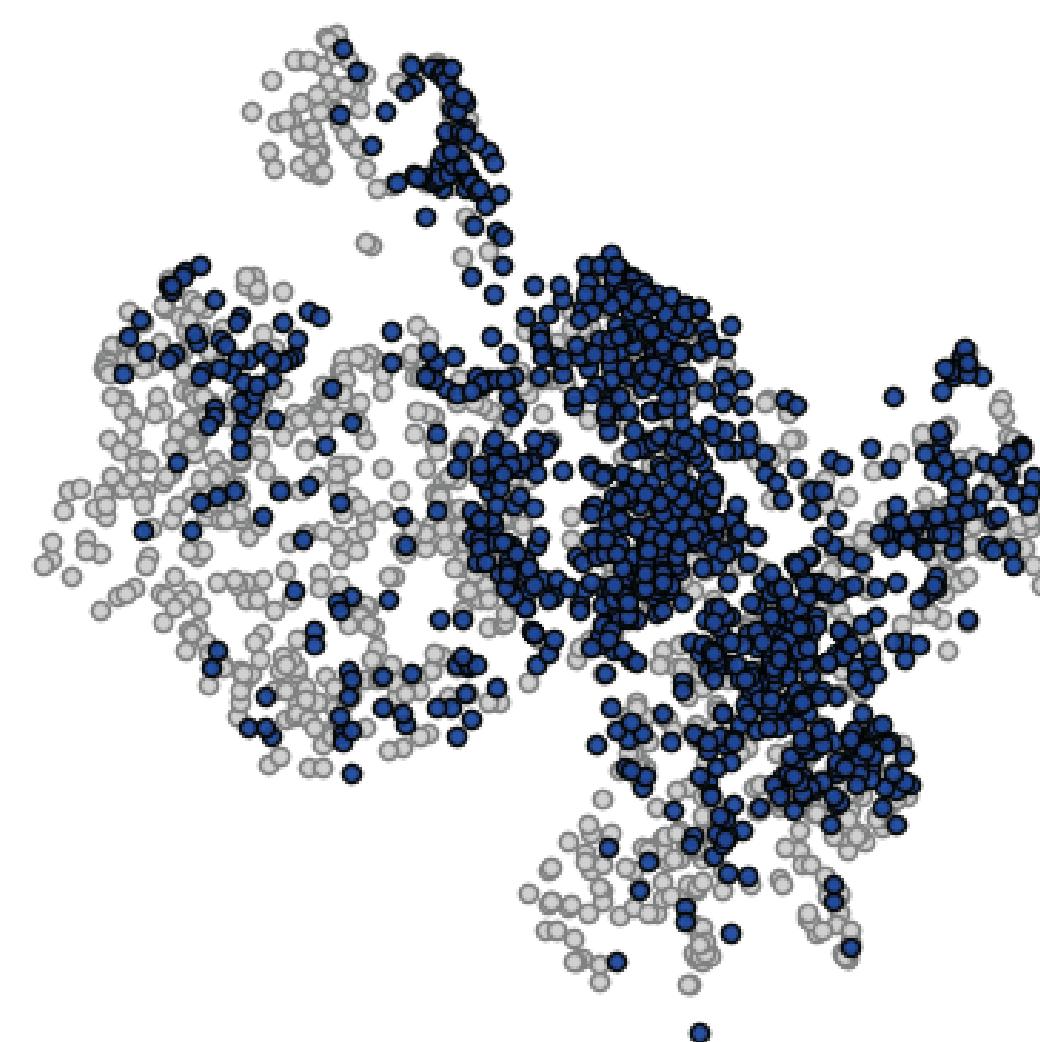
EvoDiff-Seq
FPD = 0.88



diverse

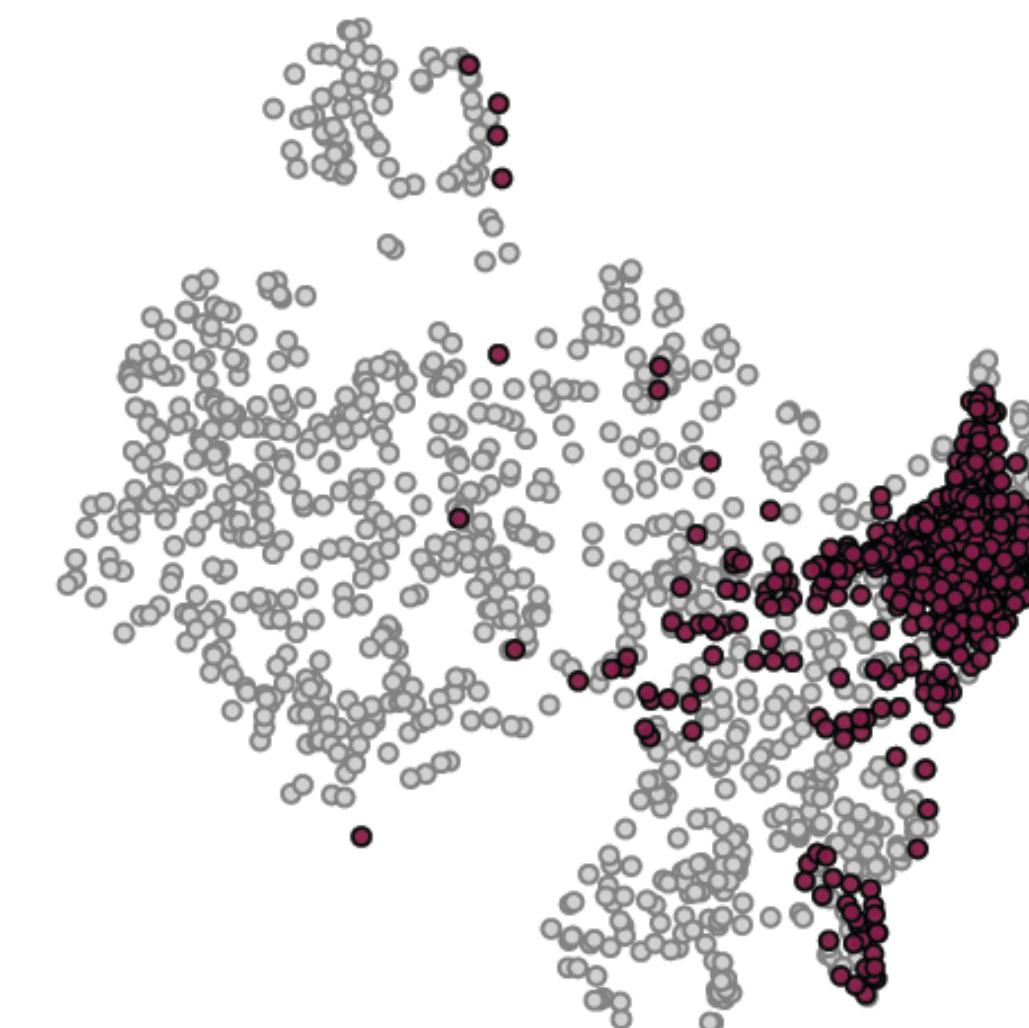
Left-to-right Language Model

LRAR
FPD = 0.63



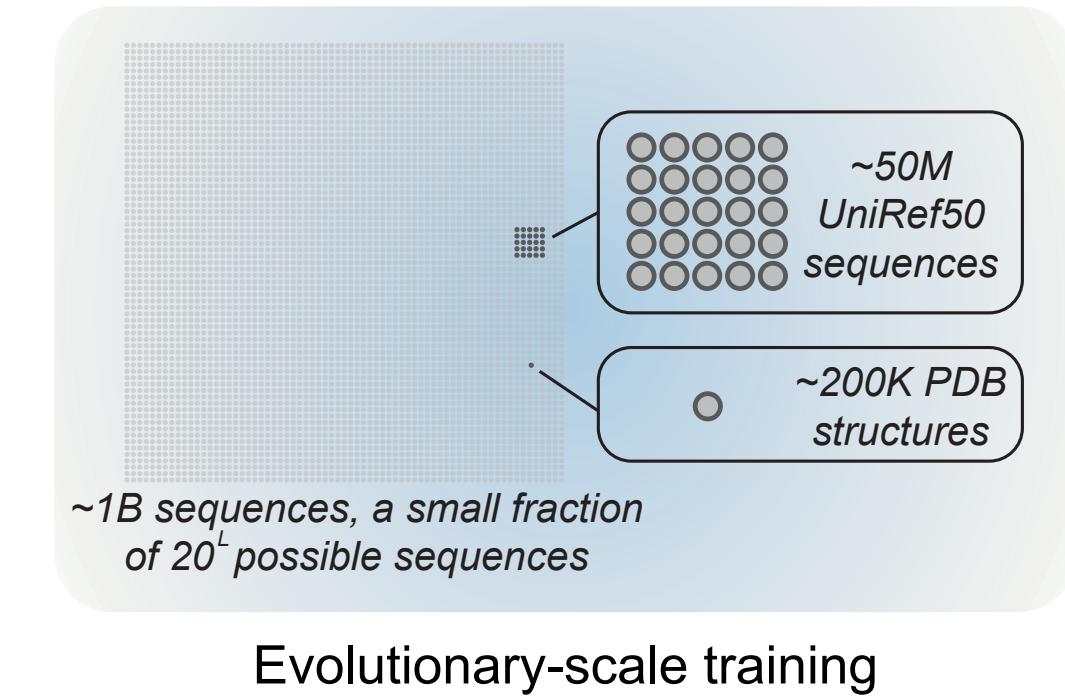
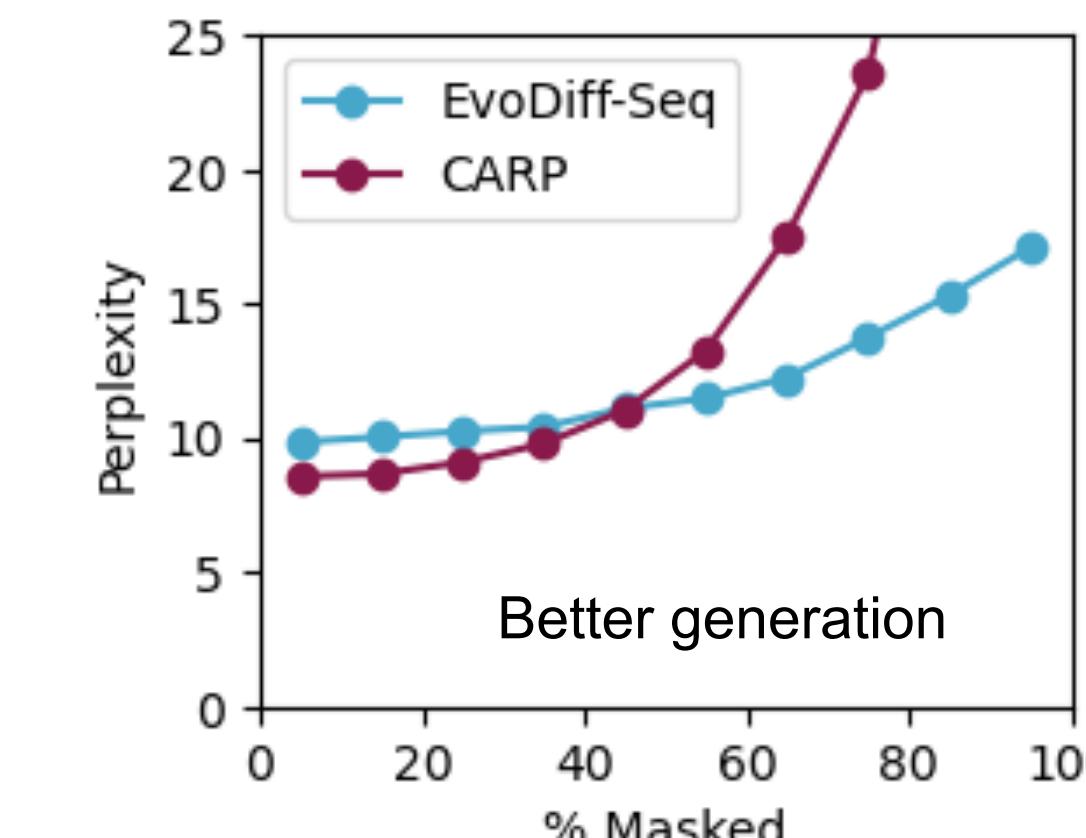
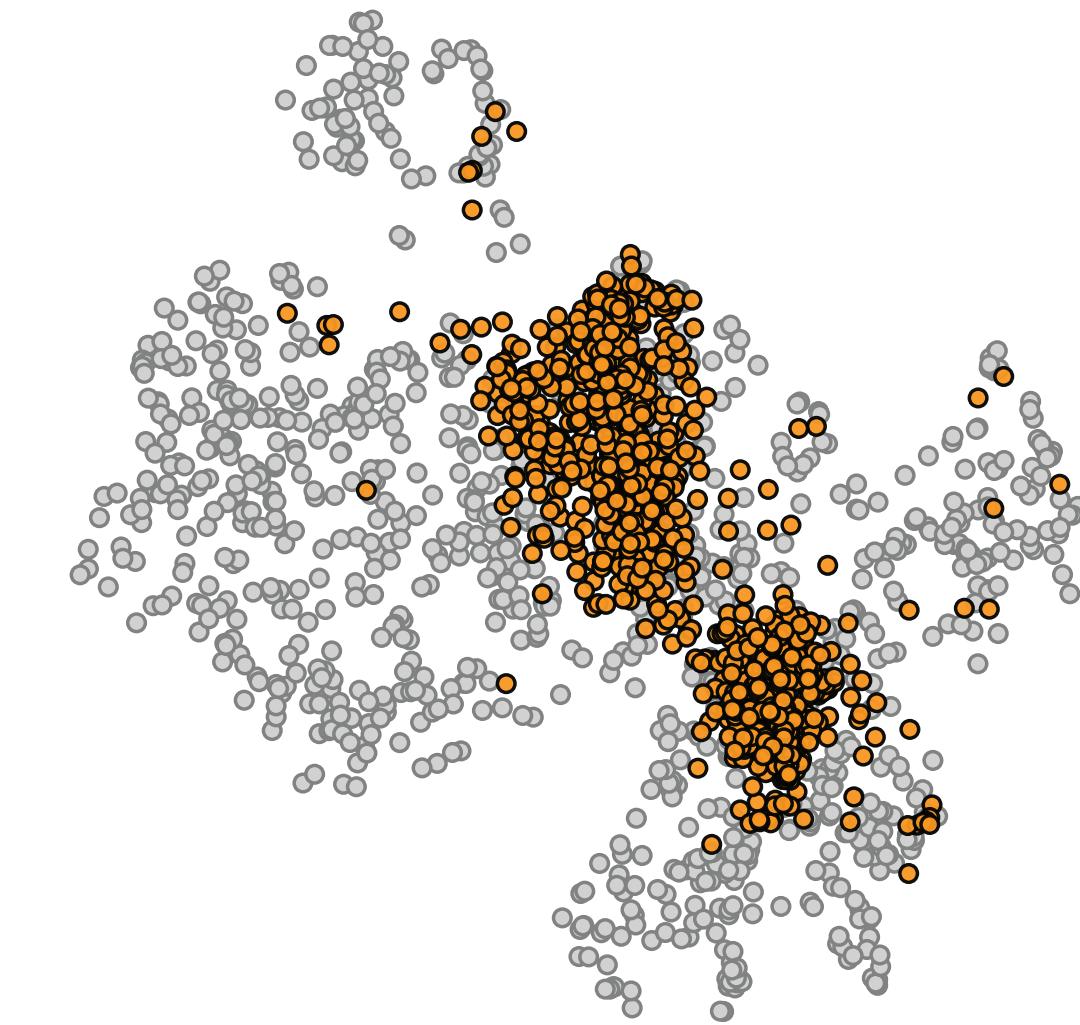
Protein Masked Language Model

ESM-2
FPD = 2.81

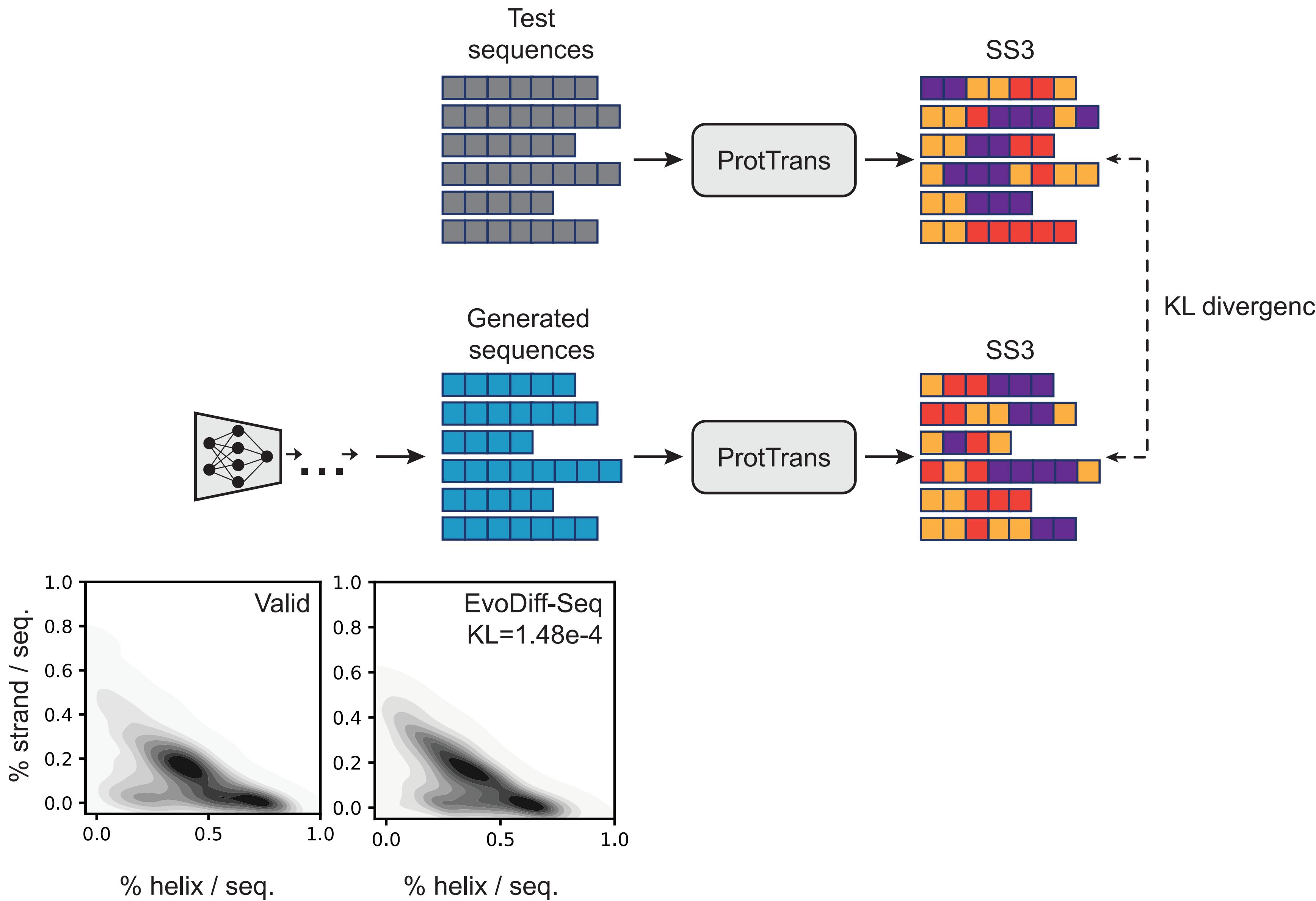


Structure-Based Diffusion

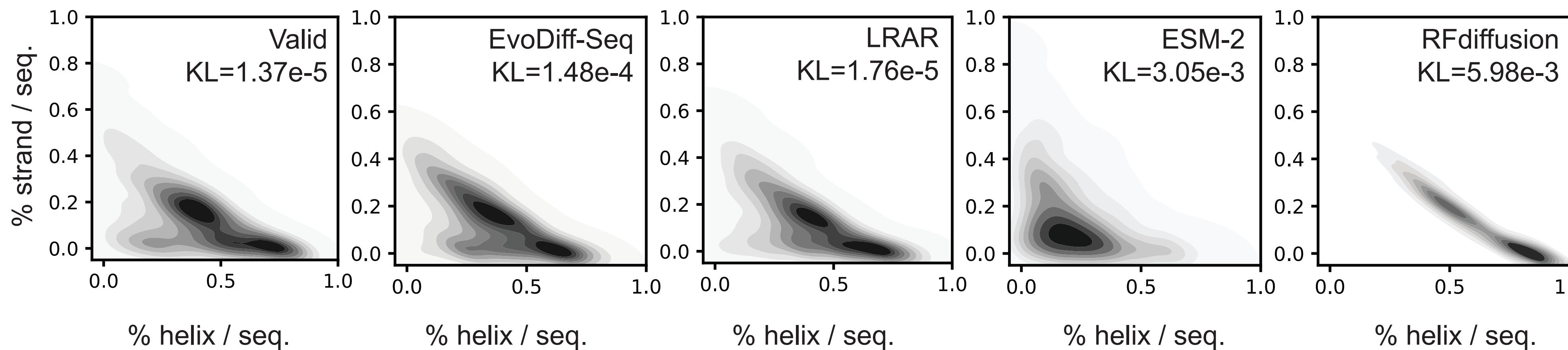
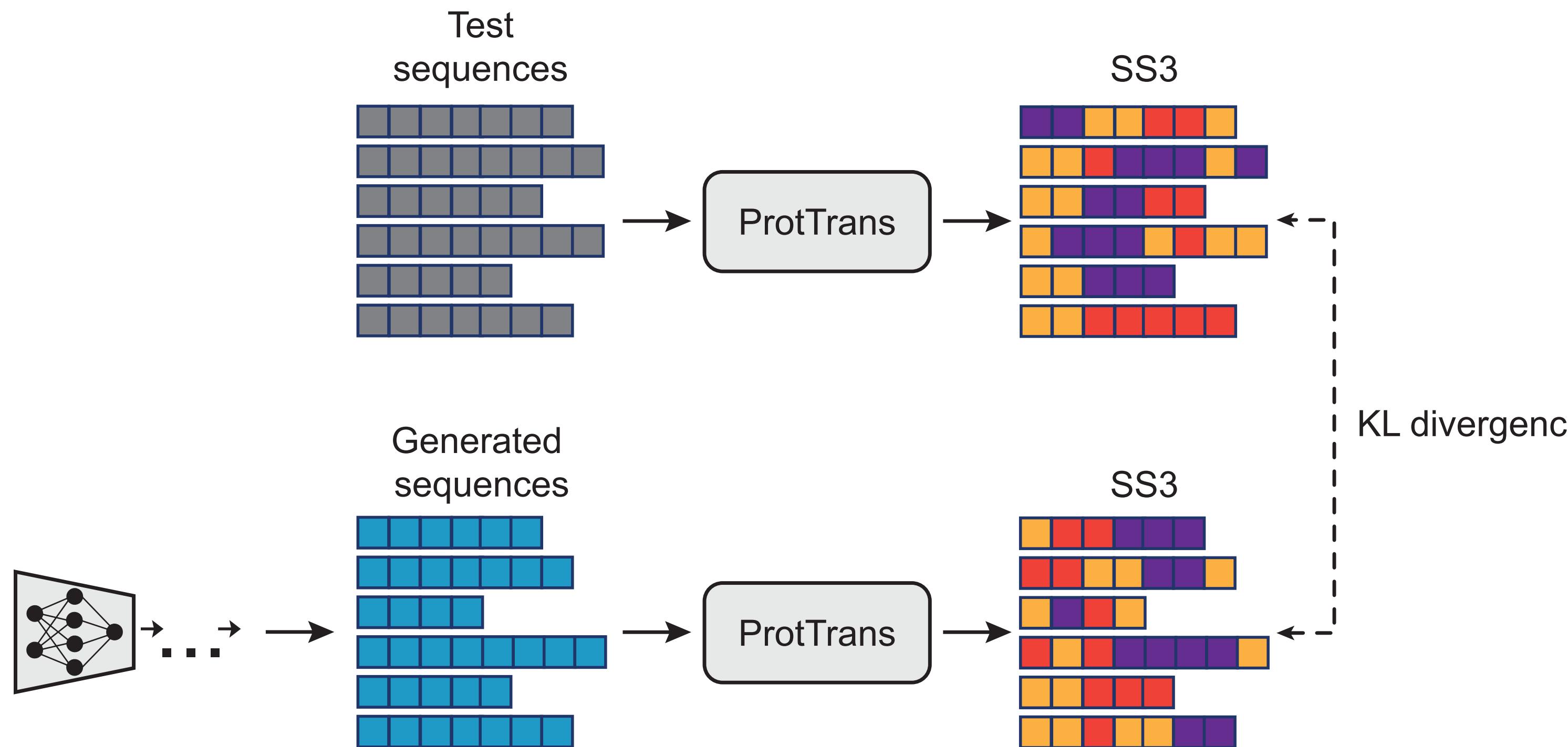
RFdiffusion
FPD = 1.96



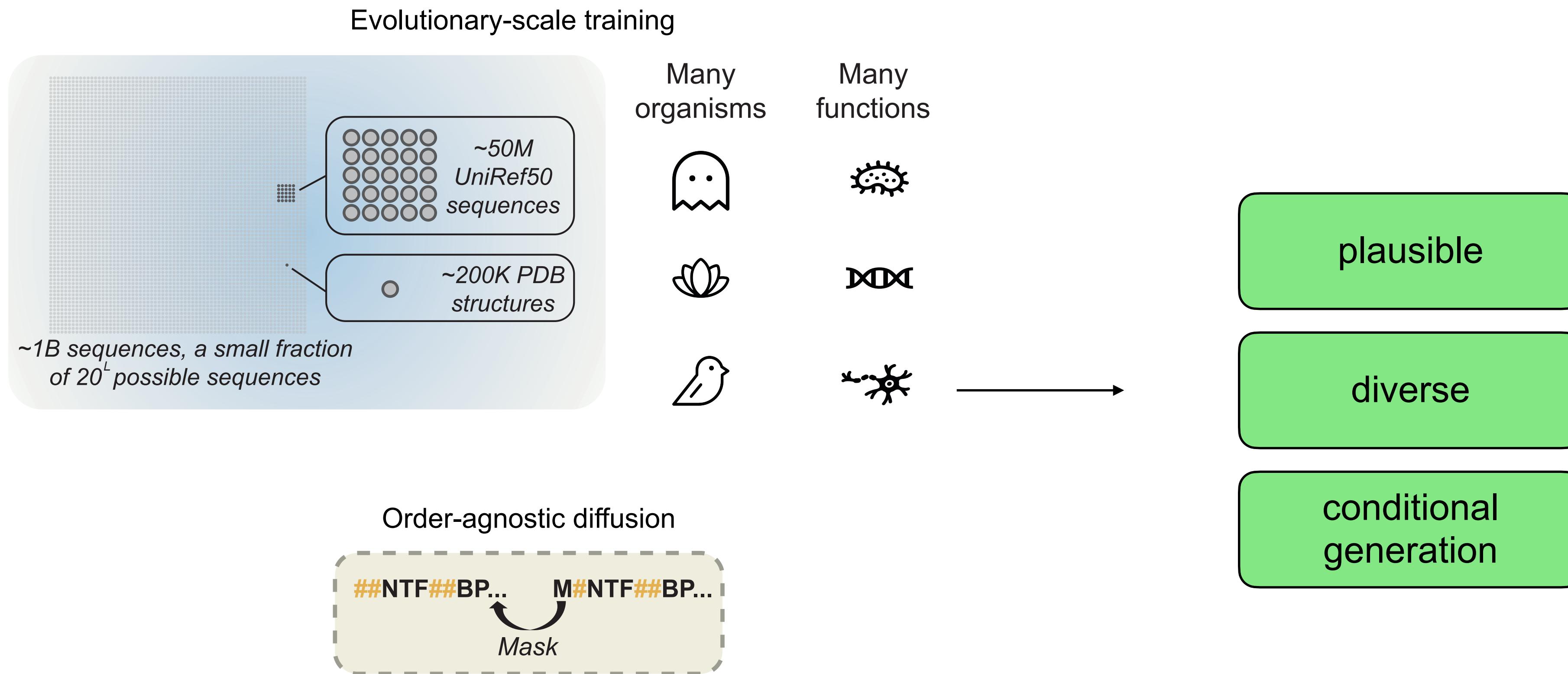
EvoDiff-Seq recapitulates natural secondary structure distribution



Evolutionary-scale diffusion improves secondary structures



EvoDiff enables controllable generation of plausible, diverse proteins



EvoDiff-MSA leverages evolutionary information

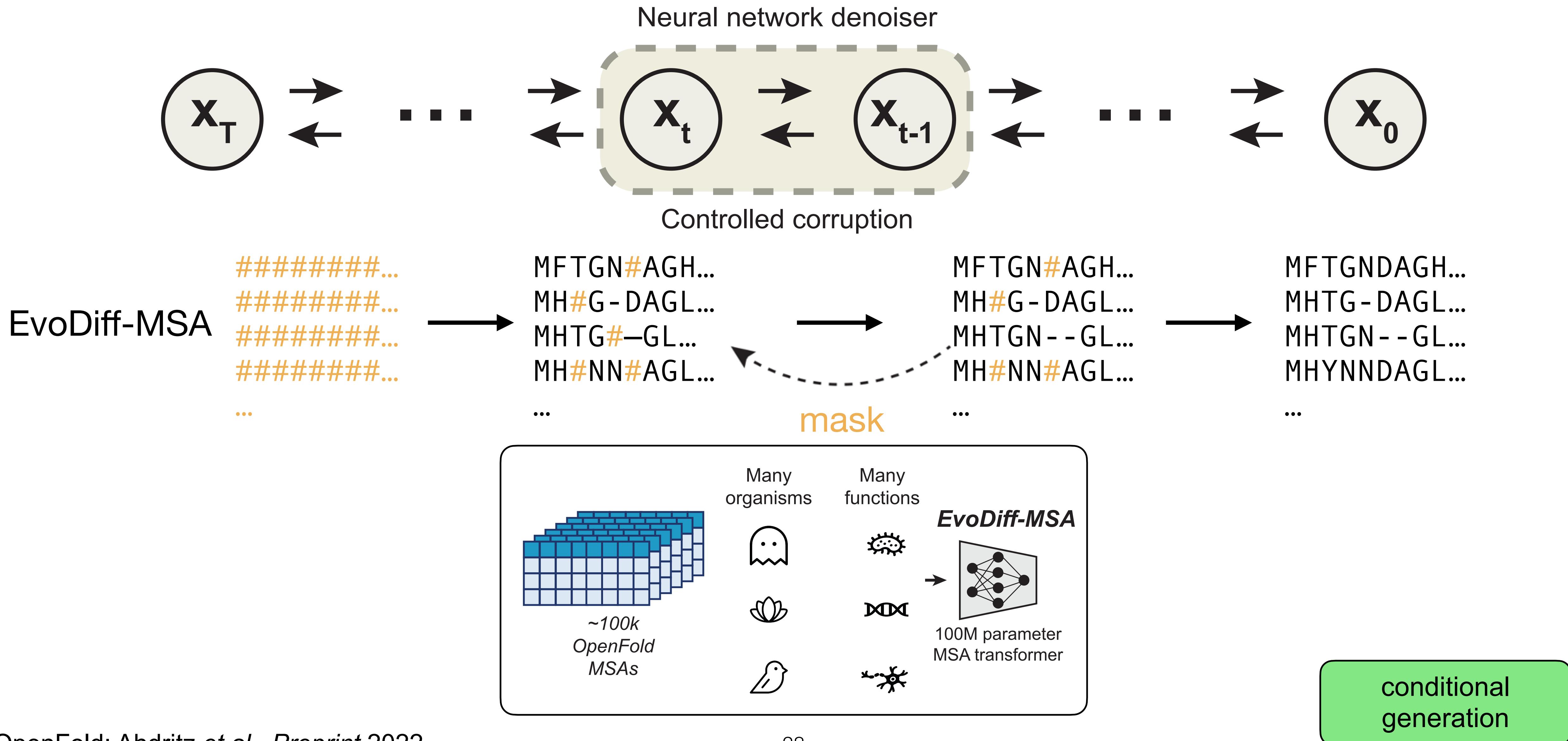
query MFTGNDAGH...

homology
search

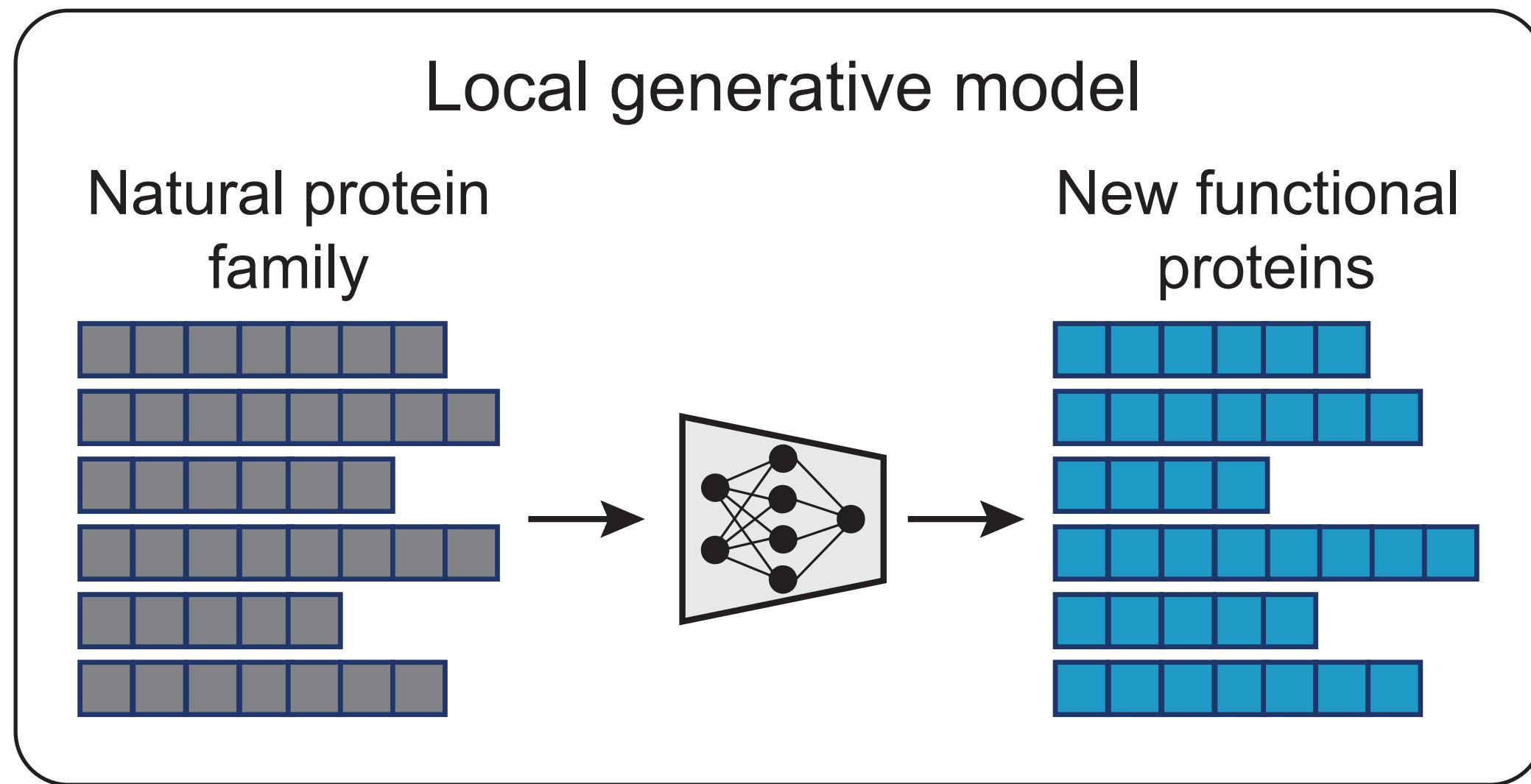
subsample to 64

conditional
generation

EvoDiff-MSA leverages evolutionary information

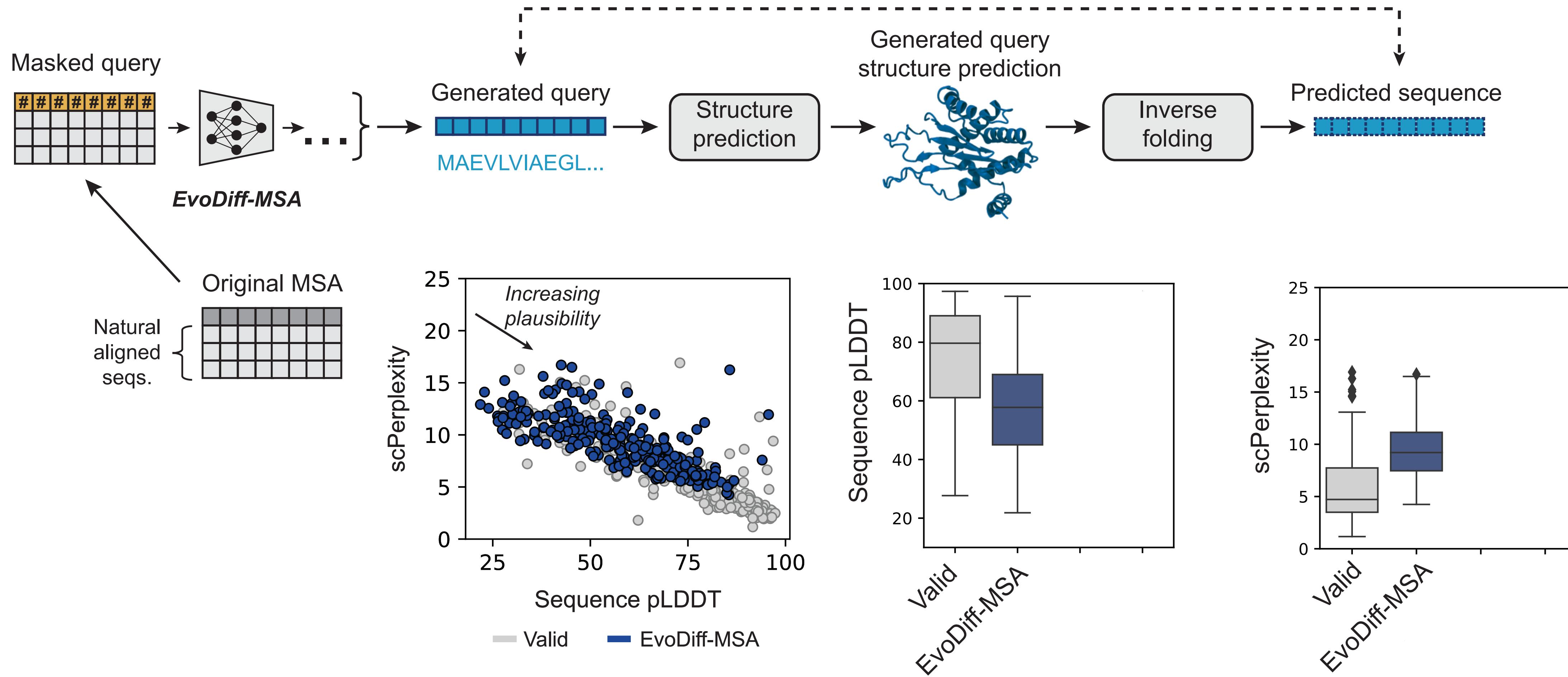


EvoDiff-MSA is an amortized local generative model



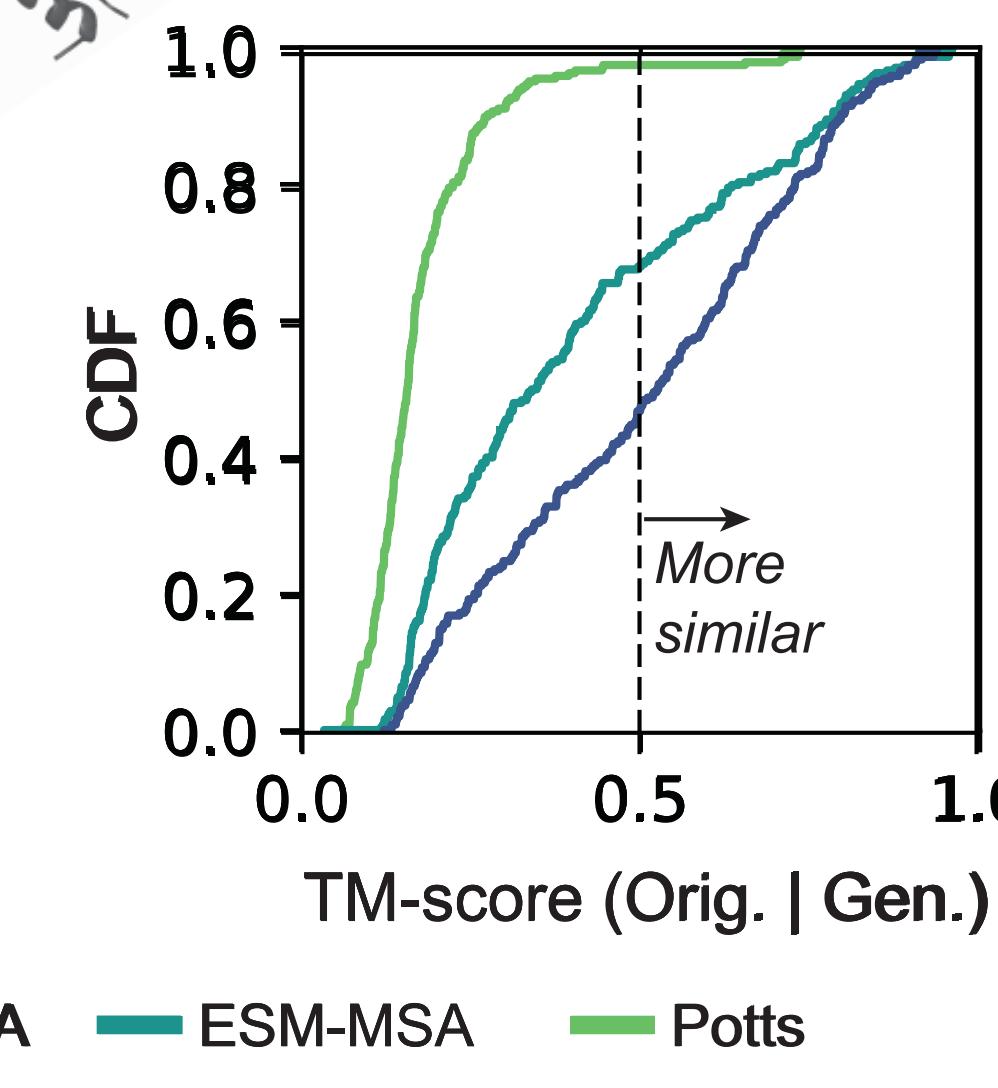
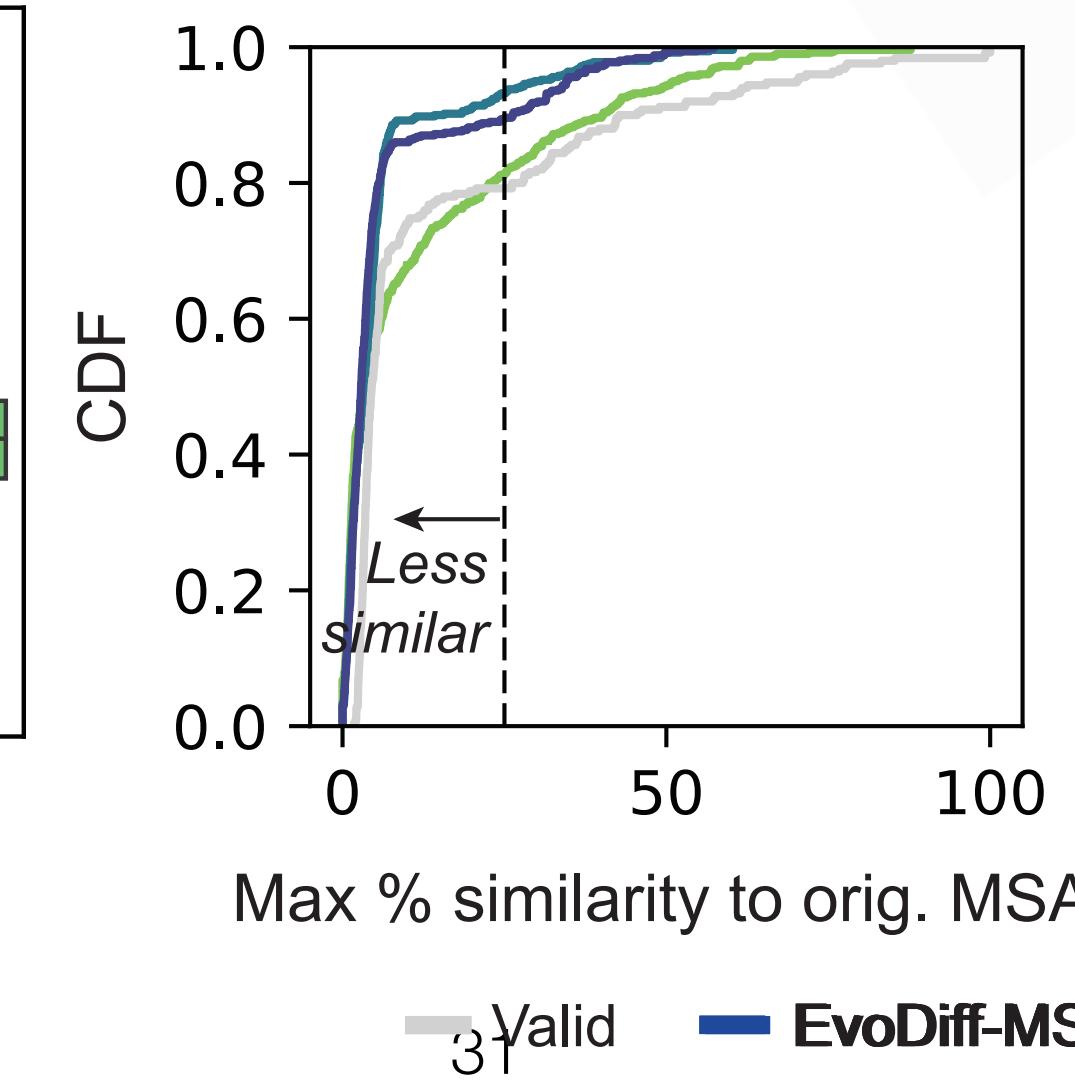
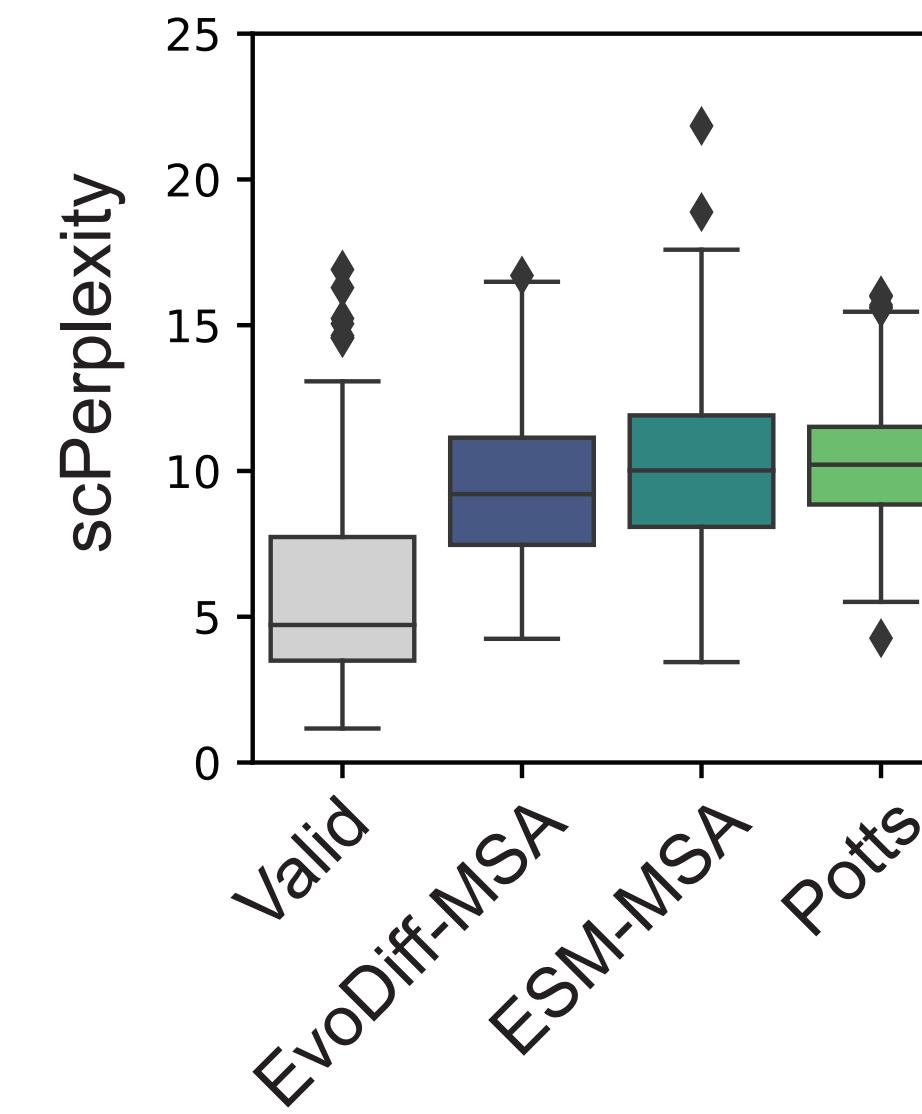
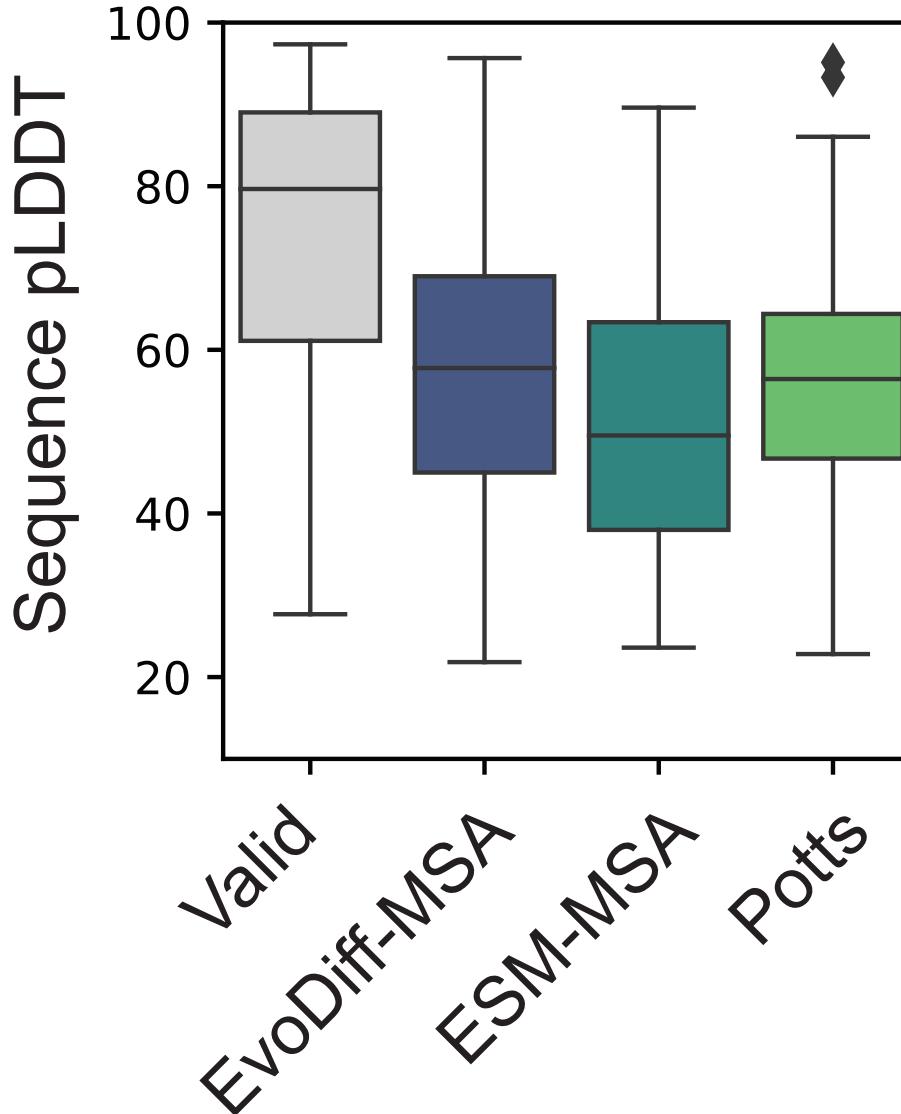
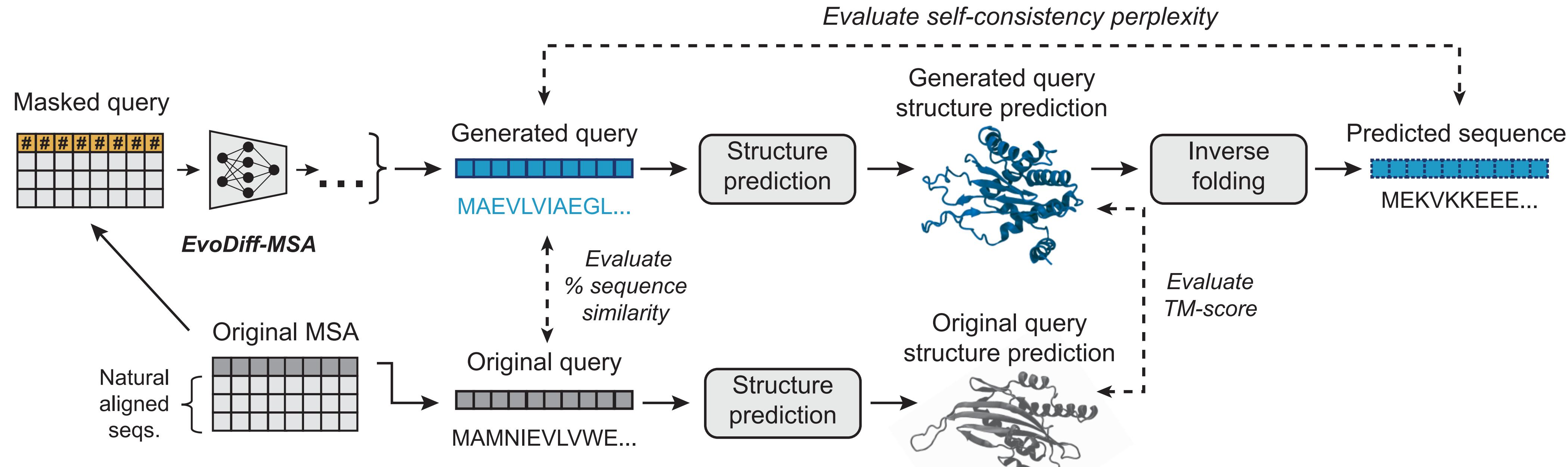
conditional
generation

Evolutionary diffusion improves local generation



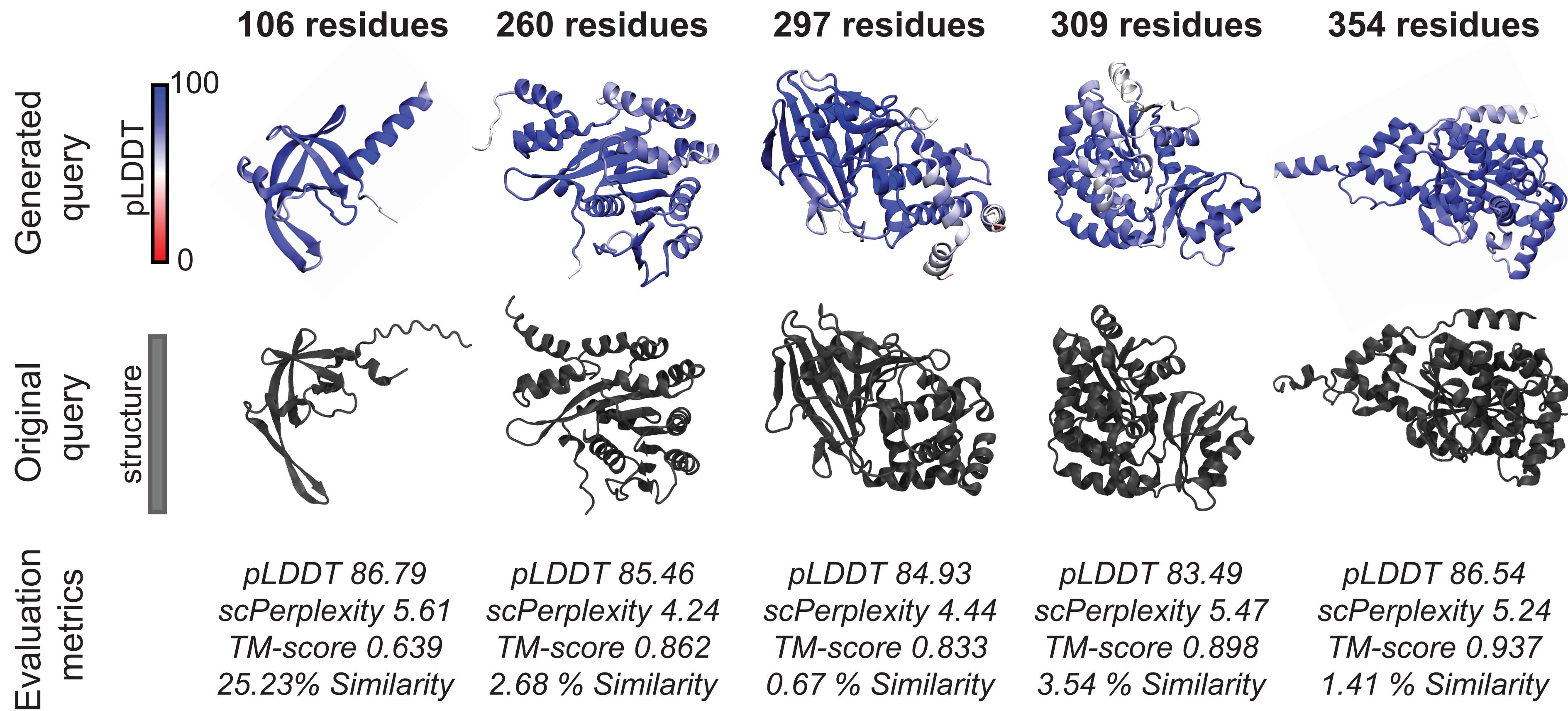
conditional
generation

Evolutionary diffusion improves local generation



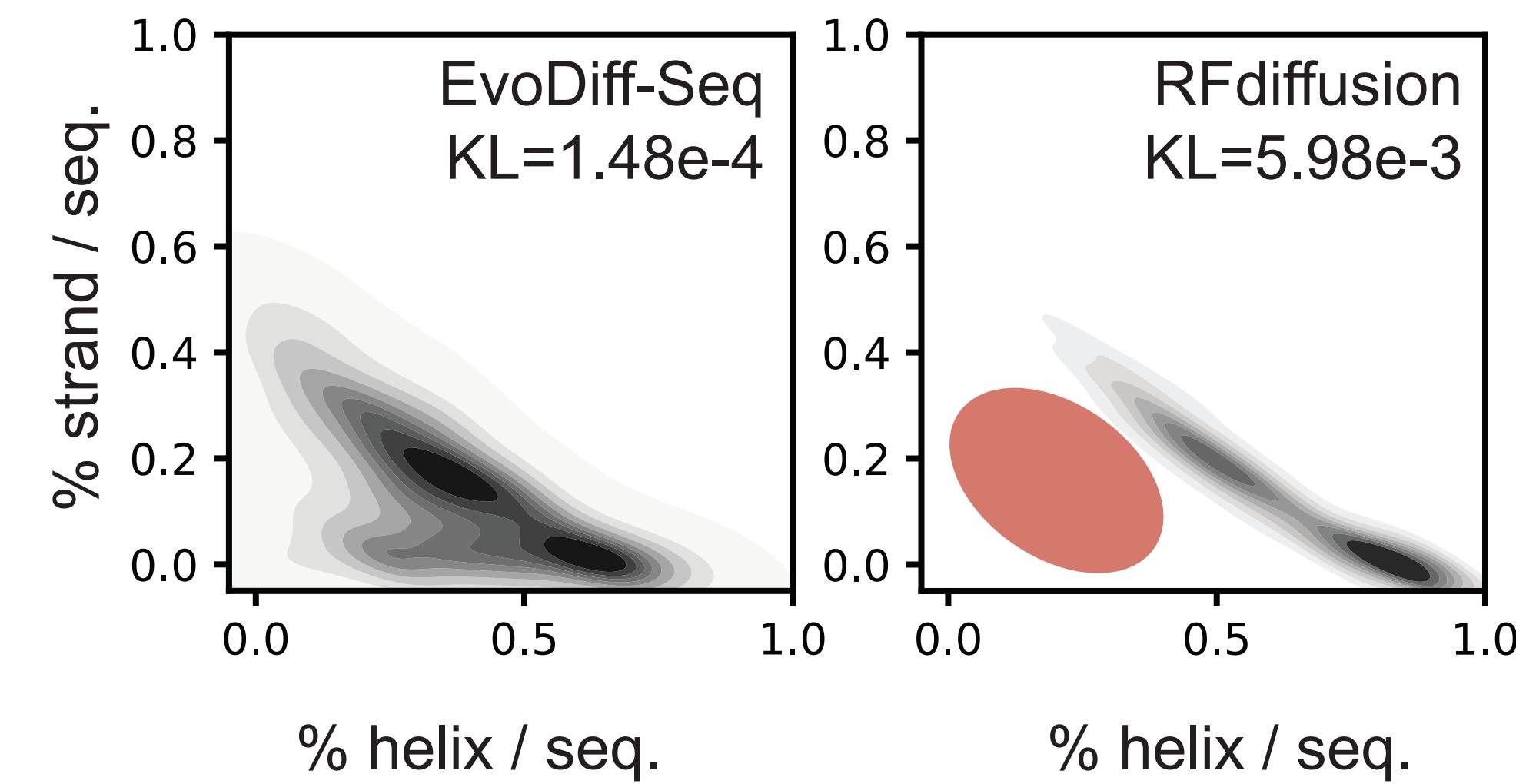
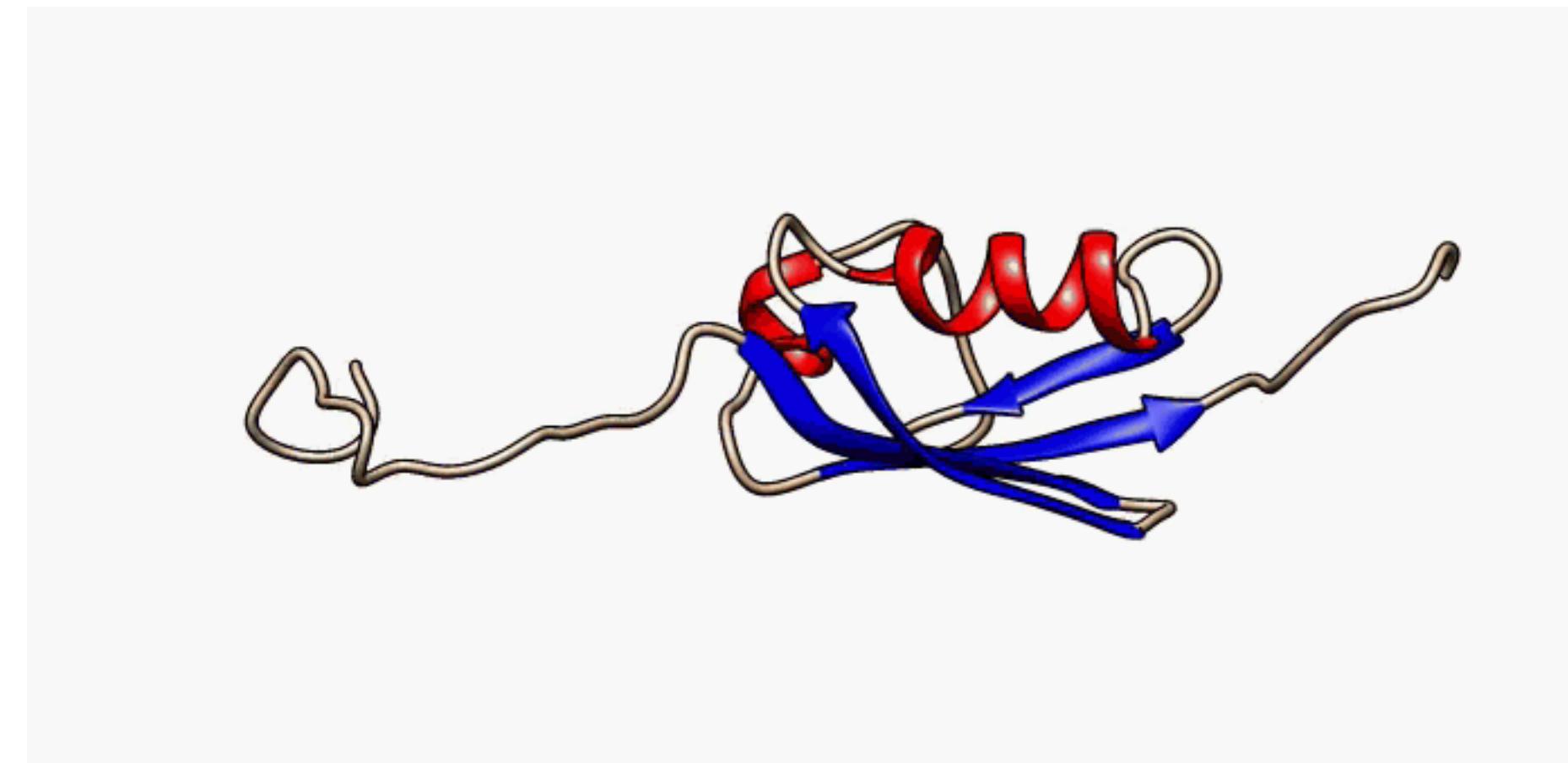
conditional
generation

Evolutionary diffusion improves local generation



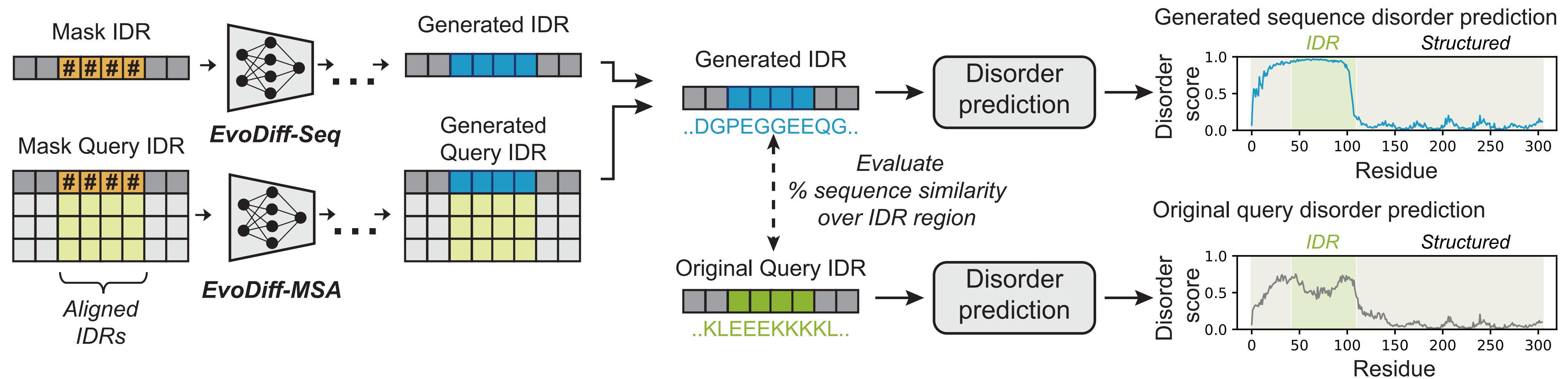
conditional
generation

Intrinsically disordered regions perform many cellular functions

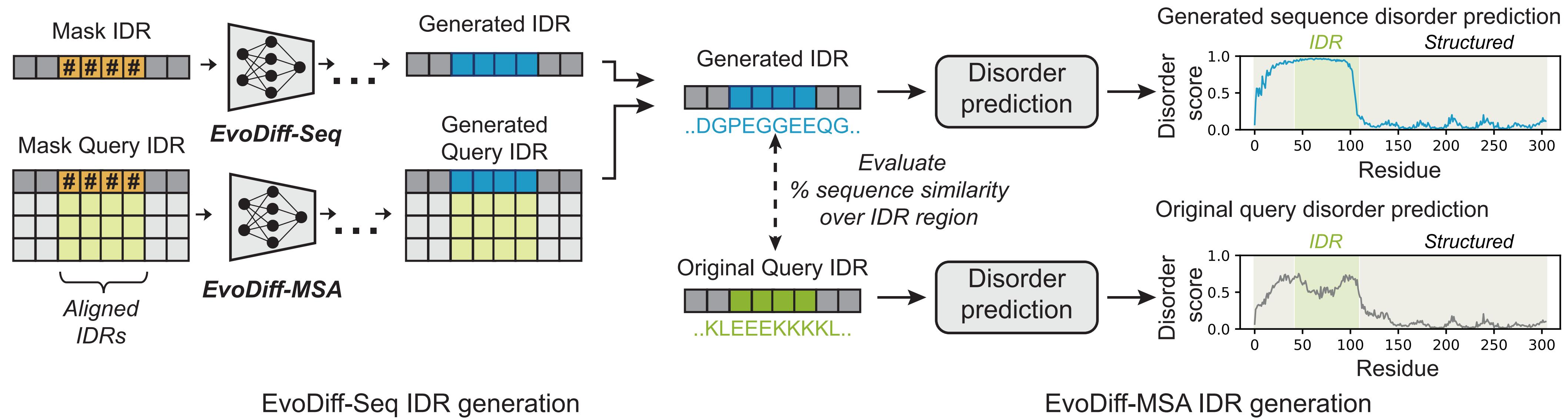


conditional
generation

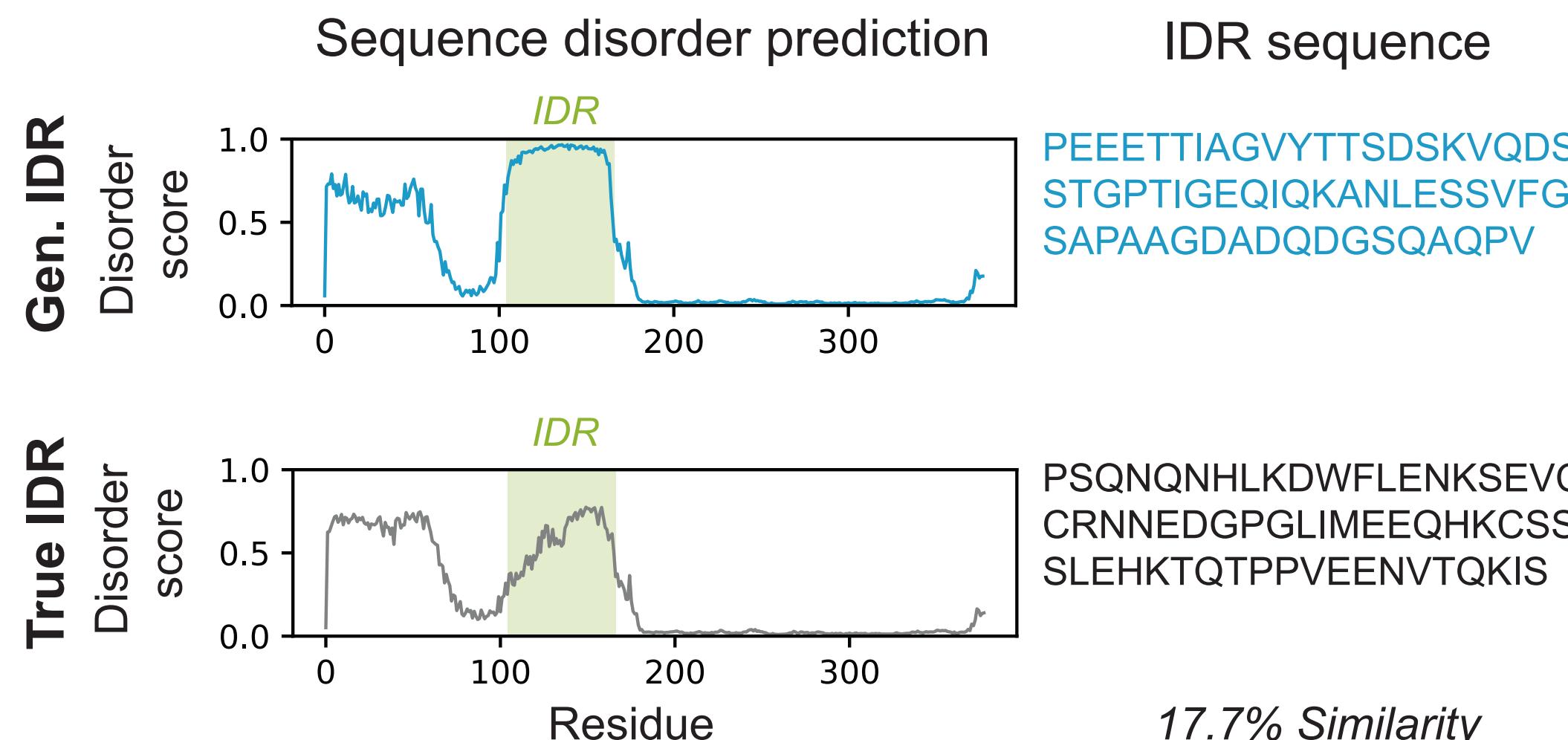
EvoDiff can generate disordered regions



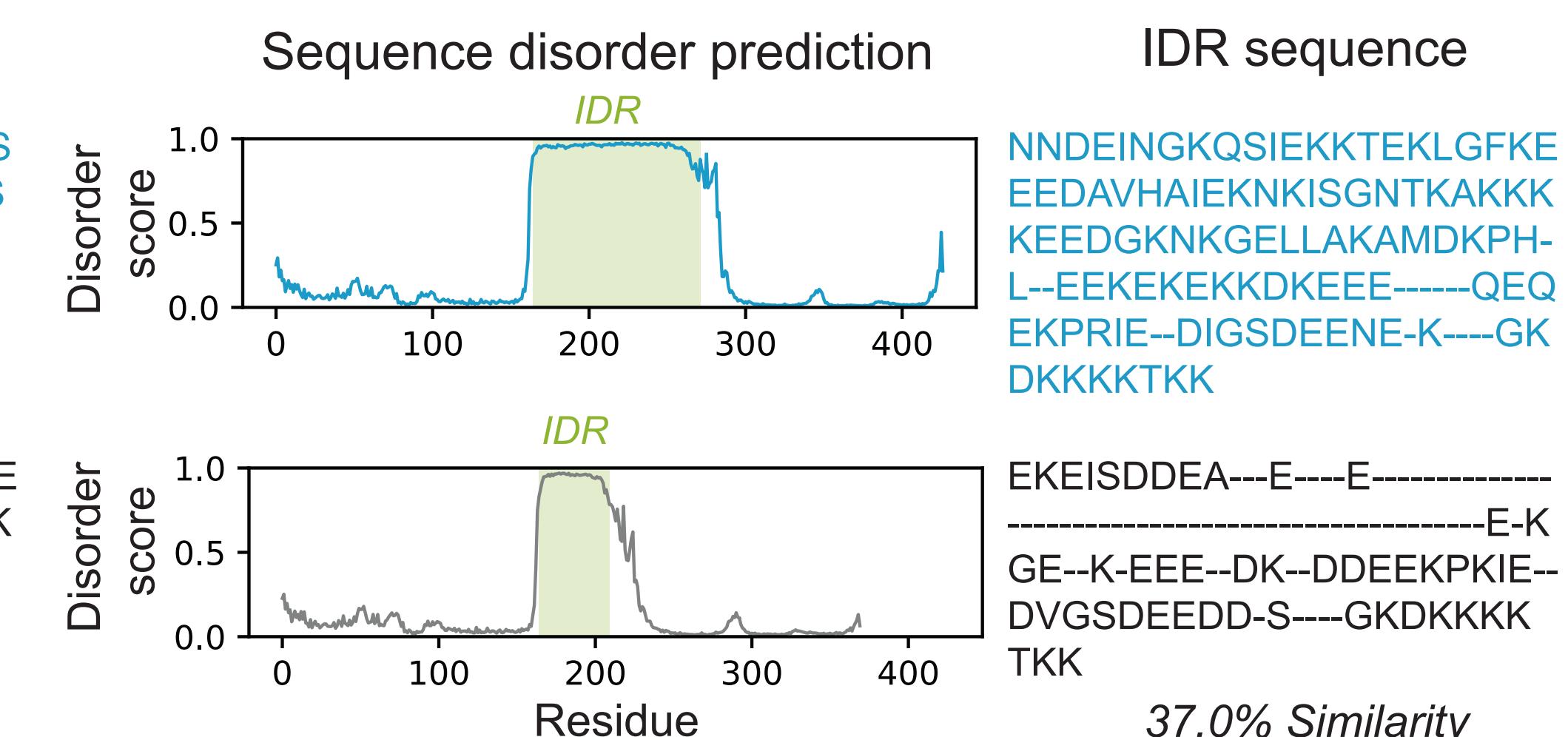
EvoDiff can generate disordered regions



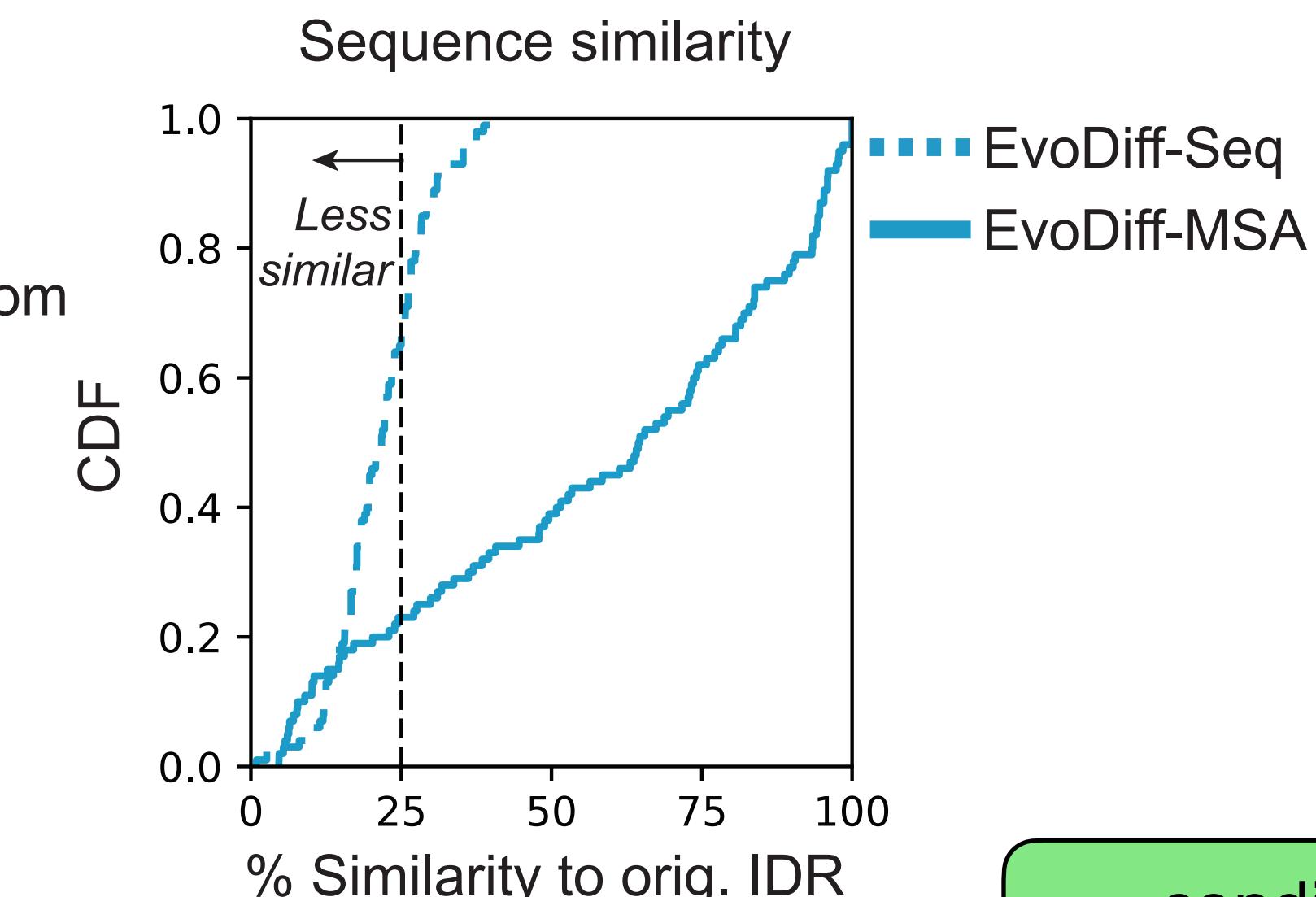
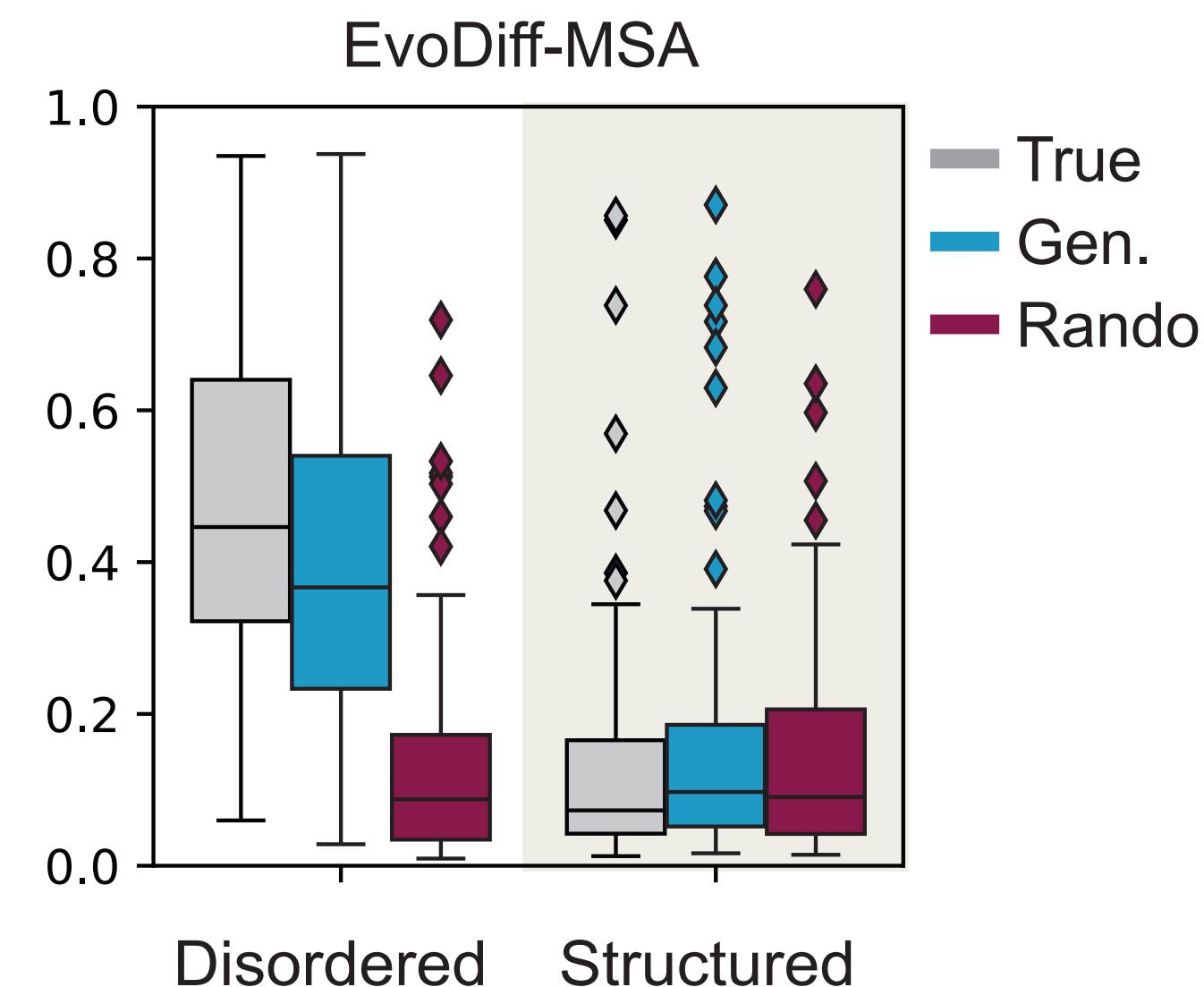
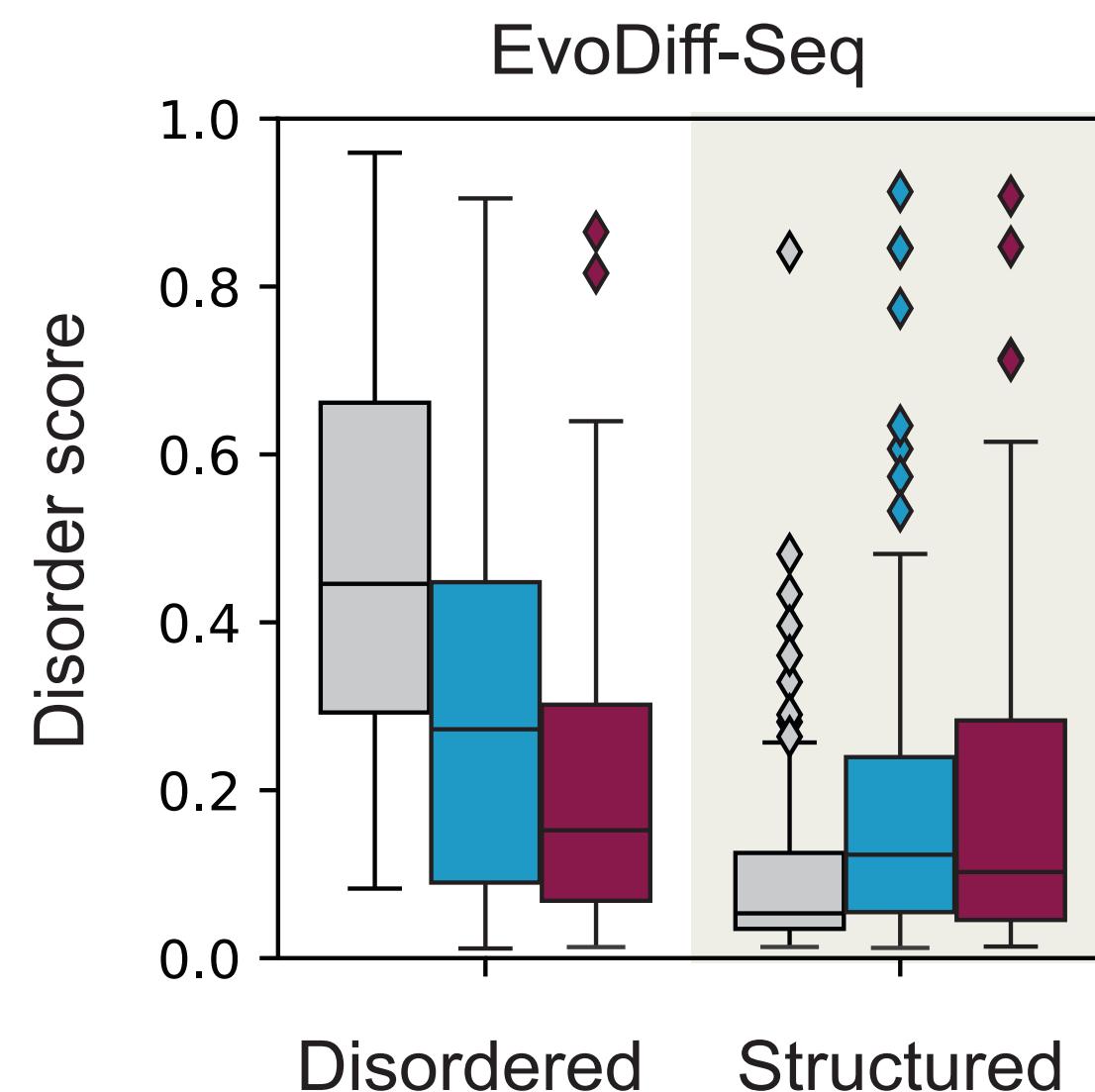
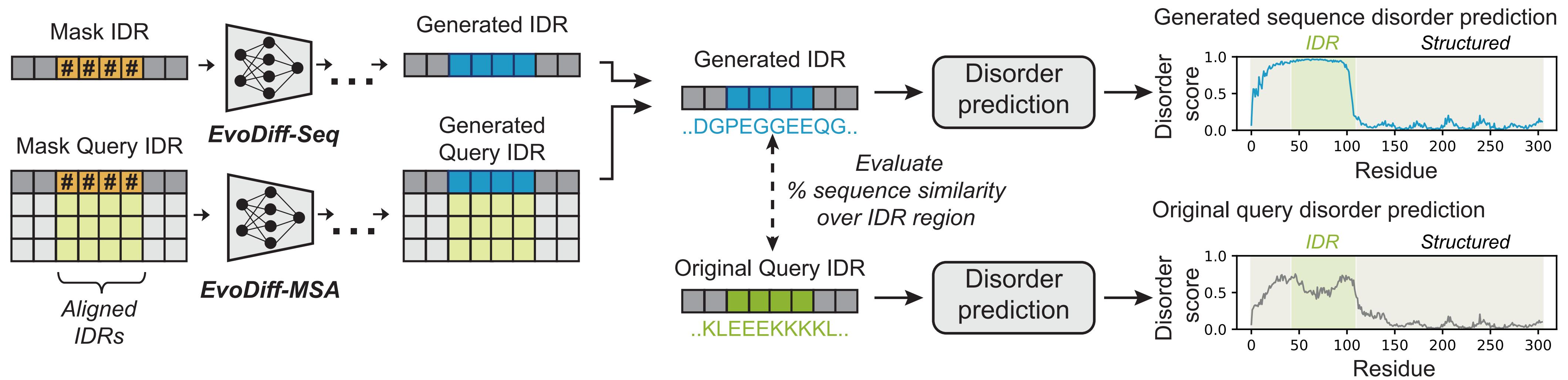
UniProtID Q96Q15: Serine/threonine protein kinase



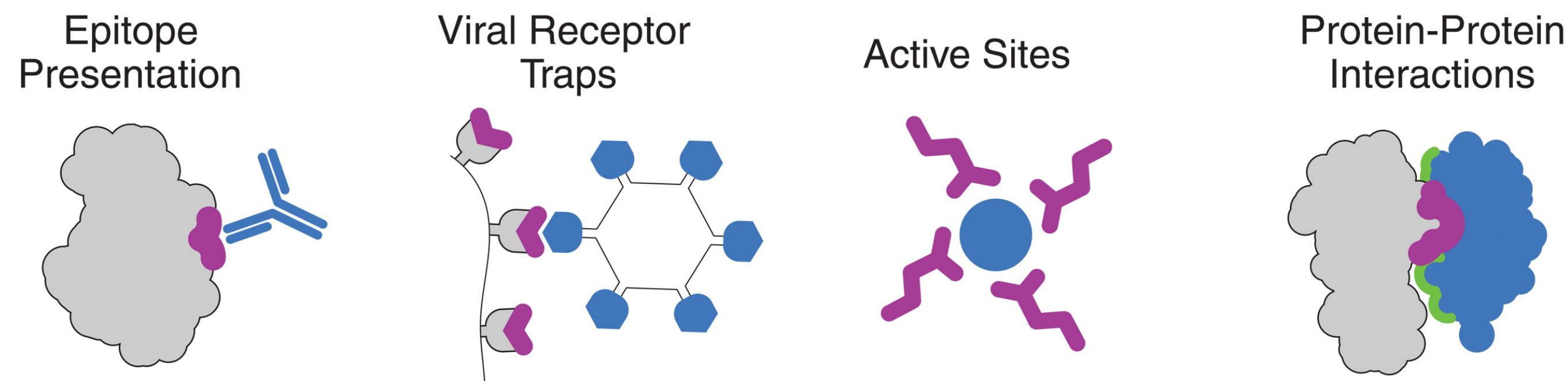
UniProtID F5H008: Endosome/lysosome vesicle associated protein



EvoDiff can generate disordered regions



Many functions are mediated by a motif stabilized by a scaffold



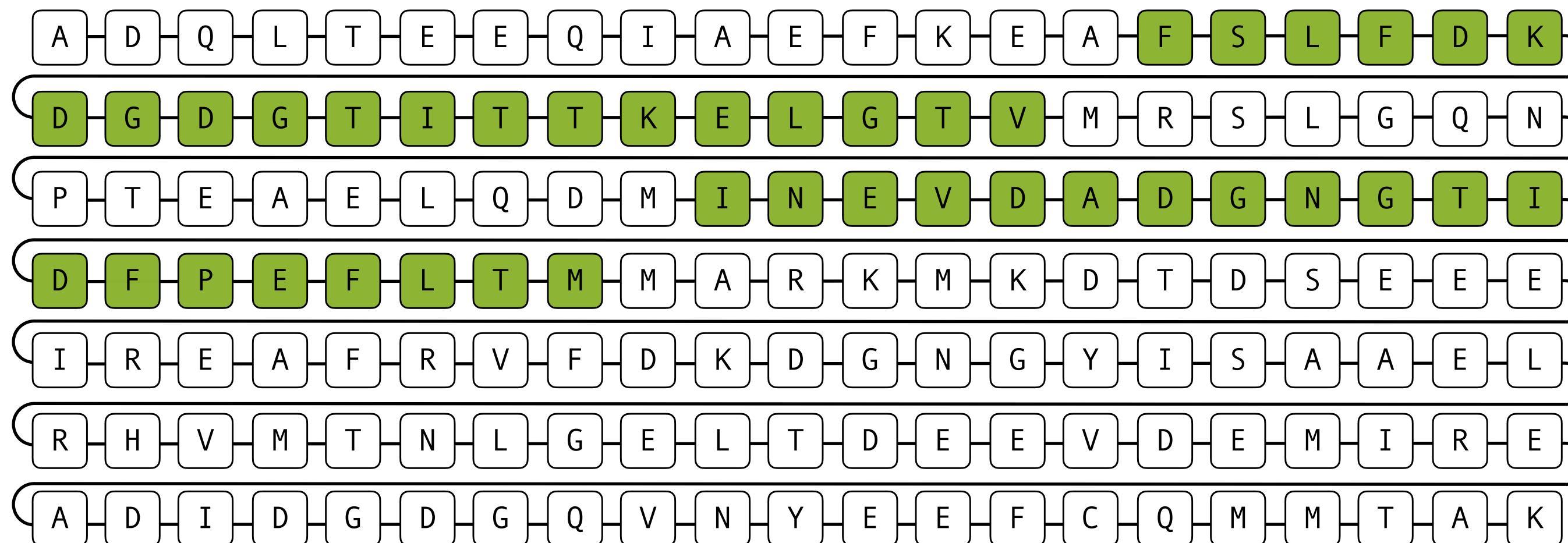
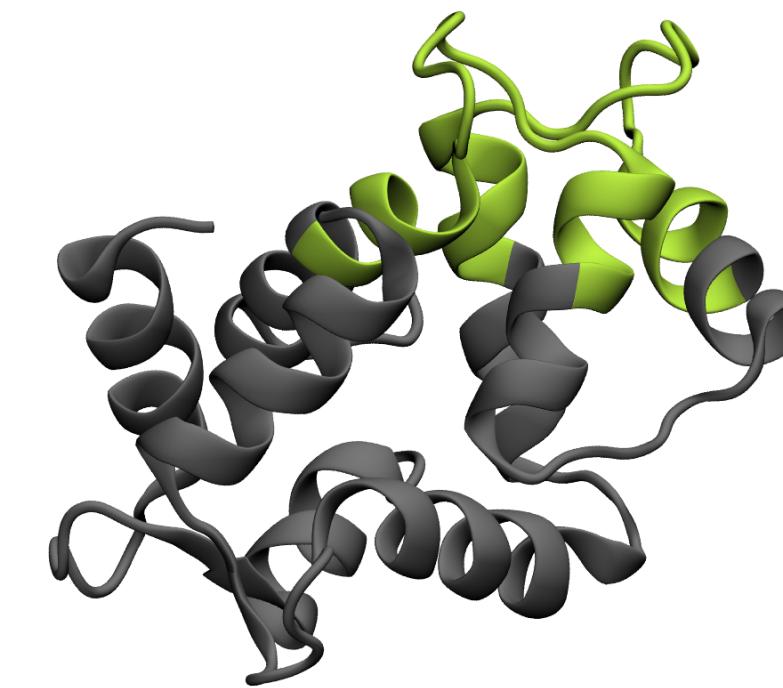
Wang *et al.*, *Science* 2022

Can we scaffold motifs in sequence space?

conditional
generation

EvoDiff can scaffold functional motifs

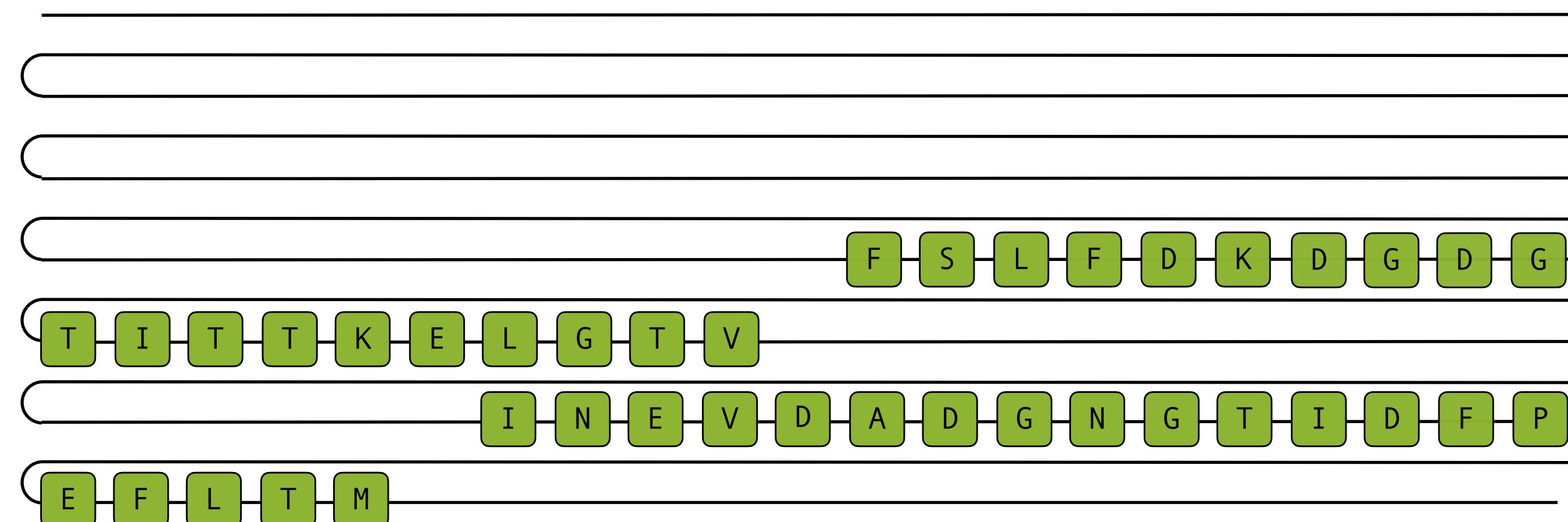
1PRW: binding site of
compact calmodulin



conditional
generation

EvoDiff can scaffold functional motifs

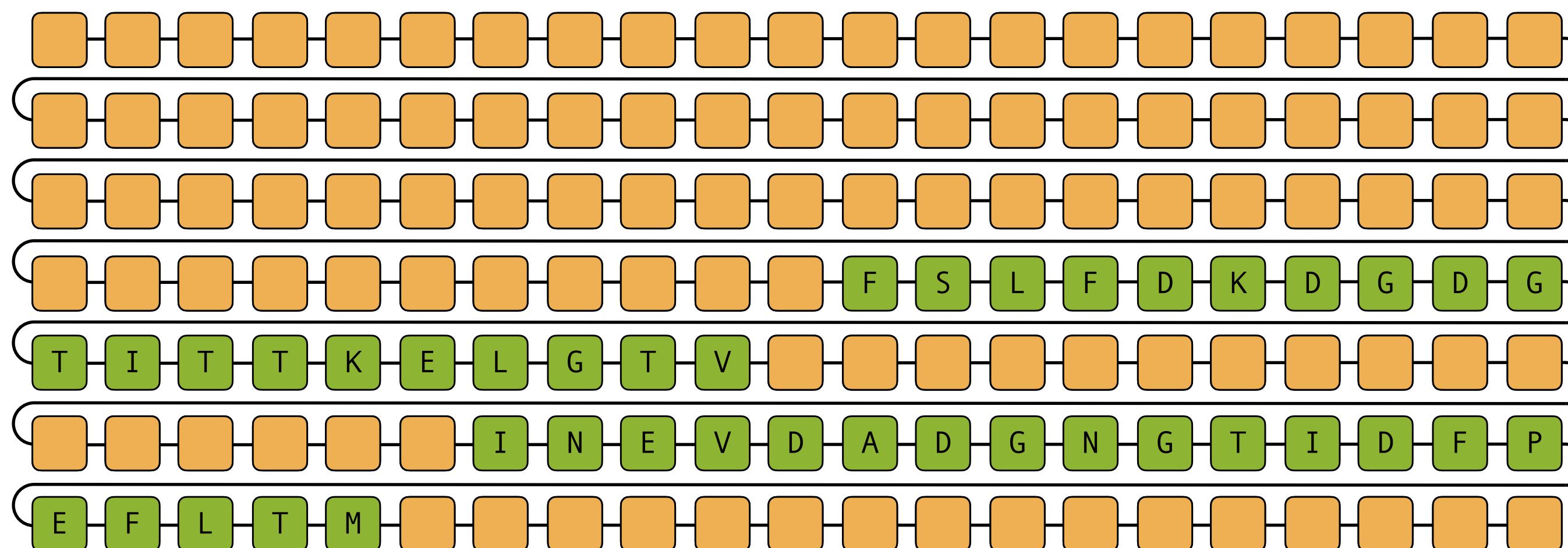
1PRW: binding site of
compact calmodulin



conditional
generation

EvoDiff can scaffold functional motifs

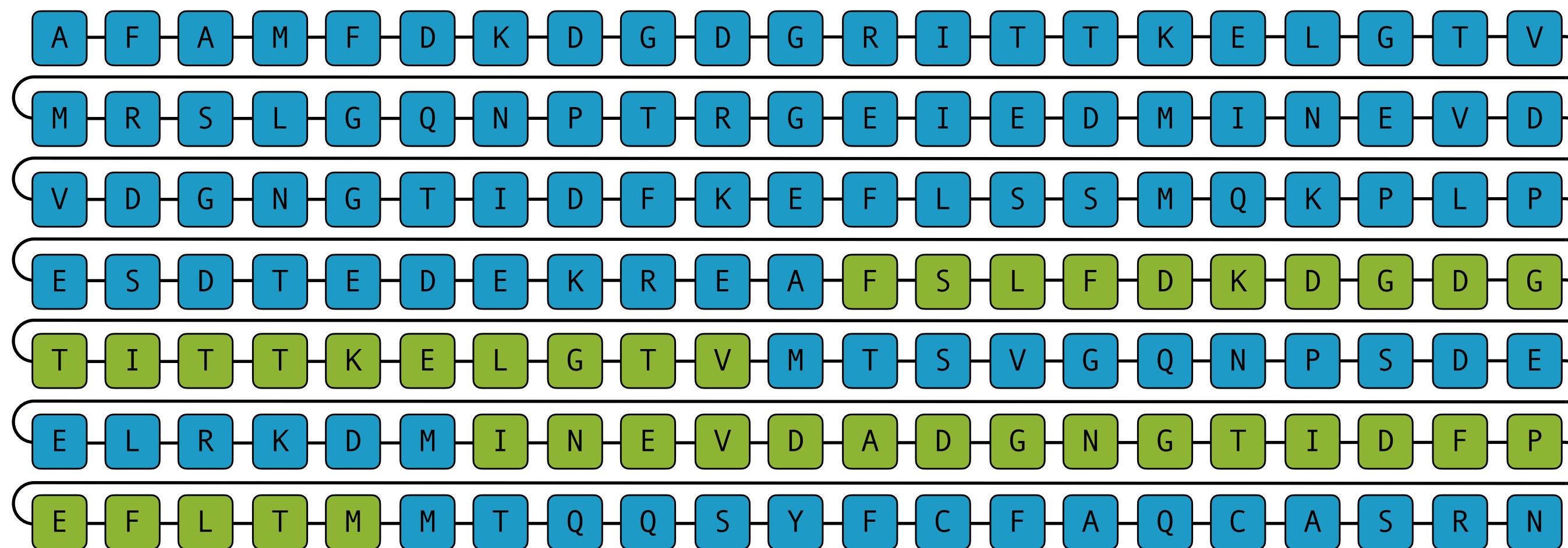
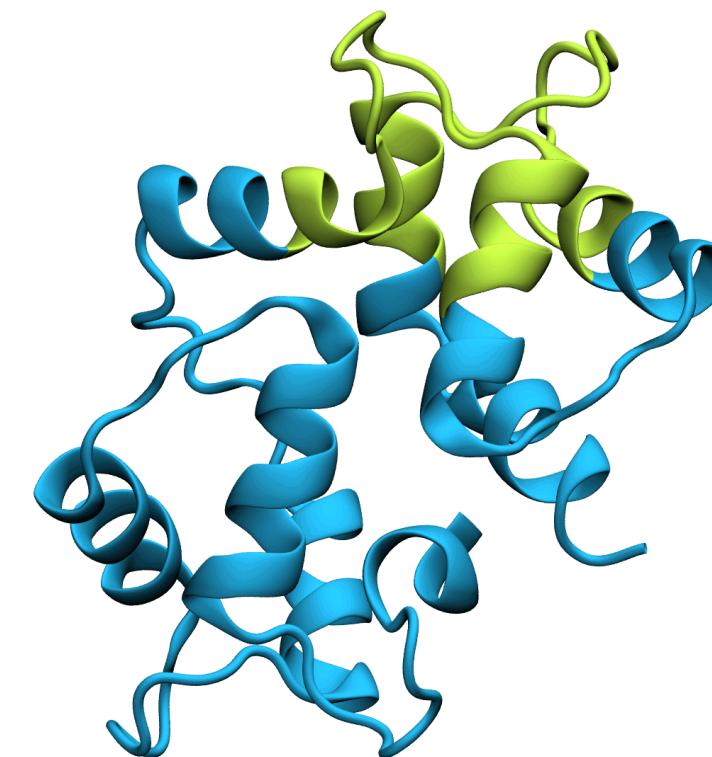
1PRW: binding site of
compact calmodulin



conditional
generation

EvoDiff can scaffold functional motifs

1PRW: binding site of
compact calmodulin



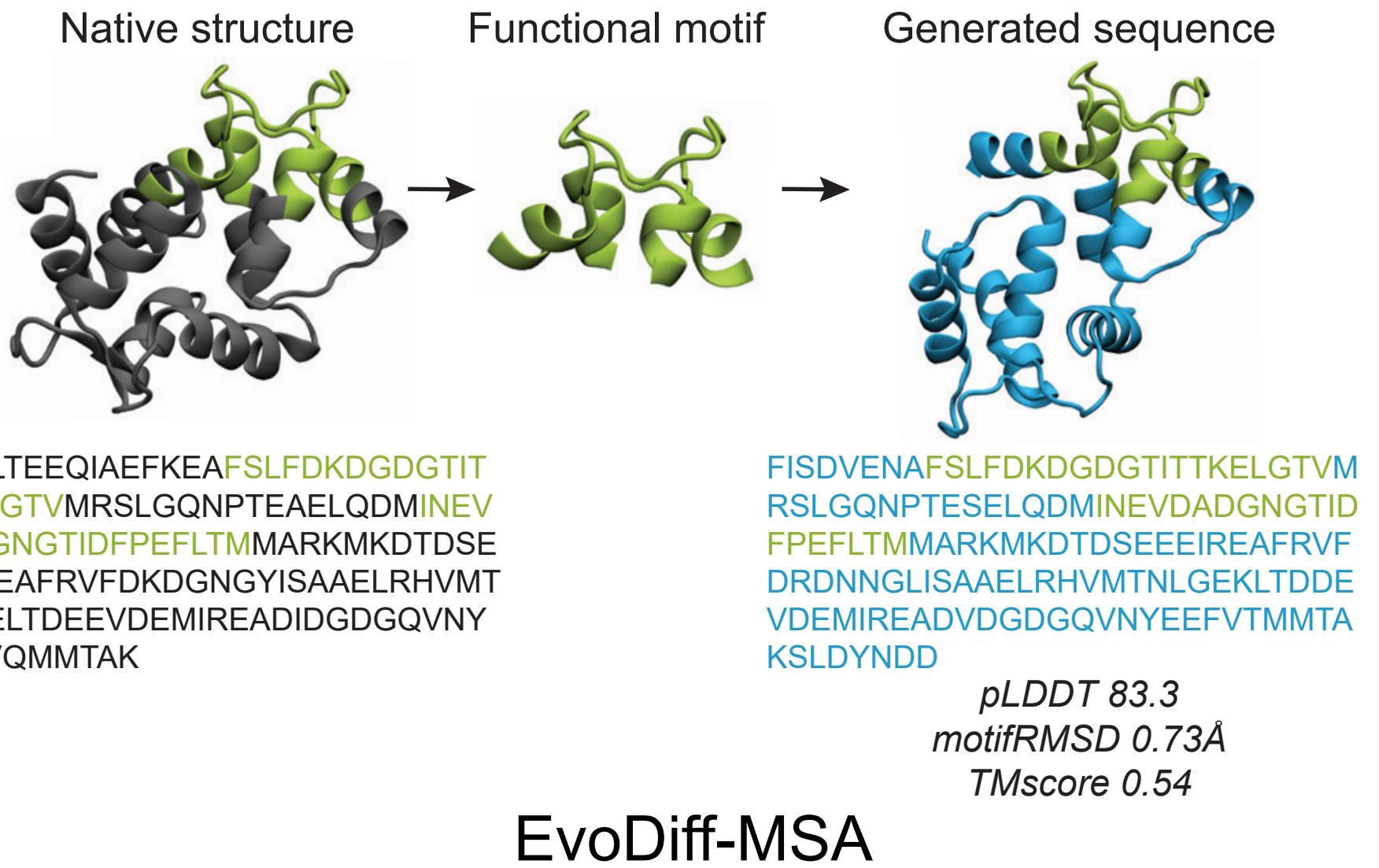
No structure needed!

conditional
generation

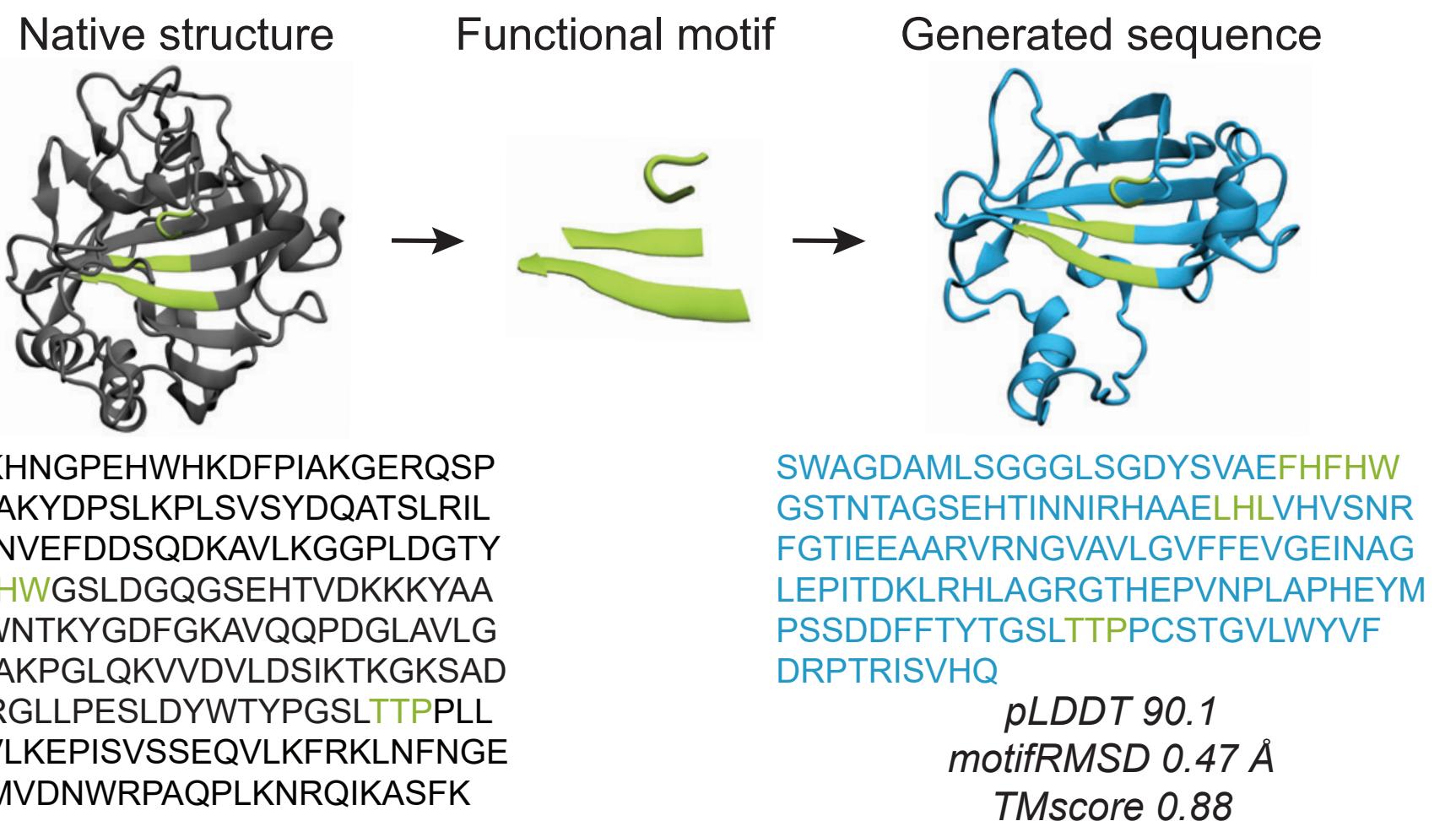
EvoDiff can scaffold functional motifs

EvoDiff-Seq

1PRW: Binding site of compact calmodulin



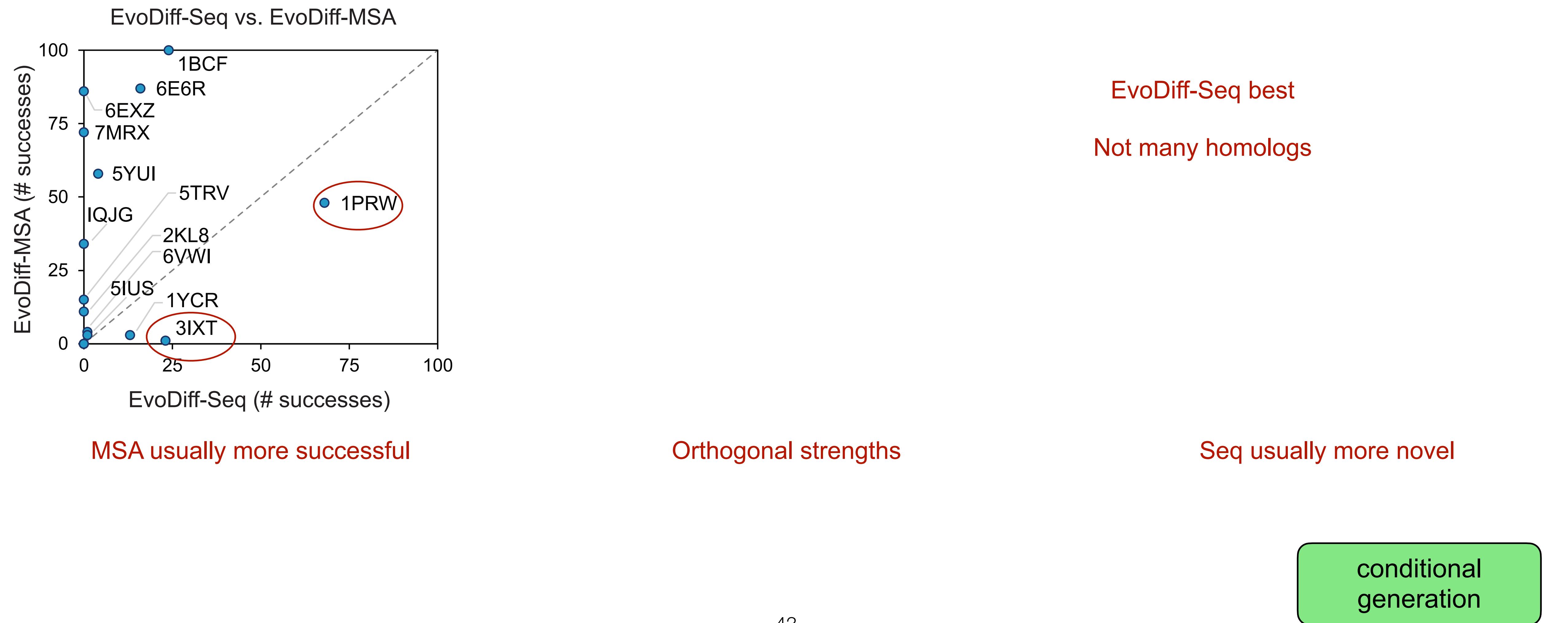
5YUI: Binding site of carbonic anhydrase metalloenzyme



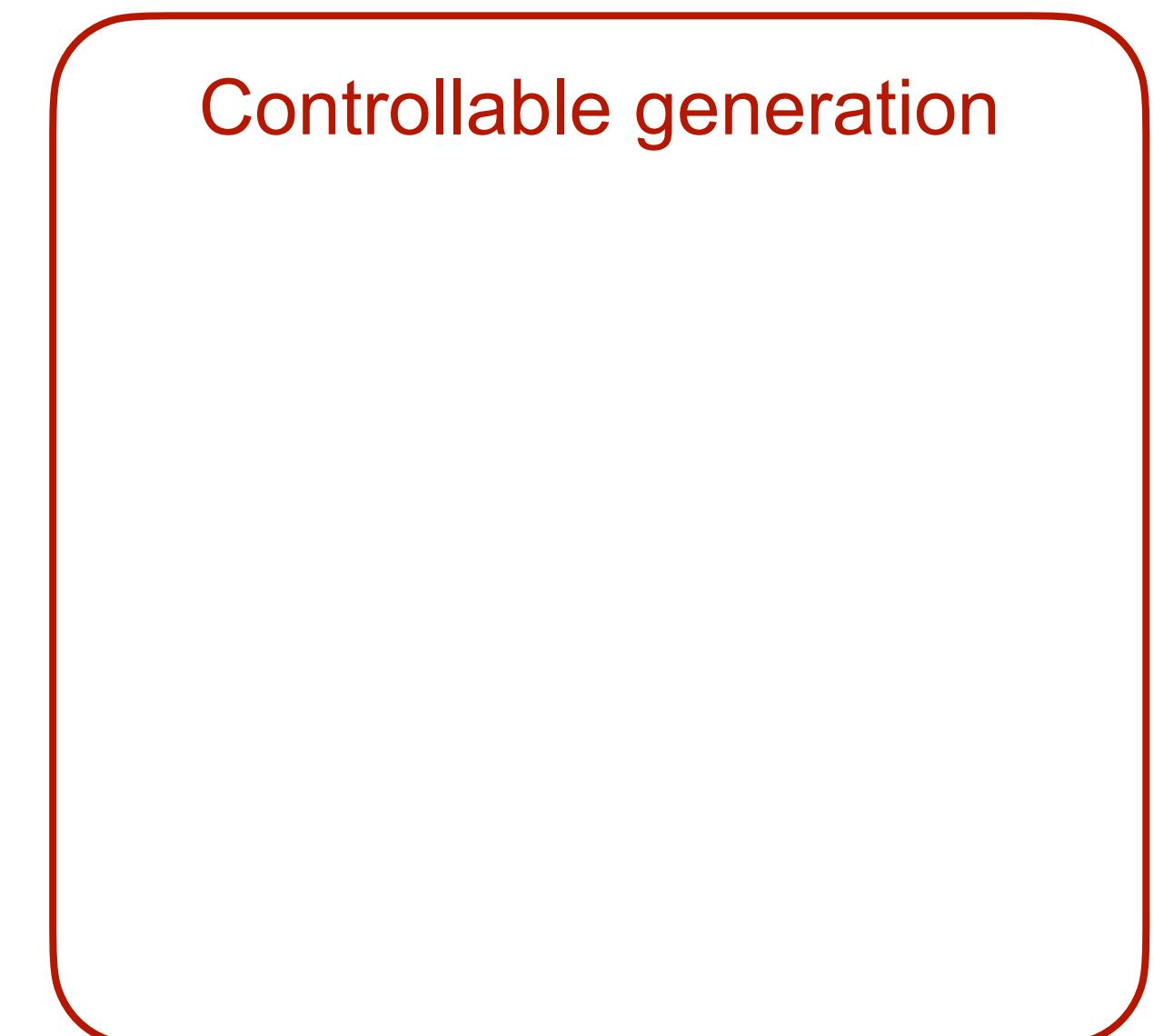
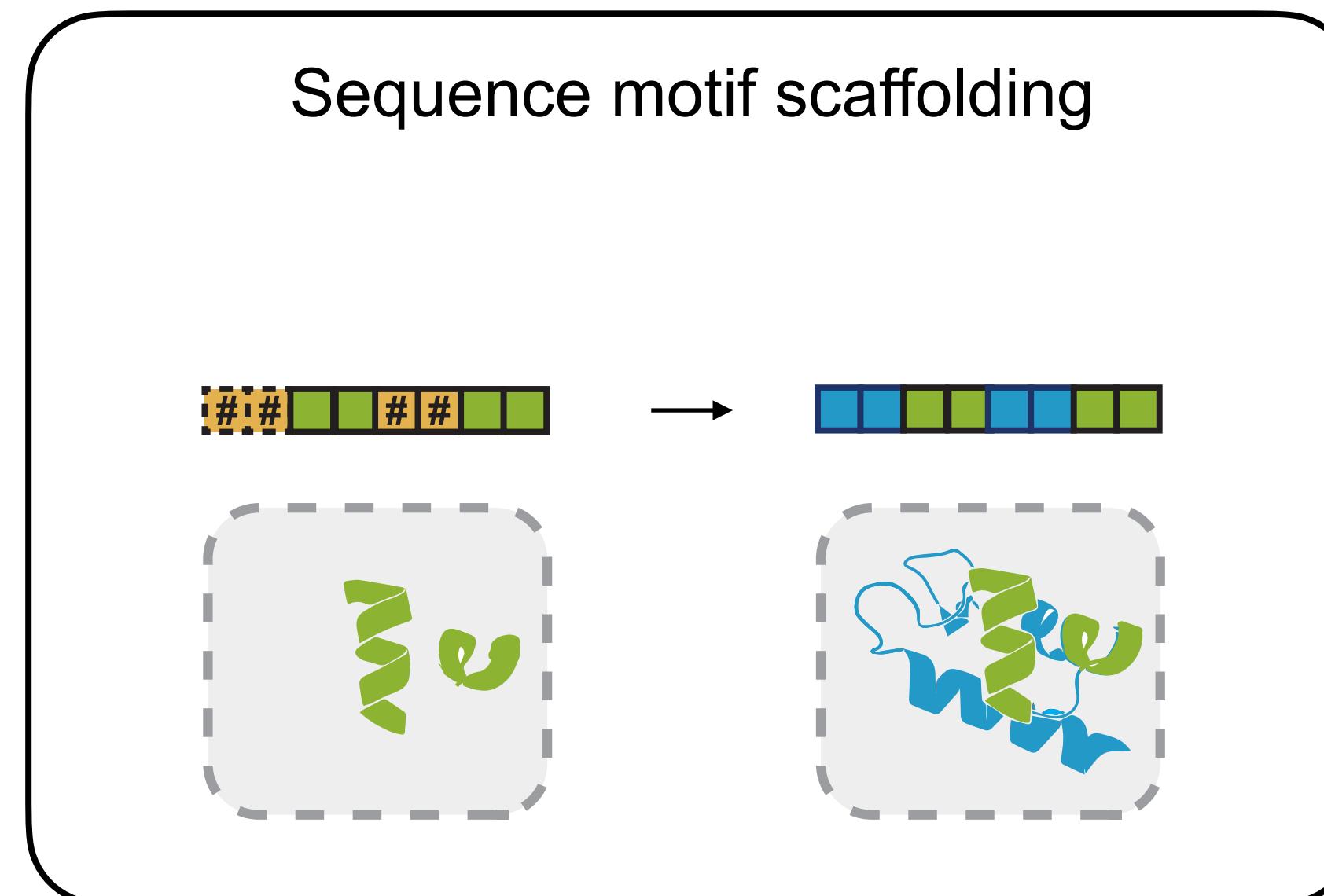
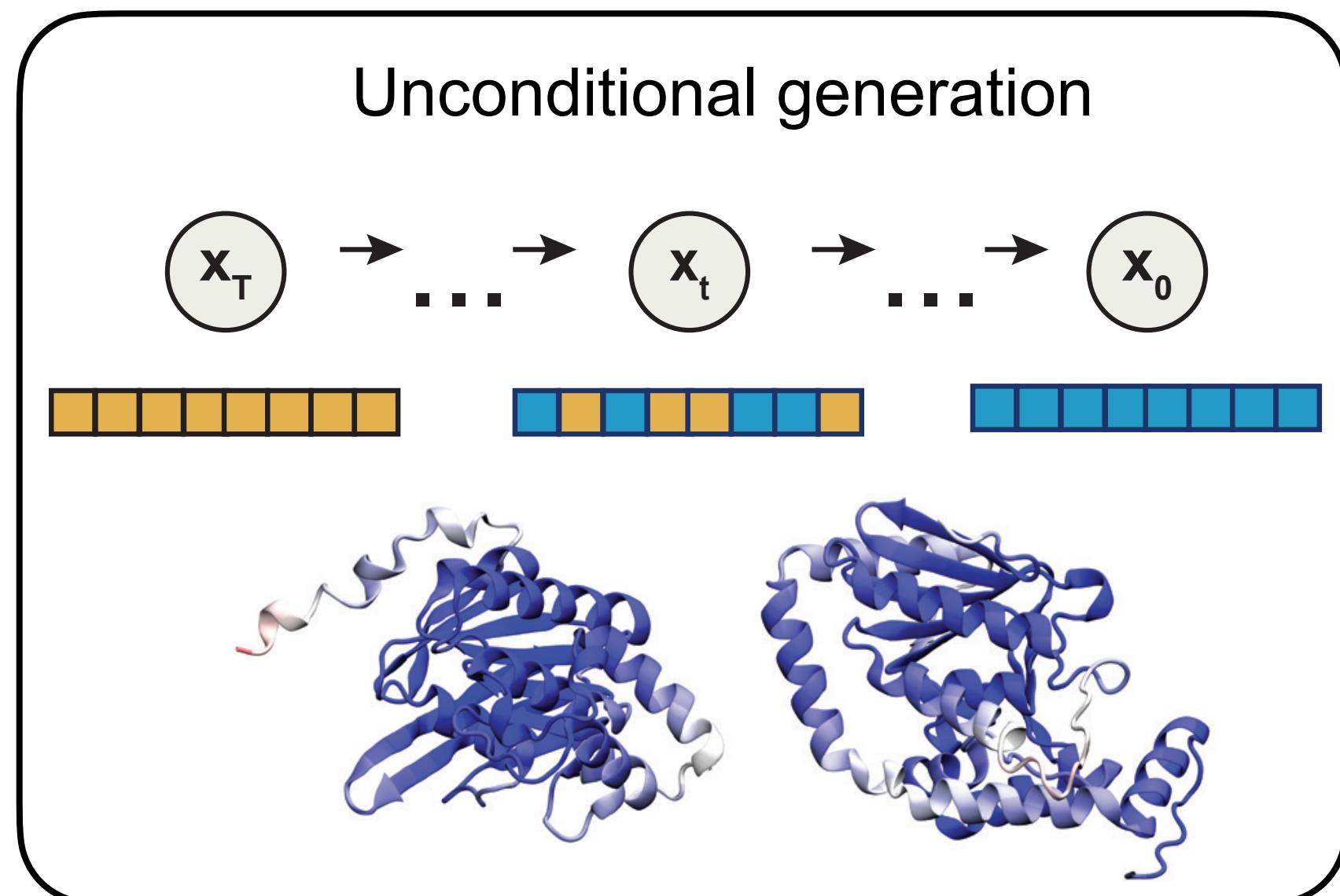
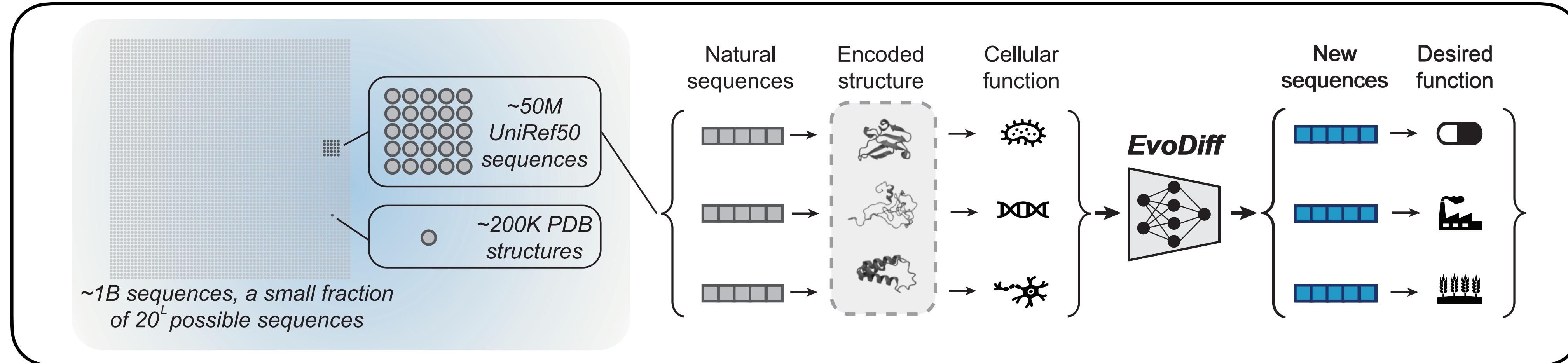
Model	# Successful (< 1Å RMSD)	# Problems solved
RFdiffusion	610	13 / 17
EvoDiff-MSA	522	13 / 17
EvoDiff-Seq	149	8 / 17

conditional
generation

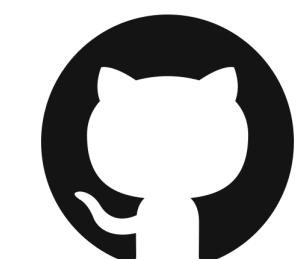
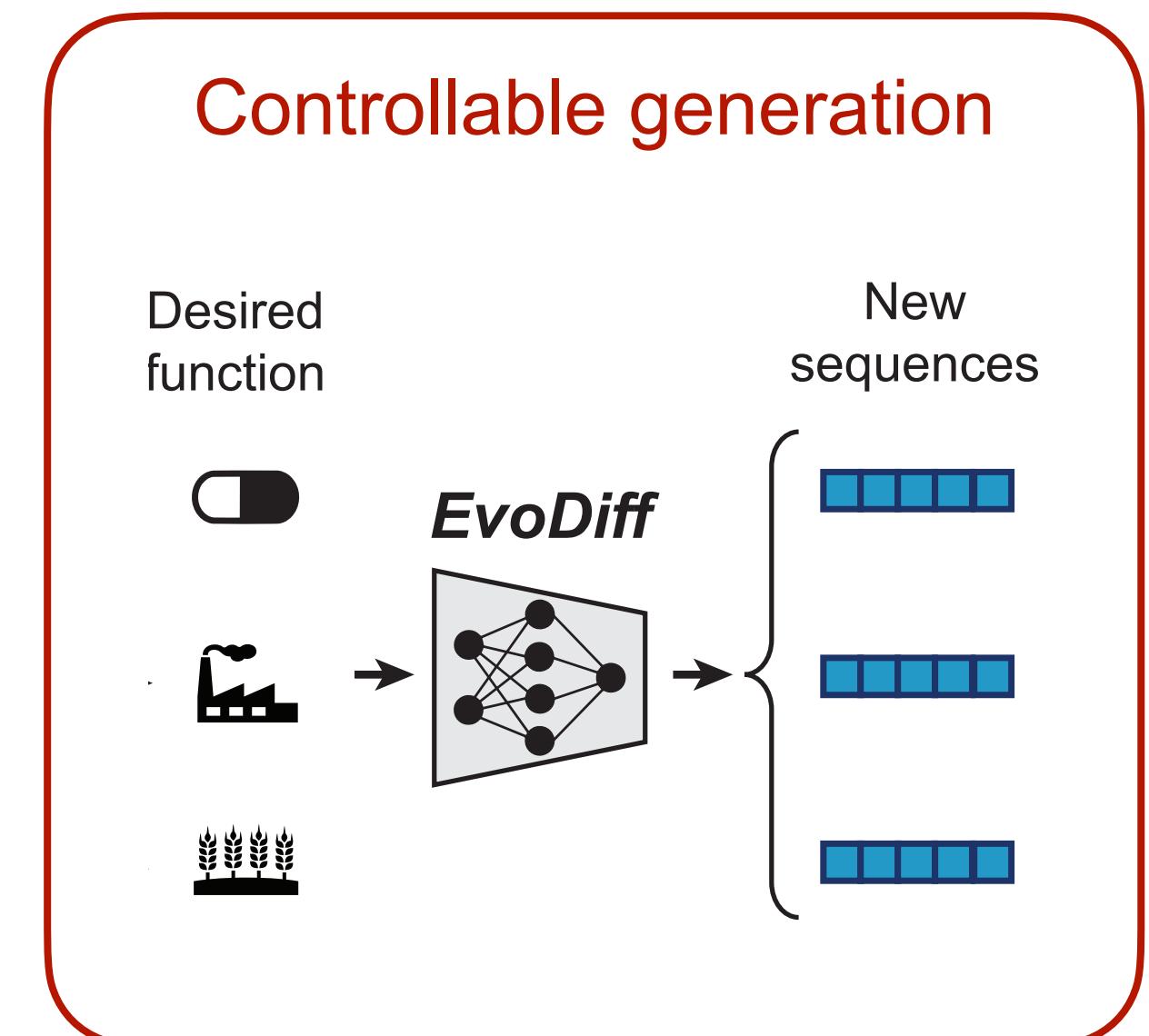
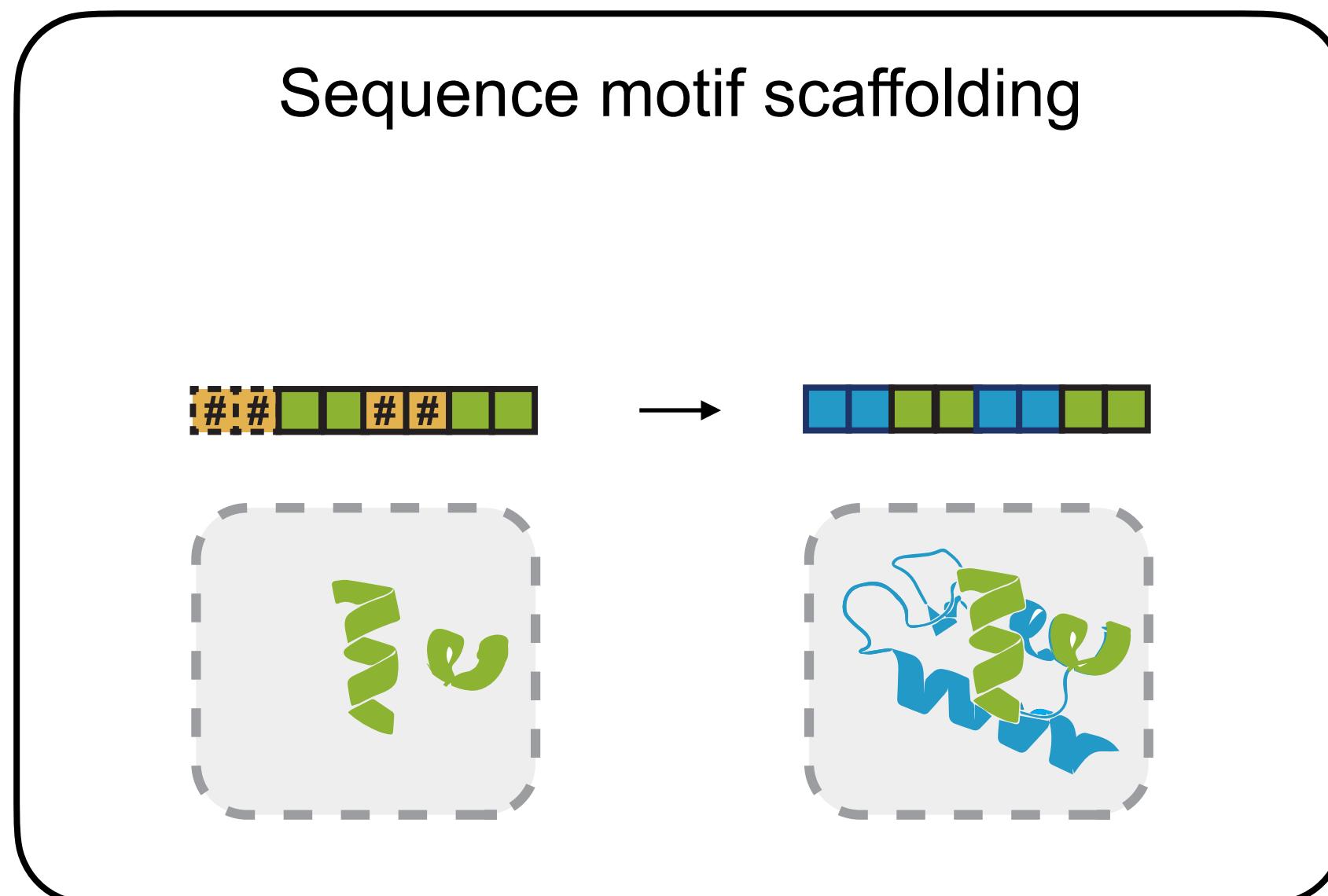
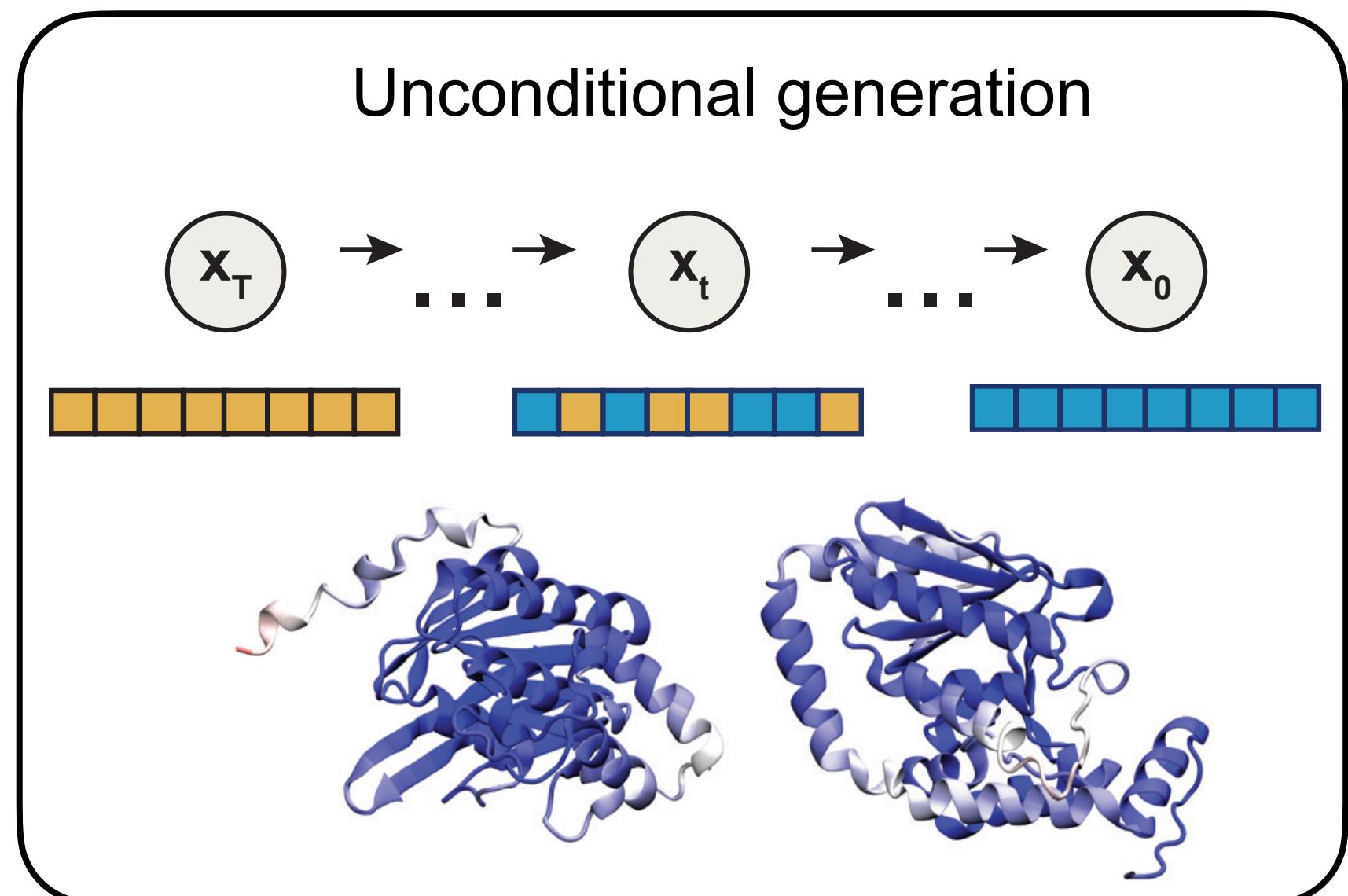
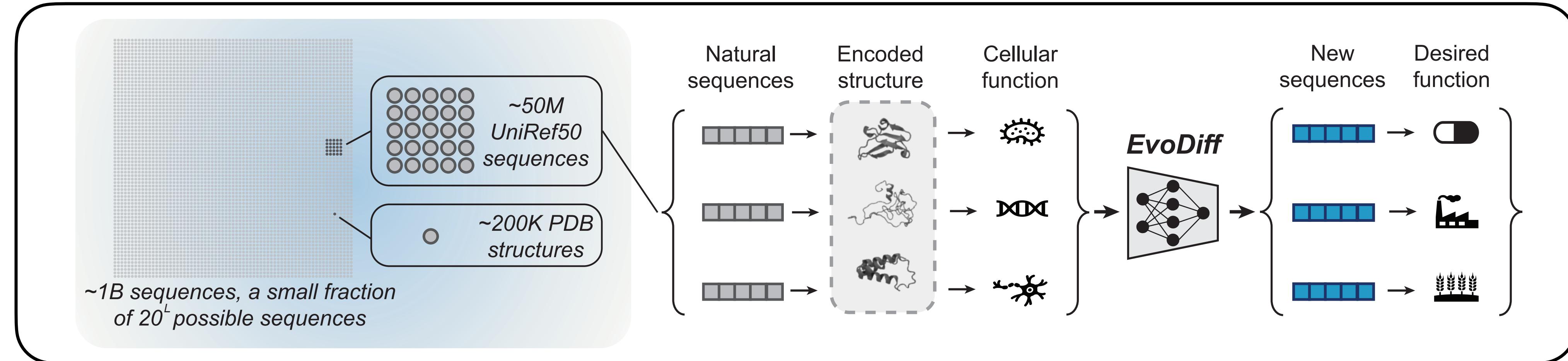
EvoDiff can scaffold functional motifs



EvoDiff: controllable protein sequence diffusion



EvoDiff: controllable protein sequence diffusion



Acknowledgments



BioML at MSR New England