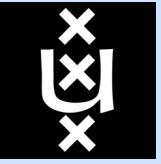


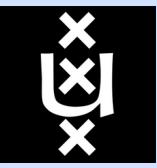
Protein Folding

Uniform Random Sampling using the HP-Model
and Randomly Generated Amino Acids



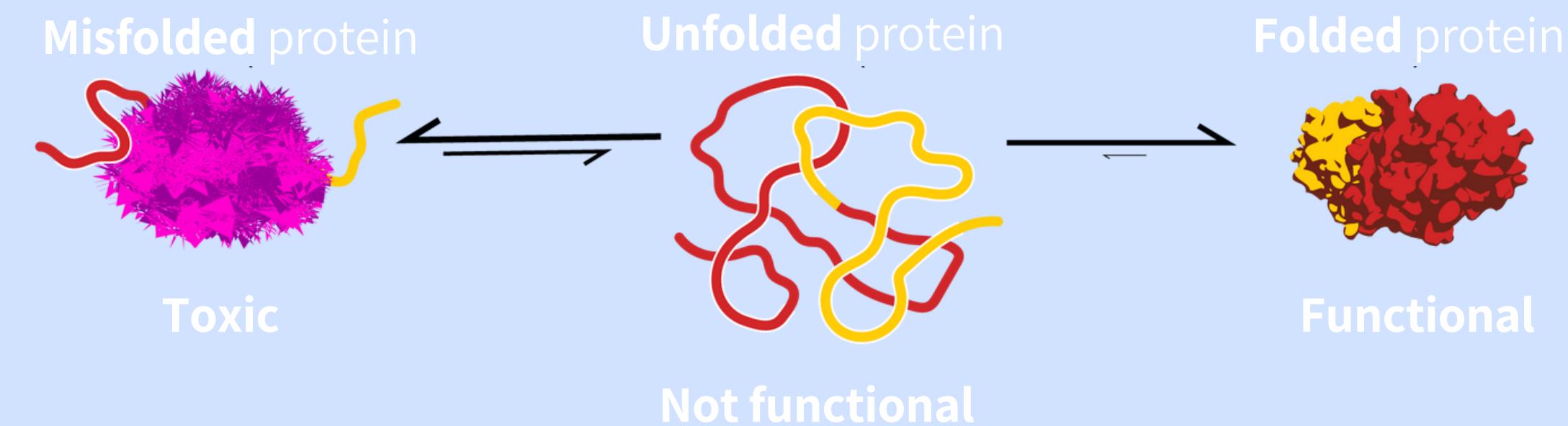
Protein Folding 101

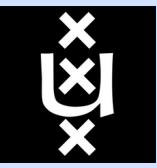




Why do proteins fold?

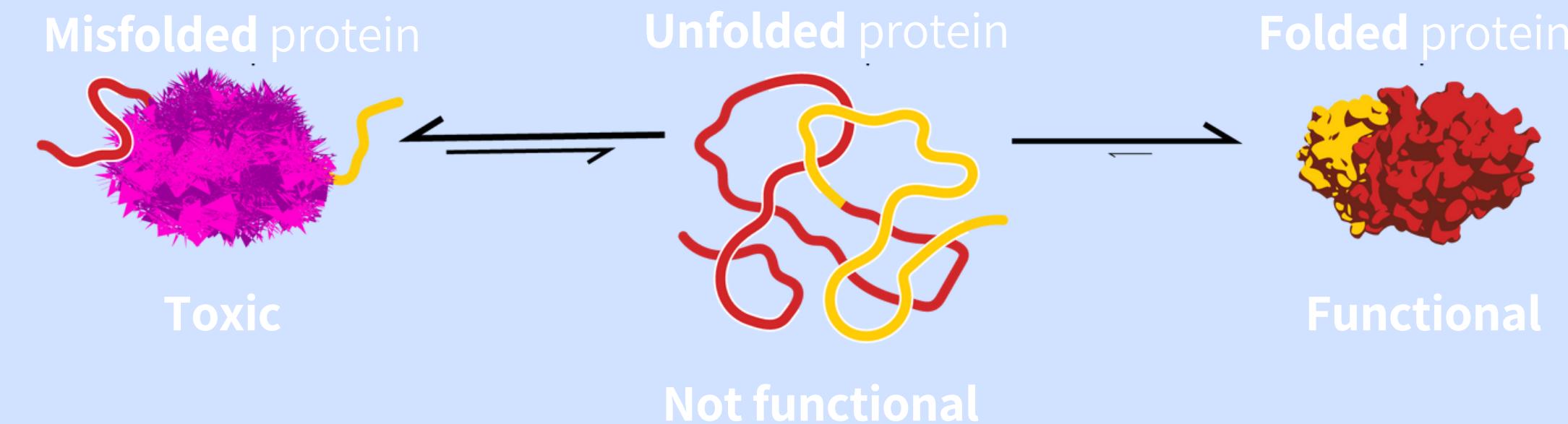
Folding is the process of taking on functional structure.

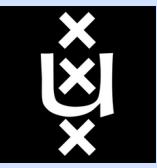




Why we fold proteins?

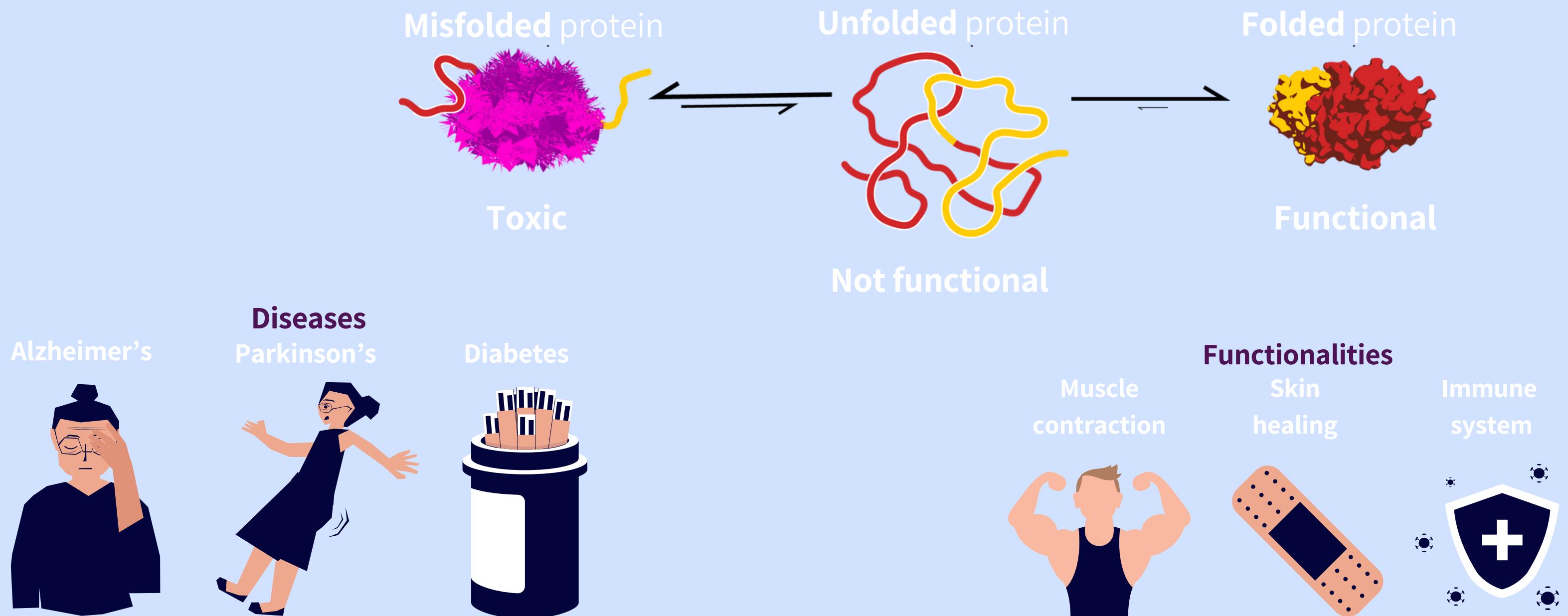
To predict the functional structure of proteins.





Why we fold proteins?

To predict the functional structure of proteins.

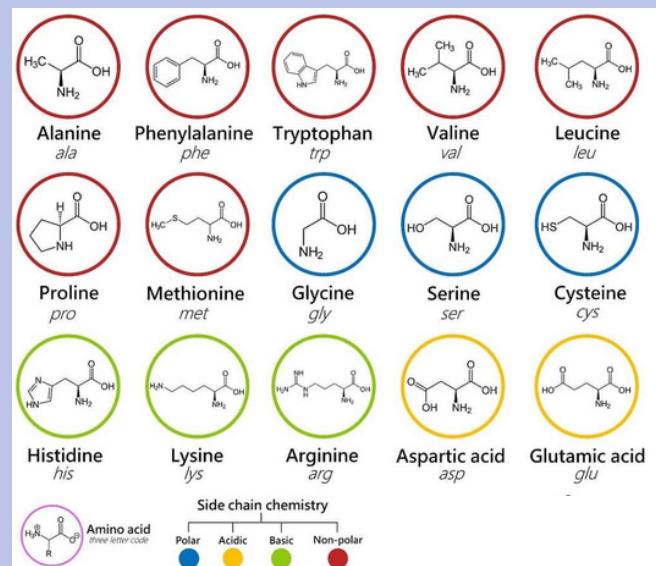




Protein and aminos

A Protein is a chain of amino acids.

20 different amino acids



Average protein
200-300 acids

Insulin - 51 amino acids

Chain A - 21 amino acids
GIVEQCCTSICSLYQLENYCN



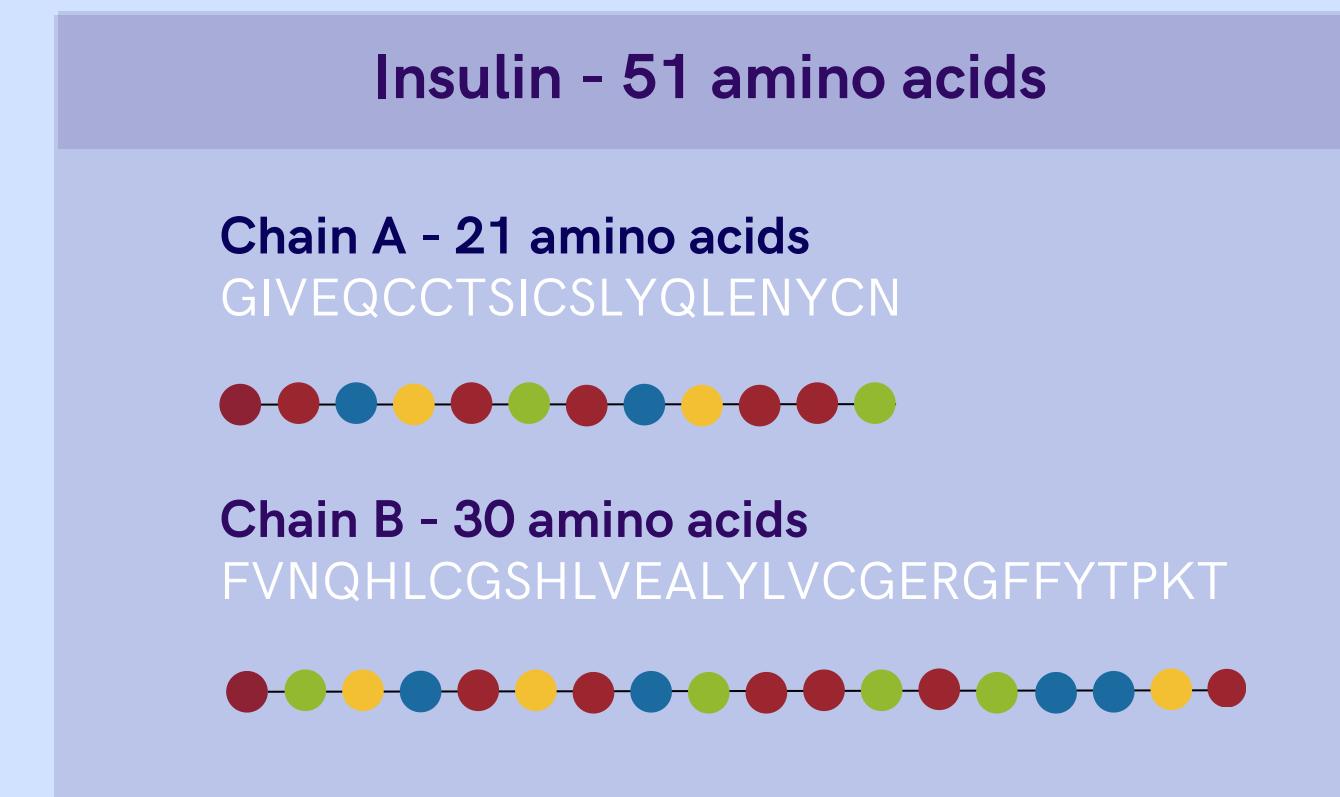
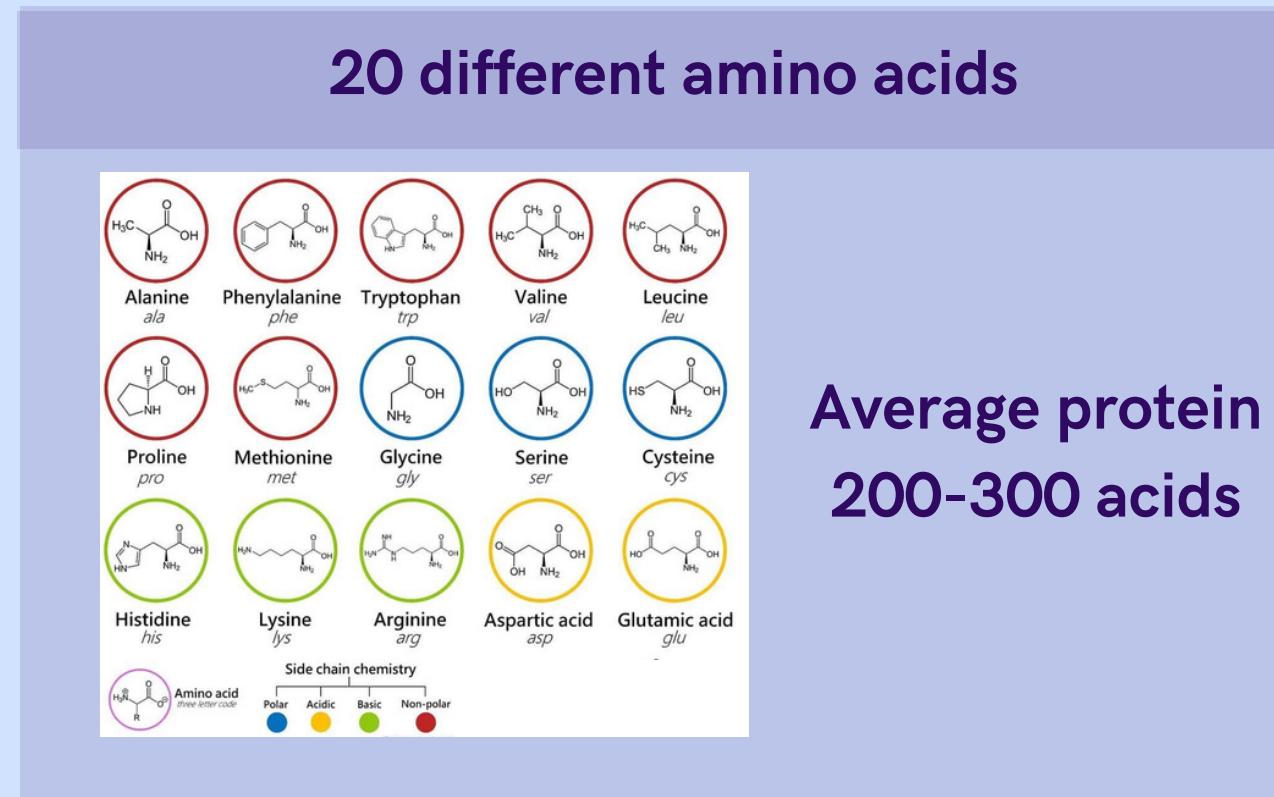
Chain B - 30 amino acids
FVNQHLCGSHLVEALYLVCGERGFFYTPKT





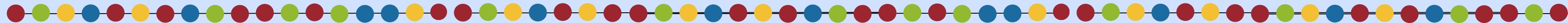
Protein and aminos

A Protein is a chain of amino acids.



Many Proteins > 2.000 amino acids

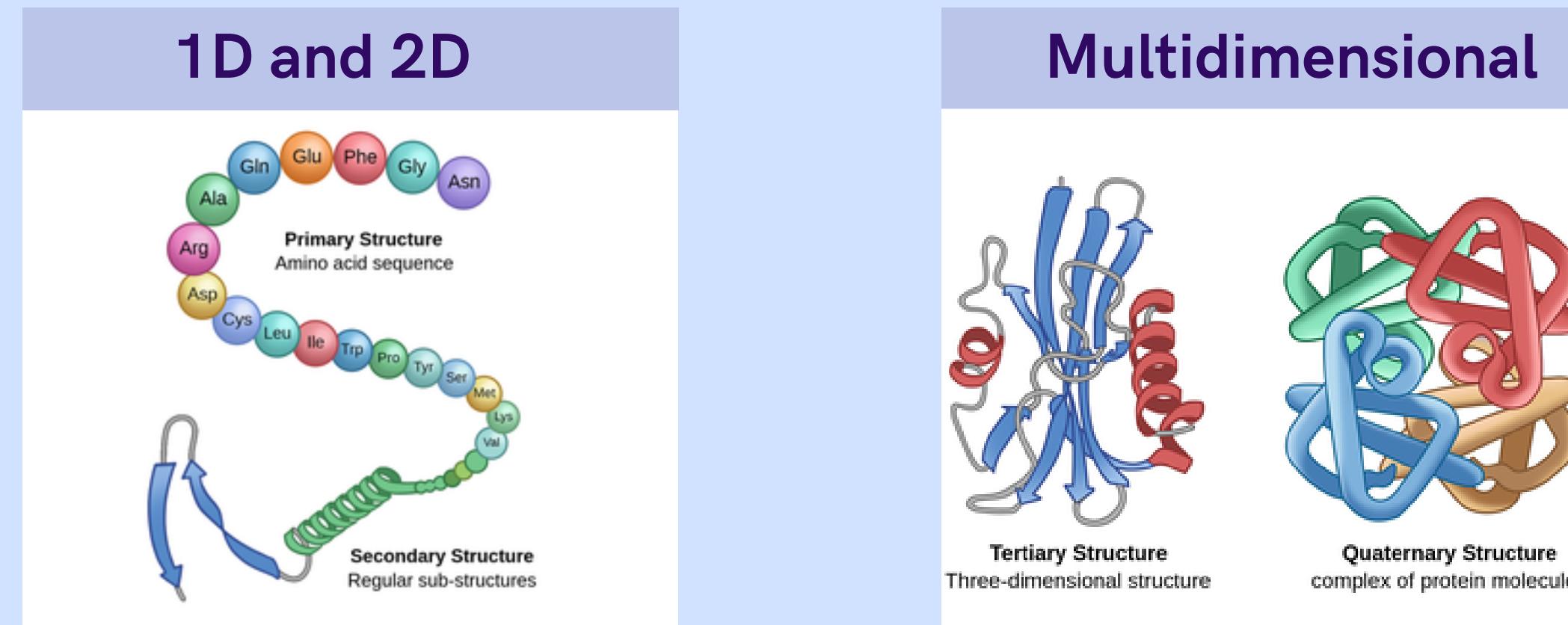
CGSHLVEALYLVCGERGFFYTPKTFVNQHLCGSHLVEALYLVCGERGFFYTPKTCGSHLVEALYLVCGERGFFYTPKTCGSHLVEALYLVCGERGFFYTPKTFVNQHL





Protein Structure

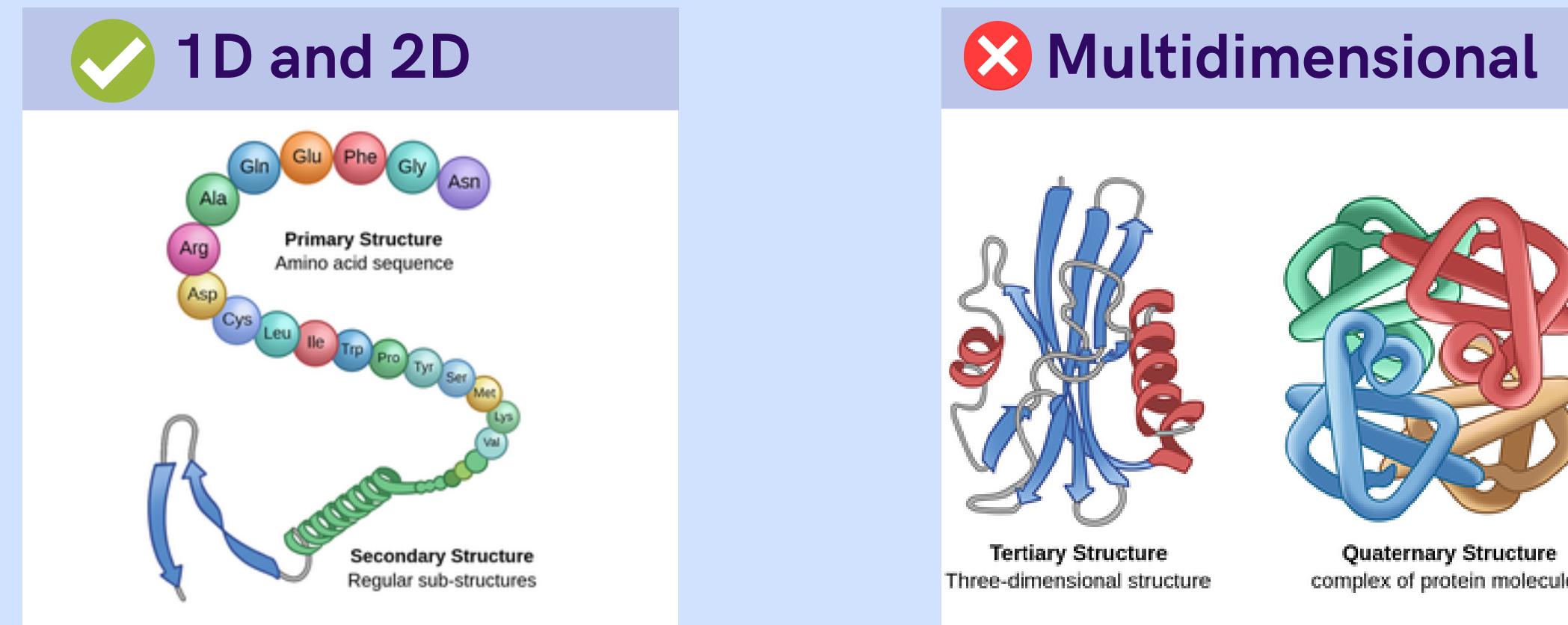
Proteins can be described and folded in terms of structural levels.





Protein Structure

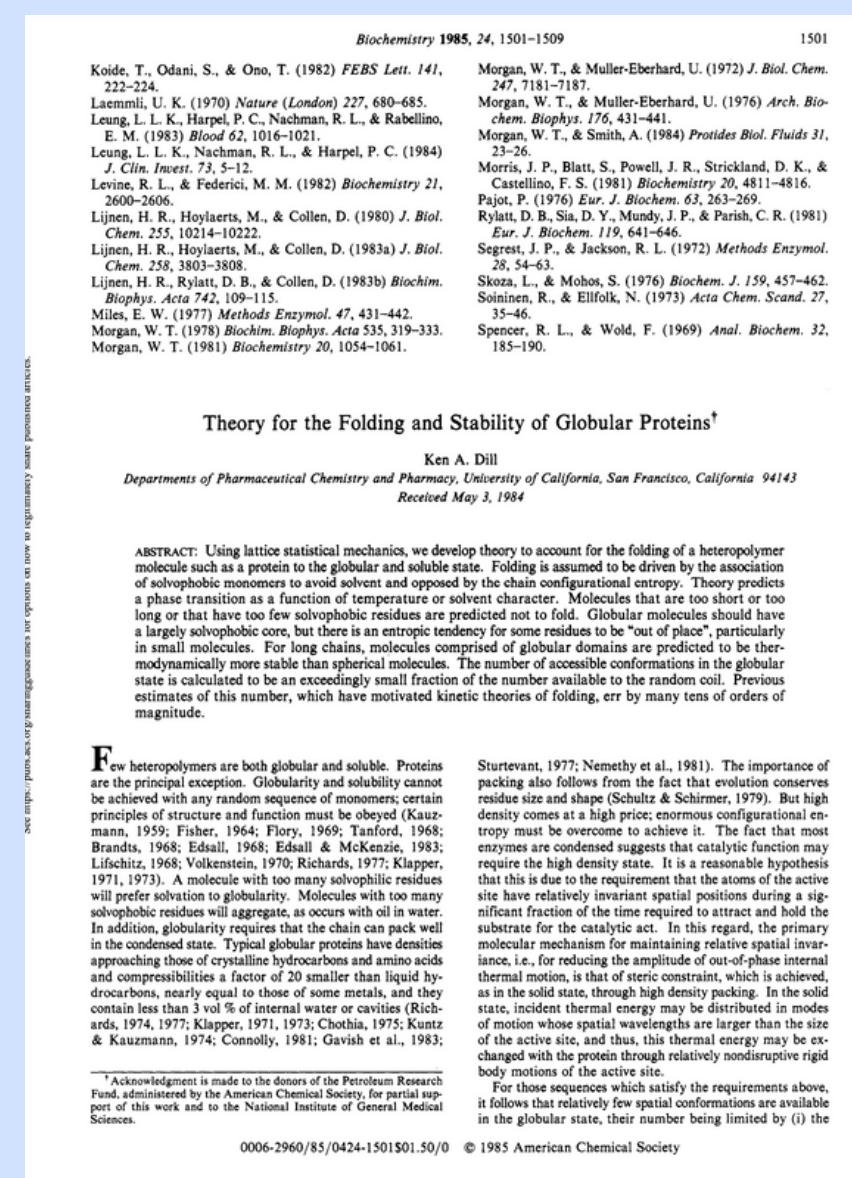
Proteins can be described and folded in terms of structural levels.



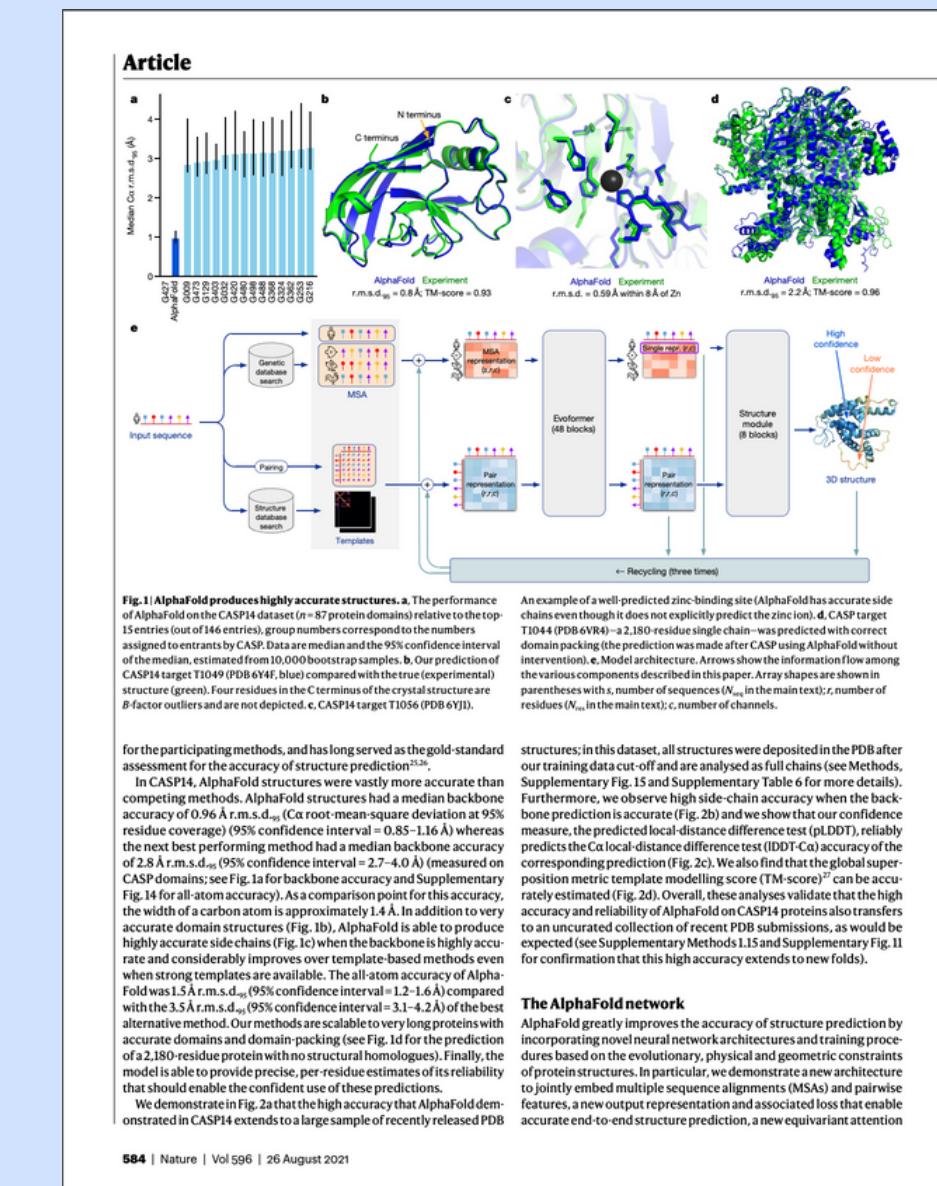


Approaches

Protein folding papers



HP-model, 1985



AlphaFold, 2023

https://doi.org/10.1038/s41586-021-03819-2
 https://doi.org/10.1021/bi00327a032

05



HP-model

Used to predict the most stable conformation of a protein.

Primary Structure (1D)

Description: Primary structure refers to the linear sequence of amino acids that make up the protein.

Importance: Sequence determines the protein's type and function. Even a small change in the sequence can significantly affect the protein's properties.

Input

Secondary Structure (2D)

Description: Folding amino acid chain into regular, repeating structures stabilized by hydrogen bonds.

Importance: Structures provide a level of organization and stability within the protein.

Output

Tertiary Structure (3D)

Description: Folding and arrangement of the secondary structural elements, stabilized by various types of bonds and interactions.

Importance: Vital for the protein's function, as it positions the protein's active sites and other functional domains in specific orientations that are necessary for activity.

Quaternary Structure (4D)

Description: Arrangement of multiple subunits into a functional protein complex. Not all proteins have a quaternary structure.

Importance: Mediates the interaction of proteins with other molecules, including other proteins



AlphaFold

AlphaFold is used to predict the 3D of proteins.

Primary Structure (1D)

Description: Primary structure refers to the linear sequence of amino acids that make up the protein.

Importance: Sequence determines the protein's type and function. Even a small change in the sequence can significantly affect the protein's properties.

Input

Secondary Structure (2D)

Description: Folding amino acid chain into regular, repeating structures stabilized by hydrogen bonds.

Importance: Structures provide a level of organization and stability within the protein.

Substructure

Tertiary Structure (3D)

Description: Folding and arrangement of the secondary structural elements, stabilized by various types of bonds and interactions.

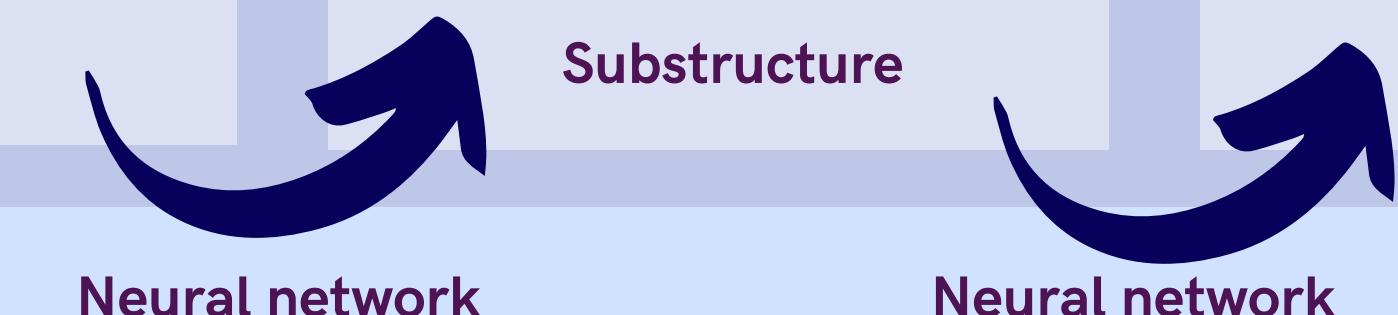
Importance: Vital for the protein's function, as it positions the protein's in specific orientations that are necessary for activity.

Output

Quaternary Structure (4D)

Description: Arrangement of multiple subunits into a functional protein complex. Not all proteins have a quaternary structure.

Importance: Mediates the interaction of proteins with other molecules, including other proteins





AlphaFold VS HP-model

The figure consists of two side-by-side panels, each divided into two vertical sections: Pros (left) and Cons (right).

HP-model Panel:

- Pros:**
 - Simplicity
 - Computational efficiency
 - Educational Tool
- Cons:**
 - Over Simplification
 - Limited Predictive Power

AlphaFold Panel:

- Pros:**
 - High Accuracy
 - Complexity Handling
- Cons:**
 - Computational expensive
 - Dependent on training data
 - Black box

Inset Figures:

- HP-model Inset:** A small image of a scientific article titled "Theory for the Folding and Stability of Globular Proteins" by K. A. Dill, published in *Biochemistry*, Vol. 24, 1501-1509. The inset shows the abstract and part of the introduction.
- AlphaFold Inset:** A small image of a scientific article titled "AlphaFold: Protein Structure Prediction via Deep Learning" by J. J. Hou et al., published in *Nature*, Vol. 596, 582-587. The inset shows the abstract and part of the introduction.

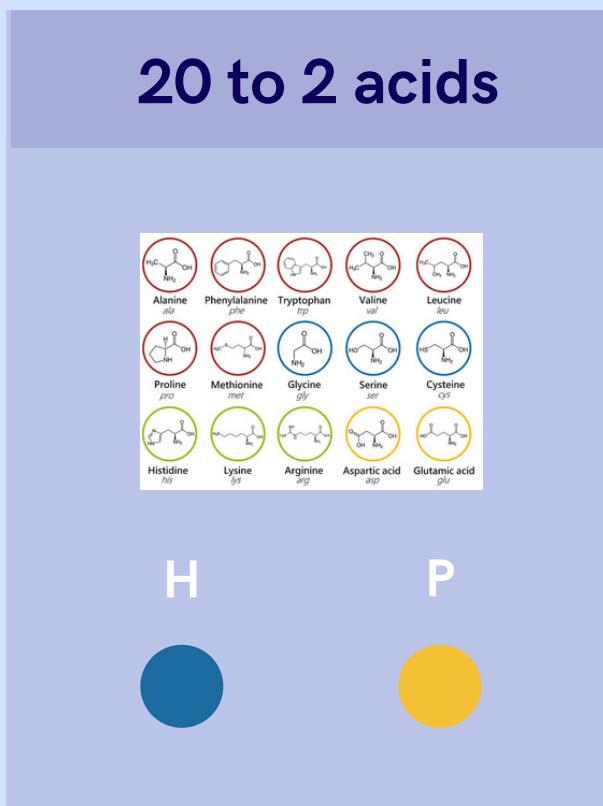


AlphaFold VS HP-model



HP-model

A binary perspective

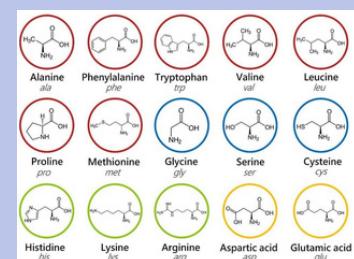




HP-model

A binary perspective

20 to 2 acids



H



P



Insulin - 51 amino acids

Chain A - 21 amino acids

GIVEQCCTSICSLYQLENYCN



Chain B - 30 amino acids

FVNQHLCGSHLVEALYLVCGERGFFYTPKT

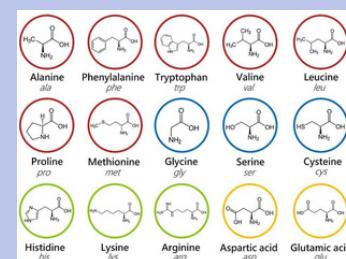




HP-model

A binary perspective

20 to 2 acids



H



P



Insulin - 51 amino acids

Chain A - 21 amino acids
GIVEQCCTSICSLYQLENYCN



Chain B - 30 amino acids
FVNQHLCGSHLVEALYLVCGERGFFYTPKT

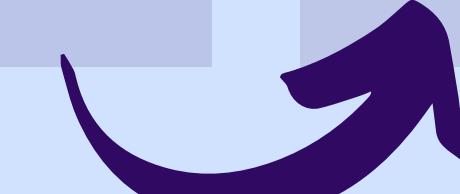


Insulin - 51 amino acids

Chain A - 21 amino acids
HPPPHHHPPPPHHHHPPPH



Chain B - 30 amino acids
PHHHPPPHPPPPHHPPPHHHPPPH

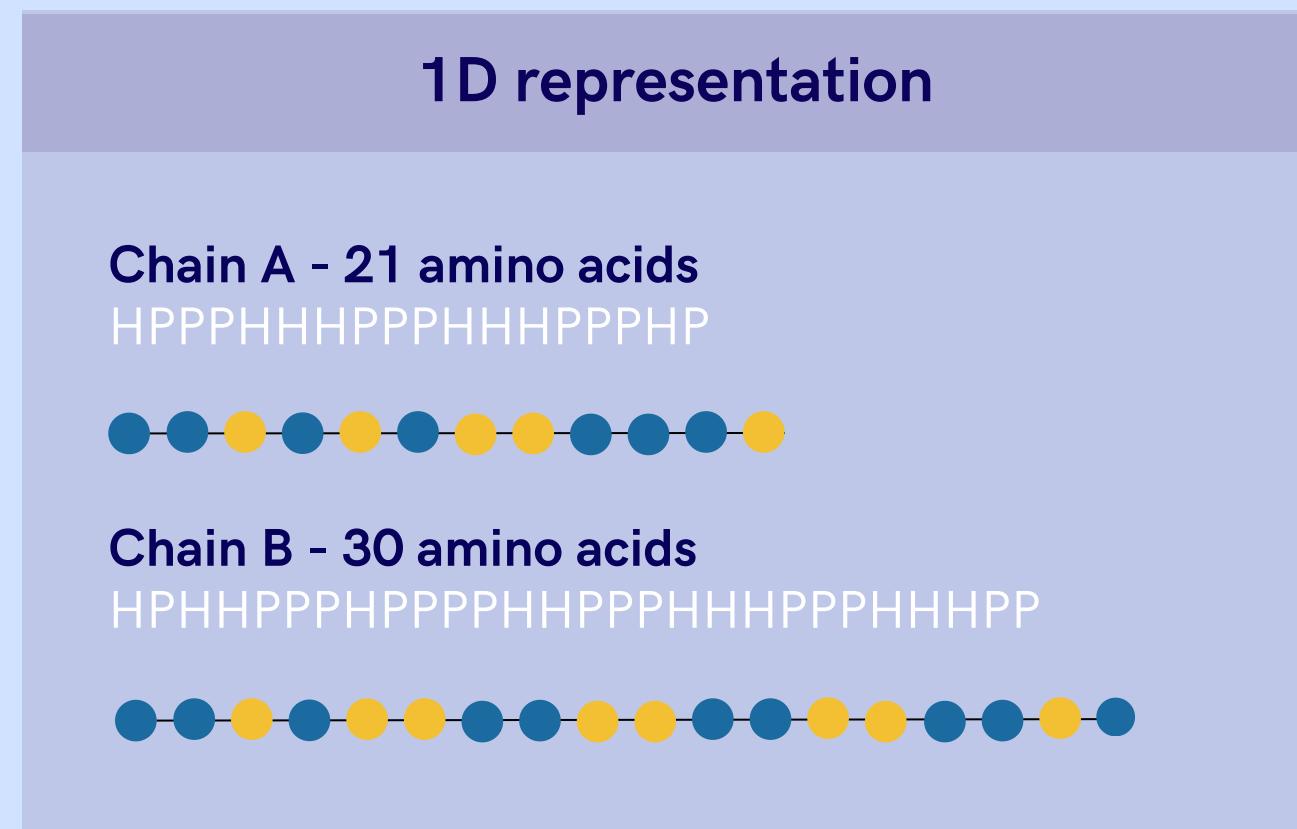


Simplification



HP-model

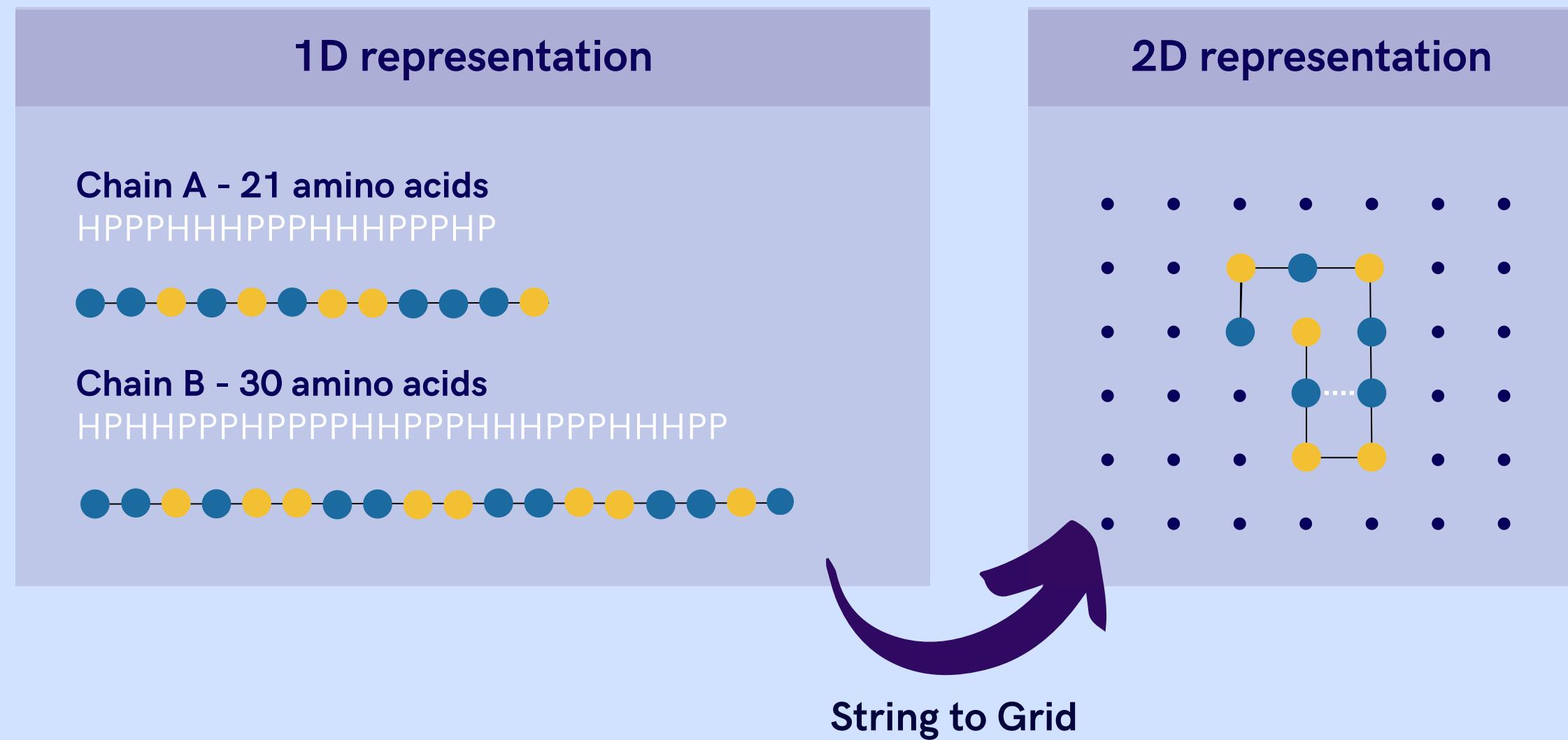
From 1D to 2D





HP-model

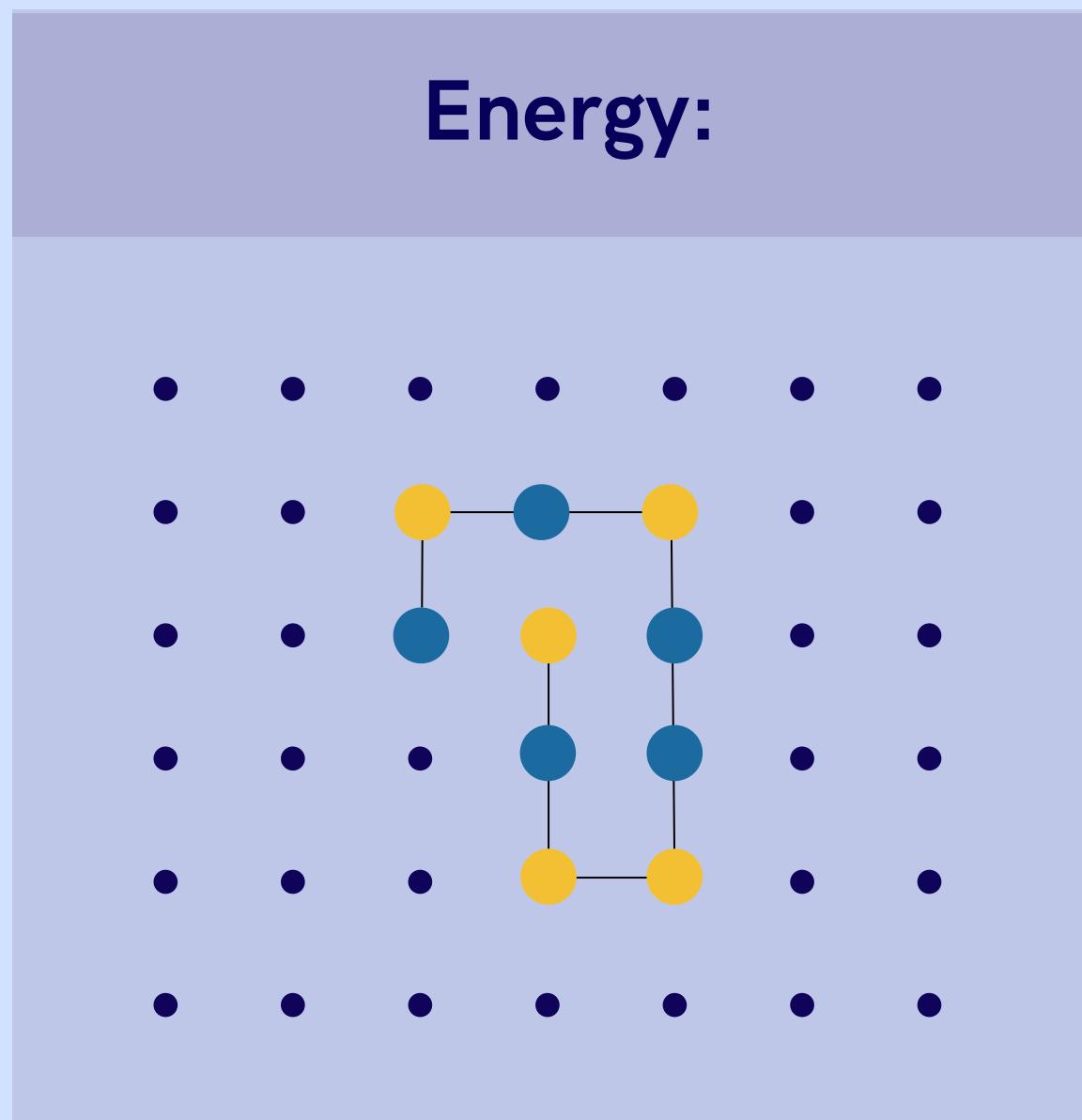
From 1D to 2D

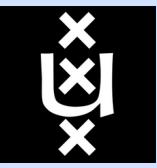




HP-model

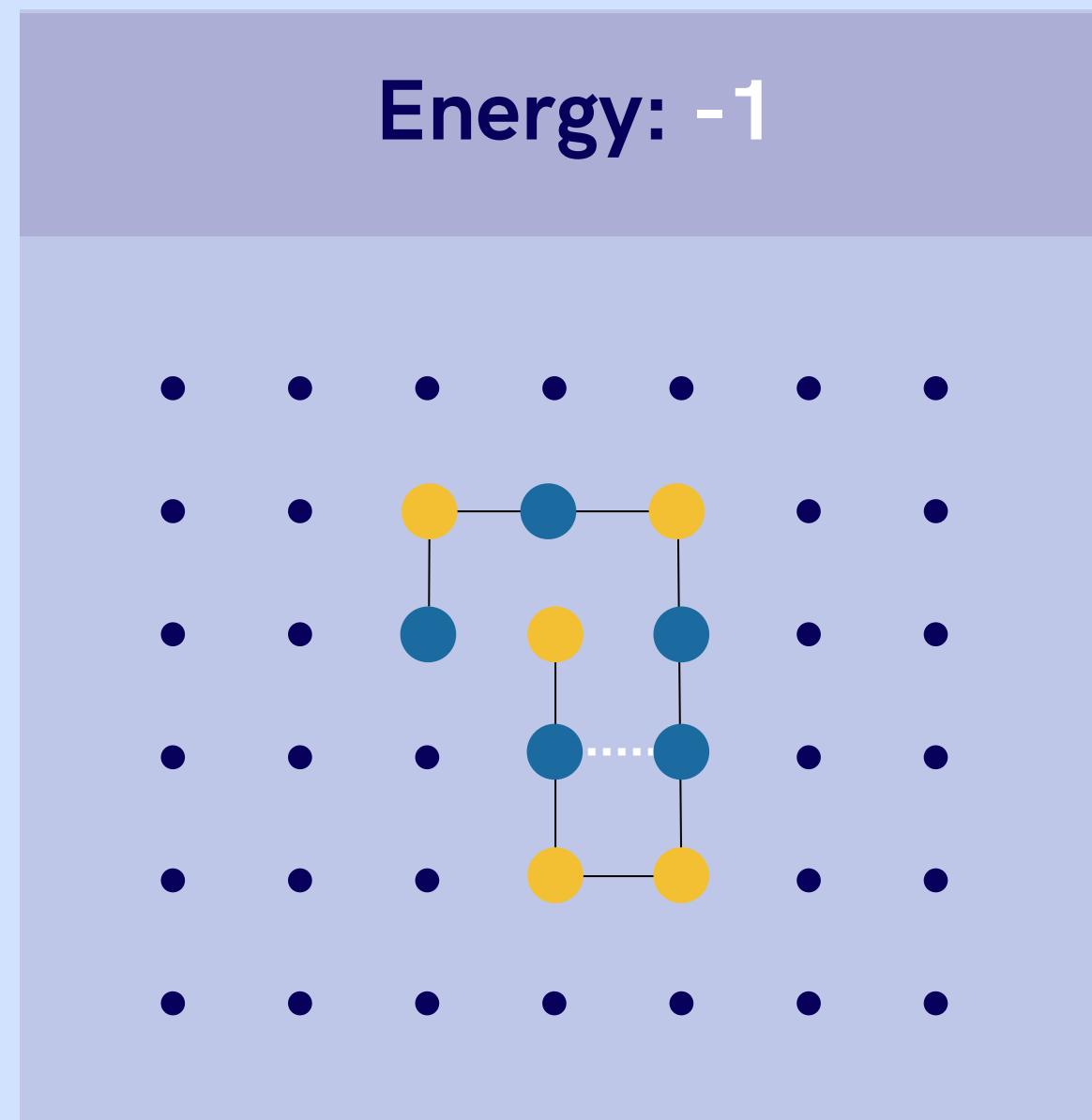
Calculating the stability with H-bonds





HP-model

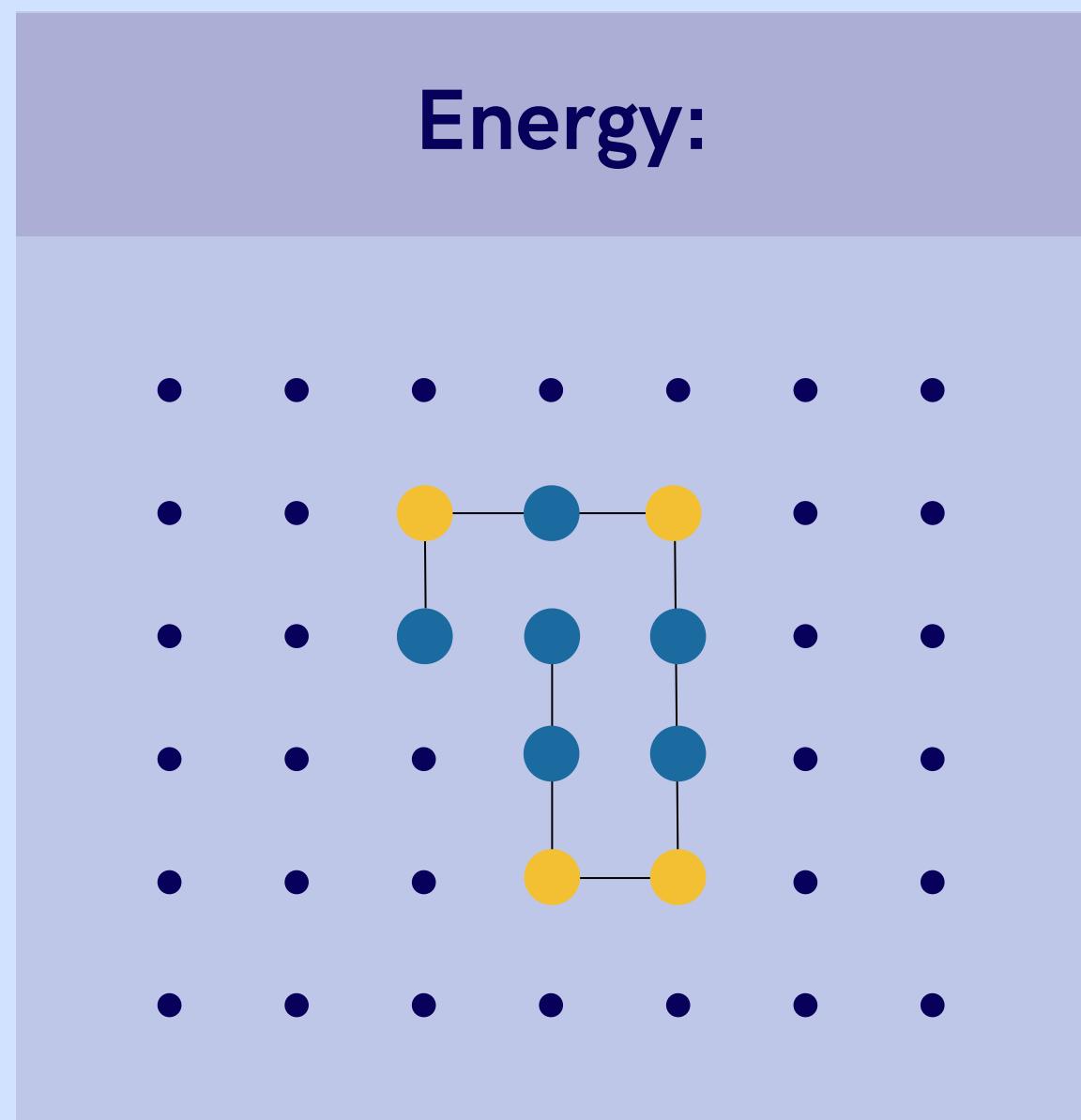
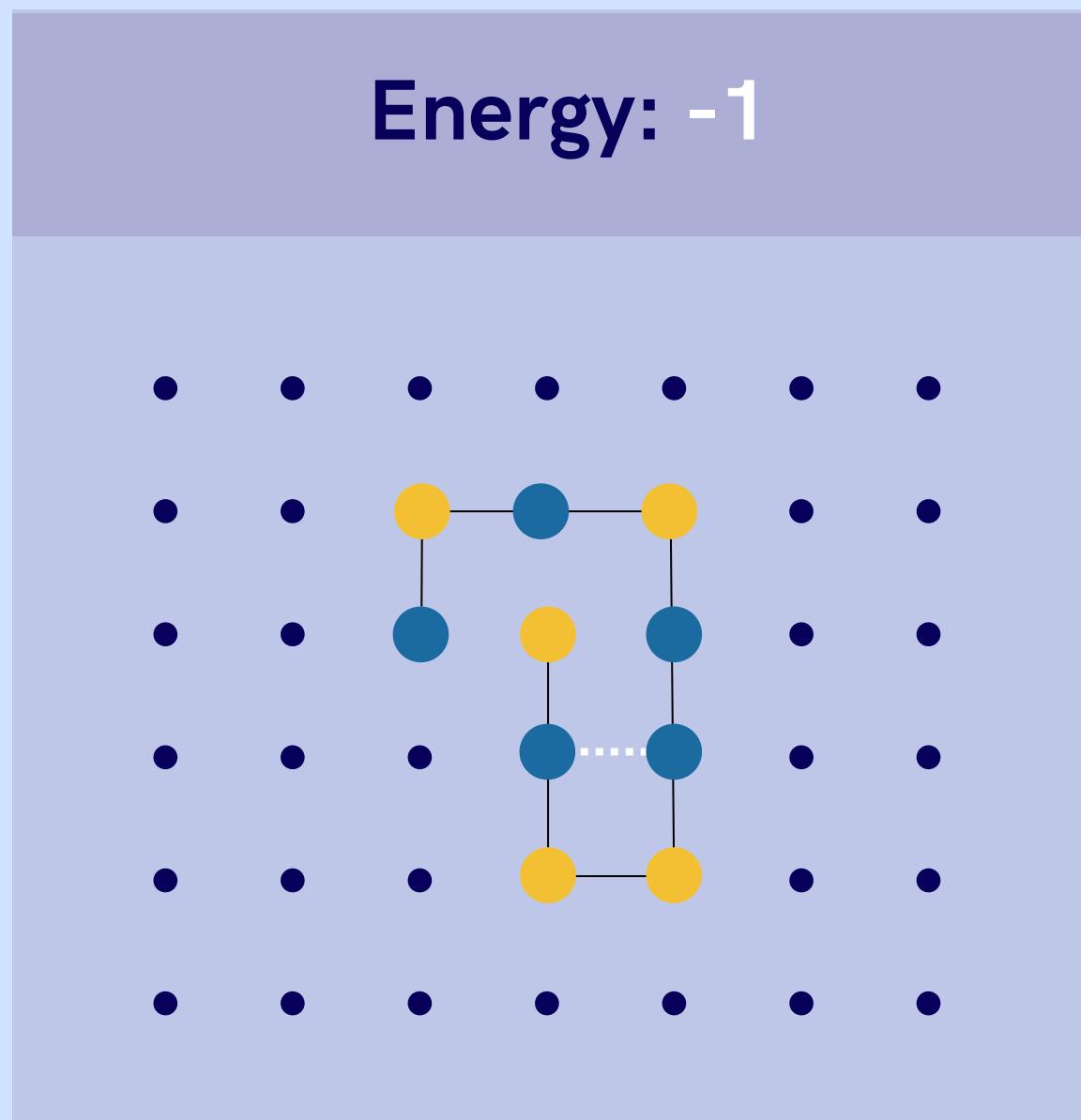
Calculating the stability with H-bonds





HP-model

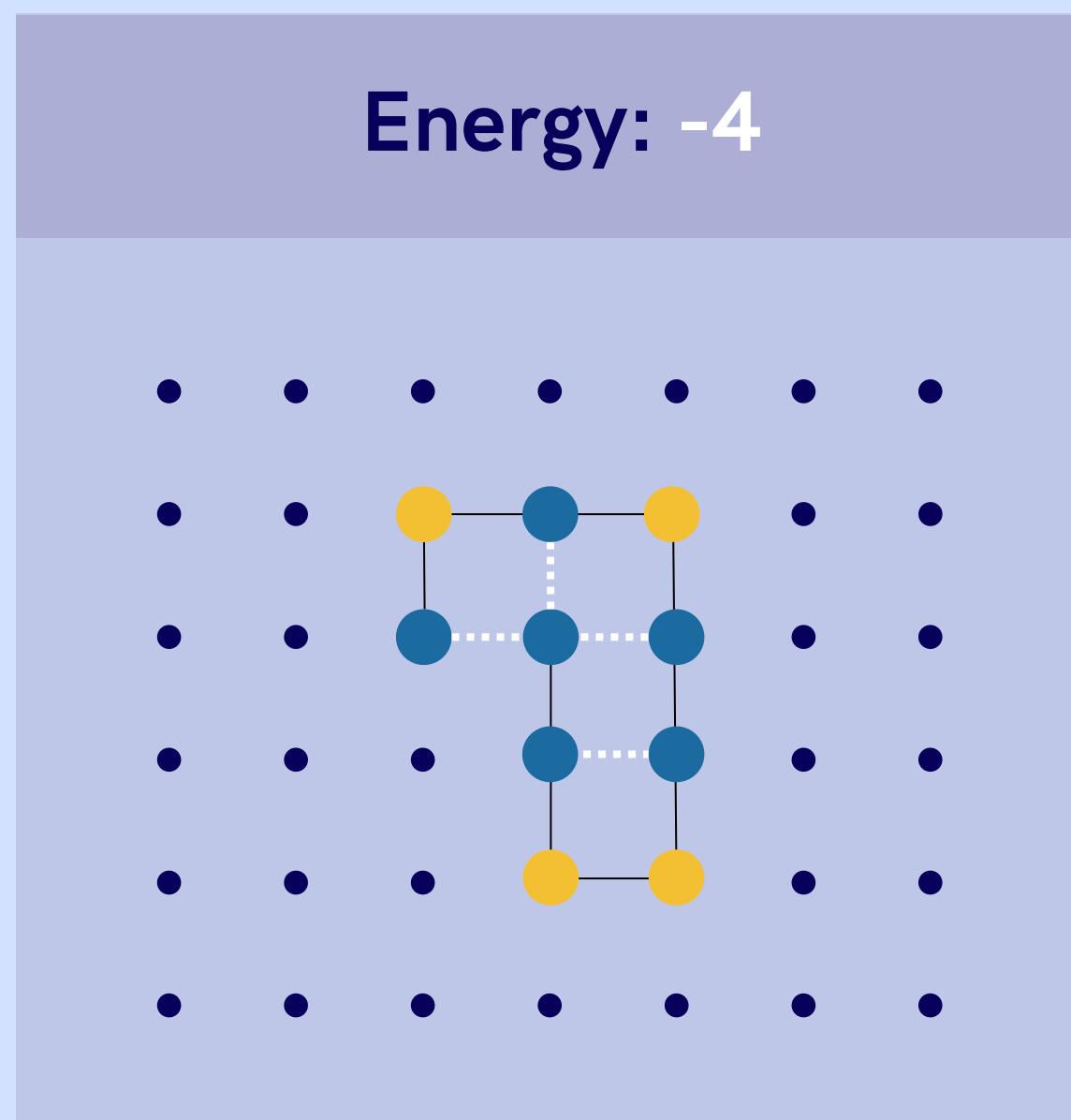
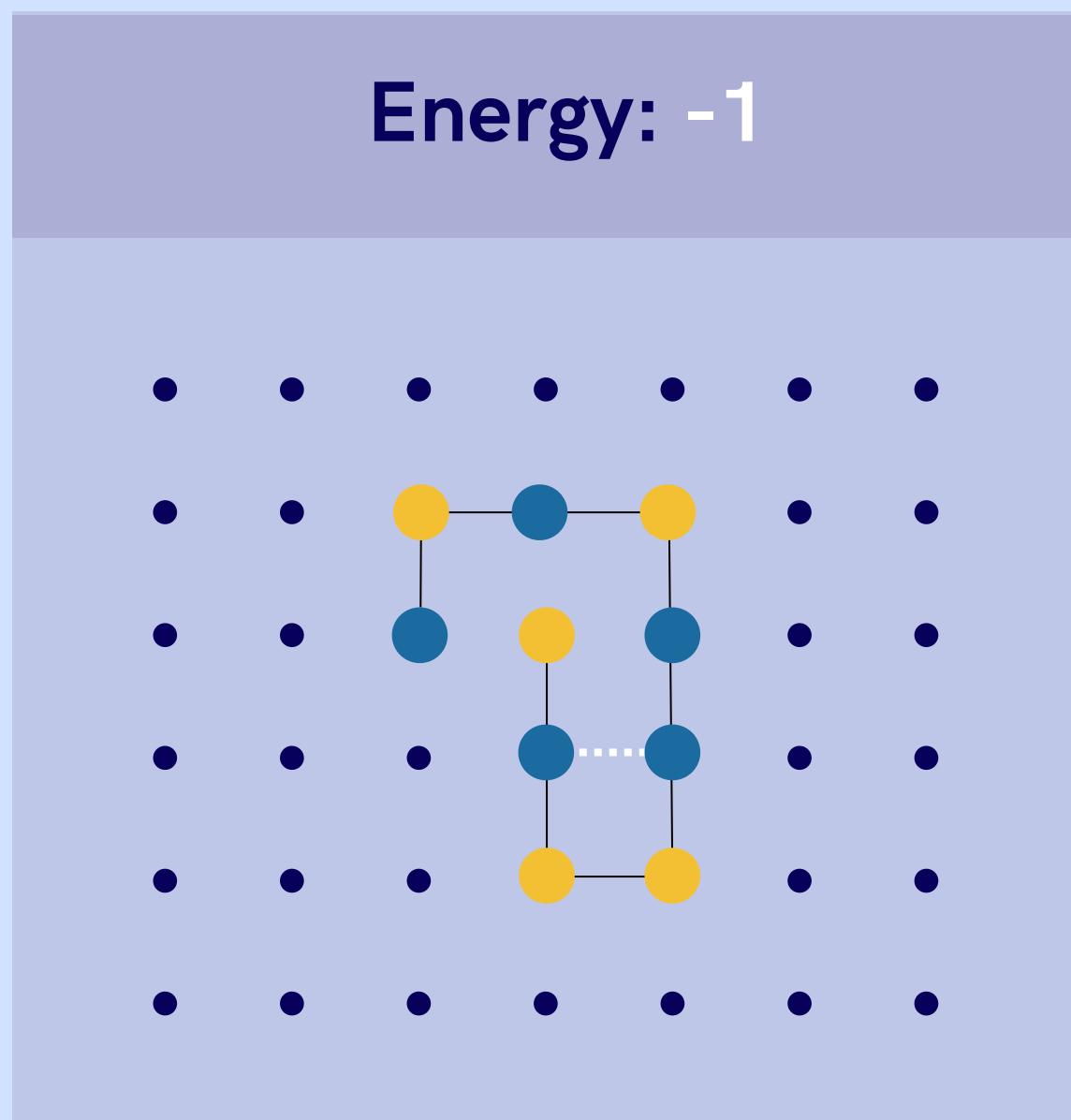
Calculating the stability with H-bonds

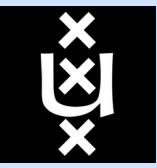




HP-model

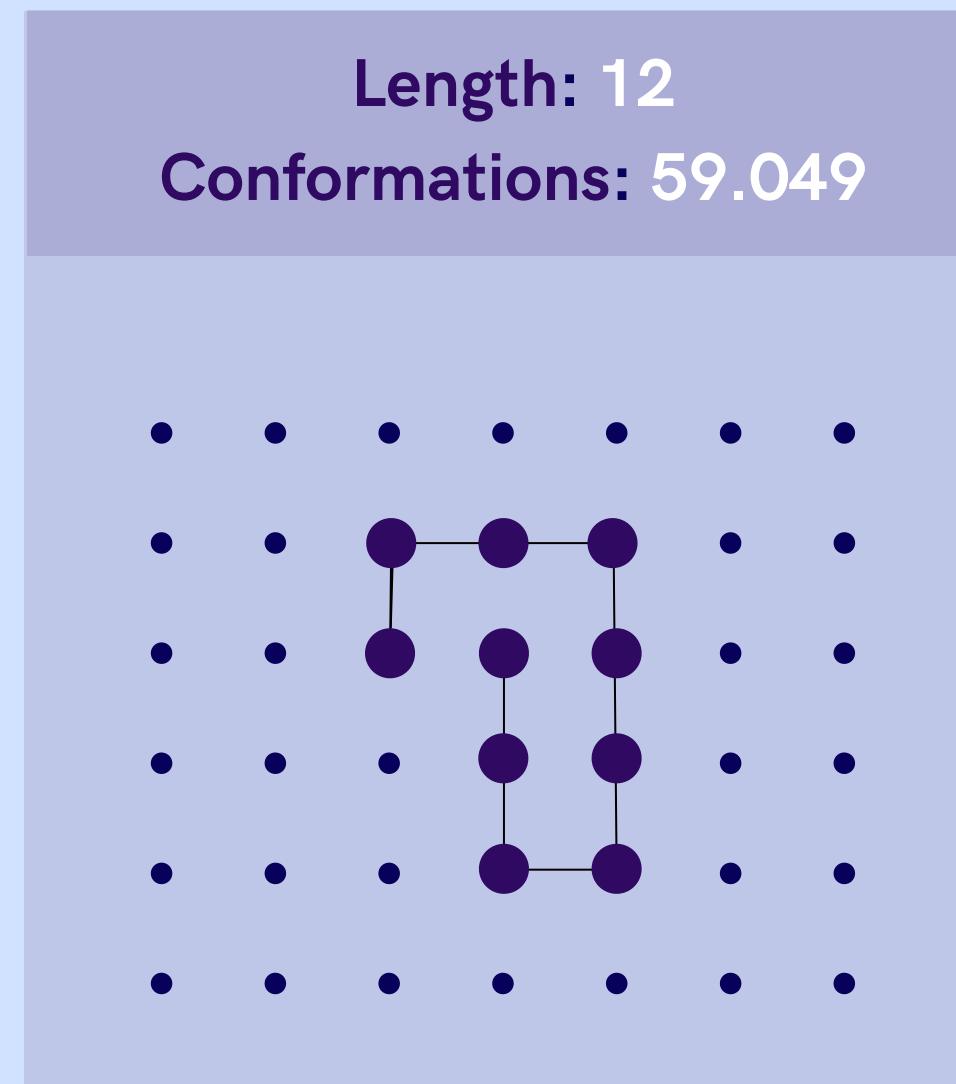
Calculating the stability with H-bonds





Protein Folding, NP-hard?

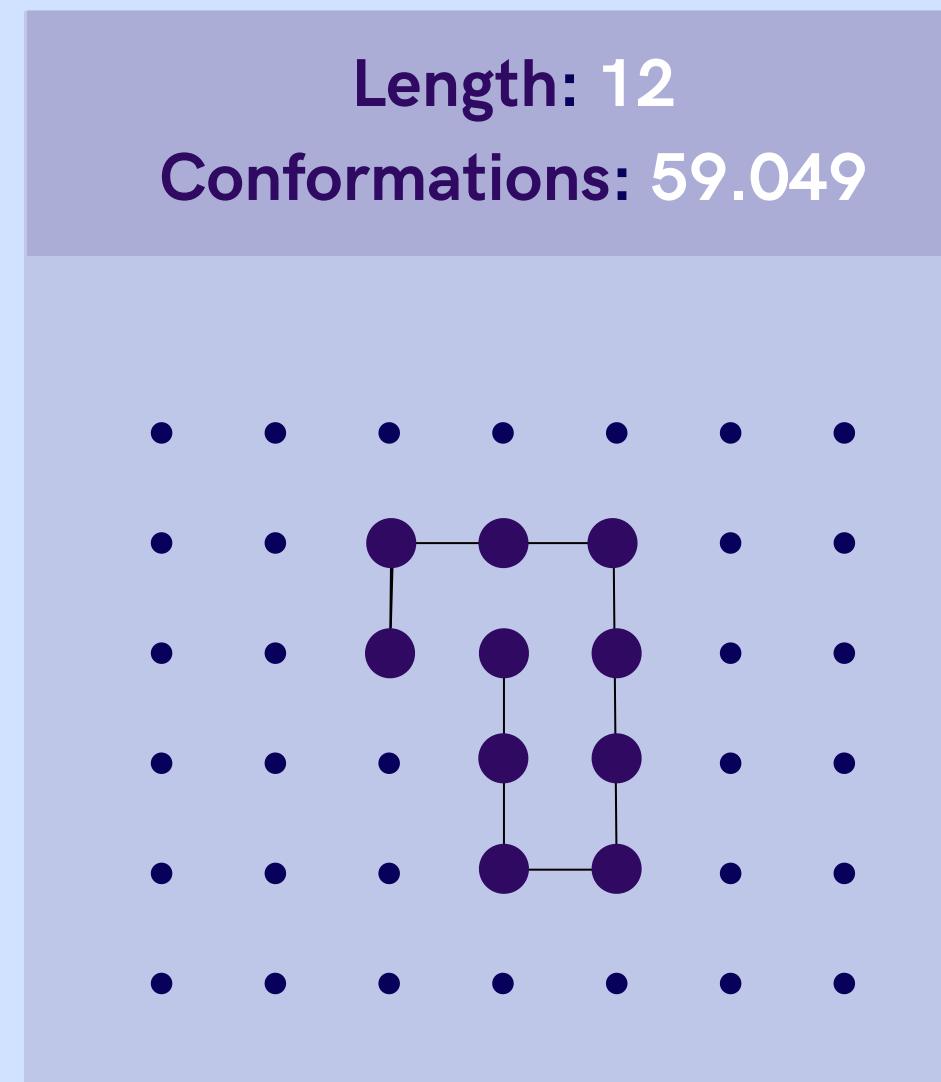
Finding most stable conformation optimal fold is considered NP-hard





Protein Folding, NP-hard?

Finding most stable conformation optimal fold is considered NP-hard

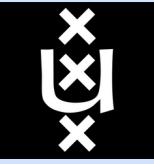


Protein Length	Number of solutions
5	27
10	6561
15	1.594.323

The problem scales exponentially with the length of the protein: $1 \times 1 \times 3(n-2)$



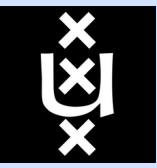
Scope



Uniform Random Sampling

Can we ensure the **uniform randomness** in initial populations before folding proteins (e.g. with Genetic algorithms)?

What is the impact of **different sampling techniques** on the quality and diversity of solutions?

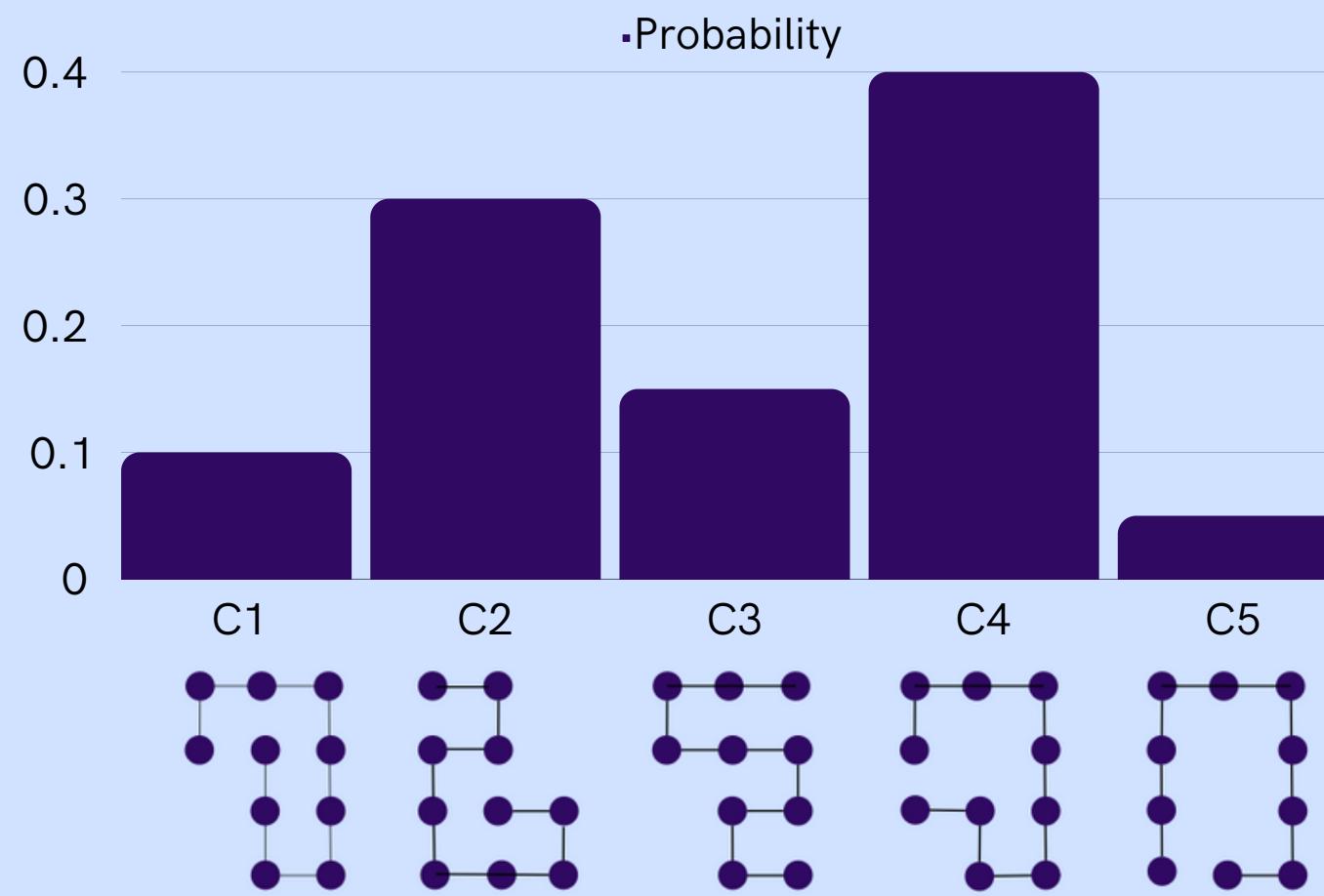


Uniform Random Sampling

What does uniform random mean?



Not Uniform Random



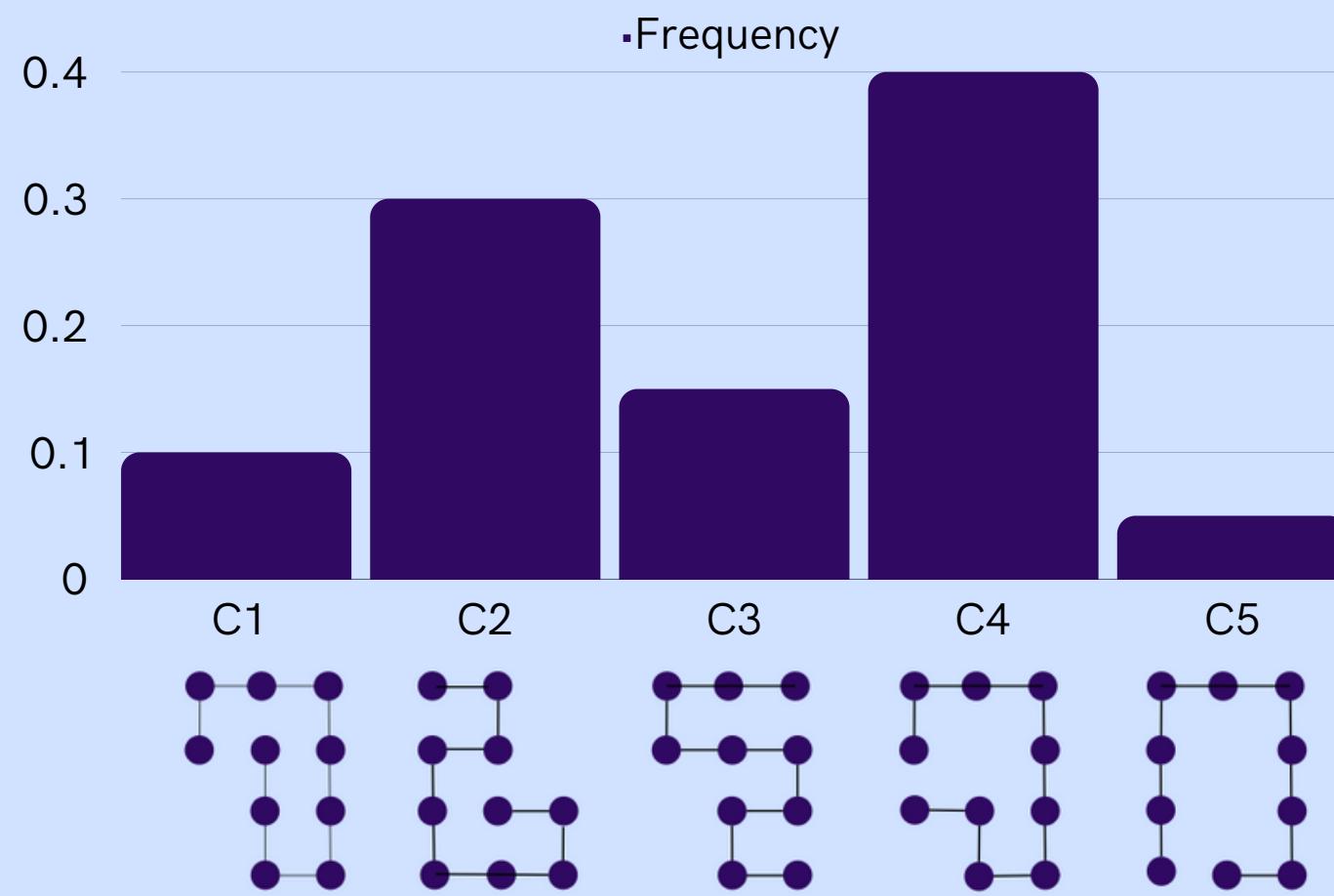


Uniform Random Sampling

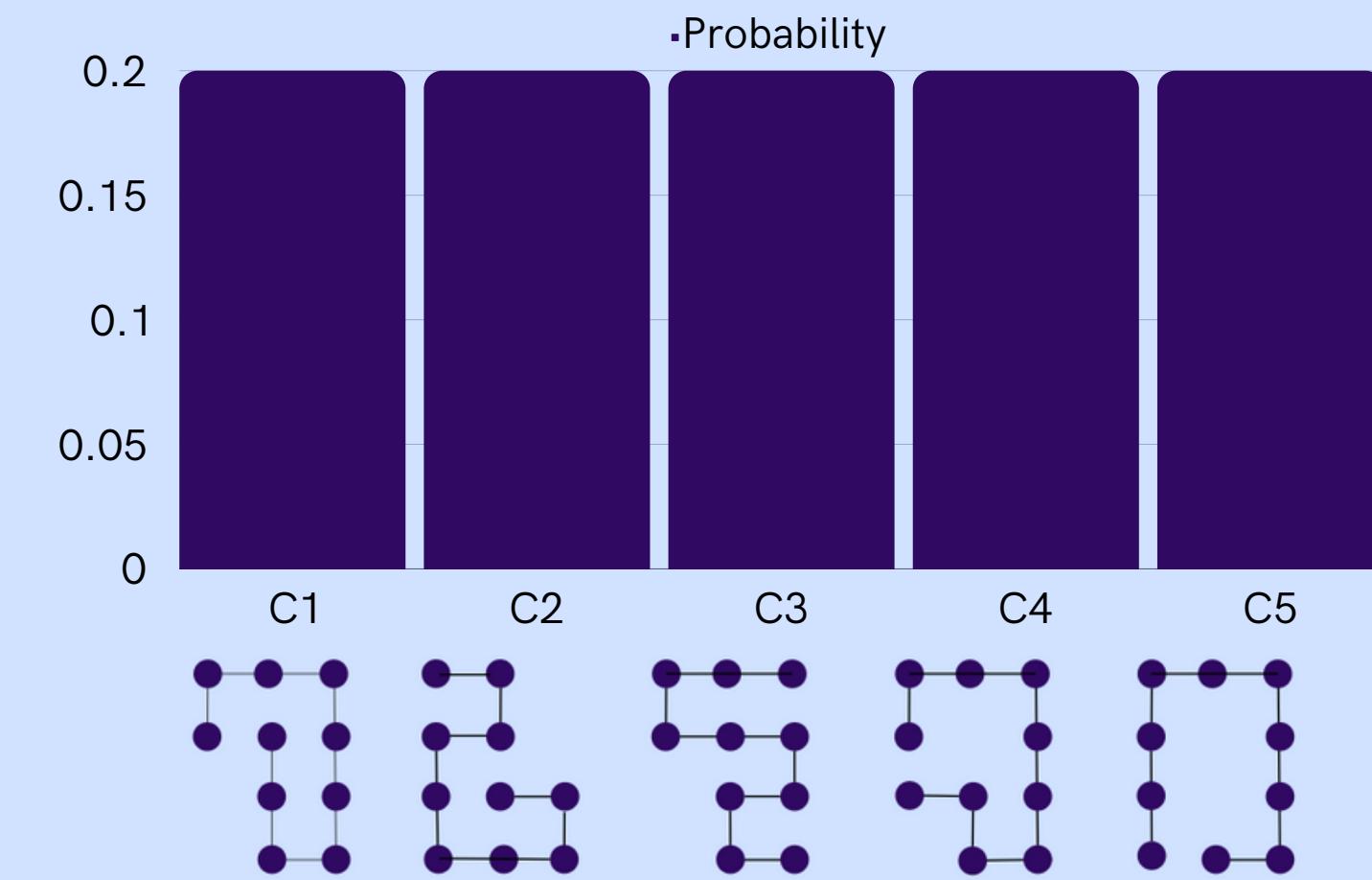
What does uniform random mean?

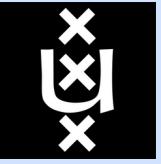


Not Uniform Random



Uniform Random



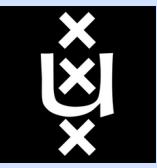


Uniform Random Sampling

Why do we want uniform random samples?

Limited Diversity

- **Problem:** Non-uniform random sampling leads to a population that is not diverse.
- **Result:** The algorithm converges to a local optimum. It doesn't have the diversity needed to explore other potentially superior solutions.



Uniform Random Sampling

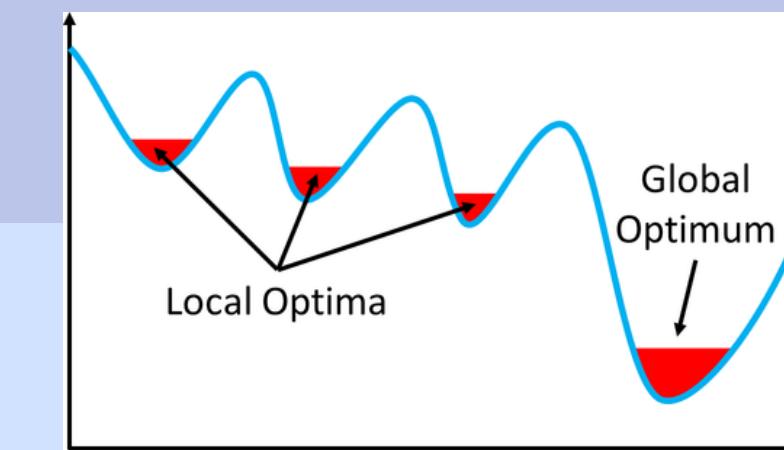
Why do we want uniform random samples?

Limited Diversity

- **Problem:** Non-uniform random sampling leads to a population that is not diverse.
- **Result:** The algorithm converges to a local optimum. It doesn't have the diversity needed to explore other potentially superior solutions.

Premature Convergence

- **Problem:** The algorithm settles on a suboptimal solution too quickly.
- **Result:** It doesn't explore enough of the solution space to find the global optimum.



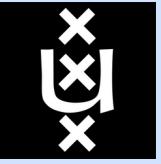


Experiment



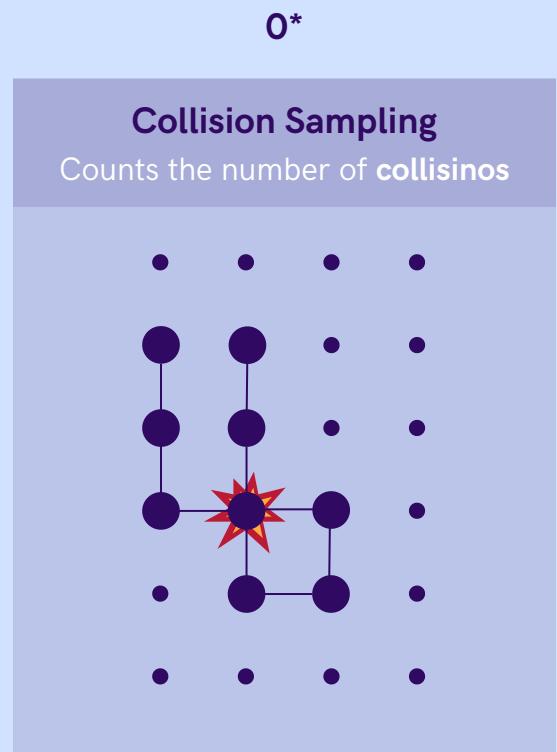
Experiment

6 Different sampling methods



Experiment

6 Different sampling methods



Lengths $\{5, 10, 15 \dots 200\}$

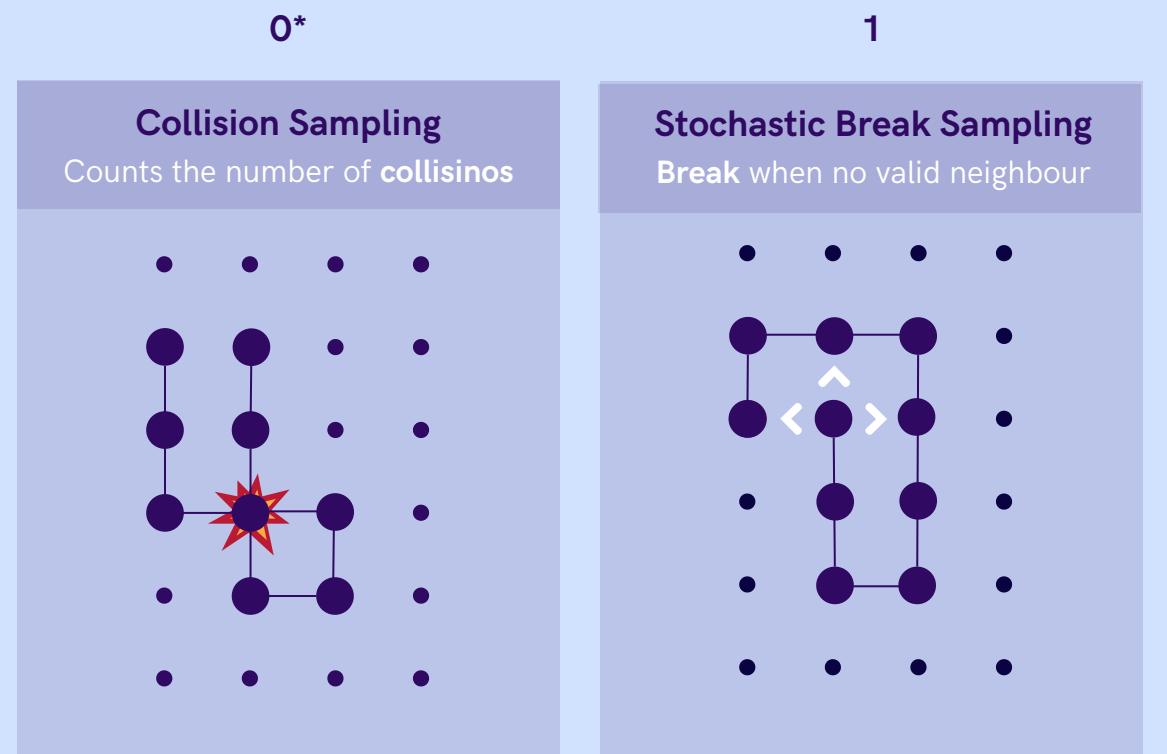
Placement Stochastic

Backtracking ---



Experiment

6 Different sampling methods



Lengths $\{5, 10, 15 \dots 200\}$

Placement Stochastic

Backtracking ---

Lengths $\{5, 10, 15 \dots 200\}$

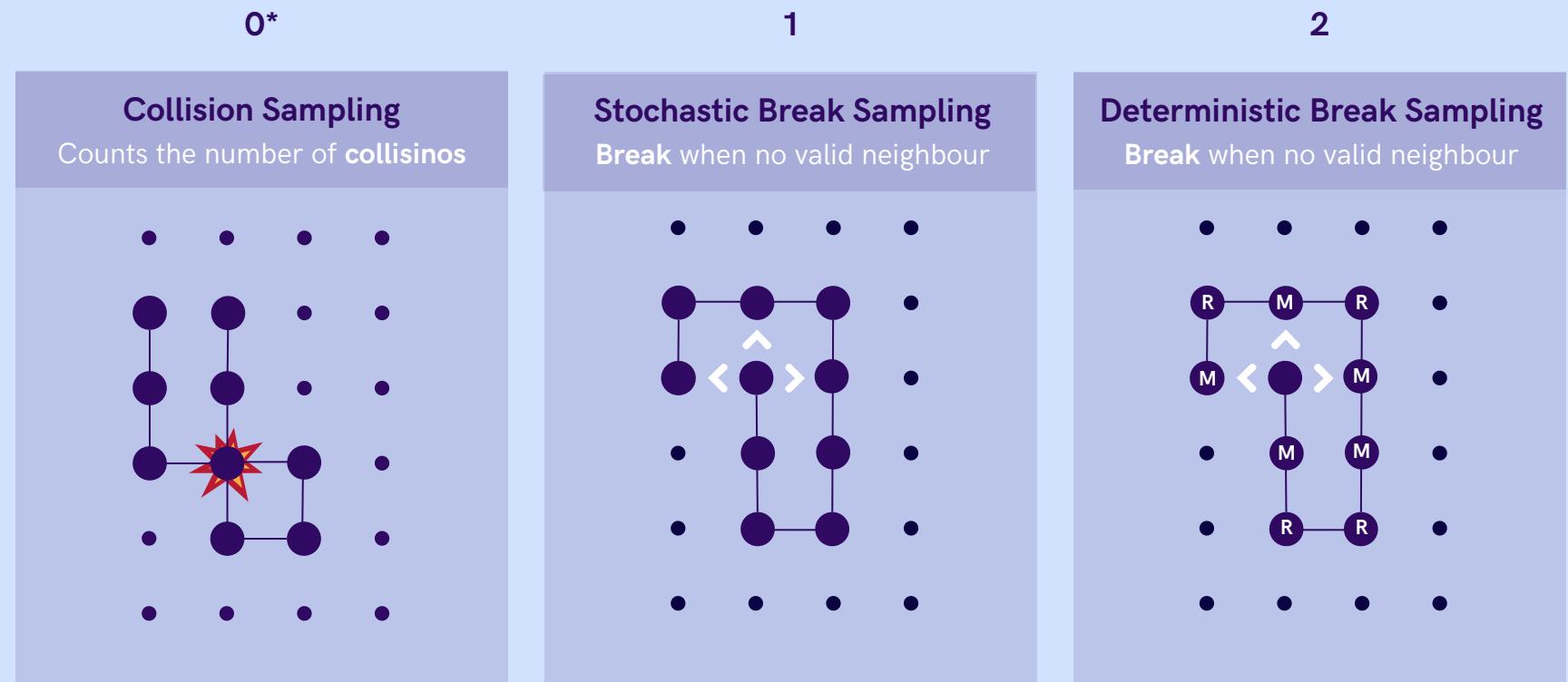
Placement Stochastic

Backtracking ---



Experiment

6 Different sampling methods



Lengths {5, 10, 15 ... 200}

Placement Stochastic

Backtracking ---

Lengths {5, 10, 15 ... 200}

Placement Stochastic

Backtracking ---

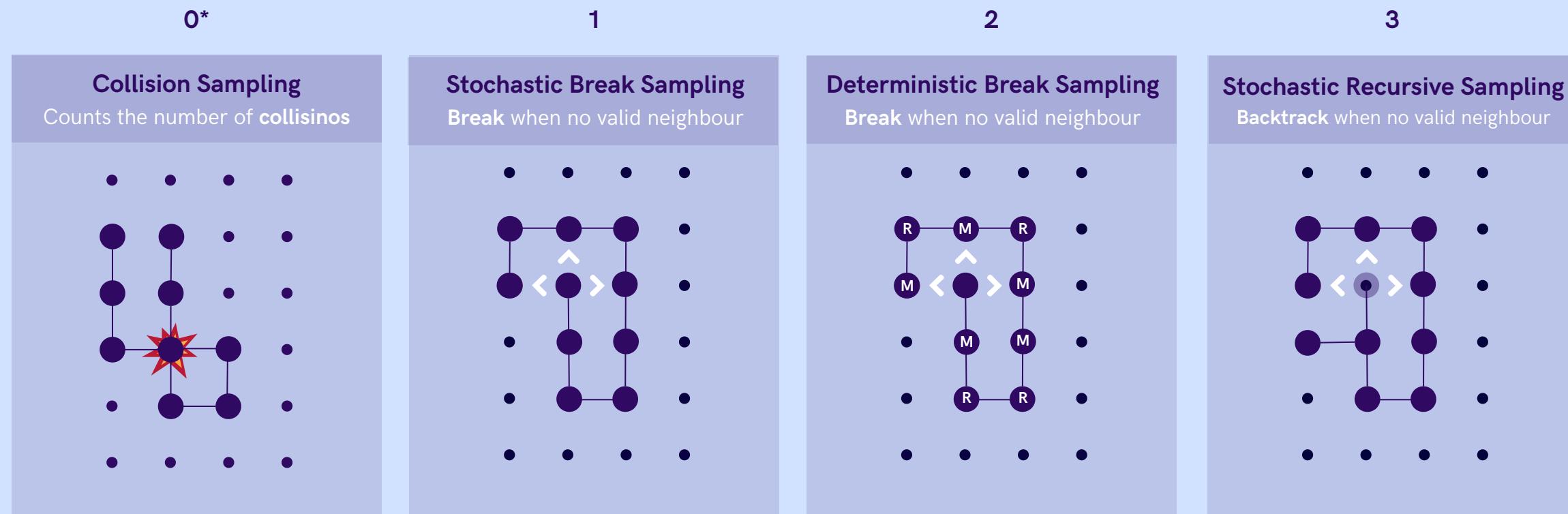
Lengths {5, 10, 15 ... 200}

Deterministic



Experiment

6 Different sampling methods



Lengths {5, 10, 15 ... 200}

Placement Stochastic

Backtracking ---

Lengths {5, 10, 15 ... 200}

Placement Stochastic

Backtracking ---

Lengths {5, 10, 15 ... 200}

Placement Deterministic

Backtracking ---

Lengths {5, 10, 15 ... 100}

Placement Stochastic

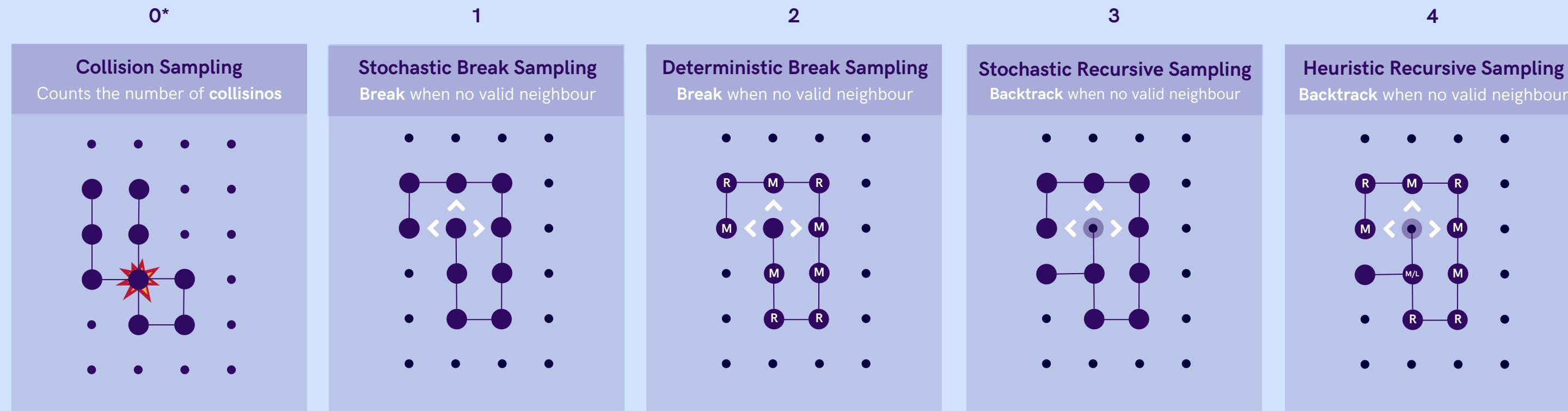
Backtracking Stochastic

[M R L M R L
L M R L M R
R L M R L M]



Experiment

6 Different sampling methods



Lengths $\{5, 10, 15 \dots 200\}$

$\{5, 10, 15 \dots 200\}$

$\{5, 10, 15 \dots 200\}$

$\{5, 10, 15 \dots 100\}$

$\{5, 10, 15 \dots 30\}$

Placement Stochastic

Stochastic

Deterministic

Stochastic

Stochastic

Backtracking ---

Stochastic

Heuristic

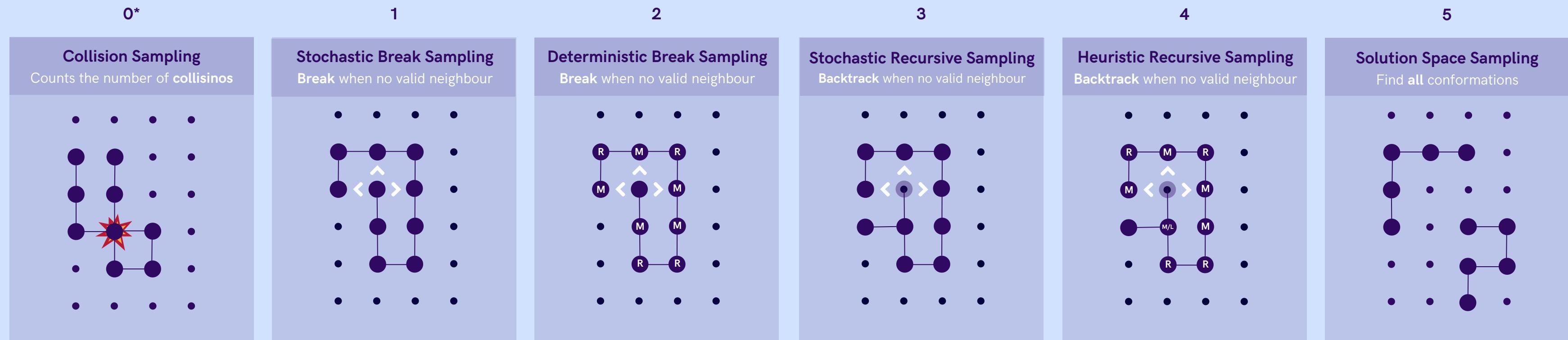
$$\begin{bmatrix} M & R & L & M & R & L \\ L & M & R & L & M & R \\ R & L & M & R & L & M \end{bmatrix}$$

$$\begin{bmatrix} L & L & L & L & L & L \\ M & M & M & M & M & M \\ R & R & R & R & R & R \end{bmatrix}$$



Experiment

6 Different sampling methods



Lengths {5, 10, 15 ... 200}

Placement Stochastic

Backtracking ---

{5, 10, 15 ... 200}

Stochastic

{5, 10, 15 ... 200}

Deterministic

{5, 10, 15 ... 100}

Stochastic

Stochastic

[
M R L M R L
L M R L M R
R L M R L M
]

{5, 10, 15 ... 30}

Stochastic

Heuristic

[
L L L L L L
M M M M M M
R R R R R R
]



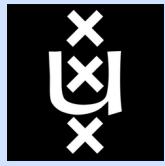
Experiment

Sampling information and parameters

X = number of folded proteins per length

N = protein lengths ∈ {5, 10, 15... N}.

	Amino Acid Length	Num Hydrophobic	Num Polar	1D protein	2D protein	Amino Acids on Grid	Trimmed 2D protein	Shape 2D protein	Amino Acid Order	H-Bonds	H-Ratio	Recursions	Time Taken (s)
19995	100	29	71	['P', 'P', 'P', 'P', 'P', 'P', 'P', 'P', 'P', ...]	[[.....]]	100	[[.....]]	(16, 18)	[('P', (100, 100)), ('P', (100, 101)), ('P', (...	4	0.04	6	7.152557e-07
19996	100	67	33	['H', 'P', 'H', 'H', 'H', 'H', 'H', 'H', 'H', ...]	[[.....]]	100	[[.....]]	(12, 17)	[('H', (100, 100)), ('P', (101, 100)), ('H', (...	19	0.19	0	1.192093e-06
19997	100	80	20	['H', 'H', 'P', 'H', 'H', 'H', 'H', 'H', 'H', ...]	[[.....]]	100	[[.....]]	(11, 23)	[('H', (100, 100)), ('H', (101, 100)), ('P', (...	24	0.24	0	0.000000e+00
19998	100	96	4	['H', 'H', 'H', 'H', 'H', 'H', 'H', 'H', 'H', ...]	[[.....]]	100	[[.....]]	(13, 14)	[('H', (100, 100)), ('H', (100, 99)), ('H', (1...	56	0.56	1	0.000000e+00
19999	100	55	45	['H', 'H', 'P', 'H', 'P', 'H', 'H', 'P', 'H', ...]	[[.....]]	100	[[.....]]	(24, 13)	[('H', (100, 100)), ('H', (101, 100)), ('P', (...	9	0.09	216331	0.000000e+00

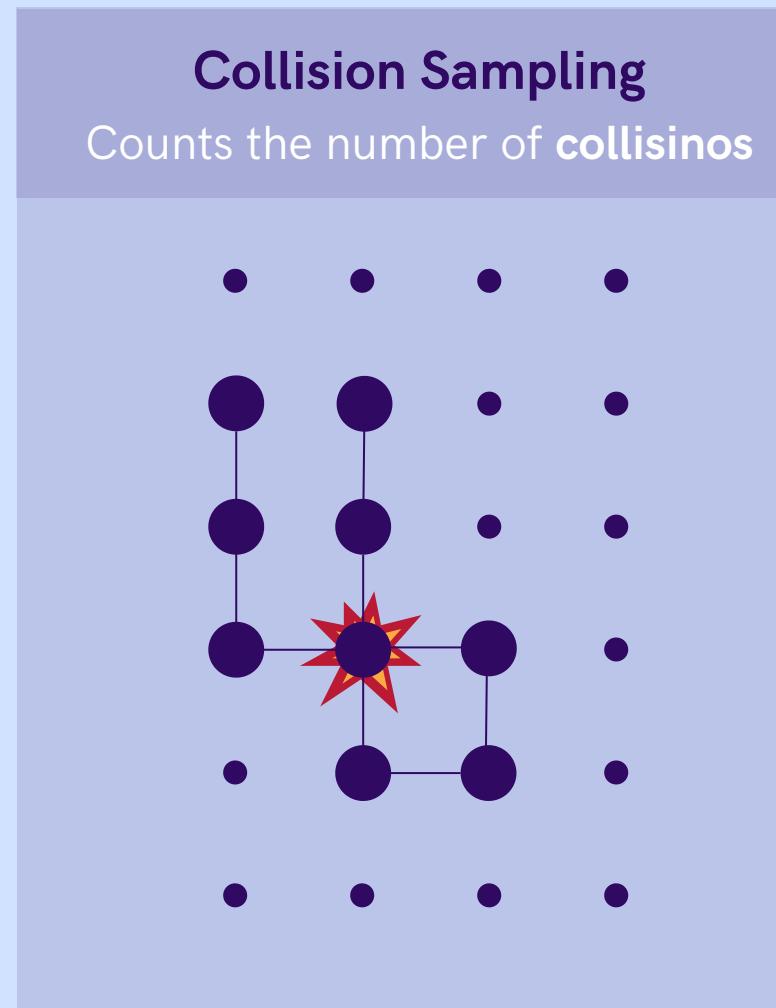


Results

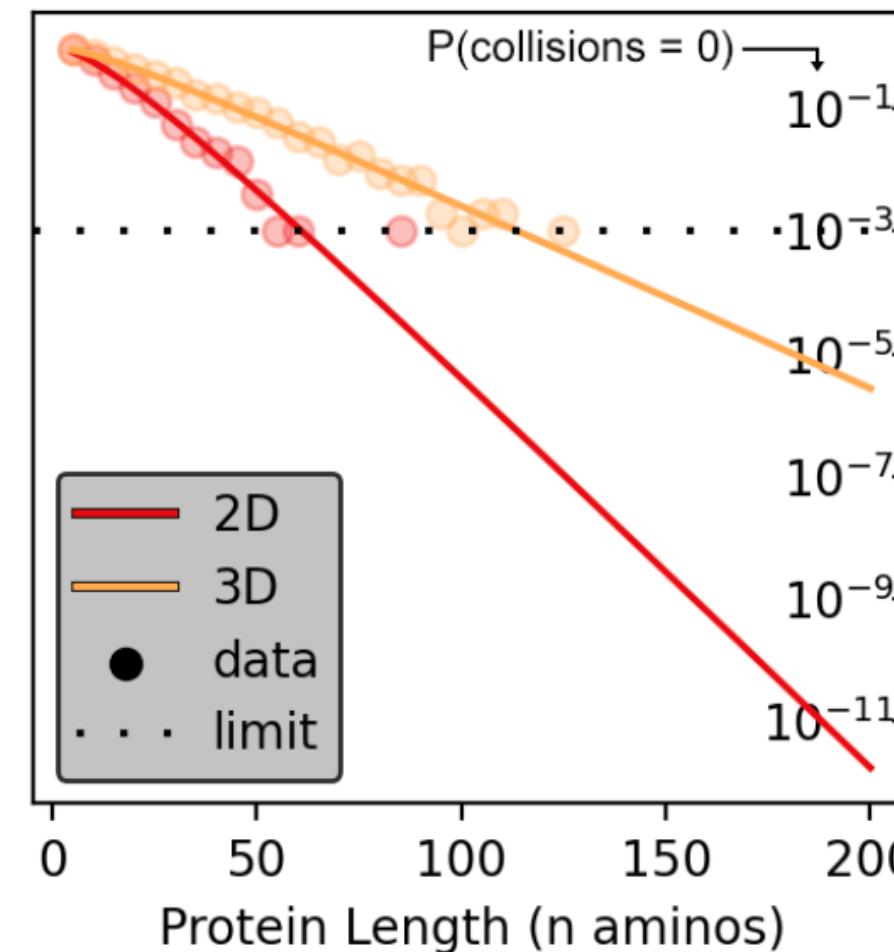


Results

Experiment 0*



Valid conformations per amino acid length



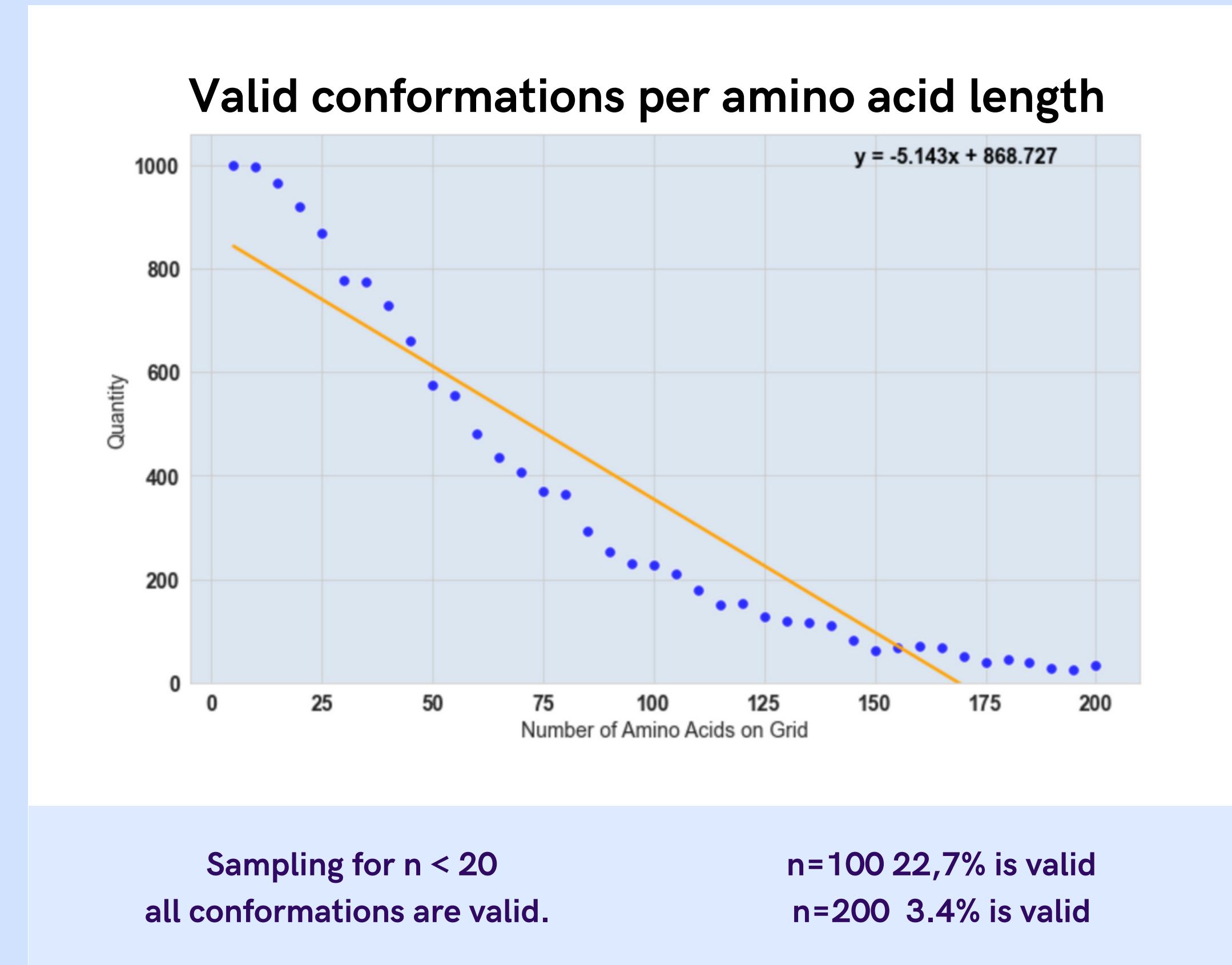
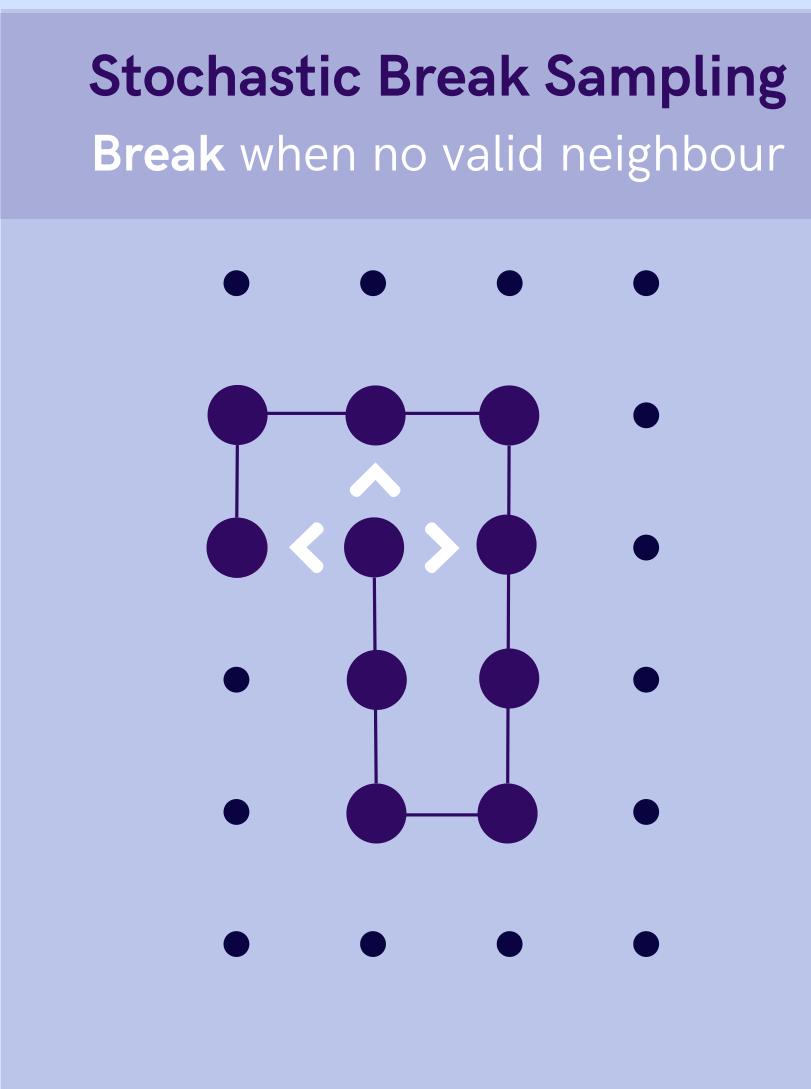
Collisions increase as protein get longer

Chance of zero-collision conformation drops exponentially in n.



Results

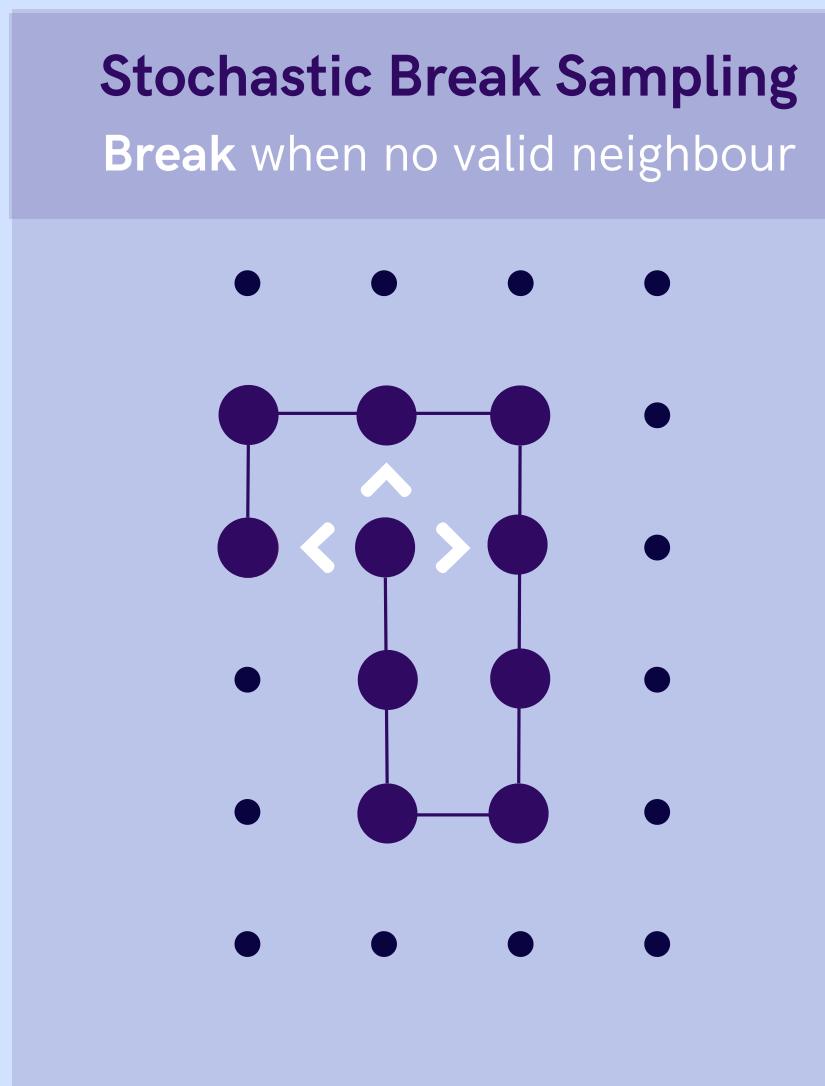
Experiment 1



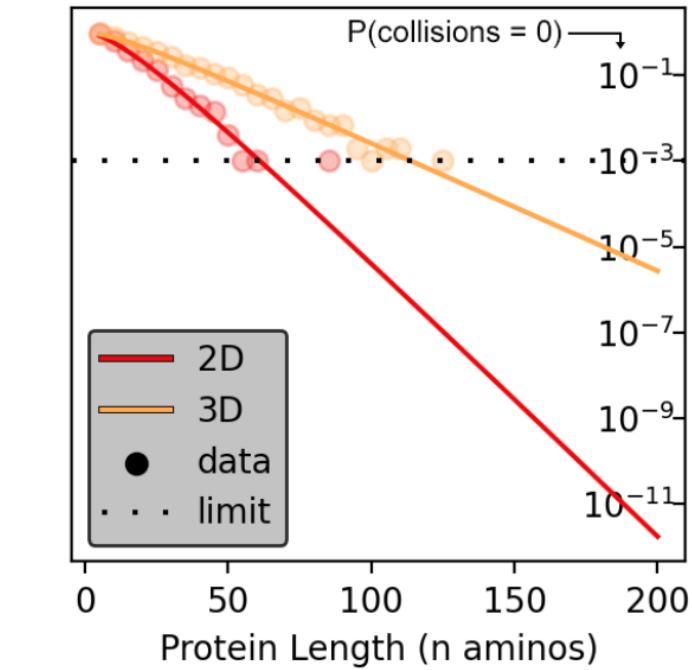
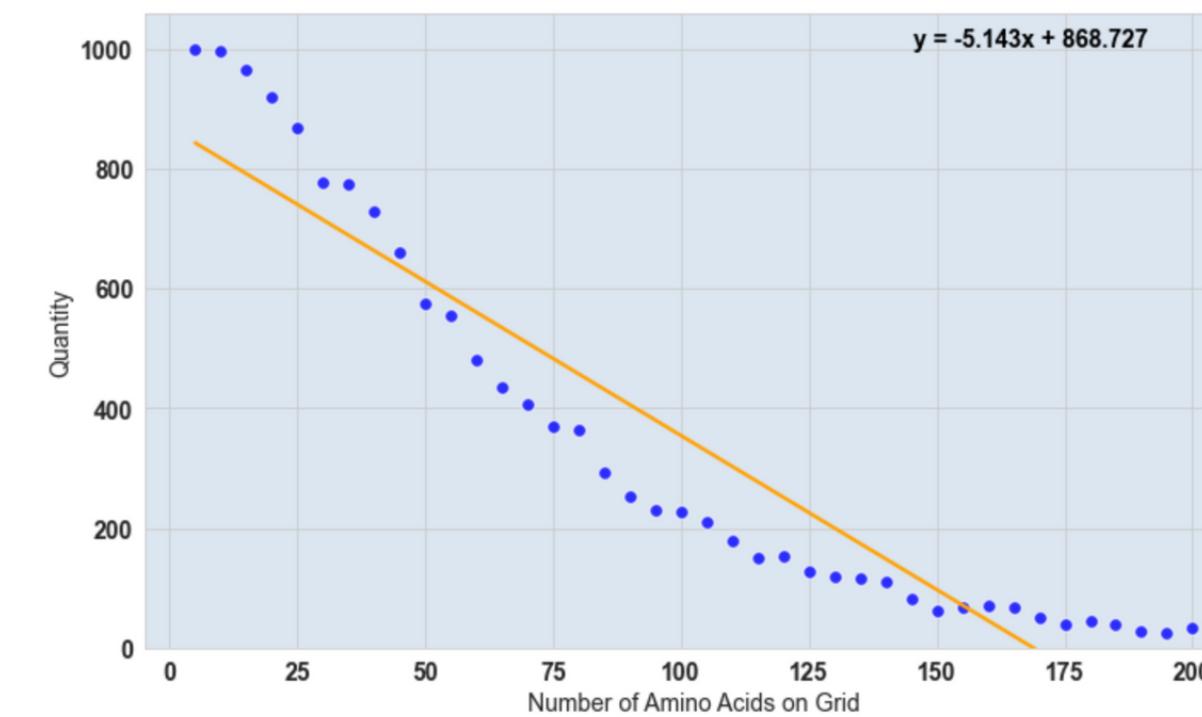


Results

Experiment 1



Valid conformations per amino acid length



Experiments show similar results

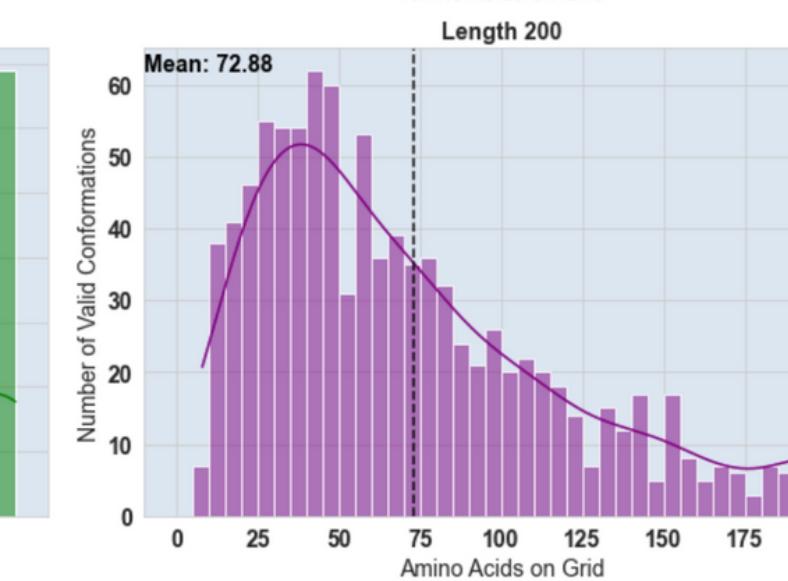
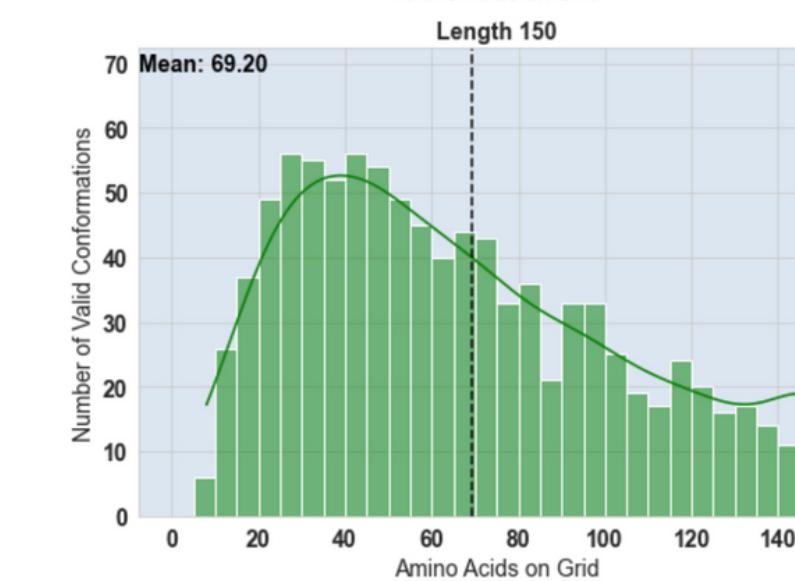
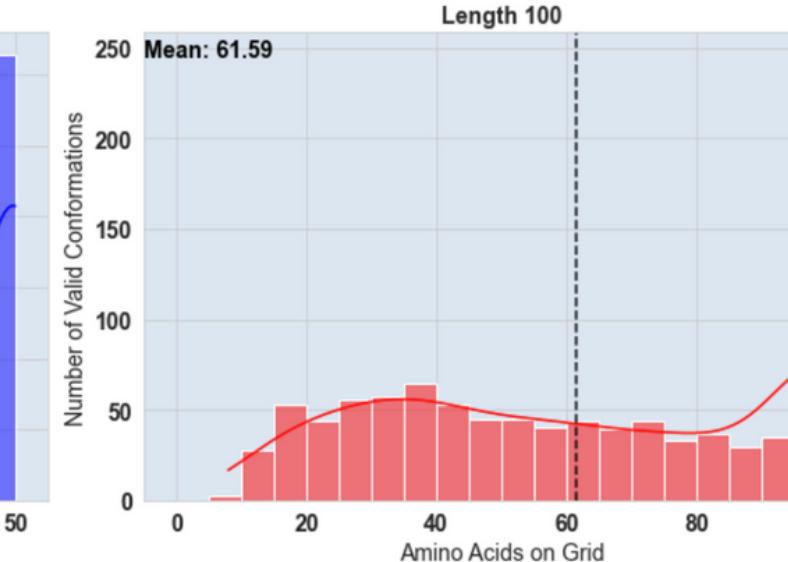
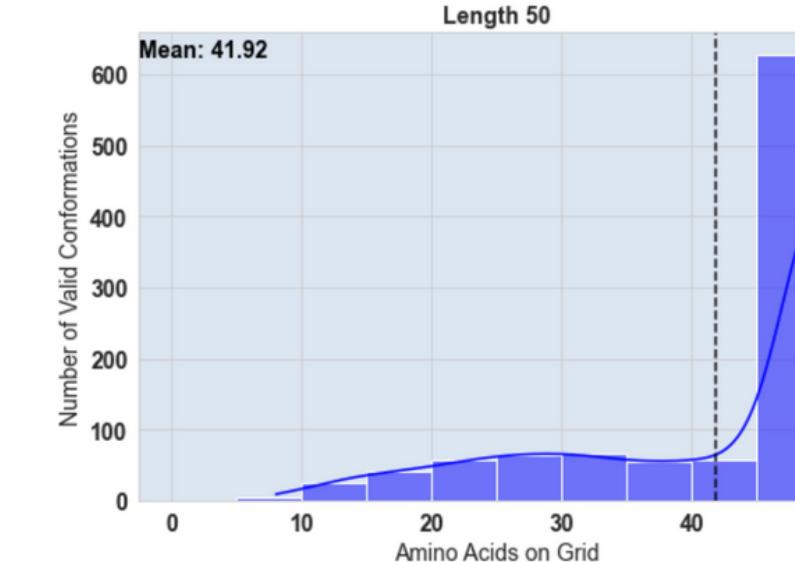
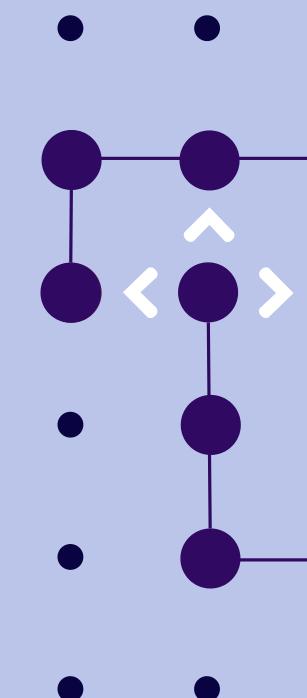
Chance of zero-collision conformation drops exponentially in n.



Results

Stochastic Break Sampling

Break when no valid neighbour



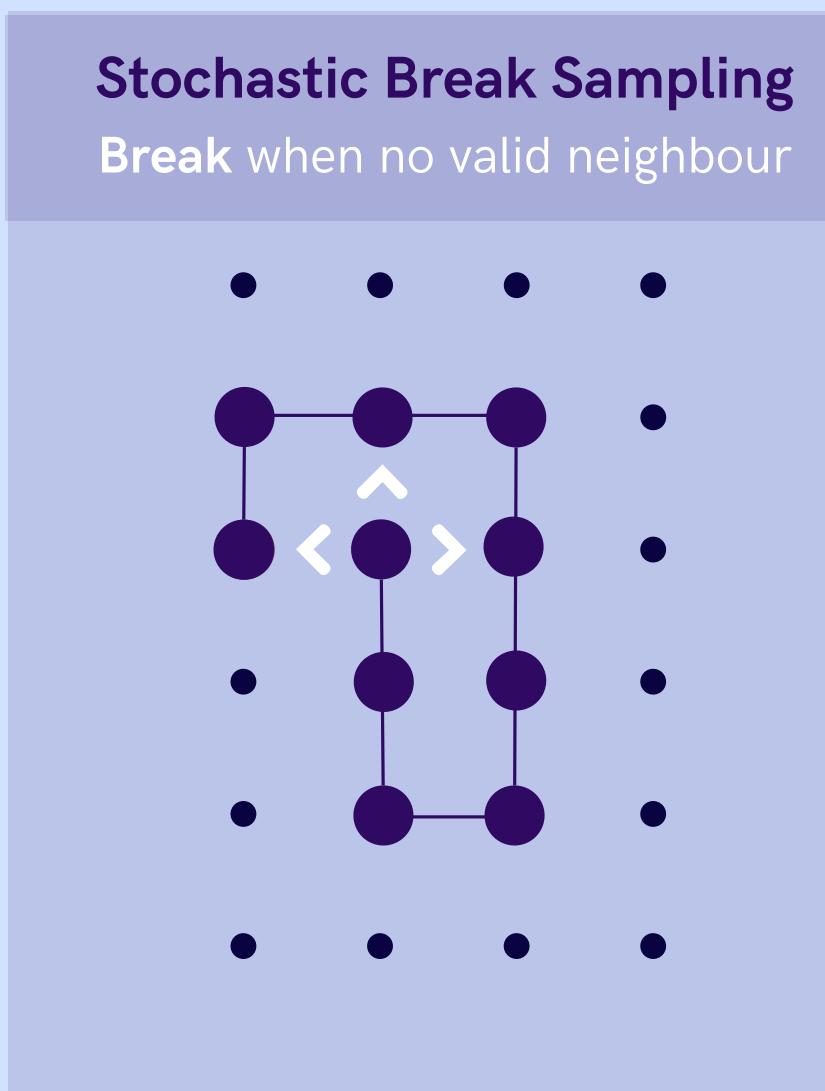
Shorter proteins exhibit peaks near maximums

Longer proteins show a broader distribution

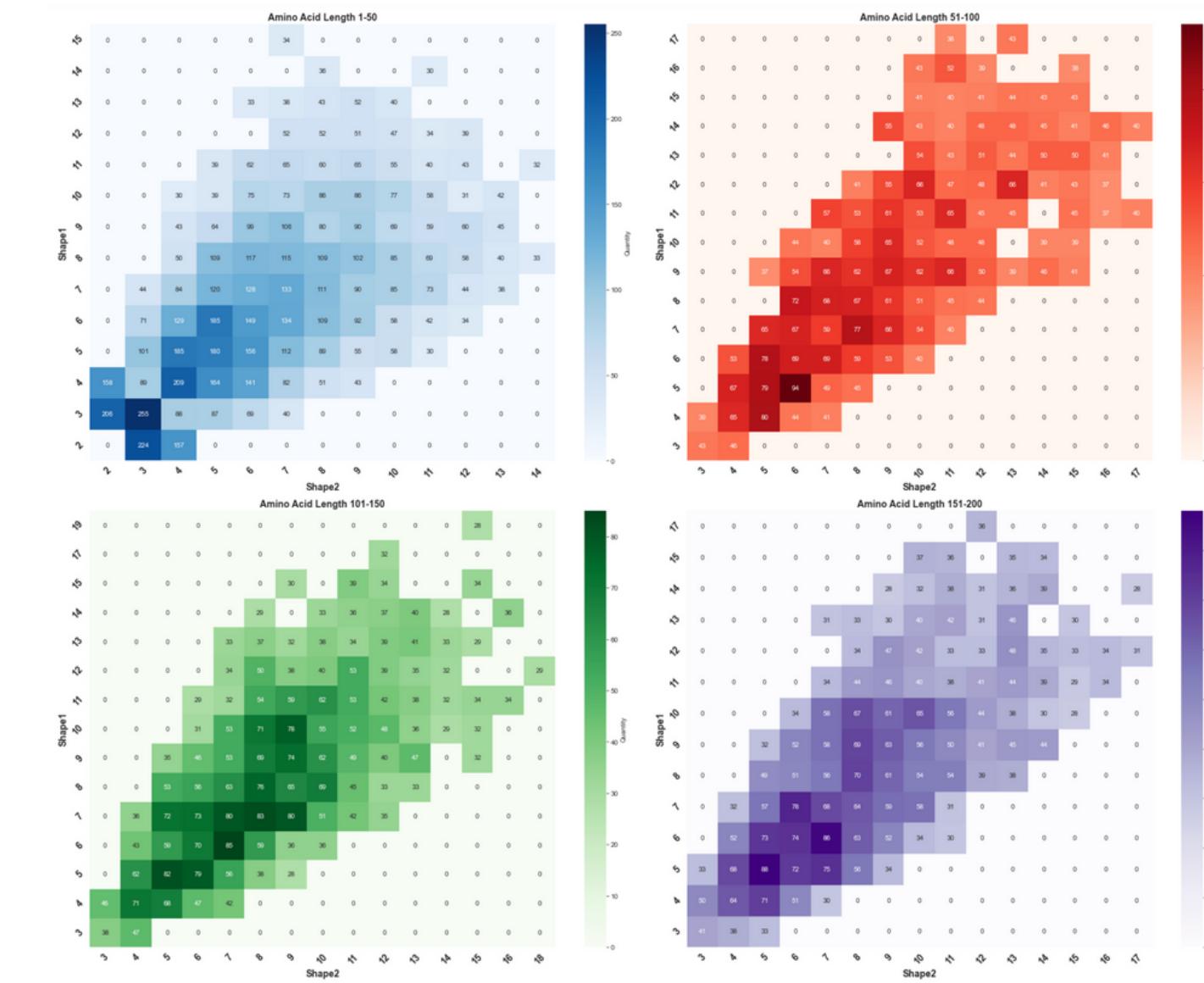


Results

Experiment 1



Heatmap distribution of shape combinations



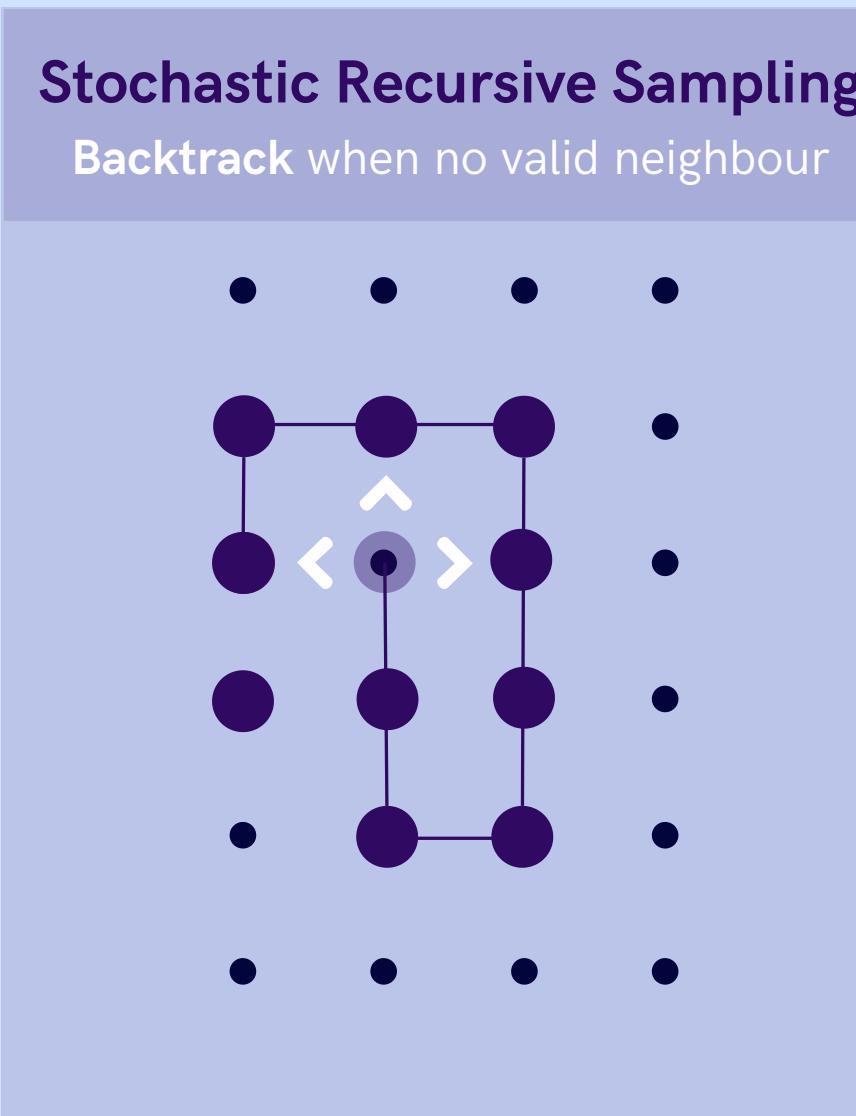
Shorter proteins manifest concentrated shape patterns

Longer proteins exhibit show more varied combinations

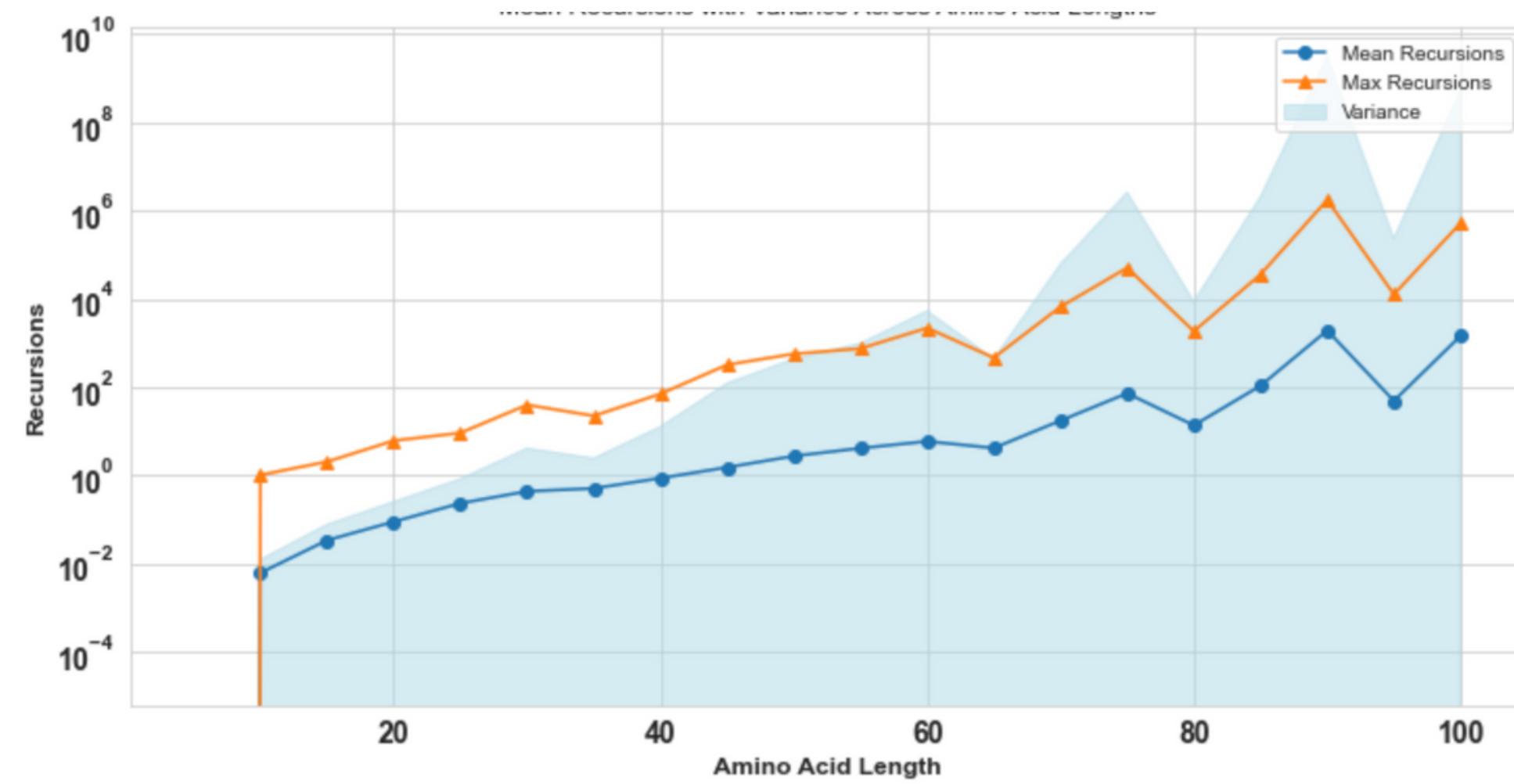


Results

Experiment 3



Distribution of mean and max recursion



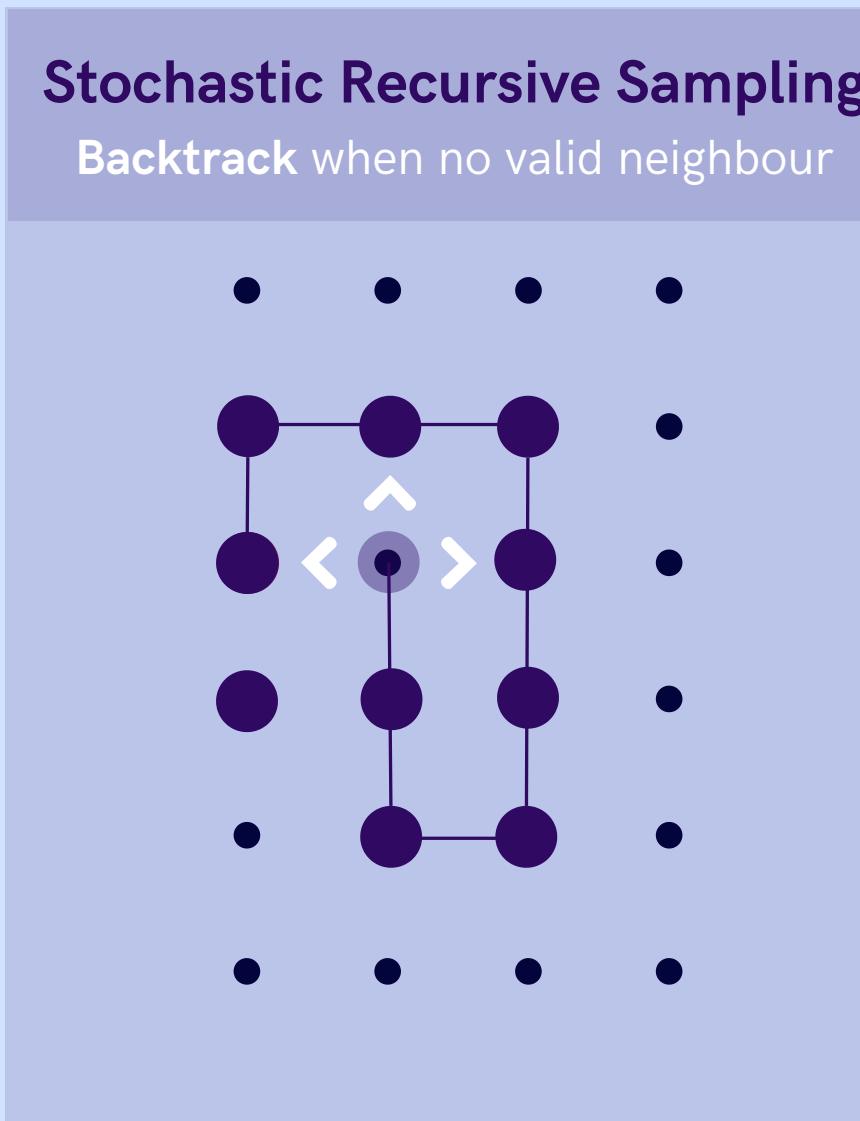
Short proteins have few recursions

Longer proteins longer show exponentially more recursions

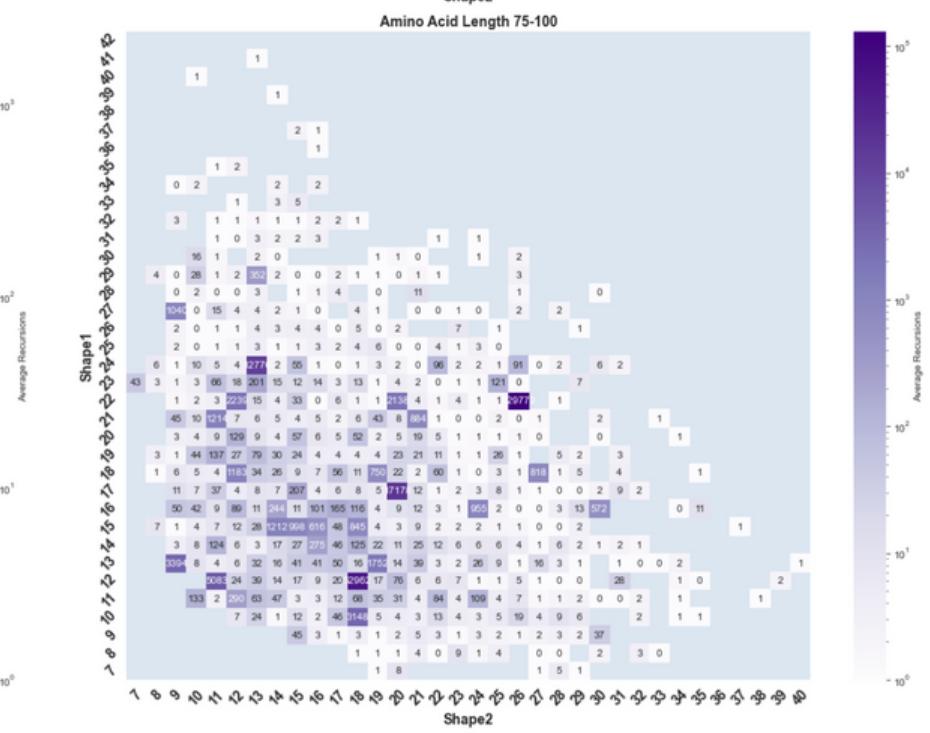
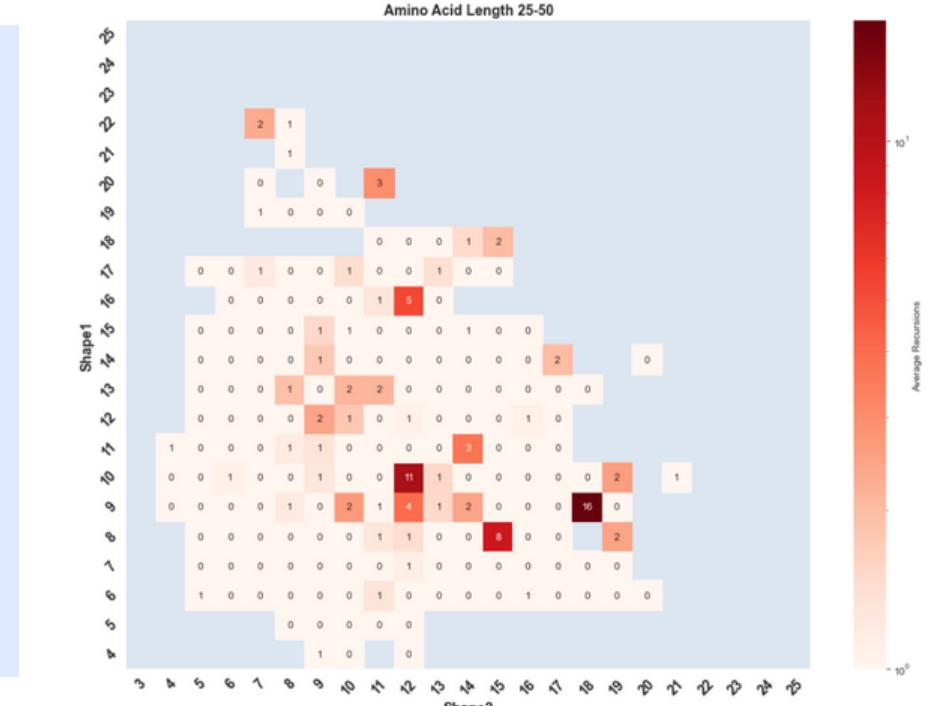
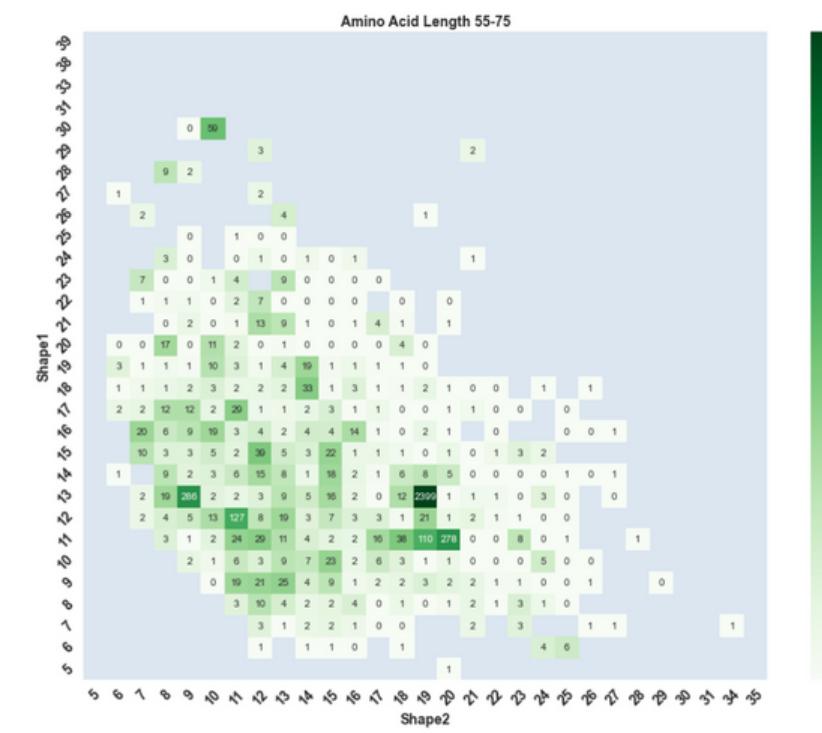
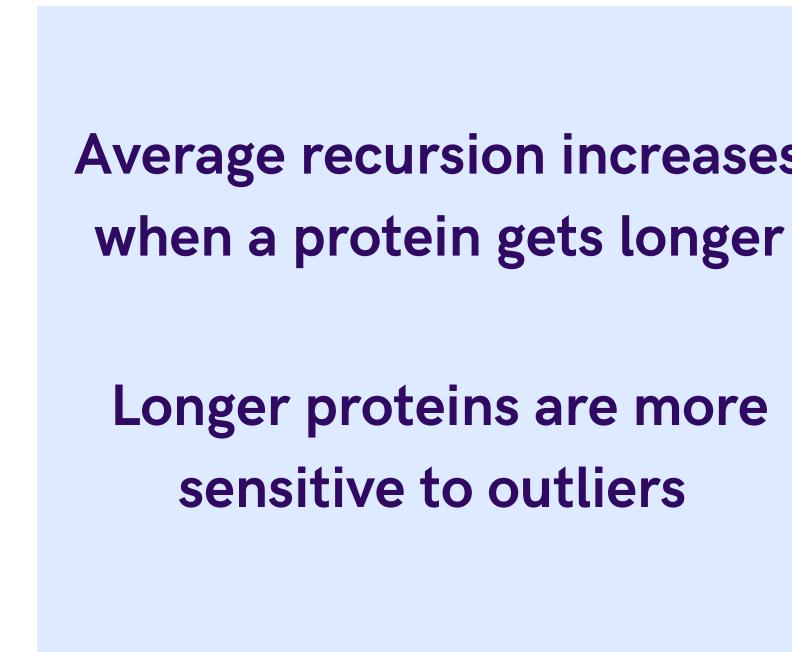


Results

Experiment 3



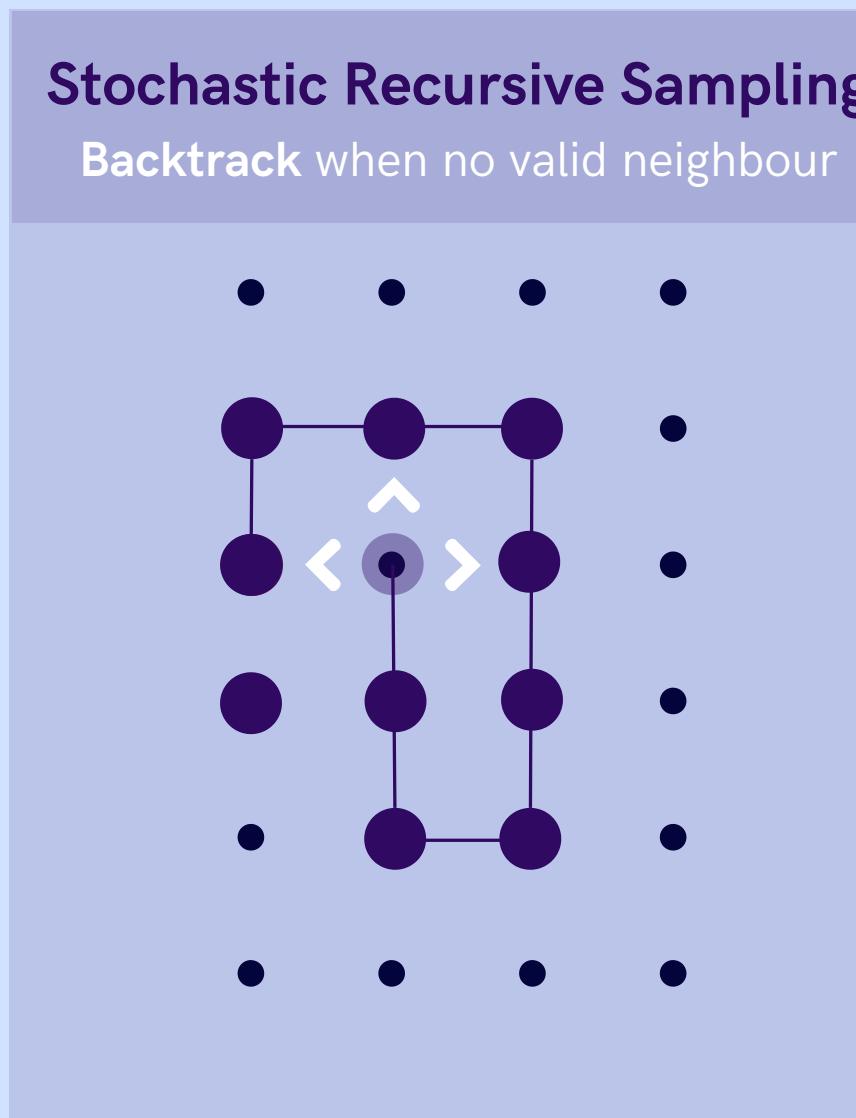
Heatmap distribution of mean recursion per shape





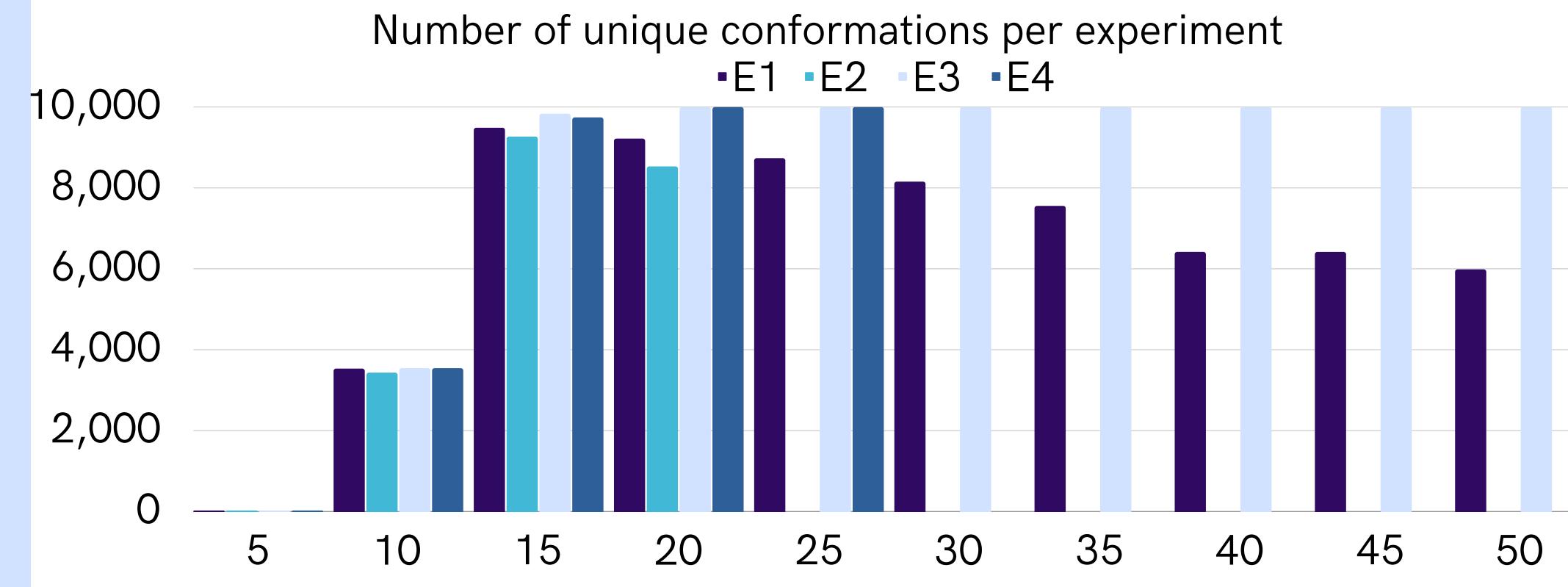
Results

Experiment 5



Benchmarking against sampling space

n	S_n	S_v	S_u
5	27	25	13
10	6561	4067	2034
15	1,594,323	593,611	296,806



Diversity of sample space is explored for short proteins

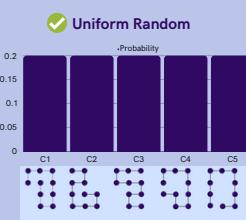
For longer proteins the unique conformations drop for non-backtracking



Conclusion

What do the results say?

Can we ensure the uniform randomness in initial populations before folding proteins?



For **short proteins**, it seems possible. However, the average protein is 200-300 acids. No real life implications.



For **longer proteins**, Probably only with backtracking, but this has many limitations and is hard(er) to verify.

- Scaling
- Recursion depth

What is the impact of different sampling techniques on the quality and diversity of solutions?

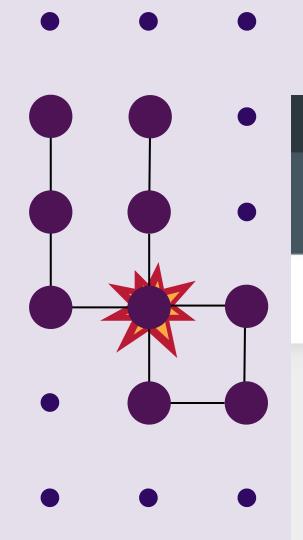
Break sampling seems to work for shorter proteins. Details have to be worked out further.

Recursive sampling seems to improve the quality (valid samples) and diversity (unique samples) of the initial population landscape. Details have to be worked out further.

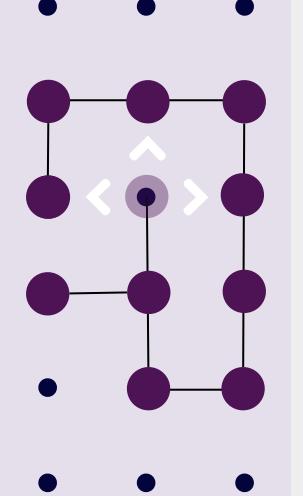


Thanks!

Collision Sampling



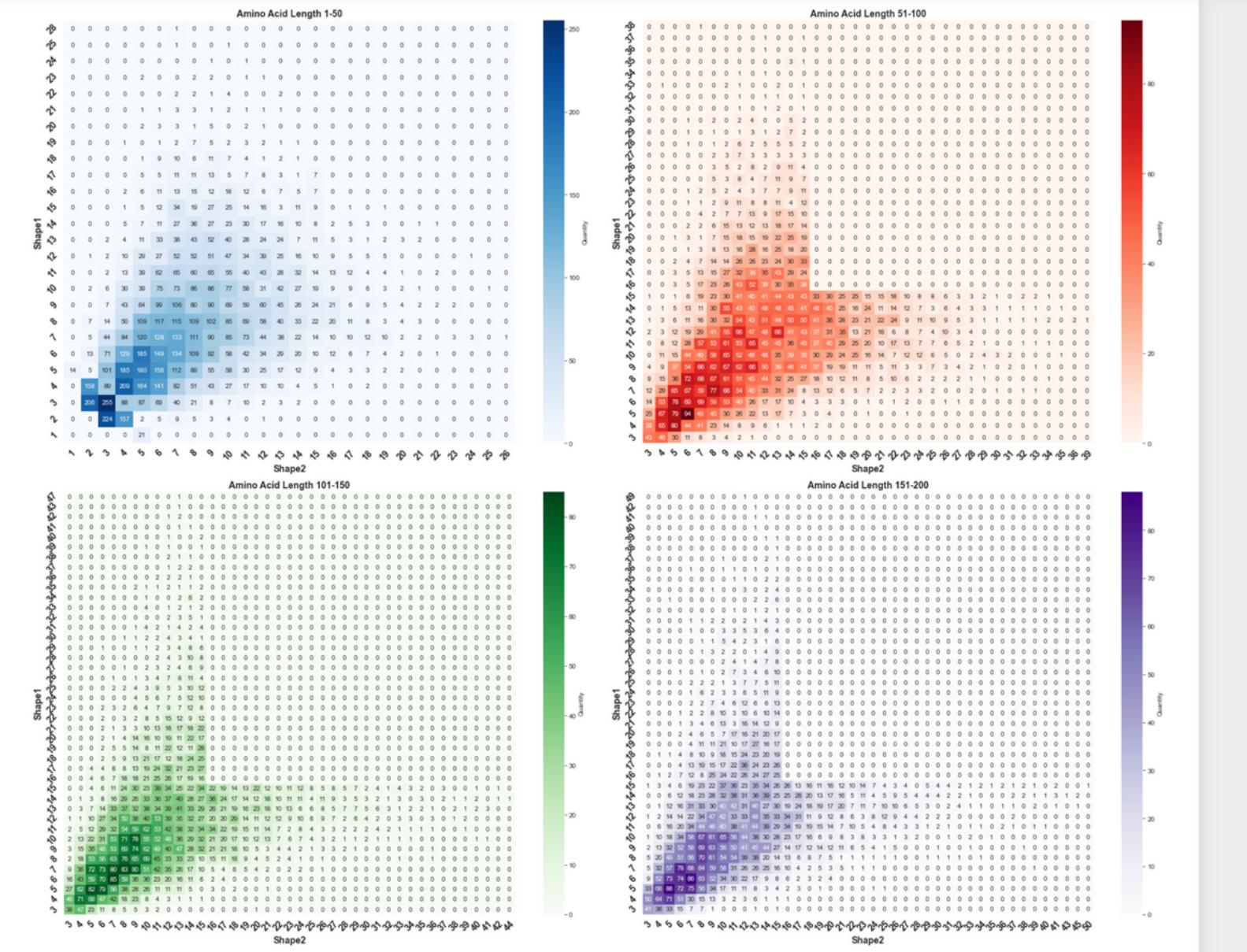
Stochastic Recursive Backtrack when no valid



Stochastic Break Sampling



Deterministic Break Sampling



In []