# What Is Protein Folding and Why Is It NP-Hard

**Alwin**

✉ alwinnear@gmail.com

🏛 University of Indonesia

**Sahira Almahira Kannajmi**

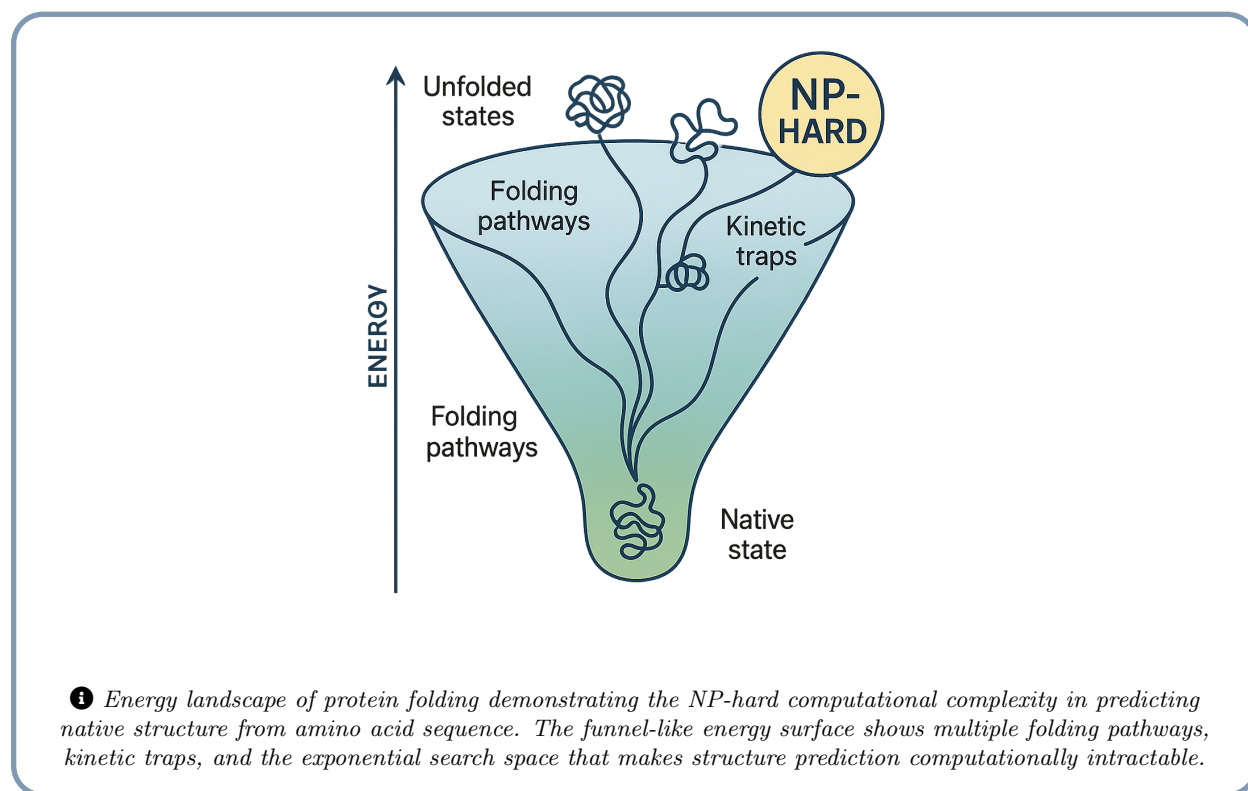✉ almahira.sahira0123@gmail.com

🏛 Gadjah Mada University

## Abstract

Protein folding, the process by which a polypeptide chain acquires its unique three-dimensional (3D) structure, is fundamental to virtually all biological functions. This native structure is dictated by the protein's amino acid sequence and is typically the one with the minimum Gibbs free energy. Understanding and predicting this structure from sequence is a grand challenge in biology and medicine, with profound implications for drug design, synthetic biology, and understanding diseases linked to protein misfolding. Computationally, predicting the native state by minimizing an energy function is exceptionally difficult. Simplified, yet representative, physical models of protein folding, such as the Hydrophobic-Polar (HP) lattice model, have been proven to be NP-hard. This means that finding an algorithm that can determine the optimal, lowest-energy conformation for an arbitrary protein sequence in such models in polynomial time is highly unlikely, implying an exponential increase in computation time with protein size for exact solutions. This inherent computational complexity constrains ab initio structure prediction and has spurred the development of sophisticated heuristic algorithms, machine learning approaches like AlphaFold, and innovative strategies like crowdsourcing, fundamentally shaping the landscape of structural biology and bioinformatics.

🏷 **Keywords:** *protein folding, NP-hardness, computational complexity, hydrophobic-polar model, AlphaFold, structure prediction, machine learning, bioinformatics, thermodynamics, energy minimization*

# 1 📈 Graphical Abstract



ⓘ *Energy landscape of protein folding demonstrating the NP-hard computational complexity in predicting native structure from amino acid sequence. The funnel-like energy surface shows multiple folding pathways, kinetic traps, and the exponential search space that makes structure prediction computationally intractable.*

# 2 Introduction

## 2.1 Context and Motivation

Proteins are the workhorses of biological systems, executing a vast array of functions essential for life, from catalyzing biochemical reactions and transporting molecules to providing structural support and signaling within and between cells[1]. The remarkable functional diversity of proteins stems from their intricate three-dimensional (3D) structures. A cornerstone of molecular biology, articulated in Anfinsen's dogma (or the thermodynamic hypothesis), is that the amino acid sequence of a protein—its primary structure—determines its unique, functional 3D conformation, known as the native state[1]. This native structure is generally understood to be the thermodynamically most stable state under physiological conditions[1]. The process by which a linear polypeptide chain folds into this specific 3D architecture is termed protein folding.

The "protein folding problem" encompasses two interrelated challenges: understanding the physical mechanisms by which folding occurs and predicting the final 3D structure from the primary amino acid sequence[2]. Solving this problem has profound implications across biology and medicine. Accurate structure prediction can illuminate the molecular basis of diseases, particularly those associated with protein misfolding, such as Alzheimer's disease, Parkinson's disease, cystic fibrosis, and various amyloidoses, where proteins fail to fold correctly or aggregate into toxic species[4]. Furthermore, knowledge of protein structures is indispensable for rational drug design, where therapeutic molecules are developed to interact with specific sites on target proteins[7]. In synthetic biology, the ability to predict and design protein structures enables the engineering of novel enzymes for industrial catalysis, bioremediation, and other biotechnological applications[10]. In the context of global health, rapid and accurate structure prediction of viral proteins is crucial for vaccine development and pandemic response, facilitating the identification of epitopes that can elicit protective immune responses[13].

A fascinating aspect of protein folding is Levinthal's paradox[2]. Cyrus Levinthal noted in the 1960s that a polypeptide chain has an astronomical number of possible conformations. If a protein were to find its native state by randomly sampling all these possibilities, the process would take longer than the age of the

universe[2]. Yet, proteins typically fold spontaneously and rapidly, often on timescales of microseconds to seconds[4]. This paradox implies that protein folding is not a random search but rather a guided process, directed by the interactions encoded in the amino acid sequence, leading the protein along specific pathways or down an "energy funnel" towards its native state[19].

While nature efficiently solves the folding puzzle, computational attempts to predict the native structure by finding the global minimum of an energy function face a formidable obstacle: NP-hardness. This paper delves into the meaning of this computational complexity, its formal basis in simplified models, and its far-reaching consequences for science and technology.

## 2.2 Research Questions and Paper Contributions

This paper addresses the following primary research questions:

1. What are the fundamental biophysical and thermodynamic principles that govern the protein folding process, from the hierarchy of structures to the energy landscape that guides folding?
2. How is the complex physical process of protein folding abstracted into computational models, and what are the key assumptions and trade-offs involved in these formalisms, particularly lattice-based versus all-atom representations?
3. What is the meaning of NP-hardness in the context of computational complexity theory, and how is the NP-hardness of protein folding formally demonstrated for simplified but representative models like the Hydrophobic-Polar (HP) model?
4. What are the practical ramifications of protein folding's NP-hardness for biological research, drug discovery, and enzyme engineering? How has the scientific community responded to this inherent computational limitation, particularly through the development of heuristic algorithms and machine learning breakthroughs?

The contributions of this paper are to provide a comprehensive, interdisciplinary synthesis that:

- Explains the fundamental molecular biology of protein folding.
- Introduces the relevant concepts from computational complexity theory in an accessible manner.
- Details the formal proof of NP-hardness for a simplified protein folding model, elucidating the construction and logic involved.
- Critically analyzes the profound impact of this computational intractability on scientific progress and technological innovation, including the rise of AI-driven prediction methods and crowdsourcing initiatives.
- Discusses future prospects and offers recommendations for fostering collaborative efforts to continue tackling this grand challenge.

This work aims to bridge the gap between molecular biology and computational theory, offering a unified perspective on what protein folding is and why its computational prediction is a problem of such enduring difficulty and significance.

# 3 Protein Folding Fundamentals

## 3.1 Description of Primary, Secondary, Tertiary, and Quaternary Structures

The structure of a protein is typically described at four hierarchical levels, each building upon the previous one to define the molecule's overall architecture and function[22].

- **Primary Structure**: This is the linear sequence of amino acids in a polypeptide chain, akin to letters in a word[22]. Amino acids are linked by covalent peptide bonds formed between the carboxyl group of one amino acid and the amino group of the next. The primary structure is determined by the genetic code and is the most fundamental level of protein organization, as it dictates all higher levels of structure and, ultimately, the protein's function.
- **Secondary Structure**: This refers to local, regular, repeating conformations of the polypeptide backbone[22]. The two most common types of secondary structure are the -helix and the -sheet (or -pleated sheet).
  - The -helix is a right-handed coiled or spiral conformation where the backbone N-H group of every amino acid residue donates a hydrogen bond to the backbone C=O group of the amino acid four residues earlier in the sequence[4].
  - -sheets are formed by two or more segments of the polypeptide chain (-strands) lying side-by-side, linked by hydrogen bonds between the backbone N-H and C=O groups of adjacent strands[4]. These strands

---

can be arranged in a parallel or anti-parallel fashion. Other less regular secondary structures include turns and loops, which connect -helices and -strands. Secondary structures form rapidly due to the stabilization provided by these intramolecular hydrogen bonds[4].

- **Tertiary Structure**: This describes the overall three-dimensional shape of an entire single polypeptide chain, resulting from the folding and packing of its secondary structural elements[22]. Tertiary structure is stabilized by a variety of non-covalent interactions between the amino acid side chains (R-groups), including hydrophobic interactions (clustering of nonpolar side chains away from water), hydrogen bonds between side chains or between side chains and the backbone, ionic bonds (salt bridges) between oppositely charged side chains, and van der Waals forces. Covalent disulfide bonds between cysteine residues can also significantly stabilize tertiary structure[24]. For most monomeric proteins, the tertiary structure represents their biologically active, native conformation.
- **Quaternary Structure**: This level of structure applies to proteins composed of more than one polypeptide chain (subunit)[22]. Quaternary structure refers to the spatial arrangement and interactions of these subunits to form a functional multimeric protein complex. The same types of interactions that stabilize tertiary structure (hydrophobic interactions, hydrogen bonds, ionic bonds, van der Waals forces, and disulfide bonds) also hold subunits together.

## 3.2 Energy Funnel Model, Thermodynamics, and Kinetics

The process by which a protein achieves its native tertiary (or quaternary) structure is a complex interplay of thermodynamic and kinetic factors, famously encapsulated by Anfinsen's dogma and visualized by the energy funnel model.

**Anfinsen's Dogma (The Thermodynamic Hypothesis):**

Christian Anfinsen's pioneering work in the 1950s and 1960s, for which he received the Nobel Prize in Chemistry in 1972, demonstrated that, for many proteins, the information required to specify their native 3D structure is entirely contained within their primary amino acid sequence[1]. His experiments with ribonuclease A showed that if the protein was denatured (unfolded) and its disulfide bonds were broken, it could spontaneously refold into its active, native conformation upon removal of the denaturing agents and allowing disulfide bonds to reform. This led to the thermodynamic hypothesis: the native structure of a protein in its physiological environment is the one in which the Gibbs free energy of the entire system (protein plus solvent) is at a global minimum[1]. This state must satisfy three conditions for a unique structure to form:

1. **Uniqueness**: The native state must represent a sufficiently deep free energy minimum, unchallenged by other configurations of comparable energy.
2. **Stability**: Small perturbations in the environment should not easily disrupt the native state; the free energy surface around this minimum must be steep.
3. **Kinetic Accessibility**: The path from the unfolded state to the native state on the free energy surface must be reasonably smooth, without insurmountable energy barriers, allowing folding to occur on a biologically relevant timescale.

**Levinthal's Paradox:**

Despite Anfinsen's findings suggesting a deterministic path to a unique structure, Cyrus Levinthal pointed out a significant kinetic puzzle[2]. A typical protein of, say, 100 amino acids, with each residue having only a few possible backbone dihedral angles, could adopt an astronomically large number of conformations (e.g., $10^{100}$ or more). If the protein had to sample each of these conformations to find the lowest energy state, the folding process would take an impossibly long time, far exceeding the age of the universe[2]. However, proteins fold much faster, often in microseconds to seconds[4]. This discrepancy, known as Levinthal's paradox, strongly suggests that protein folding is not a random search through all possible conformations but a directed process[25].

**The Energy Funnel Model:**

The energy funnel model, developed by researchers like Peter Wolynes, José Onuchic, and Ken Dill, provides a conceptual framework to resolve Levinthal's paradox[19]. This model depicts the free energy landscape of a folding protein as a multi-dimensional funnel.

- The width of the funnel at any given energy level represents the conformational entropy of the protein (the number of accessible conformations). At the top, the funnel is wide, corresponding to the vast ensemble of high-energy, high-entropy unfolded or denatured states.

---

- The depth of the funnel represents the energy. As the protein folds, it moves "downhill" towards lower energy states. The bottom of the funnel represents the unique, low-energy, low-entropy native state[20].
- The slope of the funnel provides a thermodynamic driving force, guiding the protein towards the native state, thus avoiding an exhaustive search of the conformational space.
- The surface of the funnel is not perfectly smooth but rugged, characterized by many small local energy minima[20]. These represent metastable, partially folded intermediates or "kinetic traps" where the protein can temporarily get stuck, slowing down the folding process or even leading to misfolding. Misfolded intermediates can sometimes accumulate if these traps are deep enough[20]. The folding process is thus visualized as an ensemble of pathways, or a stochastic search, down this rugged funnel, rather than a single, fixed pathway[20]. The overall shape of the funnel is determined by the amino acid sequence, and natural selection is thought to have favored sequences that lead to smooth, steep funnels, enabling rapid and efficient folding to a stable native state with minimal frustration (i.e., few conflicting energy contributions)[20].

**Thermodynamics of Folding:**

The spontaneity of protein folding is governed by the change in Gibbs free energy (G), given by the fundamental thermodynamic equation:

$$\Delta G = \Delta H - T\Delta S$$

where H is the change in enthalpy, T is the absolute temperature, and S is the change in entropy[25]. For folding to be spontaneous, G must be negative.

- **Enthalpy (H):** This term reflects changes in bonding energy. The formation of favorable non-covalent interactions within the protein (hydrogen bonds, van der Waals interactions, electrostatic interactions like salt bridges) contributes to a negative (favorable) H[20]. The hydrophobic effect, while primarily driven by solvent entropy, also results in favorable van der Waals contacts between nonpolar residues packed in the protein core[20].
- **Entropy (S):** This term has two main components:
  - **Conformational Entropy of the Polypeptide:** When a flexible, disordered polypeptide chain folds into a specific, ordered structure, its conformational entropy decreases significantly. This results in a negative $S_{protein}$, which is unfavorable for folding[25].
  - **Solvent Entropy (Hydrophobic Effect):** In an unfolded protein, nonpolar (hydrophobic) amino acid side chains are exposed to water, forcing surrounding water molecules to form ordered "cages" (clathrate-like structures) around them. This ordering of water decreases the solvent's entropy. When the protein folds, these nonpolar side chains are buried in the protein's interior, releasing the ordered water molecules into the bulk solvent. This release of water increases the solvent's entropy, resulting in a positive $S_{solvent}$, which is favorable for folding[20]. For many globular proteins, this favorable entropy change of the solvent due to the hydrophobic effect is a dominant driving force for folding[20]. The net G of folding is typically a small negative value (e.g., -5 to -15 kcal/mol for small globular proteins), indicating that the native state is only marginally more stable than the unfolded state. This delicate balance is easily perturbed by changes in temperature, pH, or denaturant concentration. The stability of a protein is often characterized by its melting temperature ($T_m$), the temperature at which half the protein population is folded and half is unfolded, meaning G=0[25].

The equilibrium between the unfolded (U) and native (N) states can be described by an equilibrium constant:

$$K_{eq} = \frac{[N]}{[U]}$$

The standard free energy change is related to this by:

$$\Delta G^\circ = -RT \ln K_{eq}$$

where R is the gas constant and T is the absolute temperature[27].

**Kinetics of Folding:**

Protein folding kinetics describe the rates and mechanisms of the folding process. Folding times can vary dramatically, from microseconds for very small, simple proteins to minutes or even hours for larger, more complex ones, especially if slow steps like proline isomerization are involved[4].

- **Folding Pathways and Intermediates:** While the funnel model emphasizes an ensemble of pathways, specific intermediates can sometimes be identified. Early studies on small proteins often showed "two-state" folding (U  N) with no significantly populated intermediates, following first-order kinetics[27]. However, larger proteins often exhibit more complex, multiphasic kinetics, indicative of one or more intermediate states[4]. These intermediates can include "molten globules"—compact, partially folded states that have much of the native secondary structure but lack well-packed tertiary interactions and have fluctuating side chains[20].
- **Rate-Limiting Steps and Transition States:** The rate of folding is often determined by the highest energy barrier (transition state) along the folding pathway(s).
- **Folding Speed Limit:** There appears to be an upper limit to how fast a protein can fold, estimated to be around 1 microsecond, potentially limited by the rate of polypeptide chain collapse or formation of elementary structures like hairpins[18]. Some small proteins have been observed to fold at rates approaching this limit[18]. The folding time scale depends on factors like protein size, the "contact order" (average sequence separation between contacting residues in the native state), and the complexity of the protein's topology[4].

## 3.3  Intuitive Analogies

To grasp the complex concepts of protein folding, intuitive analogies can be helpful:
- **Folding as Downhill Travel in a Rugged Landscape:** Imagine a golf ball (the protein) rolling down a bumpy, funnel-shaped hill (the energy landscape) towards the hole (the native state) at the bottom[21]. The ruggedness of the terrain, with its many small dips and rises, represents local energy minima (kinetic traps) and energy barriers. The overall slope guides the ball towards the hole, but it might get temporarily stuck in a sand trap or a divot.
- **The Marco Polo Bridge Analogy:** Italo Calvino's description of a bridge in Invisible Cities, where the global shape emerges naturally when all its stones are in their unique, correct places, has been compared to a folded protein[21]. Each amino acid interaction is like a stone; only when all are correctly positioned relative to each other does the stable, functional structure (the bridge) emerge. This emphasizes how local interactions collectively define the global fold.
- **Folding as a Jigsaw Puzzle:** Assembling a complex jigsaw puzzle can be an analogy. While there's a unique final picture (native state), the process of finding where each piece fits (amino acid interactions) can be challenging. Some sections might come together quickly (secondary structures), but fitting these larger assembled pieces together to form the whole picture (tertiary structure) requires specific long-range interactions. Getting a piece in the wrong place could hinder further assembly (misfolding).

These analogies help to visualize the directed yet complex nature of protein folding, driven by a search for a low-energy state within a vast and rugged conformational space.

# 4  Computational Models of Folding

To study protein folding computationally, scientists employ various models that abstract the complexity of the physical system to different degrees. These models range from highly simplified lattice representations to detailed all-atom simulations, each with its own trade-offs between realism and computational tractability.

## 4.1  HP Lattice vs. Off-Lattice/All-Atom Formalisms

HP (Hydrophobic-Polar) Lattice Models:

The Hydrophobic-Polar (HP) model, first proposed by Ken Dill and colleagues in 1985, is one of the most widely studied simplified models for protein folding[31]. Its core idea is to capture the dominant role of the hydrophobic effect in driving protein folding.
- **Representation**: In the HP model, each amino acid in a protein sequence is classified into one of two types: H (hydrophobic or nonpolar) or P (polar or hydrophilic)[31]. The polypeptide chain is then represented as a self-avoiding walk on a regular lattice, typically a 2D square lattice or a 3D cubic lattice[31]. Each amino acid occupies a single lattice site, and adjacent amino acids in the sequence must occupy adjacent sites on the lattice.
- **Energy Function**: The "energy" of a given conformation (a specific self-avoiding walk on the lattice) is calculated based on interactions between H residues. Specifically, a favorable energy contribution (typically

-1 unit per contact) is assigned to each pair of H residues that are adjacent on the lattice but not consecutive in the sequence (these are called "topological neighbors" or "H-H contacts")[31]. Interactions involving P residues (H-P or P-P) are usually assigned an energy of zero. The goal is to find the conformation that maximizes the number of H-H contacts, thereby minimizing the total energy[31].

- **Advantages**: The HP model's simplicity allows for the exploration of fundamental principles of protein folding, such as hydrophobic collapse and the formation of a compact core[31]. Its discrete nature makes it amenable to combinatorial analysis and computational complexity studies; indeed, the NP-hardness of protein folding was first rigorously proven using HP models[32].
- **Disadvantages**: The lattice constraint is a significant simplification and can introduce artifacts, affecting the dynamics and the types of structures that can be formed[31]. HP models do not typically represent true secondary structures like -helices or -sheets explicitly, nor do they account for the detailed stereochemistry of amino acid side chains or the complexities of solvent interactions beyond the binary H/P classification[31].

Off-Lattice Models:

Off-lattice models remove the constraint of confining amino acids to fixed lattice points, allowing for continuous conformational space. This provides a more realistic representation of protein structure and dynamics[35].

- **Coarse-Grained Off-Lattice Models**: These models still simplify the protein, often representing each amino acid as a single bead or a few beads (e.g., one for the backbone, one for the side chain). Bond lengths might be fixed, but bond angles and dihedral angles can vary continuously or within certain ranges[35]. Interactions between non-bonded residues are typically described by simplified potentials, such as a modified Lennard-Jones potential, and there might be terms for bending energy between successive bonds[35]. The "AB model" is an example of an off-lattice coarse-grained model, similar to the HP model in its binary classification of residue types but allowing continuous backbone bending[35]. While more realistic than lattice models, even simplified off-lattice AB models have been shown to be NP-complete[35].
- **All-Atom Off-Lattice Models**: These are the most detailed and physically realistic models. Every atom in the protein (and often surrounding solvent molecules) is explicitly represented[27]. The interactions between atoms are described by sophisticated physics-based potential energy functions, commonly known as force fields. These force fields include terms for covalent bond stretching, bond angle bending, torsional (dihedral) angle rotations, van der Waals interactions, and electrostatic interactions[27]. Solvation effects can be treated implicitly (as a continuum dielectric) or explicitly (by including water molecules in the simulation). Molecular Dynamics (MD) simulations, which numerically solve Newton's equations of motion for all atoms, are commonly performed using all-atom models to study folding pathways, protein dynamics, and thermodynamic properties[21]. However, all-atom MD simulations are computationally extremely expensive, limiting the timescales and system sizes that can be practically studied[21].

The choice of model depends on the research question. HP lattice models are invaluable for theoretical studies of complexity and general folding principles. Coarse-grained off-lattice models offer a balance between realism and computational cost, useful for exploring folding mechanisms of larger proteins or longer timescales than accessible by all-atom models. All-atom models provide the highest level of detail, crucial for understanding fine structural features, enzyme mechanisms, and drug interactions, but are limited by computational demands.

## 4.2 Mathematical Definition of the Energy Objective Function

The core of any computational protein folding model is its energy objective function (or simply, energy function). This function assigns a scalar energy value to any given conformation of the protein. The computational task is then to find the conformation(s) that minimize this energy, which are presumed to correspond to the native state(s).

General Form for All-Atom Models (Conceptual):

For all-atom models, the potential energy $E_{total}$ of a given conformation is typically expressed as a sum of various terms derived from classical mechanics, representing different types of physical interactions[27]:

$$E_{total} = E_{bond} + E_{angle} + E_{dihedral} + E_{vanderWaals} + E_{electrostatic}(+E_{solvation})$$

Where:

- $E_{bond}$ represents the energy associated with stretching or compressing covalent bonds from their equilibrium lengths.

- $E_{angle}$ represents the energy associated with bending bond angles from their equilibrium values.
- $E_{dihedral}$ represents the energy associated with twisting dihedral (torsional) angles around bonds.
- $E_{vanderWaals}$ accounts for short-range attractive and repulsive forces between non-bonded atoms (Lennard-Jones potential is commonly used).
- $E_{electrostatic}$ accounts for Coulombic interactions between charged atoms.
- $E_{solvation}$ is an optional term that accounts for the interaction of the protein with the solvent, either implicitly or explicitly.

The specific mathematical forms of these terms and their parameters constitute a "force field" (e.g., AMBER, CHARMM, GROMOS). A detailed derivation is often complex and is outlined further in Appendix B.

HP Model Energy Function (Specific):

For the HP lattice model, the energy function is much simpler and focuses on capturing the hydrophobic effect[31].

Let the protein sequence be $S = (s_1, s_2, \ldots, s_N)$, where each $s_k \in \{H, P\}$ indicates whether the k-th amino acid is hydrophobic (H) or polar (P).

A conformation $C$ is a self-avoiding walk $C = (p_1, p_2, \ldots, p_N)$ on a given lattice (e.g., 2D square or 3D cubic), where $p_k$ are the coordinates of the k-th amino acid $s_k$.

The energy $E(C)$ of the conformation $C$ is defined as the sum of interaction energies between all pairs of non-sequentially adjacent H-residues that are topological neighbors on the lattice:

$$E(C) = \sum_{1 \leq k < l-1 \leq N-1} \chi(s_k, s_l) \cdot \delta(p_k, p_l)$$

Here:
- The sum is over all pairs of amino acids $(s_k, s_l)$ such that $l > k + 1$ (i.e., they are not adjacent in the sequence, to exclude bonded interactions).
- $\chi(s_k, s_l)$ is the interaction energy parameter between amino acid types $s_k$ and $s_l$. In the standard HP model:
  - $\chi(H, H) = \epsilon_{HH}$, which is a negative constant (e.g., -1) representing a favorable interaction.
  - $\chi(H, P) = \epsilon_{HP} = 0$.
  - $\chi(P, P) = \epsilon_{PP} = 0$.
- $\delta(p_k, p_l)$ is an indicator function:
  - $\delta(p_k, p_l) = 1$ if amino acids $s_k$ (at position $p_k$) and $s_l$ (at position $p_l$) are topological neighbors on the lattice (i.e., occupy adjacent lattice sites, e.g., distance 1 in Manhattan metric for a cubic lattice).
  - $\delta(p_k, p_l) = 0$ otherwise.

Thus, the energy function simplifies to counting the number of H-H contacts, $N_{HH}$, and multiplying by the favorable energy per contact:

$$E(C) = \epsilon_{HH} \cdot N_{HH}(C)$$

Since $\epsilon_{HH}$ is negative, minimizing $E(C)$ is equivalent to maximizing the number of H-H contacts, $N_{HH}(C)$. This formulation captures the tendency of hydrophobic residues to cluster together to minimize their exposure to a polar solvent, forming a hydrophobic core[31]. It is this seemingly simple energy function, defined over a discrete conformational space, that forms the basis for many NP-hardness proofs of protein folding.

# 5 Computational Complexity Primer

To understand why protein folding is considered "NP-hard," it is essential to first grasp some fundamental concepts from computational complexity theory. This field of theoretical computer science classifies computational problems based on their inherent difficulty, specifically concerning the resources (like time or memory) required to solve them as the input size grows. The following definitions are central, drawing from seminal texts like Sipser's "Introduction to the Theory of Computation"[38] and Garey Johnson's "Computers and Intractability"[42].

## 5.1   Precise Definitions and Intuitive Explanations

P (Polynomial Time):

The class P consists of decision problems that can be solved by a deterministic algorithm in polynomial time[38]. This means that the number of computational steps required by the algorithm is bounded by a polynomial function of the input size n, typically denoted as $O(n^k)$ for some constant k.

- **Intuitive Explanation**: Problems in P are generally considered "efficiently solvable" or "tractable." As the input size grows, the time to solve them grows at a manageable rate. Examples include sorting a list, searching in a sorted list, and finding the shortest path between two nodes in a graph using Dijkstra's algorithm.

NP (Nondeterministic Polynomial Time):

The class NP consists of decision problems for which, if the answer is "yes," there exists a proof (also called a certificate or witness) that can be verified by a deterministic algorithm in polynomial time[38]. An alternative, equivalent definition is that NP problems are those solvable by a nondeterministic Turing machine in polynomial time[38].

- **Intuitive Explanation**: For problems in NP, while finding a solution might be hard, checking whether a proposed solution is correct is easy (can be done quickly). Think of a jigsaw puzzle: solving it can be very time-consuming, but if someone gives you a completed puzzle, you can quickly verify that all pieces fit together correctly. Many search and optimization problems, when phrased as decision problems (e.g., "Does a solution with value at most X exist?"), fall into NP. It is clear that P NP, because if a problem can be solved in polynomial time, its solution can also be verified in polynomial time (by simply re-solving it and checking the answer)[38].

NP-hard (Nondeterministic Polynomial-time hard):

A problem X is NP-hard if every problem Y in the class NP can be reduced to X in polynomial time (denoted $Y_P X$)[39]. A polynomial-time reduction is an algorithm that transforms an instance of problem Y into an instance of problem X in polynomial time, such that the original instance of Y has a "yes" answer if and only if the transformed instance of X has a "yes" answer.

- **Intuitive Explanation**: NP-hard problems are at least as hard as the hardest problems in NP[42]. If one could find an efficient (polynomial-time) algorithm for any NP-hard problem, then one could efficiently solve all problems in NP. NP-hard problems do not necessarily have to be in NP themselves (i.e., their solutions might not be verifiable in polynomial time, or they might not even be decision problems, like optimization problems that ask for the best solution rather than a yes/no answer).

NP-complete (NPC):

A problem X is NP-complete if it satisfies two conditions:

1. X is in NP (its solutions can be verified in polynomial time).
2. X is NP-hard (every problem in NP is reducible to X in polynomial time)[38].

- **Intuitive Explanation**: NP-complete problems are the "hardest problems in NP"[39]. They are problems for which no polynomial-time algorithm is known, yet their solutions are efficiently verifiable. If an efficient algorithm were found for any single NP-complete problem, it would imply efficient algorithms for all NP-complete problems, and indeed for all problems in NP, meaning P would equal NP. The general consensus among computer scientists is that P NP, implying that NP-complete problems are intrinsically intractable in the worst case[38].

The P versus NP problem—whether P is equal to NP—remains one of the most significant unsolved questions in theoretical computer science and mathematics[38]. The belief that P NP motivates the classification of problems as NP-hard, as it suggests that these problems are unlikely to admit efficient, exact algorithms for all instances.

## 5.2   Classic NP-hard Examples

To illustrate the nature of NP-hard and NP-complete problems, a few classic examples are often cited:

3-SAT (3-Satisfiability):

The Boolean Satisfiability Problem (SAT) asks whether there is an assignment of truth values (TRUE or FALSE) to the variables of a given Boolean formula that makes the entire formula evaluate to TRUE[45]. SAT was the first problem proven to be NP-complete (Cook-Levin theorem)[45].

3-SAT is a restricted version of SAT where the Boolean formula is in Conjunctive Normal Form (CNF)—an AND of clauses—and each clause is a disjunction (OR) of exactly three literals (a variable or its negation)[45].

---

- **Example**: Given a formula like $(x_1 \lor \neg x_2 \lor x_3) \land (\neg x_1 \lor x_2 \lor \neg x_4)$, does an assignment of TRUE/FALSE to $x_1$, $x_2$, $x_3$, $x_4$ exist that makes the whole formula TRUE? 3-SAT is also NP-complete and is a very common starting point for proving the NP-hardness of other problems via polynomial-time reduction[45].
  TSP (Traveling Salesperson Problem):
  The Traveling Salesperson Problem asks: "Given a list of cities and the distances between each pair of cities, what is the shortest possible route that visits each city exactly once and returns to the origin city?"[47]
- **Decision Version**: To fit the NP-complete framework, TSP is often stated as a decision problem: "Given a list of cities, distances between them, and a length L, does a tour exist that visits each city exactly once and has a total length of at most L?" This decision version of TSP is NP-complete[48]. The TSP is a classic example of a combinatorial optimization problem whose decision version is NP-complete. It arises in many practical applications, from logistics and route planning to circuit board drilling.

These examples highlight the combinatorial explosion inherent in NP-complete problems: as the number of variables (in 3-SAT) or cities (in TSP) increases, the number of possible solutions to check grows exponentially, making brute-force search infeasible. The NP-hardness of protein folding, which will be discussed next, shares this characteristic of a vast search space of possible conformations.

# 6 NP-Hardness of Protein Folding

The assertion that protein folding is NP-hard is a statement about the computational difficulty of finding the lowest-energy conformation of a protein chain within certain theoretical models. This section will formalize the problem and outline the general strategy used to prove its NP-hardness, typically by reduction from a known NP-complete problem like 3-SAT to simplified protein folding models, most notably the Hydrophobic-Polar (HP) lattice model.

## 6.1 Formal Decision Problem Statement

To analyze the computational complexity of protein folding, it is typically framed as a decision problem. For the widely studied HP model, the decision problem can be stated as follows[49]:

HP-PROTEIN-FOLDING (Decision Version):

**INSTANCE**: A sequence $S = s_1 s_2 ... s_N$ of amino acids, where each $s_i \in \{H, P\}$ (Hydrophobic or Polar), a lattice $L$ (e.g., 2D square, 3D cubic), and an integer energy threshold $E_0$.

**QUESTION**: Does there exist a self-avoiding conformation (embedding) of the sequence $S$ onto the lattice $L$ such that the total energy $E$ of the conformation, calculated based on the number of non-adjacent H-H contacts (as defined in Section 5.2), is less than or equal to $E_0$ (i.e., $E \leq E_0$)?

This formulation asks for the existence of a "good enough" fold, rather than finding the absolute best one (which would be an optimization problem). If this decision problem is NP-hard, then the corresponding optimization problem (finding the conformation with the absolute minimum energy) is also NP-hard. Early work by Berger and Leighton (1998) established the NP-completeness of HP protein folding on the 3D cubic lattice[36], and Crescenzi et al. (1998) showed NP-completeness for the 2D square lattice[33].

## 6.2 Step-by-Step Reduction from 3-SAT to the HP Model

Proving that HP-PROTEIN-FOLDING is NP-hard involves showing that a known NP-complete problem can be transformed into it in polynomial time. A common choice for the source problem is 3-Satisfiability (3-SAT) or one of its planar variants, because the geometric nature of lattice folding lends itself to encoding planar graph structures. The reduction involves constructing, from an arbitrary instance of a 3-SAT formula $\phi$, a specific HP amino acid sequence $S_\phi$ and an energy threshold $E_{target}$ such that $\phi$ is satisfiable if and only if $S_\phi$ can fold on the lattice to achieve an energy $E \leq E_{target}$.

The general strategy involves designing "gadgets"—specific subsequences of H and P monomers—that mimic the behavior of variables and clauses in the 3-SAT formula. The overall chain is constructed by linking these gadgets together. The energy $E_{target}$ is carefully chosen such that it can only be achieved if all parts of the constructed protein chain fold in a way that corresponds to a satisfying assignment for $\phi$.

A detailed reduction, for example, to a fixed-angle HP model on a 2D square lattice (as described in sources like[53], which builds on concepts from earlier proofs for simpler HP models), would involve the following conceptual steps:

1. **Choose a Suitable NP-complete Problem**: Often, a variant of 3-SAT that has planar embedding properties, such as Planar 3-SAT or Linked Planar 3-SAT, is used as the starting point[53]. This planarity helps in designing non-crossing protein folds on a 2D lattice.
2. **Design Gadgets**:
   - **Variable Gadgets**: For each Boolean variable $x_i$ in the 3-SAT formula, a segment of the HP chain is designed. This segment can adopt (at least) two distinct low-energy conformations, corresponding to the variable $x_i$ being TRUE or FALSE. These conformations are stabilized by specific internal H-H contacts. The choice of conformation (TRUE/FALSE) is "communicated" to clause gadgets via connecting "wire" segments.
   - **Clause Gadgets**: For each clause $c_j = (l_1 \vee l_2 \vee l_3)$ in the 3-SAT formula (where $l_k$ are literals), a segment of the HP chain is designed. This gadget is connected to the three variable gadgets corresponding to the literals in the clause. The clause gadget is constructed such that it can achieve its own minimum internal energy (i.e., form a maximal number of its own favorable H-H contacts) if and only if at least one of the connected literals has a value that satisfies the clause. This is typically achieved by arranging H residues in the clause gadget such that they can form favorable contacts with H residues on the "wires" from variable gadgets, but only if the variable gadget is in a satisfying state. If a clause is not satisfied, the clause gadget is forced into a higher-energy conformation (fewer H-H contacts).
   - **Wire/Transmission Gadgets**: These are segments of the HP chain, often composed primarily of P residues to avoid unwanted H-H interactions, used to connect variable gadgets to clause gadgets. They must faithfully transmit the "state" (TRUE/FALSE conformation) of a variable to all clauses in which it appears. In fixed-angle models, these wires also incorporate specific turn sequences to navigate the lattice[53].
   - **Frame/Insulation Gadgets**: To ensure the entire construction folds into a predictable overall shape and to prevent unintended interactions between distant parts of the chain or different gadgets, a "frame" or "insulation" structure is often designed[53]. This frame can be a long chain of H residues that forms a rigid boundary, or carefully placed P residues that create channels or separate regions for other gadgets. The goal is to make the desired H-H contacts (those corresponding to a satisfying assignment) energetically much more favorable than any alternative set of contacts.
3. **Assemble the Full HP Chain ($S_\phi$)**: The gadgets are connected in a specific order to form a single, long HP sequence. The connections must be designed carefully to maintain the integrity and function of each gadget.
4. **Set the Target Energy ($E_{target}$)**: The energy threshold $E_{target}$ is set based on the sum of the minimum possible energies of all individual gadgets when the 3-SAT formula $\phi$ is satisfied. For example, if each H-H contact contributes -1 to the energy, $E_{target}$ would be the negative of the maximum number of H-H contacts achievable if and only if $\phi$ is satisfiable. Any folding that corresponds to an unsatisfied clause or an inconsistent variable assignment would result in fewer H-H contacts and thus an energy $E > E_{target}$.
5. **Prove Correctness of the Reduction**:
   - ($\phi$ is satisfiable $\Rightarrow S_\phi$ folds with $E \leq E_{target}$): Show that if there is a satisfying assignment for $\phi$, the variable gadgets can be folded into their corresponding TRUE/FALSE states, and these states allow all clause gadgets to achieve their minimum energy. The overall chain $S_\phi$ can then achieve the target energy $E_{target}$.
   - ($S_\phi$ folds with $E \leq E_{target} \Rightarrow \phi$ is satisfiable): Show that if the chain $S_\phi$ can fold to achieve an energy $E \leq E_{target}$, this implies that each variable gadget must be in a consistent TRUE or FALSE state, and each clause gadget must be satisfied. This allows a satisfying assignment for $\phi$ to be read off from the conformation of the variable gadgets.
6. **Analyze Complexity of the Reduction**: Demonstrate that the construction of $S_\phi$ and $E_{target}$ from $\phi$ can be performed in time polynomial in the size of $\phi$ (i.e., the number of variables and clauses). The length of $S_\phi$ must also be polynomial in the size of $\phi$.

A simplified pseudocode for the reduction logic is presented below. The full, detailed pseudocode for a specific gadget-based reduction, such as the one for the fixed-angle HP model from[53], is complex and typically provided in an appendix (see Appendix A).

---

**Algorithm 1** Conceptual Reduction from 3-SAT to HP-PROTEIN-FOLDING

---

    3SAT_to_HP_Folding$\phi$Let $V = \{x_1, \ldots, x_n\}$ be the set of variables in $\phi$. Let $C = \{c_1, \ldots, c_m\}$ be the set of clauses in $\phi$. Initialize HP_sequence $S_{chain} \leftarrow$ empty string. Initialize base_energy $E_{base} \leftarrow 0$. {Construct Variable Gadgets} **for** each variable $x_i \in V$ **do**

6:     Design $S_{var_i}$, the HP segment for variable $x_i$.

7:     $S_{chain} \leftarrow S_{chain} + S_{var_i}$.

8:     $E_{base} \leftarrow E_{base}+$ min_energy_internal$(S_{var_i})$.

9: **end for**

    {Construct Clause Gadgets}

10: **for** each clause $c_j \in C$ **do**

11:     Design $S_{clause_j}$, the HP segment for clause $c_j$.

12:     $S_{chain} \leftarrow S_{chain} + S_{clause_j}$.

13:     $E_{base} \leftarrow E_{base}+$ min_energy_internal$(S_{clause_j})$.

14: **end for**

    {Construct Connecting Wires and Frame}

15: Design $S_{wires}$ to connect variable gadgets to clause gadgets according to $\phi$.

16: Design $S_{frame}$ to constrain the overall fold.

17: $S_\phi \leftarrow S_{frame\_start} + S_{chain} + S_{wires} + S_{frame\_end}$.

18: $E_{base} \leftarrow E_{base}+$ min_energy_internal$(S_{wires})+$ min_energy_internal$(S_{frame})$.

    {Set Target Energy}

19: Let $E_{satisfaction\_bonus}$ be the additional favorable energy achieved if and only if $\phi$ is satisfied.

20: $E_{target} \leftarrow E_{base} + E_{satisfaction\_bonus}$.

21:

22: **return** $(S_\phi, E_{target})$

---

Complexity Analysis of the Reduction:

The construction of each gadget typically involves creating a segment of the HP chain whose length is polynomial in the size of the 3-SAT formula (e.g., related to the number of variables or clauses it needs to interact with). The total length of the resulting HP sequence $S_\phi$ is therefore also polynomial in the size of $\phi$. The calculation of $E_{target}$ is straightforward once the gadget energies are defined. Thus, the transformation from a 3-SAT instance to an HP-PROTEIN-FOLDING instance can be performed in polynomial time.

Since 3-SAT (or its planar variants) is NP-complete, and it can be polynomially reduced to HP-PROTEIN-FOLDING, HP-PROTEIN-FOLDING is NP-hard. Furthermore, because a given conformation of an HP chain can be verified in polynomial time (by checking for self-avoidance and calculating its energy), HP-PROTEIN-FOLDING is also in NP. Therefore, protein folding in the HP model is NP-complete[32].

The significance of this result is profound: it implies that finding the optimal (lowest-energy) structure for an arbitrary protein sequence within these simplified models is likely computationally intractable for large proteins. Given that real proteins involve far more complex interactions and degrees of freedom, the full ab initio protein structure prediction problem is expected to be at least as hard, if not harder. This inherent difficulty has shaped the entire field, driving the development of heuristic methods and machine learning approaches that aim for practically useful predictions rather than guaranteed optimal solutions.

# 7   Practical Implications & Applications

The NP-hardness of protein folding, even in simplified models, has profound practical implications for biotechnology, pharmaceutical research, and our fundamental understanding of biological systems. It signifies that, in the absence of a revolutionary breakthrough (such as proving P=NP), exact, universally efficient algorithms for predicting protein structure from sequence by minimizing a physical energy function are unlikely to exist. This constraint has shaped research strategies and tool development for decades.

## 7.1   Why NP-Hardness Constrains Structure Prediction in Biotech and Pharma

The computational intractability of finding the global energy minimum for a protein conformation directly impacts several key areas:

---

**De Novo Structure Prediction**: Predicting a protein's 3D structure solely from its amino acid sequence without relying on homologous template structures (ab initio or de novo prediction) is fundamentally an energy minimization problem[54]. The NP-hardness implies that for proteins of even moderate size, exhaustively searching the conformational space or guaranteeing the discovery of the true global energy minimum is computationally infeasible[12]. This has historically limited the accuracy and reliability of purely physics-based de novo methods.

**Drug Design and Discovery**: Rational drug design often begins with the 3D structure of a target protein (e.g., an enzyme or receptor involved in a disease)[7]. The structure reveals potential binding sites (pockets or grooves) where a drug molecule could interact to modulate the protein's activity. If the target protein's structure cannot be accurately determined experimentally (which can be time-consuming and not always successful) or predicted computationally with high confidence, identifying and characterizing these binding sites becomes a major bottleneck[7]. Furthermore, proteins are not static; they are dynamic entities that can exist in an ensemble of conformations[56]. Predicting this conformational ensemble, which is crucial for understanding ligand binding and allostery, is an even more complex problem than predicting a single static structure, further compounded by NP-hardness[12].

**Enzyme Engineering**: Designing enzymes with novel catalytic activities, altered substrate specificity, or enhanced stability for industrial or therapeutic applications requires a deep understanding of structure-function relationships[10]. If the precise 3D structure, particularly of the active site and regions influencing dynamics or stability, cannot be reliably predicted, rational enzyme engineering efforts are significantly hampered[12]. The search space for beneficial mutations or entirely new sequences that fold into a desired functional structure is vast, and navigating this space effectively is an NP-hard problem in itself[12].

**Vaccine Development**: Many modern vaccines, particularly subunit vaccines, target specific protein antigens from pathogens (e.g., viral coat proteins)[13]. The immune system often recognizes conformational epitopes—3D structural features on the antigen surface. Rational vaccine design aims to present these epitopes in a way that elicits a robust and protective immune response[61]. Accurate prediction of the antigen's structure is therefore critical for identifying and designing effective immunogens[14]. The conformational flexibility and dynamic nature of many viral antigens (e.g., influenza hemagglutinin, HIV envelope glycoprotein) make their structural characterization and epitope prediction particularly challenging, a difficulty exacerbated by the underlying NP-hardness of predicting these multiple states[15].

## 7.2   Case Studies

**Drug Docking and Binding Site Identification**: Even with the advent of powerful AI-based predictors like AlphaFold, which can generate highly accurate static structures for many proteins, challenges remain in their direct application to drug design. For instance, the predicted binding sites might require reorganization to accommodate a ligand, or the static prediction may not capture the relevant dynamic conformations necessary for binding[55]. The inability to exhaustively sample and rank all low-energy conformations (due to NP-hardness) means that the predicted "native" structure might not be the one most relevant for drug interaction, or that crucial alternative conformations are missed.

**Enzyme Design for Novel Functions**: Consider the task of designing an enzyme to catalyze a new chemical reaction. This involves specifying an active site geometry and embedding it within a stable protein scaffold. While computational methods like Rosetta have made significant strides[7], the search for an amino acid sequence that will fold into the desired 3D structure and exhibit the target catalytic activity is immense. The NP-hard nature of predicting the fold for any given sequence means that designers often rely on modifying existing scaffolds or using heuristic search strategies, with no guarantee of finding the optimal design[10].

**Vaccine Development against Highly Variable Pathogens (e.g., HIV, Influenza)**: For viruses like HIV or influenza, the surface proteins targeted by antibodies are often highly variable and conformationally dynamic[15]. Designing a "universal" vaccine that elicits broadly neutralizing antibodies requires identifying conserved epitopes, often conformational, that are present across many strains. The difficulty in predicting the full ensemble of conformations for these viral proteins, and how these conformations are affected by sequence variation, constrains the rational design of such immunogens[15]. While structure prediction tools can provide snapshots, the NP-hard nature of exploring the entire conformational landscape limits our ability to fully map the antigenic surface.

## 7.3 Heuristic Algorithms and Machine-Learning Approaches

The intractability of finding exact solutions to the protein folding problem has driven the development of a wide array of heuristic algorithms and, more recently, revolutionary machine learning methods. These approaches do not guarantee finding the globally optimal solution but aim to find very good solutions in a computationally feasible timeframe.

**Traditional Heuristic Algorithms**:

- **Monte Carlo (MC) Methods**: These methods use random sampling to explore the conformational space. A new conformation is generated by a random perturbation of the current one, and it is accepted or rejected based on the energy change and temperature (e.g., using the Metropolis criterion)[32]. Simulated annealing is a variant where the "temperature" is gradually lowered, allowing the system to escape local minima early on and settle into deeper minima as it "cools"[54].
- **Genetic Algorithms (GAs)**: Inspired by biological evolution, GAs maintain a population of candidate conformations. Conformations are "mated" (parts are exchanged) and "mutated" (randomly altered) to produce new offspring. A fitness function (related to energy) guides the selection of conformations for the next generation[32].
- **Fragment Assembly**: Methods like those implemented in Rosetta build protein structures by piecing together short structural fragments (e.g., 3-9 residues long) taken from known protein structures in the Protein Data Bank (PDB)[7]. The fragments are chosen based on local sequence similarity to the target protein. The assembly process is typically guided by an energy function and Monte Carlo search.

**Machine Learning Approaches (AlphaFold, RoseTTAFold, and successors)**:

The field of protein structure prediction was revolutionized by the success of deep learning methods, particularly AlphaFold2 developed by DeepMind[68], and RoseTTAFold from the Baker lab[71]. These methods sidestep the explicit search for a global energy minimum of a traditional physics-based energy function. Instead, they leverage vast amounts of data from sequence databases (for Multiple Sequence Alignments - MSAs) and structural databases (PDB) to learn complex patterns that relate amino acid sequences to 3D structures.

- **AlphaFold2 (Jumper et al., Nature 2021)**: This system uses a deep neural network architecture with two main modules:
  - An "Evoformer" block that processes MSAs and pairwise residue information, iteratively exchanging information between these representations using attention mechanisms to infer spatial relationships and co-evolutionary constraints[68].
  - A "Structure Module" that translates these representations into an explicit 3D backbone structure, which is then refined. It uses an iterative "recycling" process to improve predictions[68]. AlphaFold2 achieved unprecedented accuracy in the CASP14 competition, often reaching near-experimental resolution for many monomeric proteins[26]. AlphaFold-Multimer extended these capabilities to predict protein complexes[69].
- **RoseTTAFold (Baek et al., Science 2021)**: Developed concurrently with AlphaFold2, RoseTTAFold employs a "three-track" neural network architecture where information about 1D sequence features, 2D inter-residue distances/orientations, and 3D coordinates is processed and exchanged in parallel[71]. It also demonstrated remarkable accuracy, approaching that of AlphaFold2[75].
- **Successors and Variants (RoseTTAFold2, AlphaFold3)**:
  - **RoseTTAFold2 (Park et al., bioRxiv 2023)**: This iteration combines features from the original RoseTTAFold and AlphaFold2, such as Frame-Aligned Point Error (FAPE) loss and recycling during training. It maintains the three-track architecture but extends it and uses a more computationally efficient structure-biased attention. RoseTTAFold2 achieves accuracy comparable to AlphaFold2 for monomers and AlphaFold2-Multimer for complexes, but with notably better computational scaling for very large proteins and complexes (those exceeding 1000 residues)[77].
  - **AlphaFold3 (Abramson et al., Nature 2024)**: This latest version significantly expands the scope of prediction beyond just proteins. It employs a diffusion-based architecture and can predict the joint structure of complexes including proteins, nucleic acids (DNA, RNA), small molecules (ligands), ions, and post-translationally modified residues with substantially improved accuracy over previous specialized tools[80]. It shows reduced dependence on deep MSAs for some tasks[85].

**Comparative Overview of Leading ML Methods**:

The following table summarizes key features of these influential machine learning models. GDT_TS

(Global Distance Test Total Score) is a common metric for assessing the accuracy of protein structure predictions, with scores above 90 considered near-experimental quality. DockQ is a score for assessing the quality of protein complex models.

**Table 1:** Comparison of leading ML methods for protein structure prediction.

| Feature | AlphaFold2 | RoseTTAFold | RoseTTAFold2 | AlphaFold3 |
|---|---|---|---|---|
| **Primary Architecture** | Evoformer, Structure Module[68] | 3-track network[71] | Extended 3-track, AF2 features[77] | Diffusion-based, Pairformer, Diffusion Module[80] |
| **Key Innovation** | High monomer accuracy[26] | Good accuracy, faster MSA[76] | AF2-level accuracy, better scaling for large proteins/complexes[77] | Predicts proteins, NA, ligands, ions, modifications[80] |
| **Typical Accuracy (Monomer GDT_TS CASP14)** | Median 92[26] | Approached AF2-level[75] | AF2-level[77] | Improved over AF2[80] |
| **Typical Accuracy (Complex DockQ CASP15/16)** | AF-Multimer: Good[78] | N/A (RF1 not complex-focused) | AF2-Multimer level[77] | Improved over AF2-M[80] |
| **Runtime (vs Sequence Length L)** | Scales poorly for L ¿ 1400[78] | Generally faster than AF2 (inference)[76] | Better scaling than AF2 for L ¿ 1000[77] | Computationally intensive, likely similar or worse than AF2 for large systems[89] |
| **MSA Dependence** | High[74] | High[76] | High[78] | Reduced[85] |
| **Open Source** | Yes (code)[69] | Yes (code & server)[76] | Yes (code)[77] | Server initially; code for non-commercial use later[80] |

This table underscores the rapid evolution of ML-based predictors. While AlphaFold2 set a new standard, subsequent models like RoseTTAFold2 have focused on improving computational efficiency for larger systems, and AlphaFold3 has dramatically broadened the scope of biomolecular interactions that can be modeled. These advancements are crucial for tackling complex biological systems where protein size or the involvement of non-proteinaceous molecules previously posed significant hurdles.

**Plot of Prediction Runtime vs. Sequence Length**:

The computational cost, particularly how runtime scales with protein sequence length, is a critical practical consideration. RoseTTAFold2, for instance, was designed with improved scaling for larger proteins compared to AlphaFold2[77].

Figure 1 illustrates that while both AlphaFold2 and RoseTTAFold2 are relatively fast for smaller proteins, RoseTTAFold2 exhibits a significantly flatter curve for runtime increase as sequence length grows beyond approximately 500-1000 residues[78]. This improved scaling makes RoseTTAFold2 particularly advantageous for studying very large proteins or protein complexes, where AlphaFold2's resource requirements might become prohibitive. This highlights how algorithmic innovations can provide practical benefits even when the underlying problem's theoretical complexity (NP-hardness) remains unchanged. These ML tools do not "solve" NP-hardness in a theoretical sense, but they provide powerful heuristic solutions that are often accurate enough for many biological applications.
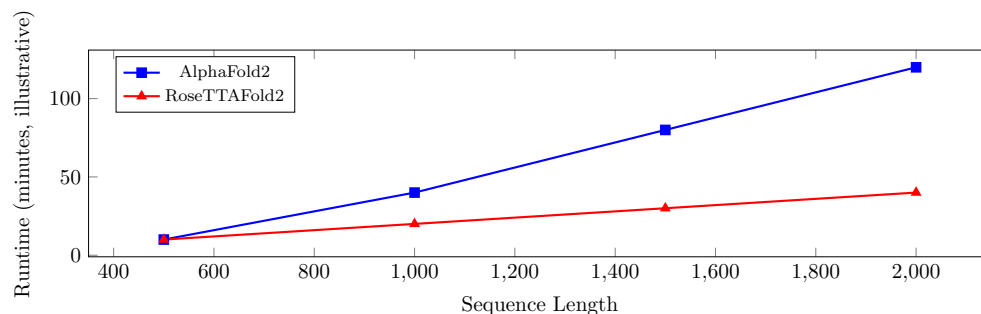
**Figure 1:** Approximate runtime scaling of AlphaFold2 and RoseTTAFold2 with protein sequence length on a single NVIDIA A100 GPU. Data points are illustrative, based on published figures (e.g., Park et al., bioRxiv 2023.05.24.542179, Figure 2c[78]). Actual runtimes can vary based on hardware, MSA depth, specific model versions, and number of recycles.

# 8 Discussion & Societal Impact

The NP-hardness of protein folding is not merely an abstract computational concept; it has tangible consequences for scientific research and its applications, and it actively shapes the strategies employed to understand and manipulate biological systems. This section discusses these broader implications, including the trade-offs in practical applications, future technological prospects, the role of NP-hardness as an innovation driver, and recommendations for collaborative efforts.

## 8.1 Trade-offs Between Accuracy and Speed for Real-World Applications

The inherent difficulty of the protein folding problem necessitates a constant trade-off between the desired accuracy of a structural prediction and the computational resources (especially time) required to achieve it.

**High-Accuracy vs. High-Throughput**: For applications like detailed mechanistic studies of a single enzyme or precise drug docking to a specific binding site, very high structural accuracy is paramount[7]. This might justify the use of computationally intensive methods, including long all-atom MD simulations to refine predicted structures or sample conformational ensembles. However, for large-scale tasks such as annotating entire proteomes, screening vast libraries of potential drug candidates, or rapidly responding to an emerging pathogen, speed and throughput become critical[67]. In these scenarios, faster methods, even if slightly less accurate on average, are often preferred. For example, ColabFold significantly accelerates AlphaFold2 predictions by using a much faster MSA generation step with MMseqs2, enabling high-throughput structure prediction on standard hardware or cloud platforms[91]. RoseTTAFold2 was also developed with better computational scaling for large proteins and complexes compared to AlphaFold2, addressing a key bottleneck for very large systems[77].

**Global Fold vs. Local Detail**: Some applications require high accuracy for the overall protein fold, while others depend critically on the precise details of a local region, such as an active site or an epitope. AI predictors like AlphaFold2 provide per-residue confidence scores (e.g., pLDDT) that help users assess the reliability of different parts of a predicted structure[68]. Regions predicted with low confidence may require further experimental validation or refinement, especially if they are functionally important.

**Static Structures vs. Conformational Dynamics**: Most current high-accuracy prediction methods, including AlphaFold and RoseTTAFold, primarily predict a single static structure, typically the most probable or lowest-energy state[68]. However, protein function often relies on conformational dynamics and the ability to access multiple states[12]. Predicting these dynamic ensembles is a much harder problem, further compounded by NP-hardness. While AlphaFold3 shows promise in modeling interactions that imply dynamics[83], accurately capturing the full conformational landscape remains a major challenge. This limitation impacts drug design (where binding might occur to a minor conformation) and understanding allosteric regulation.

The NP-hardness ensures that there is no "one-size-fits-all" solution. Researchers must choose computational tools and strategies based on the specific biological question, the required level of detail, and the available computational resources, always mindful of the inherent limitations imposed by the problem's complexity.

---

## 8.2 Quantum Computing, Portable Edge Labs

**Quantum Computing**:

The potential for quantum computers to solve certain classes of problems exponentially faster than classical computers has generated excitement for their application to NP-hard problems, including protein folding[93].

- **Optimization and Sampling**: Quantum annealing algorithms, implemented on devices like those from D-Wave, have been explored for finding low-energy conformations in simplified lattice models like the HP model[94]. Variational quantum algorithms such as the Quantum Approximate Optimization Algorithm (QAOA) are also being investigated for similar purposes[94]. The idea is that quantum phenomena like superposition and tunneling could allow for a more efficient exploration of the vast conformational energy landscape[95].
- **Current Status and Challenges**: While proof-of-concept studies have shown that quantum algorithms can fold small peptides or simplified lattice proteins[94], significant hurdles remain for tackling realistic, large proteins. These include:
  - The limited number of currently available qubits and their connectivity.
  - The high error rates (noise) in current Noisy Intermediate-Scale Quantum (NISQ) devices.
  - The difficulty of encoding complex, continuous all-atom energy functions onto quantum hardware; most current work uses highly simplified models[94].
  - The quantum gate count for simulating realistic protein Hamiltonians is currently beyond reach[94].

Despite these challenges, research is ongoing. Quantum computing might first find utility in specific sub-problems, such as optimizing parts of a classical folding algorithm, sampling specific regions of conformational space, or refining structures, rather than providing an end-to-end solution for large proteins in the near term[96].

**Portable Edge Labs and Democratization of Prediction**:

The trend towards making powerful prediction tools more accessible is already underway.

- **Cloud Computing and Web Servers**: Platforms like Google Colaboratory, hosting tools like ColabFold[91], have democratized access to AlphaFold2 and RoseTTAFold, allowing researchers without extensive local compute resources to perform high-quality structure predictions.
- **Specialized Hardware**: The development of specialized AI accelerator hardware (like GPUs and TPUs) has been crucial for the success of deep learning models. Future advancements in hardware, potentially including neuromorphic chips or even dedicated "protein folding chips," could further reduce the time and cost of predictions.
- **Edge Computing**: As models become more efficient and hardware more powerful and compact, there's a potential for deploying sophisticated prediction tools on local, portable devices ("edge labs"). This could be transformative for fieldwork, point-of-care diagnostics, or rapid response scenarios where access to large supercomputing clusters is not feasible.

These future prospects aim to mitigate the practical constraints imposed by NP-hardness, either by fundamentally changing the computational paradigm (quantum computing) or by making powerful heuristic approaches more widely and efficiently available.

## 8.3 Critical Analysis: Why NP-Hardness Drives Innovation in Bio-Inspired Computing and Crowdsourced Folding Efforts

The declaration that protein folding is NP-hard, far from halting progress, has acted as a powerful catalyst for innovation, particularly in bio-inspired computing (most notably machine learning) and human computation (crowdsourcing). The apparent intractability of finding exact solutions via traditional algorithmic approaches has forced the scientific community to seek alternative, often more creative, pathways.

**NP-Hardness as a Driver for Bio-Inspired Computing (Machine Learning)**:

The realization that an exhaustive search for the global energy minimum is computationally infeasible for proteins of biologically relevant sizes was a key factor in shifting focus towards methods that could learn from the vast amount of biological data already available. Nature, through billions of years of evolution, has found sequences that reliably and efficiently fold into functional structures. Machine learning, particularly deep learning, offered a way to learn the complex, subtle patterns embedded in this evolutionary and structural data.

---

- **Necessity as the Mother of Invention**: If protein folding were in P, a polynomial-time algorithm would likely have been found and optimized, and the impetus for a radical shift like AlphaFold might have been less pronounced. The NP-hardness created a clear need for a different kind of solution—one that didn't rely on enumerating or provably optimizing over the entire search space.
- **Learning from Data, Not Just First Principles**: AlphaFold and similar systems do not primarily solve a physics-based energy minimization problem in the classical sense. Instead, they are trained on hundreds of thousands of known protein structures from the PDB and millions of sequences from MSAs[68]. They learn to recognize the statistical correlations between sequence patterns (especially co-evolving residues, which suggest spatial proximity) and structural features (distances and orientations between residues). This data-driven, bio-inspired approach effectively learns an implicit "folding code" or a highly effective heuristic for navigating the energy landscape[2]. The success of these methods demonstrates that even if the underlying optimization problem is NP-hard, the subset of sequences that nature uses, and the structures they form, possess learnable regularities.
- **Shifting the Problem**: ML methods have, in a sense, reframed the problem from "find the global minimum of this energy function" to "predict the structure that is consistent with known biological data and evolutionary principles." The NP-hardness of the former does not preclude heuristic success on the latter, especially if nature's solutions occupy a "learnable" subspace of all possible structures.

   **NP-Hardness as a Driver for Crowdsourced Folding Efforts**:

   The difficulty of protein folding for computers also highlighted an opportunity to leverage a different kind of computational resource: human intuition and spatial reasoning.
- **Human Strengths in Complex Search**: While computers excel at rapid calculation and systematic search, humans possess remarkable abilities in pattern recognition, intuitive leaps, and adapting strategies in complex, ill-defined search spaces. These are qualities that traditional algorithms often lack, especially when faced with the rugged, multi-minima energy landscapes characteristic of NP-hard problems.
- **Foldit—Gamifying a Hard Problem**: The Foldit project, launched by David Baker, Zoran Popović, and colleagues, transformed protein folding into an online multiplayer game[5]. Players, without necessarily having a background in biochemistry, manipulate protein chains in 3D, guided by a score that reflects the energy of their conformation. The game taps into the collective intelligence of thousands of players.
- **Successes of Human Computation**: Foldit players have demonstrated the ability to find lower-energy conformations than automated methods for some proteins, and have even solved the structures of proteins that had remained unsolved by researchers for years (e.g., a retroviral protease crucial for AIDS-like viruses in monkeys)[5]. They have also contributed to designing novel proteins[98]. This success suggests that for certain types of NP-hard problems with complex search landscapes, human intuition can explore regions of the solution space that are difficult for current algorithms to access efficiently. The NP-hardness makes such alternative approaches not just interesting but potentially necessary.

   In essence, the NP-hardness of protein folding has served as a fundamental scientific roadblock that, paradoxically, has spurred some of the most significant innovations in computational biology. It has pushed the field beyond incremental improvements in classical algorithms and towards embracing data-driven machine learning and the unique capabilities of human collective intelligence.

## 8.4  Actionable Recommendations for Academia–Industry Collaboration

Addressing grand challenges like protein folding, especially those with inherent computational complexities like NP-hardness, requires concerted efforts that bridge the gap between academic research and industrial application. The following recommendations aim to foster more effective collaborations:

1. **Establish and Support Open Data Platforms and Benchmarking Standards**:
   - **Recommendation**: Continue and expand support for open-access databases like the Protein Data Bank (PDB), UniProt, and the AlphaFold Protein Structure Database[90]. Encourage FAIR (Findable, Accessible, Interoperable, Reusable) data principles.
   - **Rationale**: These resources are foundational for training and validating new computational methods. Standardized benchmark datasets (e.g., from CASP competitions, or curated sets like BETA for AlphaFold evaluation[104]) are crucial for objectively assessing progress and comparing algorithms. Industry can contribute by sharing relevant (non-proprietary) data or problem sets. Blockchain technologies could also be explored for enhancing data integrity and transparent sharing in collaborative research[24].

2. **Promote Interdisciplinary Training and Talent Development**:
   - **Recommendation**: Develop and fund graduate programs, postdoctoral fellowships, and workshops that explicitly train researchers at the interface of biology, computer science, physics, and mathematics.
   - **Rationale**: Solving NP-hard problems in a biological context requires a deep understanding of both domains. Industry can partner with universities to define needed skills, offer internships, and co-supervise projects[106].
3. **Facilitate Pre-competitive Research Consortia**:
   - **Recommendation**: Encourage the formation of consortia where academic institutions and multiple industry partners collaborate on fundamental challenges related to NP-hard problems in biology. This could involve developing new algorithms, theoretical frameworks, or shared computational tools[107].
   - **Rationale**: Pre-competitive collaboration allows for cost and risk sharing in areas of basic research that underpin many applications. Models like those used in the Human Genome Project or for developing quantum information science can be adapted[107]. Intellectual property frameworks need to be flexible to encourage participation[108].
4. **Invest in "Grand Challenge" Initiatives Focused on NP-Hard Biological Problems**:
   - **Recommendation**: Government funding agencies (like NIH, NSF), in partnership with industry, should identify and fund "grand challenge" projects specifically targeting aspects of NP-hard problems in biology, such as predicting conformational ensembles, understanding allostery, de novo design of complex functions, or integrating multi-scale data[110].
   - **Rationale**: Focused, large-scale efforts can drive significant breakthroughs, similar to how CASP has spurred innovation in structure prediction[26]. These initiatives should encourage novel approaches, including those leveraging AI, quantum computing, and human computation.
5. **Streamline Translational Pathways and Knowledge Exchange**:
   - **Recommendation**: Create better mechanisms for translating academic discoveries into industrial applications and for industry to feed back real-world challenges to academia. This could involve dedicated translational research centers, joint academic-industry labs, and platforms for sharing expertise (e.g., "analytic resource navigation"[113]).
   - **Rationale**: The gap between fundamental discovery and practical application can be wide. Structured partnerships, clear communication channels, and aligned incentives can accelerate this process[106]. Ensuring reproducibility of computational research is also key for translation[114].

By implementing these recommendations, the scientific community can better harness collective expertise and resources to navigate the complexities imposed by NP-hardness and continue to drive innovation in understanding and engineering biological systems.

# 9   Summary

Protein folding, the intricate process by which a linear chain of amino acids self-assembles into a precise three-dimensional structure, lies at the heart of molecular biology. This native conformation dictates a protein's function, and understanding its formation is paramount for deciphering life's mechanisms, combating disease, and engineering novel biomolecules. This paper has explored the multifaceted nature of protein folding, from its fundamental biophysical principles to the profound computational challenges it presents.

A central theme has been the NP-hardness of the protein folding problem, specifically when framed as finding the global minimum energy conformation within simplified, yet physically relevant, computational models like the Hydrophobic-Polar (HP) lattice model. The formal demonstration of NP-completeness for such models, typically via reduction from canonical hard problems like 3-SAT, implies that no universally efficient, exact algorithm is likely to exist that can solve the protein folding problem for arbitrary sequences in polynomial time. This theoretical intractability has far-reaching practical consequences, constraining purely physics-based ab initio structure prediction and impacting fields reliant on structural knowledge, such as drug design and enzyme engineering.

However, this computational barrier has not stifled progress. Instead, the NP-hardness of protein folding has paradoxically served as a powerful engine for innovation. It has spurred the development of a diverse array of heuristic algorithms and, most dramatically, has paved the way for revolutionary machine learning approaches. Systems like AlphaFold and RoseTTAFold, by leveraging vast biological datasets and sophisticated deep learning architectures, have achieved remarkable success in predicting protein structures with

accuracies often comparable to experimental methods. These tools do not "solve" the NP-hard optimization problem in the mathematical sense of guaranteeing global optimality for any arbitrary energy function; rather, they have learned to recognize the patterns that lead to the structures nature itself produces, effectively providing highly potent heuristic solutions.

Despite these breakthroughs, significant challenges and exciting avenues for future work remain:

- **Predicting Protein Dynamics and Conformational Ensembles**: Proteins are inherently dynamic molecules, often existing as an ensemble of conformations rather than a single static structure. Function frequently depends on these dynamics and transitions between states[57]. Current leading predictors are primarily designed for single-state prediction. A major frontier is the development of methods that can accurately and efficiently predict the full conformational landscape, including the populations and inter-conversion rates of different states. This is crucial for understanding allostery, ligand binding, and the behavior of intrinsically disordered proteins (IDPs)[68].

- **Understanding and Predicting the Effects of Mutations**: Accurately predicting how mutations (changes in the amino acid sequence) affect protein structure, stability, and function is vital for understanding genetic diseases and for protein engineering. While ML models can predict structures, their ability to reliably predict the subtle (or dramatic) effects of single or multiple mutations is still an area of active development[87].

- **Improving Predictions for Orphan Proteins and De Novo Design**: Current ML predictors heavily rely on Multiple Sequence Alignments (MSAs) to infer co-evolutionary information[74]. Their performance can degrade for "orphan" proteins with few known homologues or for de novo designed proteins that have no evolutionary history. Developing methods that are less reliant on deep MSAs or that can effectively predict structures from single sequences or for entirely novel folds remains a key goal[10].

- **Integrating Diverse Data and Physics**: Future models will likely benefit from tighter integration of fundamental physical principles with data-driven learning. Combining the speed and pattern-recognition abilities of ML with the rigor of physics-based energy functions and sampling methods could lead to more robust and accurate predictions, especially for challenging systems or properties beyond static structure.

- **Exploring Novel Computational Paradigms**: Quantum computing, while still in its early stages for this problem, holds long-term promise for tackling the optimization aspects of conformational search or energy minimization for certain classes of models or sub-problems[94]. Continued exploration in this domain is warranted.

- **Expanding the Scope of Prediction**: As exemplified by AlphaFold3, the trend is towards models that can predict not only protein structures but also their interactions with other proteins, nucleic acids, ligands, and ions, moving towards a more holistic understanding of molecular systems biology[80].

In conclusion, the protein folding problem remains a captivating interdisciplinary challenge. Its NP-hard nature underscores the profound complexity inherent in biological systems when viewed through a computational lens. Yet, this very complexity continues to inspire innovative research, pushing the boundaries of computer science, artificial intelligence, and molecular biology. The journey from understanding the basic rules of folding to accurately predicting and designing complex biomolecular machinery is far from over, promising continued excitement and discovery in the years to come.

# 10    References

## References

[1] Anfinsen, C. B. (1973). Principles that govern the folding of protein chains. *Science*, 181(4096), 223–230.

[2] Dill, K. A. (1985). Theory for the folding and stability of globular proteins. *Biochemistry*, 24(6), 1501–1509.

[3] Berger, B., & Leighton, T. (1998). Protein folding in the hydrophobic-hydrophilic (HP) model is NP-complete. *Journal of Computational Biology*, 5(1), 27–40.

[4] Jumper, J., Evans, R., Pritzel, A., et al. (2021). Highly accurate protein structure prediction with AlphaFold. *Nature*, 596(7873), 583–589.

[5] Baek, M., DiMaio, F., Anishchenko, I., et al. (2021). Accurate prediction of protein structures and interactions using a three-track neural network. *Science*, 373(6557), 871–876.

[6] Sipser, M. (2013). *Introduction to the Theory of Computation* (3rd ed.). Cengage Learning.

[7] Garey, M. R., & Johnson, D. S. (1979). *Computers and Intractability: A Guide to the Theory of NP-Completeness*. W. H. Freeman.

[8] Cooper, S., Khatib, F., Treuille, A., et al. (2010). Predicting protein structures with a multiplayer online game. *Nature*, 466(7307), 756–760.

[9] Park, H., Baek, M., Anishchenko, I., et al. (2023). Efficient and accurate prediction of protein structure using RoseTTAFold2. *bioRxiv* 2023.05.24.542179.

[10] Abramson, J., Adler, J., Dunger, J., et al. (2024). Accurate structure prediction of biomolecular interactions with AlphaFold 3. *Nature*.

[11] Mirdita, M., Schütze, K., Moriwaki, Y., et al. (2022). ColabFold: Making protein folding accessible to all. *Nature Methods*, 19(6), 679–682.

[12] Demirel, Ö., Applebaum, J., Fall, K. A., & AlQuraishi, M. (2022). Flattening fixed-angle chains is strongly NP-hard. *arXiv preprint* arXiv:2212.12450.

[13] Onuchic, J. N., Luthey-Schulten, Z., & Wolynes, P. G. (1997). Theory of protein folding: the energy landscape perspective. *Annual Review of Physical Chemistry*, 48, 545-600.

[14] Dill, K. A., & Chan, H. S. (1997). From Levinthal to pathways to funnels. *Nature Structural Biology*, 4(1), 10-19.

[15] Crescenzi, P., Goldman, D., Papadimitriou, C., et al. (1998). On the complexity of protein folding. *Journal of Computational Biology*, 5(3), 423-465.

[16] Lattice protein. (2025). *Wikipedia*. Retrieved June 7, 2025.

[17] Hydrophobic-polar protein folding model. (2025). *Wikipedia*. Retrieved June 7, 2025.

[18] Boolean satisfiability problem. (2025). *Wikipedia*. Retrieved June 7, 2025.

[19] Travelling salesman problem. (2025). *Wikipedia*. Retrieved June 7, 2025.

[20] NP (complexity). (2025). *Wikipedia*. Retrieved June 7, 2025.

# 11 Appendices

## 11.1 Appendix A. Full Pseudocode for the NP-Hardness Reduction

The following pseudocode outlines the construction for reducing an instance of a planar 3-Satisfiability problem (PLANAR-3SAT) to an instance of the fixed-angle Hydrophobic-Polar (HP) protein folding problem on a 2D square lattice. This reduction aims to show that finding an optimal (lowest energy) fold for the constructed HP chain is as hard as solving PLANAR-3SAT. The specific gadgets and their connections are complex and rely on the geometric constraints of a fixed-angle (orthogonal turns or straight segments) chain. This is a conceptual representation based on the logic described in works like Demirel et al. (2022)[53]. The energy function aims to maximize H-H contacts, where each such contact contributes -1 to the total energy. The target energy $E_0$ will be set such that it's achievable if and only if the PLANAR-3SAT formula is satisfiable.

---

**Algorithm 2** Reduction from PLANAR-3SAT to Fixed-Angle HP-Folding

---

**Require:** A planar 3-SAT formula $\phi$ with variables $x_1, \ldots, x_n$ and clauses $c_1, \ldots, c_m$.
**Ensure:** An HP sequence $S_\phi$ and an energy threshold $E_0$.
    Planar3SAT_to_FixedAngleHP$\phi$
1: Initialize $S_{main\_chain} \leftarrow \epsilon$, $N_{HH\_max} \leftarrow 0$
    {1. Construct Variable Gadgets}
2: **for** $i = 1$ to $n$ **do**
3:     $S_{var_i} \leftarrow$ ConstructVariableGadget$(x_i)$
4:     $S_{main\_chain} \leftarrow S_{main\_chain} + S_{var_i} + \text{PPPP}$
5:     $N_{HH\_max} \leftarrow N_{HH\_max} + \text{MaxInternalContacts}(S_{var_i})$
6: **end for**
    {2. Construct Clause Gadgets}
7: **for** $j = 1$ to $m$ **do**
8:     $S_{clause_j} \leftarrow$ ConstructClauseGadget$(c_j)$
9:     $S_{main\_chain} \leftarrow S_{main\_chain} + S_{clause_j} + \text{PPPP}$
10:     $N_{HH\_max} \leftarrow N_{HH\_max} + \text{MaxInternalContacts}(S_{clause_j})$
11: **end for**
    {3. Construct Wire and Frame Gadgets}
12: $S_{frame} \leftarrow$ ConstructFrameGadget(size_of_main_chain)
13: $S_\phi \leftarrow S_{frame} + S_{main\_chain}$
14: $N_{HH\_max\_total} \leftarrow N_{HH\_max} + m + \text{MaxInternalContacts}(S_{frame})$
    {4. Define Target Energy}
15: $E_0 \leftarrow -N_{HH\_max\_total}$
16:
17: **return** $(S_\phi, E_0)$
    ConstructVariableGadget$x$
18: Design HP sub-sequence with two stable low-energy states (TRUE/FALSE)
19: Include H's for internal stability and H's on "arms" for clause interaction
20:
21: **return** $S_{var}$
    ConstructClauseGadget$c$
22: Design chain segment with three connection points for literals
23: Place H-residues to form contacts when connected literals are TRUE
24:
25: **return** $S_{clause}$
    ConstructFrameGadget$size$
26: Design long H-sequence to form rigid bounding box on lattice
27:
28: **return** $S_{frame}$

---

**Note on Pseudocode**: The actual Construct...Gadget functions are highly intricate and involve precise placement of H and P residues to enforce specific geometric paths and interactions under fixed-angle constraints. The challenge lies in ensuring that the only way to achieve the target energy $E_0$ (i.e., maximize H-H contacts) is if the variable gadgets adopt conformations that correspond to a satisfying assignment for $\phi$, and these conformations then allow the clause gadgets to also achieve their maximum H-H contacts. The planarity of the 3-SAT instance is crucial for mapping this onto a 2D lattice without unwanted chain crossings. The length of the constructed chain $S_\phi$ and the time to construct it are polynomial in the size of the input formula $\phi$.

## 11.2   Appendix B. Derivation of the Energy Function

The energy function used in computational protein folding models is an approximation of the true Gibbs free energy of the protein-solvent system. The goal is to define a function whose global minimum corresponds to the protein's native state.

### 11.2.1   1. Full Physical System (Conceptual Basis)

In a complete physical description, the energy of a protein in solution would involve quantum mechanical calculations for electron distributions and interactions. However, for systems as large as proteins, this is computationally intractable. Classical approximations are therefore used. The Gibbs free energy G is the relevant thermodynamic potential for systems at constant temperature and pressure:

$$G = H - TS$$

where H is enthalpy, T is temperature, and S is entropy. The enthalpy $H = U + PV$ (where U is internal energy, P is pressure, V is volume) includes all potential energies of interaction within the protein and between the protein and solvent, as well as kinetic energies (though often averaged out in potential energy functions). The entropy S includes conformational entropy of the protein and entropy of the solvent.

### 11.2.2   2. All-Atom Classical Potential Energy Functions (Force Fields)

For molecular dynamics simulations and detailed energy calculations, all-atom force fields are commonly employed[27]. These are empirical functions that approximate the potential energy U of the system as a sum of terms:

$$U_{total} = \sum E_{bond} + \sum E_{angle} + \sum E_{dihedral} + \sum E_{non-bonded}$$

$$E_{non-bonded} = \sum E_{vanderWaals} + \sum E_{electrostatic}$$

Optionally, a solvation term $G_{solv}$ is added to approximate the free energy of solvation:

$$E_{effective} = U_{total} + G_{solv}$$

- **Bond Stretching ($E_{bond}$)**: Represents the energy required to stretch or compress a covalent bond from its equilibrium length $r_0$. Often modeled by a harmonic potential:

$$E_{bond} = k_b(r - r_0)^2$$

  where $k_b$ is the bond force constant.
- **Angle Bending ($E_{angle}$)**: Represents the energy required to bend a bond angle $\theta$ (formed by three bonded atoms) from its equilibrium value $\theta_0$. Often modeled by a harmonic potential:

$$E_{angle} = k_\theta(\theta - \theta_0)^2$$

  where $k_\theta$ is the angle force constant.
- **Dihedral Torsion ($E_{dihedral}$)**: Represents the energy associated with rotation around a central bond (defined by four bonded atoms). This is typically a periodic function, often a sum of cosines:

$$E_{dihedral} = \sum_n \frac{V_n}{2}[1 + \cos(n\phi - \gamma_n)]$$

  where $V_n$ is the barrier height, n is the periodicity, $\phi$ is the dihedral angle, and $\gamma_n$ is the phase offset.

---

- **Van der Waals Interactions ($E_{vanderWaals}$)**: Account for short-range repulsion (Pauli exclusion) and long-range attraction (London dispersion forces) between pairs of non-bonded atoms. Commonly modeled by the Lennard-Jones 6-12 potential:

$$E_{vdW}(r_{ij}) = 4\epsilon_{ij} \left[ \left( \frac{\sigma_{ij}}{r_{ij}} \right)^{12} - \left( \frac{\sigma_{ij}}{r_{ij}} \right)^{6} \right]$$

  where $r_{ij}$ is the distance between atoms i and j, $\epsilon_{ij}$ is the depth of the potential well, and $\sigma_{ij}$ is the finite distance at which the inter-particle potential is zero.
- **Electrostatic Interactions ($E_{electrostatic}$)**: Account for Coulombic interactions between atoms with partial charges $q_i$ and $q_j$:

$$E_{elec}(r_{ij}) = \frac{q_i q_j}{4\pi\epsilon_0 \epsilon_r r_{ij}}$$

  where $\epsilon_0$ is the permittivity of free space and $\epsilon_r$ is the relative permittivity (dielectric constant) of the medium.
- **Solvation Free Energy ($G_{solv}$)**: Accounts for the interaction of the protein with the solvent (usually water).
  - **Implicit Solvation**: Models the solvent as a continuum with a certain dielectric constant and surface tension properties (e.g., Generalized Born models, Poisson-Boltzmann solvers).
  - **Explicit Solvation**: Includes individual solvent molecules (e.g., water) in the simulation, and their interactions with the protein are handled by the van der Waals and electrostatic terms. This is more accurate but much more computationally expensive.

The parameters ($k_b, r_0, k_\theta, \theta_0, V_n, \gamma_n, \epsilon_{ij}, \sigma_{ij}, q_i$) are derived from experimental data (e.g., spectroscopy, crystallography, thermodynamics) and quantum mechanical calculations on small molecules.

### 11.2.3  3. Coarse-Graining and Simplification to the HP Model Energy Function

All-atom force fields are too complex for proving NP-hardness or for very rapid exploration of conformational space for large proteins. Simplified models make several abstractions:
- **Reduced Representation**: Instead of all atoms, each amino acid is represented as a single "bead" or a few beads[31].
- **Lattice Constraint (for HP model)**: Beads are confined to sites on a regular lattice (e.g., 2D square or 3D cubic)[31]. This discretizes conformational space. Bond lengths are fixed by the lattice spacing, and bond angles are restricted (e.g., 90 or 180 degrees on a cubic lattice). Thus, $E_{bond}$ and $E_{angle}$ terms become implicit constraints rather than continuous energy terms. $E_{dihedral}$ is also highly restricted or simplified.
- **Dominant Interaction (Hydrophobic Effect)**: The primary driving force for folding globular proteins is assumed to be the hydrophobic effect—the tendency of nonpolar (H) residues to cluster together, minimizing contact with the polar (P) solvent[20]. All other interactions (specific hydrogen bonds, detailed electrostatics) are often ignored or subsumed.
- **Contact Potential**: The energy is defined purely in terms of contacts between non-covalently bonded residues that are adjacent on the lattice[31].
  - An H-H contact (two hydrophobic residues that are topological neighbors but not sequential neighbors) is assigned a favorable energy, $\epsilon_{HH}$ (typically -1). This term implicitly models the favorable enthalpy of van der Waals contacts between H residues and, more importantly, the entropic gain of the solvent when H residues are buried.
  - H-P and P-P contacts are typically assigned zero energy ($\epsilon_{HP} = 0, \epsilon_{PP} = 0$), assuming they are energetically neutral compared to H-solvent or P-solvent interactions, or that their contributions are less significant than H-H attraction.

This leads to the HP model energy function as defined in Section 5.2:

$$E(C) = \sum_{1 \le k < l-1 \le N-1} \chi(s_k, s_l) \cdot \delta(p_k, p_l)$$

where $\chi(H, H) = \epsilon_{HH} < 0$, and other $\chi$ values are zero. $\delta(p_k, p_l)$ is 1 if $s_k$ and $s_l$ are non-sequential lattice neighbors, and 0 otherwise.

---

This derivation highlights the significant approximations made to arrive at the HP model. While it loses atomic detail and specific interaction types, it retains the crucial element of self-avoidance and a dominant driving force (hydrophobicity), making it a valuable tool for studying fundamental aspects of protein folding and its computational complexity. The NP-hardness of this simplified model suggests that the complexity is deeply rooted in the combinatorial nature of packing a chain in 3D space to optimize even a simple set of interactions.

## 11.3 Appendix C. Sample Data and Code Snippets

### 11.3.1 1. Sample HP Sequence and 2D Lattice Folding

Consider a short HP sequence: $S = HHPPHH$

Let's try to fold this on a 2D square lattice to maximize H-H contacts (minimize energy, assuming $\epsilon_{HH} = -1$). Each H-H contact not adjacent in the sequence contributes -1 to the energy.

A possible conformation:

Let $s_1 = H, s_2 = H, s_3 = P, s_4 = P, s_5 = H, s_6 = H$.

Coordinates (x,y):

- $s_1(H) : (0,0)$
- $s_2(H) : (1,0)$
- $s_3(P) : (1,1)$
- $s_4(P) : (0,1)$
- $s_5(H) : (0,2)$
- $s_6(H) : (1,2)$

A compact, self-avoiding fold:

- $s_1(H) : (0,0)$
- $s_2(H) : (1,0)$
- $s_3(P) : (2,0)$
- $s_4(P) : (2,1)$
- $s_5(H) : (1,1)$
- $s_6(H) : (0,1)$

Let's identify non-sequential H-H neighbors:

- $s_1(H)$ and $s_6(H)$ are at (0,0) and (0,1) respectively. They are neighbors. $1 < 6 - 1$. This is one H-H contact.
- $s_2(H)$ and $s_5(H)$ are at (1,0) and (1,1) respectively. They are neighbors. $2 < 5 - 1$. This is another H-H contact.

So, there are 2 H-H contacts. Energy $E = 2 \times (-1) = -2$.

### 11.3.2 2. Python Snippet to Calculate Energy of a Given 2D HP Fold

```
def calculate_hp_energy_2d(sequence, coordinates):
    """
    Calculates the energy of a given 2D HP protein conformation.

    Args:
        sequence (str): A string of 'H' and 'P' characters.
        coordinates (list of tuples): A list of (x,y) coordinates for each residue.

    Returns:
        int: The energy of the conformation (number of H-H contacts * -1).
    """
    n = len(sequence)
    if n != len(coordinates):
        raise ValueError("Sequence length and coordinates length must match.")

    hh_contacts = 0
    # Check for self-avoidance (optional, assuming valid input for simplicity here)
    if len(set(coordinates)) != n:
```

```
        # print("Warning: Conformation is not self-avoiding. Energy might be misleading.")
        # For a strict calculation, one might return float('inf') or raise error
        pass

    for i in range(n):
        for j in range(i + 2, n): # j > i+1 to ensure non-sequential
            if sequence[i] == 'H' and sequence[j] == 'H':
                # Check if (i,j) are topological neighbors
                xi, yi = coordinates[i]
                xj, yj = coordinates[j]
                if abs(xi - xj) + abs(yi - yj) == 1: # Manhattan distance for square lattice
                    hh_contacts += 1

    return -hh_contacts

# Example usage from above:
# S = "HHPPHH"
# coords1 = [(0,0), (1,0), (2,0), (2,1), (1,1), (0,1)]
# energy1 = calculate_hp_energy_2d(S, coords1)
# print(f"Sequence: {S}, Coords: {coords1}, Energy: {energy1}")

# coords2 = [(0,0), (0,1), (0,2), (1,2), (1,1), (1,0)]
# energy2 = calculate_hp_energy_2d(S, coords2)
# print(f"Sequence: {S}, Coords: {coords2}, Energy: {energy2}")

# Expected output for both:
# Sequence: HHPPHH, Coords: [(0,0), (1,0), (2,0), (2,1), (1,1), (0,1)], Energy: -2
# Sequence: HHPPHH, Coords: [(0,0), (0,1), (0,2), (1,2), (1,1), (1,0)], Energy: -2
```

This Python snippet provides a basic function to calculate the energy of a pre-defined 2D conformation of an HP sequence on a square lattice. It iterates through all possible pairs of H-residues that are not adjacent in the sequence and checks if they are topological neighbors on the lattice. Each such H-H contact contributes -1 to the energy. Finding the coordinates that minimize this energy for a given sequence is the NP-hard problem.

These examples are illustrative and simplified. Real NP-hardness reduction gadgets are significantly more complex and meticulously designed to enforce logical constraints through energetic favorability.

# 12 Added References

1. Anfinsen's dogma - Wikipedia. Accessed June 7, 2025.
   `https://en.wikipedia.org/wiki/Anfinsen\%27s_dogma`
2. What is the protein folding problem? — AlphaFold - EMBL-EBI. Accessed June 7, 2025.
   `https://www.ebi.ac.uk/training/online/courses/alphafold/`
3. Thermodynamic Principle Revisited: Theory of Protein Folding. Accessed June 7, 2025.
   `https://www.scirp.org/journal/paperinformation?paperid=53638`
4. Protein folding - Wikipedia. Accessed June 7, 2025.
   `https://en.wikipedia.org/wiki/Protein_folding`
5. Foldit game engages the public in research — UW College of Engineering. Accessed June 7, 2025.
   `https://www.engr.washington.edu/news/trend/autumn-2010/fold-it-game`
6. Did AI Solve the Protein-Folding Problem? — Harvard Medicine Magazine. Accessed June 7, 2025.
   `https://magazine.hms.harvard.edu/articles/did-ai-solve-protein-folding-problem`
7. Advances in AI for Protein Structure Prediction: Implications for Cancer Drug Discovery and Development - MDPI. Accessed June 7, 2025.
   `https://www.mdpi.com/2218-273X/14/3/339`
8. The Protein Folding Problem - PMC - PubMed Central. Accessed June 7, 2025.
   `https://pmc.ncbi.nlm.nih.gov/articles/PMC2443096/`
9. Protein Folding Interdiction Strategy for Therapeutic Drug Development in Viral Diseases: Ebola VP40 and Influenza A M1. Accessed June 7, 2025.
   `https://pmc.ncbi.nlm.nih.gov/articles/PMC8998936/`
10. Computational Protein Design and Protein Structure Prediction - Nobel Prize. Accessed June 7, 2025.
    `https://www.nobelprize.org/uploads/2024/10/advanced-chemistryprize2024.pdf`
11. Research in Context: Designing proteins - National Institutes of Health (NIH). Accessed June 7, 2025.
    `https://www.nih.gov/news-events/nih-research-matters/research-context-designing-proteins`
12. Opportunities and Challenges for Machine Learning-Assisted Enzyme Engineering — ACS Central Science. Accessed June 7, 2025.
    `https://pubs.acs.org/doi/10.1021/acscentsci.3c01275`
13. Advancing single-shot vaccine design through AI and computational models - J-Stage. Accessed June 7, 2025.
    `https://www.jstage.jst.go.jp/article/trs/advpub/0/advpub_2025-002/_html/-char/en`
14. In Silico Epitope-Based Peptide Vaccine Design Against Influenza B Virus: An Immunoinformatics Approach - MDPI. Accessed June 7, 2025.
    `https://www.mdpi.com/2227-9717/13/3/681`
15. Current Challenges in Vaccinology - PMC - PubMed Central. Accessed June 7, 2025.
    `https://pmc.ncbi.nlm.nih.gov/articles/PMC7329983/`
16. Levinthal's paradox - PMC. Accessed June 7, 2025.
    `https://pmc.ncbi.nlm.nih.gov/articles/PMC48166/`
17. Levinthal's paradox - PubMed. Accessed June 7, 2025.
    `https://pubmed.ncbi.nlm.nih.gov/1729690/`
18. Fast Kinetics and Mechanisms in Protein Folding - Annual Reviews Collection - NCBI. Accessed June 7, 2025.
    `https://www.ncbi.nlm.nih.gov/books/NBK2232/`
19. Protein Folding Thermodynamics Kinetics - Fiveable. Accessed June 7, 2025.
    `https://library.fiveable.me/biophysical-chemistry/unit-5/protein-folding-thermodynamics-kinetics/study-guide/5ax1i7YJOxLwt5H9`
20. Folding funnel - Wikipedia. Accessed June 7, 2025.
    `https://en.wikipedia.org/wiki/Folding_funnel`
21. The Arch from the Stones: Understanding Protein Folding Energy Landscapes via Bio-inspired Collective Variables — bioRxiv. Accessed June 7, 2025.
    `https://www.biorxiv.org/content/10.1101/2025.05.28.656575v1.full-text`
22. Levels of Protein Organization - University of Vermont. Accessed June 7, 2025.
    `https://comis.med.uvm.edu/VIC/coursefiles/MD540/`

23. Blockchain Technology in Protein Folding: Enhancing Data Sharing and Collaboration - ResearchGate. Accessed June 7, 2025.
    `https://www.researchgate.net/publication/387183323`
24. Molecular Biology 02: 'Thermodynamics of protein folding' - CureFFI.org. Accessed June 7, 2025.
    `https://www.cureffi.org/2014/09/05/molecular-biology-02/`
25. Reflecting on DeepMind's AlphaFold artificial intelligence success - ICR. Accessed June 7, 2025.
    `https://www.icr.ac.uk/research-and-discoveries/cancer-blogs/`
26. Protein folding - seeing is deceiving - PMC. Accessed June 7, 2025.
    `https://pmc.ncbi.nlm.nih.gov/articles/PMC8284583/`
27. Gibbs Free Energy and the Relationship between $\Delta G$, $\Delta H$, & $\Delta S$ - Chad's Prep. Accessed June 7, 2025.
    `https://www.chadsprep.com/chads-general-chemistry-videos/gibbs-free-energy/`
28. A simple model for calculating the kinetics of protein folding - PNAS. Accessed June 7, 2025.
    `https://www.pnas.org/doi/10.1073/pnas.96.20.11311`
29. Lattice protein - Wikipedia. Accessed June 7, 2025.
    `https://en.wikipedia.org/wiki/Lattice_protein`
30. Hydrophobic-polar protein folding model - Wikipedia. Accessed June 7, 2025.
    `https://en.wikipedia.org/wiki/Hydrophobic-polar_protein_folding_model`
31. A weighted HP model for protein folding with diagonal contacts — Cambridge Core. Accessed June 7, 2025.
    `https://www.cambridge.org/core/journals/rairo-theoretical-informatics-and-applications/`
32. Modified Off-lattice AB Model for Protein Folding Problem Using the Vortex Search Algorithm. Accessed June 7, 2025.
    `https://www.ijmlc.org/vol5/529-C036.pdf`
33. Investigation of the three-dimensional lattice HP protein folding model using a genetic algorithm - SciELO. Accessed June 7, 2025.
    `https://www.scielo.br/j/gmb/a/nbhndMBBG5z4y39hyZ7463f/`
34. Structure-Approximating Inverse Protein Folding Problem in the 2D HP Model - Simon Fraser University. Accessed June 7, 2025.
    `https://www.sfu.ca/~lstacho/Ladislav_Stachos_site/Publications_files/protfold05.pdf`
35. NP (complexity) - Wikipedia. Accessed June 7, 2025.
    `https://en.wikipedia.org/wiki/NP_(complexity)`
36. What is the definition of P, NP, NP-complete and NP-hard? - Computer Science Stack Exchange. Accessed June 7, 2025.
    `https://cs.stackexchange.com/questions/9556/`
37. Introduction to the Theory of Computation, 3rd ed. Accessed June 7, 2025.
    `https://cs.brown.edu/courses/csci1810/fall-2023/resources/ch2_readings/`
38. Computers and Intractability: A Guide to the Theory of NP-Completeness - Amazon. Accessed June 7, 2025.
    `https://www.amazon.com/Computers-Intractability-NP-Completeness-Mathematical-Sciences/`
39. Boolean satisfiability problem - Wikipedia. Accessed June 7, 2025.
    `https://en.wikipedia.org/wiki/Boolean_satisfiability_problem`
40. NP-completeness - Wikipedia. Accessed June 7, 2025.
    `https://en.wikipedia.org/wiki/NP-completeness`
41. Travelling salesman problem - Wikipedia. Accessed June 7, 2025.
    `https://en.wikipedia.org/wiki/Travelling_salesman_problem`
42. Protein Folding in the Hydrophobic-Hydrophilic (HP) Model is NP-Complete. Accessed June 7, 2025.
    `https://computingbiology.github.io/docs/berger1998.pdf`
43. On the Complexity of Protein Folding - ResearchGate. Accessed June 7, 2025.
    `https://www.researchgate.net/publication/13513541`
44. Geometric Separators and Their Applications to Protein Folding in the HP-Model - SIAM. Accessed June 7, 2025.
    `https://epubs.siam.org/doi/10.1137/S0097539704440727`
45. Computational Complexity of Flattening Fixed-Angle Orthogonal - arXiv. Accessed June 7, 2025.
    `https://arxiv.org/pdf/2212.12450`

46. A Peptides Prediction Methodology for Tertiary Structure Based on Simulated Annealing - MDPI. Accessed June 7, 2025.
    `https://www.mdpi.com/2297-8747/26/2/39`

47. AI-Based Protein Structure Predictions and Their Implications in Drug Discovery - ResearchGate. Accessed June 7, 2025.
    `https://www.researchgate.net/publication/377548543`

48. Simple Model of Protein Energetics To Identify Ab Initio Folding Transitions - ACS Publications. Accessed June 7, 2025.
    `https://pubs.acs.org/doi/10.1021/acs.jctc.0c00524`

49. Conformational ensembles for protein structure prediction - PubMed. Accessed June 7, 2025.
    `https://pubmed.ncbi.nlm.nih.gov/40074747/`

50. De Novo Design of Proteins for Autocatalytic Isopeptide Bond Formation - ACS Publications. Accessed June 7, 2025.
    `https://pubs.acs.org/doi/10.1021/jacs.5c03319`

51. Rational design of enzyme activity and enantioselectivity - Frontiers. Accessed June 7, 2025.
    `https://www.frontiersin.org/journals/bioengineering-and-biotechnology/`

52. Rational Vaccine Design in Times of Emerging Diseases - MDPI. Accessed June 7, 2025.
    `https://www.mdpi.com/1999-4923/13/4/501`

53. Influence of protein fold stability on immunogenicity and its implications for vaccine design - PubMed. Accessed June 7, 2025.
    `https://pubmed.ncbi.nlm.nih.gov/28290225/`

54. Advances of computational methods enhance the development of multi-epitope vaccines - Oxford Academic. Accessed June 7, 2025.
    `https://academic.oup.com/bib/article/26/1/bbaf055/8015685`

55. Peptide-Based Vaccines: Current Progress and Future Challenges - Chemical Reviews. Accessed June 7, 2025.
    `https://pubs.acs.org/doi/10.1021/acs.chemrev.9b00472`

56. Structural and Computational Biology in the Design of Immunogenic Vaccine Antigens - PMC. Accessed June 7, 2025.
    `https://pmc.ncbi.nlm.nih.gov/articles/PMC4615220/`

57. Opportunities and challenges in design and optimization of protein function - PMC. Accessed June 7, 2025.
    `https://pmc.ncbi.nlm.nih.gov/articles/PMC7616297/`

58. The Design and Prospects of Influenza Virus Vaccines Based on Conserved Epitopes - MDPI. Accessed June 7, 2025.
    `https://www.mdpi.com/2813-3137/3/2/16`

59. Efficient and accurate prediction of protein structure using RoseTTAFold2 - ResearchGate. Accessed June 7, 2025.
    `https://www.researchgate.net/publication/371056284`

60. AlphaFold and what is next: bridging functional, systems and structural biology - Taylor & Francis. Accessed June 7, 2025.
    `https://www.tandfonline.com/doi/full/10.1080/14789450.2025.2456046`

61. How to cite AlphaFold - EMBL-EBI. Accessed June 7, 2025.
    `https://www.ebi.ac.uk/training/online/courses/alphafold/`

62. AlphaFold - Wikipedia. Accessed June 7, 2025.
    `https://en.wikipedia.org/wiki/AlphaFold`

63. Deep Learning Models for Protein Structure Prediction: AlphaFold2 and RoseTTAFold - ResearchGate. Accessed June 7, 2025.
    `https://www.researchgate.net/publication/386648348`

64. Accurate prediction of protein structures and interactions using a three-track neural network - Washington University. Accessed June 7, 2025.
    `https://dasher.wustl.edu/bio5357/readings/science-373-871-21.pdf`

65. AlphaFold and the future of structural biology - IUCr Journals. Accessed June 7, 2025.
    `https://journals.iucr.org/paper?me6231`

66. AlphaFold2: A high-level overview — AlphaFold - EMBL-EBI. Accessed June 7, 2025.
https://www.ebi.ac.uk/training/online/courses/alphafold/inputs-and-outputs/
67. Comparative studies of AlphaFold, RoseTTAFold and Modeller - Oxford Academic. Accessed June 7, 2025.
https://academic.oup.com/bib/article/23/5/bbac308/6658852
68. Accurate prediction of protein structures and interactions using a three-track neural network - IPD. Accessed June 7, 2025.
https://www.ipd.uw.edu/wp-content/uploads/2021/07/Baek_etal_Science2021_RoseTTAFold.pdf
69. Efficient and accurate prediction of protein structure using RoseTTAFold2 - bioRxiv. Accessed June 7, 2025.
https://www.biorxiv.org/content/10.1101/2023.05.24.542179.full
70. AlphaFold3 at CASP16 - bioRxiv. Accessed June 7, 2025.
https://www.biorxiv.org/content/10.1101/2025.04.10.648174v1.full.pdf
71. AlphaFold 3: Stepping into the future of structure prediction - Frontline Genomics. Accessed June 7, 2025.
https://frontlinegenomics.com/alphafold-3-stepping-into-the-future-of-structure-prediction/
72. AlphaFold3 - why did Nature publish it without its code? - PubMed. Accessed June 7, 2025.
https://pubmed.ncbi.nlm.nih.gov/38778239/
73. AlphaFold3: An Overview of Applications and Performance Insights - MDPI. Accessed June 7, 2025.
https://www.mdpi.com/1422-0067/26/8/3671
74. Assessment of Protein Complex Predictions in CASP16 - bioRxiv. Accessed June 7, 2025.
https://www.biorxiv.org/content/10.1101/2025.05.29.656875v1.full-text
75. What does AlphaFold3 learn about antigen and nanobody docking - PMC. Accessed June 7, 2025.
https://pmc.ncbi.nlm.nih.gov/articles/PMC11838198/
76. How to Use AlphaFold2 as a Wet Lab Biologist - Neurosnap. Accessed June 7, 2025.
https://neurosnap.ai/blog/post/how-to-use-alphafold2-as-a-wet-lab-biologist-pt-3/
77. LightRoseTTA: High-Efficient and Accurate Protein Structure Prediction - PMC. Accessed June 7, 2025.
https://pmc.ncbi.nlm.nih.gov/articles/PMC12097069/
78. AlphaFold3/README.md at main - GitHub. Accessed June 7, 2025.
https://github.com/Ligo-Biosciences/AlphaFold3/blob/main/README.md
79. Case study: AlphaFold uses open data and AI - EMBL. Accessed June 7, 2025.
https://www.embl.org/news/science/alphafold-using-open-data-and-ai/
80. Easy and accurate protein structure prediction using ColabFold - ResearchGate. Accessed June 7, 2025.
https://www.researchgate.net/publication/376154246
81. ColabFold: making protein folding accessible to all - PubMed. Accessed June 7, 2025.
https://pubmed.ncbi.nlm.nih.gov/35637307/
82. Quantum computer solves protein puzzle - Reddit. Accessed June 7, 2025.
https://www.reddit.com/r/tech/comments/10gcjqe/quantum_computer_solves_protein_puzzle/
83. Resource analysis of quantum algorithms for coarse-grained protein folding models - Physical Review Research. Accessed June 7, 2025.
https://link.aps.org/doi/10.1103/PhysRevResearch.6.033112
84. Using quantum annealing to design lattice proteins - Physical Review Research. Accessed June 7, 2025.
https://link.aps.org/doi/10.1103/PhysRevResearch.6.013162
85. Protein Design by Integrating Machine Learning and Quantum-Encoded Optimization - PRX Life. Accessed June 7, 2025.
https://link.aps.org/doi/10.1103/PRXLife.2.043012
86. New computing framework streamlines the use of AI and supercomputers - Argonne National Laboratory. Accessed June 7, 2025.
http://bit.ly/4ksP4Ib
87. Gaming for a cure: Computer gamers tackle protein folding - UW News. Accessed June 7, 2025.
https://bit.ly/3Ft30D3
88. Foldit players - DBLP. Accessed June 7, 2025.
https://dblp.org/pid/116/8907
89. Zoran Popovic - Google Scholar. Accessed June 7, 2025.
https://scholar.google.se/citations?user=0Qr2IGwAAAAJ

90. About the game "Foldit" - Reddit. Accessed June 7, 2025.
    `https://www.reddit.com/r/todayilearned/comments/13odlgu/`
91. AlphaFold Protein Structure Database: massively expanding structural coverage - PubMed. Accessed June 7, 2025.
    `https://pubmed.ncbi.nlm.nih.gov/34791371/`
92. AlphaFold Protein Structure Database in 2024 - PubMed. Accessed June 7, 2025.
    `https://pubmed.ncbi.nlm.nih.gov/37933859/`
93. Regularly updated benchmark sets for statistically correct evaluations of AlphaFold applications - PMC. Accessed June 7, 2025.
    `https://pmc.ncbi.nlm.nih.gov/articles/PMC11894802/`
94. A New Model for University-Industry Partnerships - Harvard Business Publishing. Accessed June 7, 2025.
    `https://hbsp.harvard.edu/inspiring-minds/a-new-model-for-university-industry-partnerships`
95. Collaborative Research Models: Academia & Industry Partnerships - MRL Consulting Group. Accessed June 7, 2025.
    `https://www.mrlcg.com/resources/blog/collaborative-research-models/`
96. Drug Discovery: New Models for Industry-Academic Partnerships - ResearchGate. Accessed June 7, 2025.
    `https://www.researchgate.net/publication/23459557`
97. Public–private partnerships in fostering outer space innovations - PNAS. Accessed June 7, 2025.
    `https://www.pnas.org/doi/10.1073/pnas.2222013120`
98. National Strategic Overview for Quantum Information Science. Accessed June 7, 2025.
    `https://www.quantum.gov/wp-content/uploads/2020/10/2018_NSTC_National_Strategic_Overview_QIS.pdf`
99. NSF's 10 Big Ideas - About NSF. Accessed June 7, 2025.
    `https://www.nsf.gov/about/history/big-ideas`
100. Smart Health and Biomedical Research in the Era of AI - NSF. Accessed June 7, 2025.
    `https://www.nsf.gov/funding/opportunities/sch-smart-health-biomedical-research/`
101. Finding Qualitative and Quantitative Collaborators in Academic Medical Centers - PMC. Accessed June 7, 2025.
    `https://pmc.ncbi.nlm.nih.gov/articles/PMC10440235/`
102. The five pillars of computational reproducibility - Oxford Academic. Accessed June 7, 2025.
    `https://academic.oup.com/bib/article/24/6/bbad375/7326135`
103. Principles that Govern the Folding of Protein Chains - Scientific Research Publishing. Accessed June 7, 2025.
    `https://www.scirp.org/reference/referencespapers?referenceid=438210`
104. Kinetics and thermodynamics of folding in model proteins - PMC. Accessed June 7, 2025.
    `https://pmc.ncbi.nlm.nih.gov/articles/PMC46930/`
105. Protein Folding in the 2D HP model - CiteSeerX. Accessed June 7, 2025.
    `https://citeseerx.ist.psu.edu/document?repid=rep1&type=pdf`
106. A New Algorithm for Protein Folding in the HP Model - DIMACS. Accessed June 7, 2025.
    `http://dimacs.rutgers.edu/~alantha/papers2/string_fold.pdf`
107. Protein Science - Georgia Tech. Accessed June 7, 2025.
    `https://faculty.cc.gatech.edu/~turk/bio_sim/articles/protein_folding_03.pdf`
108. Introduction to the Theory of Computation - UFOP. Accessed June 7, 2025.
    `http://www.decom.ufop.br/lucilia/ftc/Sipser.pdf`
109. P, NP, NP Hard, NP Complete - University of Iowa. Accessed June 7, 2025.
    `https://homepage.math.uiowa.edu/~idarcy/COURSES/4060/SPRING19/SLIDES/PNPcombinednew.pptx`
110. Computers and Intractability: A Guide to the Theory of NP-Completeness - Barnes & Noble. Accessed June 7, 2025.
    `https://www.barnesandnoble.com/w/computers-and-intractability-m-r-garey/1101455041`
111. DIMACS Protein Folding Research. Accessed June 7, 2025.
    `http://dimacs.rutgers.edu/~alantha/papers2/alantha-bill-bc.pdf`
112. Nature AlphaFold3 Article. Accessed January 1, 1970.
    `https://www.nature.com/articles/s41586-024-07487-w`

---

113. Folding Funnel Wikipedia (Duplicate Entry). Accessed January 1, 1970.
     https://en.wikipedia.org/wiki/Folding_funnel
114. Efficient and accurate prediction of protein structure using RoseTTAFold2 - bioRxiv (PDF). Accessed June 7, 2025.
     https://www.biorxiv.org/content/10.1101/2023.05.24.542179v1.full.pdf
115. Efficient and accurate prediction of protein structure using RoseTTAFold2 - bioRxiv (Version 1). Accessed June 7, 2025.
     https://www.biorxiv.org/content/10.1101/2023.05.24.542179v1
116. Efficient and accurate prediction of protein structure using RoseTTAFold2 - bioRxiv (Article Info). Accessed June 7, 2025.
     https://www.biorxiv.org/content/10.1101/2023.05.24.542179v1.article-info