# Structural barriers of the discrete Hasimoto map applied to protein backbone geometry

Yiquan Wang[*]

*College of Mathematics and System Science, Xinjiang University, Urumqi, Xinjiang, China*
*Xinjiang Key Laboratory of Biological Resources and Genetic Engineering,*
*College of Life Science and Technology, Xinjiang University, Urumqi, Xinjiang, China and*
*Shenzhen X-Institute, Shenzhen, China*

(Dated: February 16, 2026)

Determining the three-dimensional structure of a protein from its amino-acid sequence remains a fundamental problem in biophysics. The discrete Frenet geometry of the $C_\alpha$ backbone can be mapped, via a Hasimoto-type transform, onto a complex scalar field $\psi = \kappa e^{i\Sigma\tau}$ satisfying a discrete nonlinear Schrödinger equation (DNLS), whose soliton solutions reproduce observed secondary-structure motifs. Whether this mapping, which provides an elegant geometric description of folded states, can be extended to a predictive framework for protein folding remains an open question. We derive an exact closed-form decomposition of the DNLS effective potential $V_{\text{eff}} = V_{\text{re}} + iV_{\text{im}}$ in terms of curvature ratios and torsion angles, validating the result to machine precision across 856 non-redundant proteins. Our analysis identifies three structural barriers to forward prediction: (i) $V_{\text{im}}$ encodes chirality via the odd symmetry of $\sin\tau$, accounting for $\sim$31% of the total information and implying a $2^N$ degeneracy if neglected; (ii) $V_{\text{re}}$ is determined primarily ($\sim$95%) by local geometry, rendering it effectively sequence-agnostic; and (iii) self-consistent field iterations fail to recover native structures (mean RMSD = 13.1 Å) even with hydrogen-bond terms, yielding torsion correlations indistinguishable from zero. Constructively, we demonstrate that the residual of the DNLS dispersion relation serves as a geometric order parameter for $\alpha$-helices (ROC AUC = 0.72), defining them as regions of maximal integrability. These findings establish that the Hasimoto map functions as a kinematic identity rather than a dynamical governing equation, presenting fundamental obstacles to its use as a predictive framework for protein folding.

## I. INTRODUCTION

The relationship between amino-acid sequence and three-dimensional structure is a central problem in molecular biophysics. A protein's native fold is encoded in its sequence [1], yet the physical principles that govern the mapping from one-dimensional chemical information to three-dimensional geometry remain incompletely understood. Among the many theoretical frameworks proposed to address this problem, a geometric approach based on the differential geometry of space curves has attracted sustained interest. In this approach the protein backbone is treated as a discrete curve in $\mathbb{R}^3$, and its local shape is characterized by two scalar fields: the bond angle $\kappa[n]$ and the torsion angle $\tau[n]$ at each $C_\alpha$ vertex. These two fields constitute geometric order parameters that fully determine the backbone conformation up to rigid-body motion.

The idea of using curvature and torsion as dynamical variables for space curves originates in fluid mechanics. Hasimoto [2] showed that the local induction approximation for a vortex filament can be exactly transformed, via the complex scalar field $\psi = \kappa e^{i\int\tau ds}$, into the cubic nonlinear Schrödinger equation (NLS). Because the NLS is completely integrable, this transformation converts the geometric evolution of a three-dimensional curve into a one-dimensional soliton problem with exact analytical solutions. The success of the Hasimoto map in vortex dynamics naturally raises the question of whether an analogous construction can be applied to other physical filaments whose geometry is described by curvature and torsion.

Niemi and collaborators pursued this analogy systemat-

ically for protein backbones, developing a geometric program rooted in gauge field theory and nonlinear dynamics [3–6]. Working with the discrete Frenet frame of the $C_\alpha$ chain, they constructed a generalized discrete nonlinear Schrödinger equation (DNLS) whose dark-soliton solutions reproduce the characteristic $(\kappa, \tau)$ profiles of $\alpha$-helices and $\beta$-strands. Molkenthin, Hu, and Niemi [7] showed that a two-soliton configuration reproduces the villin headpiece HP35 with RMSD = 0.72 Å, and that each constituent soliton describes over 7 000 supersecondary structures in the Protein Data Bank (PDB). Krokhotin, Niemi, and Peng [8] constructed a library of 200 soliton parameter sets covering over 90% of PDB loop structures at sub-ångström accuracy. More recently, the framework has been extended to thermal dynamics and structural stability modeling: Begun *et al.* [9] simulated the folding and unfolding of the slipknotted protein AFV3-109 using multi-soliton ansätze, and Begun *et al.* [10] introduced Arnold's perestroikas to characterize topological phase transitions in myoglobin. Complementing this topological perspective, Liubimov *et al.* [11] applied the underlying lattice Abelian Higgs model to the same myoglobin structure, demonstrating that native conformations can be stabilized by introducing heterogeneous external fields to mimic environmental interactions. These studies collectively demonstrate that the DNLS and Abelian Higgs frameworks provide a compact and accurate geometric language for characterizing known protein conformations. However, a fundamental distinction must be drawn between descriptive capacity and predictive power. A critical unresolved issue is whether this formalism allows for the determination of the native structure strictly from the energy function, without reliance on *a priori* structural targets.

Concretely, if the DNLS effective potential $V_{\text{eff}}[n]$ could be determined from the amino-acid sequence alone, one could in

principle solve the DNLS forward to obtain the native $(\kappa, \tau)$ profile and reconstruct the three-dimensional structure. This leads to a fundamental theoretical question: does the Hasimoto map constitute a dynamical governing equation that drives folding, or is it merely a kinematic identity that describes the final state?

To address this, we must situate the Hasimoto framework within the broader landscape of theoretical biophysics, where predictive success has invariably relied on capturing non-local information that the nearest-neighbor structure of the DNLS cannot inherently represent. Energy landscape theory [12] and coarse-grained models like AWSEM [13] achieve accuracy by incorporating explicit non-local contact potentials that bias the free-energy surface. Similarly, from a geometric perspective, tube models [14, 15] rely on non-local excluded-volume interactions to select secondary structures, while direct coupling analysis (DCA) [16, 17] extracts long-range contacts from evolutionary covariance. Topological complexities such as knots [18, 19] further imply global constraints that exceed local curvature descriptions.

This distinction is sharpened by recent advances in deep learning. AlphaFold 2 [20] and AlphaFold 3 [21] solve the prediction problem by predicting a full rigid-body frame (rotation and translation) for every residue, effectively retaining all orientational degrees of freedom. Even single-sequence methods like ESMFold [22, 23], driven by protein language models, achieve competitive accuracy by mapping the sequence into a high-dimensional embedding space that captures non-local structural and dynamic context [24] implicitly. By integrating evolutionary covariance or contextual embeddings, these methods explicitly reconstruct the global information defining the native fold. In contrast, the Hasimoto transform projects the backbone geometry onto a single complex scalar field $\psi$, compressing the full coordinate frame into local curvature and torsion. The critical unresolved issue is whether this scalar projection retains sufficient information to determine the three-dimensional structure, or whether the reduction to a local effective potential fundamentally strips away the non-local and chiral information required for folding.

In this work, we investigate the applicability limits of the discrete Hasimoto map, positing that the widespread success of non-local and frame-based approaches suggests a structural barrier for purely local scalar theories. To demonstrate this, we derive an exact closed-form decomposition of the DNLS effective potential $V_{\mathrm{eff}} = V_{\mathrm{re}} + i V_{\mathrm{im}}$ into explicit functions of the curvature ratios $r^{\pm} = \kappa[n \pm 1]/\kappa[n]$ and torsion angles $\tau[n]$, and verify this decomposition to machine precision ($< 10^{-14}$) on 856 non-redundant proteins spanning all four SCOP structural classes. The decomposition reveals three independent structural barriers to forward prediction. First, the imaginary part $V_{\mathrm{im}}$ encodes the sign of the torsion angle through the odd symmetry of $\sin \tau$, carrying approximately 31% of the total potential information; discarding it introduces a $2^N$ chiral degeneracy that grows exponentially with chain length. Second, the real part $V_{\mathrm{re}}$ is 95% determined by local backbone geometry rather than by the amino-acid sequence, leaving the potential nearly devoid of the chemical information needed for prediction. Third, self-consistent field iterations

driven by hydrophobic and elastic potentials, with and without hydrogen-bond terms, fail uniformly on all 856 proteins (mean RMSD = 13.1 Å, torsion correlation indistinguishable from zero), and the two settings produce statistically identical outcomes.

These barriers are structural rather than algorithmic: they arise from the mathematical content of the Hasimoto map itself and are unlikely to be circumvented by parameter tuning or improved optimization alone. At the same time, the decomposition yields a constructive result. We show that the residual of the DNLS dispersion relation serves as a purely geometric helix detector with ROC AUC = 0.72 across all SCOP classes, identifying $\alpha$-helical segments as the backbone regions where the DNLS most closely approximates an integrable system. This suggests a geometric criterion for secondary-structure identification that does not require hydrogen-bond information.

The remainder of this paper is organized as follows. Section II introduces the discrete Frenet frame and the DNLS formulation. Section III presents the exact decomposition and its statistical validation. Section IV analyzes piecewise integrability and its connection to secondary structure. Section V reports the self-consistent field tests. Section VI discusses the physical origin of the barriers and the constructive applications of the framework.

## II. DISCRETE FRENET FRAME AND DNLS FORMULATION

### A. Discrete Frenet geometry of the $C_\alpha$ backbone

We represent the protein backbone as a polygonal chain of $C_\alpha$ positions $\{\mathbf{r}[n]\}_{n=1}^{N}$ in $\mathbb{R}^3$, where $N$ is the number of residues. The bond vectors are defined as

$$\mathbf{t}[n] = \mathbf{r}[n+1] - \mathbf{r}[n], \quad n = 1, \ldots, N-1. \tag{1}$$

For a $C_\alpha$ trace the bond lengths $|\mathbf{t}[n]|$ are approximately 3.8 Å [25, 26], corresponding to the virtual-bond distance between consecutive $\alpha$-carbons. We work with the unit tangent vectors $\hat{\mathbf{t}}[n] = \mathbf{t}[n]/|\mathbf{t}[n]|$.

The discrete Frenet frame at vertex $n$ is constructed from three successive $C_\alpha$ positions. The bond angle $\kappa[n]$ is defined through

$$\cos \kappa[n] = \hat{\mathbf{t}}[n] \cdot \hat{\mathbf{t}}[n-1], \quad n = 2, \ldots, N-1, \tag{2}$$

so that $\kappa[n] \in [0, \pi]$ measures the bending of the chain at vertex $n$. The torsion angle $\tau[n]$ requires four successive $C_\alpha$ positions and is defined as the dihedral angle

$$\begin{aligned}
\tau[n] = \mathrm{atan2}\big(&(\hat{\mathbf{t}}[n-1] \times \hat{\mathbf{t}}[n]) \times (\hat{\mathbf{t}}[n] \times \hat{\mathbf{t}}[n+1]) \cdot \hat{\mathbf{t}}[n], \\
&(\hat{\mathbf{t}}[n-1] \times \hat{\mathbf{t}}[n]) \cdot (\hat{\mathbf{t}}[n] \times \hat{\mathbf{t}}[n+1])\big),
\end{aligned} \tag{3}$$

with $\tau[n] \in (-\pi, \pi]$, defined for $n = 2, \ldots, N-2$. The sign of $\tau$ encodes the local handedness of the backbone: positive values correspond to right-handed twisting and negative values to left-handed twisting.

The pair $(\kappa[n], \tau[n])$ constitutes a complete set of internal coordinates for the discrete curve. Given an initial position $\mathbf{r}[1]$, an initial tangent $\hat{\mathbf{t}}[1]$, and an initial normal direction, the full three-dimensional backbone can be reconstructed from $\{\kappa[n], \tau[n]\}$ by sequential application of discrete Frenet rotations. This reconstruction is unique up to a global rigid-body transformation (three translations and three rotations), so that $\kappa$ and $\tau$ encode all shape information of the backbone.

The discrete normal and binormal vectors are defined as

$$\hat{\mathbf{n}}[n] = \frac{\hat{\mathbf{t}}[n] - \cos\kappa[n]\,\hat{\mathbf{t}}[n-1]}{\sin\kappa[n]}, \quad \hat{\mathbf{b}}[n] = \hat{\mathbf{t}}[n-1] \times \hat{\mathbf{n}}[n], \quad (4)$$

forming an orthonormal triad $\{\hat{\mathbf{t}}[n-1], \hat{\mathbf{n}}[n], \hat{\mathbf{b}}[n]\}$ at each interior vertex. The torsion angle $\tau[n]$ then governs the rotation of the normal plane from vertex $n$ to vertex $n+1$ about the tangent direction.

### B. Hasimoto transform and complex scalar field

Following the construction introduced by Hasimoto for continuous curves and adapted to the discrete setting by Niemi and collaborators, we define a complex scalar field on the backbone through

$$\psi[n] = \kappa[n]\exp\left(i\sum_{k=2}^{n}\tau[k]\right), \quad n = 2,\ldots,N-2. \quad (5)$$

The modulus $|\psi[n]| = \kappa[n]$ records the local bending, while the phase $\arg\psi[n] = \sum_{k=2}^{n}\tau[k]$ accumulates the torsion along the chain. This transform maps the two real geometric fields $(\kappa, \tau)$ into a single complex field $\psi$, reducing the description of backbone geometry to a one-component problem at the cost of entangling curvature and torsion in the amplitude and phase of $\psi$.

The inverse transform recovers the geometric variables as

$$\kappa[n] = |\psi[n]|, \quad \tau[n] = \arg\psi[n] - \arg\psi[n-1]. \quad (6)$$

Given $\psi[n]$ for all interior vertices, the backbone can therefore be reconstructed by extracting $(\kappa, \tau)$ via Eq. (6) and applying the discrete Frenet reconstruction.

### C. Discrete nonlinear Schrödinger equation

The DNLS equation for the backbone field $\psi[n]$ takes the form of a discrete eigenvalue problem

$$\beta^{+}[n]\big(\psi[n+1] - \psi[n]\big) - \beta^{-}[n]\big(\psi[n] - \psi[n-1]\big) = V_{\text{eff}}[n]\,\psi[n], \quad (7)$$

where $\beta^{+}[n]$ and $\beta^{-}[n]$ are site-dependent coupling parameters and $V_{\text{eff}}[n]$ is the effective potential. The left-hand side is a discrete second difference of $\psi$ with nonuniform coefficients, analogous to a lattice Laplacian.

The coupling parameters $\beta^{\pm}[n]$ encode the local mechanical properties of the peptide chain. In the simplest formulation

they depend on the virtual-bond lengths as

$$\beta^{+}[n] = \frac{1}{|\mathbf{t}[n]|}, \quad \beta^{-}[n] = \frac{1}{|\mathbf{t}[n-1]|}. \quad (8)$$

Because the $C_\alpha$–$C_\alpha$ virtual-bond length varies only weakly along the chain (standard deviation $\sim 0.1$ Å around the mean of 3.8 Å), the coupling parameters are nearly uniform. This near-uniformity will play a central role in the analysis of Sec. III, where we show that the sequence dependence of $\beta^{\pm}$ contributes less than 5% of the variance of $V_{\text{re}}$.

The effective potential is obtained by rearranging Eq. (7):

$$V_{\text{eff}}[n] = \frac{\beta^{+}[n]\big(\psi[n+1] - \psi[n]\big) - \beta^{-}[n]\big(\psi[n] - \psi[n-1]\big)}{\psi[n]}. \quad (9)$$

This expression is well defined wherever $\psi[n] \neq 0$, i.e., wherever the bond angle $\kappa[n] \neq 0$. For physical protein backbones $\kappa[n]$ is strictly positive at all interior vertices, so $V_{\text{eff}}$ is defined throughout the chain.

Because $\psi$ is complex, $V_{\text{eff}}$ is in general complex:

$$V_{\text{eff}}[n] = V_{\text{re}}[n] + iV_{\text{im}}[n]. \quad (10)$$

The real and imaginary parts carry distinct geometric content, as we will demonstrate through the exact decomposition in Sec. III.

Two features of Eq. (9) merit emphasis. First, the equation is an algebraic identity: for any discrete curve with $\kappa[n] > 0$, one can always compute $\psi$ via Eq. (5) and then define $V_{\text{eff}}$ via Eq. (9). No physical assumption enters this construction. The potential $V_{\text{eff}}$ is not postulated; it is read off from the known geometry. Second, the equation becomes a dynamical equation only when $V_{\text{eff}}$ is specified independently of the geometry, for example through a physical energy functional. In the vortex-filament context, the local induction approximation provides exactly such a specification, and the Hasimoto transform converts the geometric evolution into the integrable NLS. For proteins, the question is whether an analogous physical specification of $V_{\text{eff}}$ exists. The analysis of the following sections indicates that the answer is negative: the mathematical structure of $V_{\text{eff}}$ on protein backbones presents fundamental obstacles to its determination from sequence information alone.

### D. Connection to the continuum Hasimoto transform

In the continuum limit where the lattice spacing $\Delta s \to 0$ and the coupling becomes uniform ($\beta^{\pm} \to 1/\Delta s$), the discrete second difference in Eq. (7) reduces to $\partial^2\psi/\partial s^2$, and the DNLS recovers the structure of the cubic NLS

$$i\frac{\partial\psi}{\partial t} = \frac{\partial^2\psi}{\partial s^2} + \tfrac{1}{2}|\psi|^2\psi. \quad (11)$$

The soliton solutions of Eq. (11) describe localized bending excitations that propagate without dispersion along the filament. In the protein context, the discrete soliton solutions of Eq. (7) have been shown to reproduce the $(\kappa, \tau)$ profiles of

$\alpha$-helices (dark solitons with $\kappa \approx 1.5\,\mathrm{rad}$, $\tau \approx 1.0\,\mathrm{rad}$) and $\beta$-strands (bright solitons with larger $\kappa$ and $\tau \approx \pm\pi$). The key distinction is that for vortex filaments Eq. (11) is both kinematic (relating $\psi$ to geometry) and dynamic (governing time evolution), whereas for proteins only the kinematic content survives. Establishing this distinction quantitatively is the purpose of the remainder of this paper.

## III. EXACT DISCRETE DECOMPOSITION

The effective potential $V_{\mathrm{eff}}[n]$ defined by Eq. (7) encodes the full geometric content of the backbone in a single complex-valued sequence. In this section we derive a closed-form decomposition of $V_{\mathrm{eff}}$ into real and imaginary parts, each expressed entirely in terms of the local curvature ratios and torsion angles. The decomposition is an algebraic identity: it holds for any discrete space curve, independent of any physical model or approximation.

### A. Decomposition identity

*Proposition (Decomposition identity).* Let $\psi[n] = \kappa[n]\exp\left(i\sum_{k=1}^{n}\tau[k]\right)$ be the Hasimoto field constructed from the discrete Frenet curvature $\kappa[n] > 0$ and torsion $\tau[n]$, and let $V_{\mathrm{eff}}[n]$ be defined through the discrete Schrödinger equation

$$\beta_n^+\left(\psi[n+1]-\psi[n]\right)-\beta_n^-\left(\psi[n]-\psi[n-1]\right)=V_{\mathrm{eff}}[n]\,\psi[n],\tag{12}$$

where $\beta_n^\pm$ are the (possibly sequence-dependent) bond stiffness parameters. Then $V_{\mathrm{eff}}[n]=V_{\mathrm{re}}[n]+iV_{\mathrm{im}}[n]$ with

$$V_{\mathrm{re}}[n]=\beta_n^+\,r^+[n]\cos\tau[n+1]+\beta_n^-\,r^-[n]\cos\tau[n]-\left(\beta_n^++\beta_n^-\right),\tag{13}$$

$$V_{\mathrm{im}}[n]=\beta_n^+\,r^+[n]\sin\tau[n+1]-\beta_n^-\,r^-[n]\sin\tau[n],\tag{14}$$

where the curvature ratios are

$$r^+[n]\equiv\frac{\kappa[n+1]}{\kappa[n]},\qquad r^-[n]\equiv\frac{\kappa[n-1]}{\kappa[n]}.\tag{15}$$

*Proof.* We substitute the Hasimoto ansatz into Eq. (12) and divide both sides by $\psi[n]\neq 0$. The left-hand side becomes

$$\frac{V_{\mathrm{eff}}[n]\,\psi[n]}{\psi[n]}=\beta_n^+\left(\frac{\psi[n+1]}{\psi[n]}-1\right)-\beta_n^-\left(1-\frac{\psi[n-1]}{\psi[n]}\right).\tag{16}$$

The ratio of adjacent Hasimoto fields is

$$\frac{\psi[n+1]}{\psi[n]}=\frac{\kappa[n+1]}{\kappa[n]}\exp\left(i\tau[n+1]\right)=r^+[n]\,e^{i\tau[n+1]},\tag{17}$$

and similarly

$$\frac{\psi[n-1]}{\psi[n]}=\frac{\kappa[n-1]}{\kappa[n]}\exp\left(-i\tau[n]\right)=r^-[n]\,e^{-i\tau[n]}.\tag{18}$$

Collecting terms,

$$V_{\mathrm{eff}}[n]=\beta_n^+\left(r^+[n]\,e^{i\tau[n+1]}-1\right)-\beta_n^-\left(1-r^-[n]\,e^{-i\tau[n]}\right).\tag{19}$$

Separating real and imaginary parts via Euler's formula yields Eqs. (13)–(14) directly. ☐ The derivation is algebraically direct; the value of stating it explicitly lies in disentangling the roles of $r^\pm$ and $\tau$ that are obscured in the complex field $\psi$, and in enabling the systematic tests of Secs. III–V.

Several features of this decomposition merit emphasis. First, the result is exact: no continuum limit, Taylor expansion, or slowly varying envelope approximation has been invoked. Second, the decomposition holds for arbitrary bond parameters $\beta_n^\pm$ and for any discrete curve with $\kappa[n] > 0$; it is a kinematic identity rather than a dynamical equation. Third, the structure of Eqs. (13)–(14) makes transparent the distinct roles of the curvature ratios $r^\pm$ and the torsion angles $\tau$, which are entangled in the original complex field $\psi[n]$.

### B. Numerical verification

We verify the decomposition on eight representative proteins spanning the four SCOP [27, 28] structural classes (Table I). For each protein, we compute $V_{\mathrm{eff}}[n]$ in two independent ways: (i) directly from the definition Eq. (12) using the Hasimoto field $\psi[n]$, and (ii) from the analytic expressions Eqs. (13)–(14) using $\kappa[n]$ and $\tau[n]$. The maximum absolute difference over all residues is reported in the column labeled $\varepsilon$ in Table I.

In all eight cases the error $\varepsilon$ is of order $10^{-15}$ to $10^{-14}$, consistent with IEEE 754 double-precision arithmetic. This confirms that Eqs. (13)–(14) reproduce $V_{\mathrm{eff}}$ to machine precision and contain no hidden approximation. The remaining columns of Table I preview the three structural barriers analyzed in the following subsections and in Sec. IV; we defer their full statistical treatment to the 856-protein dataset.

### C. Barrier I: imaginary potential and torsion-sign ambiguity

The imaginary part $V_{\mathrm{im}}$ [Eq. (14)] depends on $\tau$ exclusively through the sine function. Because $\sin(\tau)$ is odd, $V_{\mathrm{im}}$ changes sign under $\tau \to -\tau$ while $V_{\mathrm{re}}$ [Eq. (13)], which depends on $\cos(\tau)$, remains invariant. This has a direct consequence for the inverse problem.

*Corollary 1 (Torsion-sign degeneracy).* Given only $V_{\mathrm{re}}[n]$ for $n=1,\ldots,L$, the torsion angles $\tau[n]$ are determined only up to independent sign flips $\tau[n] \to -\tau[n]$ at each site. The number of degenerate backbone configurations compatible with a given $V_{\mathrm{re}}$ profile is therefore $2^L$.

This degeneracy has significant physical implications. The sign of $\tau$ encodes the local chirality of the backbone; consequently, enantiomeric transformations at each residue leave $V_{\mathrm{re}}$ invariant while generating distinct three-dimensional topologies. Reconstruction of the backbone from $V_{\mathrm{re}}$ alone therefore confronts a $2^L$-fold ambiguity, representing an exponential loss of chiral information with chain length.

We quantify the information content of $V_{\mathrm{im}}$ across the full dataset of $N=856$ non-redundant proteins. Figure 1(a) shows the distribution of the ratio $\langle|V_{\mathrm{im}}|/|V_{\mathrm{re}}|\rangle$ per protein. The mean ratio is 0.31, indicating that the imaginary component carries

TABLE I. Numerical verification of the exact decomposition on eight representative proteins (two per SCOP class). $N$: number of C$\alpha$ atoms; $\varepsilon$: maximum absolute error between the two independent computations of $V_{\text{eff}}$; $\langle|V_{\text{im}}|/|V_{\text{re}}|\rangle$: mean ratio of imaginary to real potential magnitude; $\rho_{\text{geom}}$: Spearman correlation between $V_{\text{re}}$ computed with physical $\beta(s)$ and with uniform $\beta = 1$; $\delta_H$, $\delta_E$, $\delta_C$: dispersion-relation RMSE (in degrees) for helix, strand, and coil residues respectively; RMSD: backbone RMSD (Å) of the structure reconstructed from $V_{\text{re}}$ only (setting $V_{\text{im}} = 0$). A dash indicates that the corresponding secondary-structure type is absent.

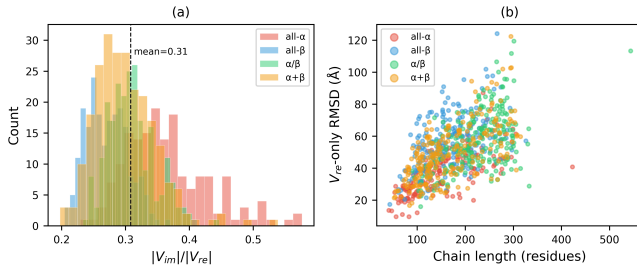| PDB | Class | $N$ | $\varepsilon$ | $\langle|V_{\text{im}}|/|V_{\text{re}}|\rangle$ | $\rho_{\text{geom}}$ | $\delta_H$ | $\delta_E$ | $\delta_C$ | RMSD |
|---|---|---|---|---|---|---|---|---|---|
| 1PA7 | all-$\alpha$ | 130 | $4.9 \times 10^{-15}$ | 0.324 | 0.966 | 20.4 | — | 39.4 | 51.2 |
| 2RH3 | all-$\alpha$ | 130 | $3.1 \times 10^{-15}$ | 0.332 | 0.963 | 27.1 | 36.6 | 35.6 | 32.5 |
| 5SV5 | all-$\beta$ | 134 | $1.4 \times 10^{-14}$ | 0.290 | 0.913 | — | 41.7 | 40.0 | 28.9 |
| 1AYO | all-$\beta$ | 130 | $1.8 \times 10^{-14}$ | 0.247 | 0.971 | 20.1 | 37.1 | 40.3 | 56.6 |
| 2B1L | $\alpha/\beta$ | 129 | $5.1 \times 10^{-15}$ | 0.300 | 0.945 | 17.8 | 38.2 | 42.5 | 55.9 |
| 1JBE | $\alpha/\beta$ | 132 | $1.8 \times 10^{-15}$ | 0.305 | 0.956 | 21.3 | 45.8 | 28.8 | 36.6 |
| 3CIP | $\alpha+\beta$ | 130 | $9.3 \times 10^{-15}$ | 0.273 | 0.949 | 22.1 | 31.1 | 36.3 | 55.4 |
| 2R4I | $\alpha+\beta$ | 130 | $3.6 \times 10^{-15}$ | 0.264 | 0.969 | 24.6 | 39.4 | 32.7 | 48.2 |



FIG. 1. Information cost of discarding the imaginary potential. (a) Distribution of the ratio $\langle|V_{\text{im}}|/|V_{\text{re}}|\rangle$ across 856 non-redundant proteins, colored by SCOP class. The mean ratio is 0.31, indicating that the imaginary component, which encodes the sign of the torsion angle through the odd symmetry of $\sin(\tau)$, carries roughly one-third of the total potential information. The distribution is largely class-independent, with all-$\alpha$ proteins showing a slightly broader tail toward higher values. (b) Backbone RMSD of structures reconstructed using $V_{\text{re}}$ only (setting $V_{\text{im}} = 0$) versus chain length. RMSD grows roughly linearly with chain length, reaching 40–120 Å for chains of 200–300 residues. This reflects the cumulative effect of torsion-sign errors: each residue contributes 1 bit of unresolved chiral ambiguity, and the resulting $2^N$ degeneracy makes $V_{\text{re}}$-only reconstruction physically meaningless beyond short peptides.

roughly one-third of the total potential magnitude. The distribution is largely independent of SCOP class, with all-$\alpha$ proteins showing a slightly broader tail toward higher values.

To demonstrate the practical consequence of discarding $V_{\text{im}}$, we reconstruct the backbone of each protein using $V_{\text{re}}$ only (setting $V_{\text{im}} = 0$ and choosing $\text{sign}(\tau)$ uniformly positive). Figure 1(b) plots the resulting C$\alpha$ RMSD against chain length. The RMSD grows approximately linearly, reaching 40–120 Å for chains of 200–300 residues. This linear growth reflects the cumulative nature of torsion-sign errors: each incorrectly assigned sign produces a local angular deviation that propagates along the chain. In information-theoretic terms, each residue contributes one bit of unresolved chiral information, for a total of $L$ bits per chain. The $V_{\text{re}}$-only reconstruction therefore yields structures that deviate substantially from the native fold for all but the shortest peptides.

## D. Barrier II: geometric dominance of $V_{\text{re}}$

The decomposition [Eq. (13)] shows that $V_{\text{re}}$ depends on the bond parameters $\beta_n^{\pm}$ only as multiplicative prefactors of the geometric terms $r^{\pm} \cos \tau$. A natural question is whether the sequence dependence encoded in $\beta(s)$ contributes significantly to the spatial profile of $V_{\text{re}}$, or whether the geometric factors dominate.

We test this by computing, for each of the 856 proteins, two versions of $V_{\text{re}}$: one with the physical, amino-acid-dependent bond parameters $\beta_n^{\pm}(s)$, and one with uniform $\beta^+ = \beta^- = 1$. The Spearman rank correlation $\rho_{\text{geom}}$ between the two profiles measures the fraction of the $V_{\text{re}}$ pattern that is determined by geometry alone.

*Corollary 2 (Geometric dominance).* Across 856 non-redundant proteins, the Spearman correlation between $V_{\text{re}}(\beta(s))$ and $V_{\text{re}}(\beta = 1)$ has mean $\bar{\rho}_{\text{geom}} = 0.951$ and minimum 0.88 (Fig. 2). The residual variance attributable to sequence-dependent $\beta(s)$ is less than $1 - \bar{\rho}^2 \approx 0.05$, i.e., below 5% on average.

This finding indicates a decoupling of sequence information from the effective potential. Although the bond parameters $\beta_n^{\pm}$ depend on amino-acid identity, their contribution to the variance of $V_{\text{re}}$ is negligible. The potential is dominated by the geometric terms $r^{\pm} \cos \tau$, the fluctuations of which exceed those induced by sequence-dependent stiffness by an order of magnitude (Fig. 4). All four SCOP classes overlap in the $\rho_{\text{geom}}$ distribution (Fig. 2), confirming that this pattern is universal across fold types. An important caveat is in order: the geometric terms $r^{\pm}$ and $\tau$ are themselves determined by the amino-acid sequence through the folding process. The Spearman test does not show that $V_{\text{re}}$ is independent of sequence in an absolute sense; rather, it establishes that the only *explicit* pathway from sequence to $V_{\text{re}}$ within the decomposition [Eqs. (13)–(14)], namely the bond stiffnesses $\beta_n^{\pm}$, is too narrow to carry significant information. The sequence dependence of $V_{\text{re}}$ is almost entirely *implicit*, mediated by the three-dimensional structure $(\kappa, \tau)$ that the folding process produces.

To test this interpretation directly, we compared $V_{\text{re}}$ profiles between protein pairs that share the same SCOP superfamily but have low sequence identity. From the 856 non-redundant
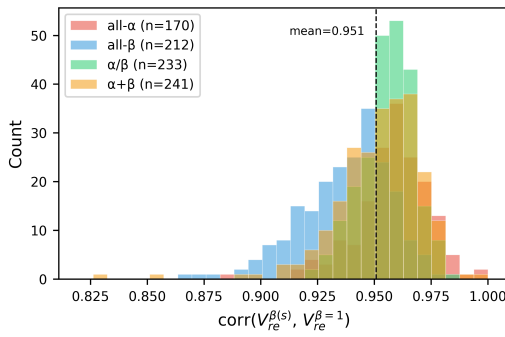
FIG. 2. Geometric dominance of the real effective potential. Distribution of Spearman rank correlation between $V_{re}$ computed with the physical, sequence-dependent bond parameters $\beta(s)$ and $V_{re}$ computed with uniform $\beta = 1$, evaluated over 856 non-redundant proteins. Colors denote SCOP structural classes: all-$\alpha$ (170), all-$\beta$ (212), $\alpha/\beta$ (233), and $\alpha+\beta$ (241). The distribution is sharply peaked near unity (mean $= 0.951$), indicating that the explicit sequence dependence carried by $\beta(s)$ accounts for less than 5% of the variance of $V_{re}$ on average. The dominant contribution comes from the geometric terms $r^{\pm}\cos\tau$, which depend on the backbone structure $(\kappa, \tau)$ rather than directly on amino-acid identity. All four SCOP classes overlap, confirming that this pattern is universal across protein folds.
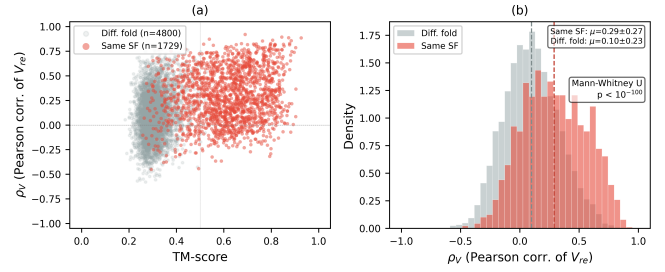


FIG. 3. $V_{re}$ tracks fold rather than sequence. (a) Pearson correlation $\rho_V$ of structurally aligned $V_{re}$ profiles versus TM-score for 1729 same-superfamily pairs (red) and 4800 different-fold pairs (gray). Among the same-superfamily pairs, 79% have TM-score $> 0.5$ and cluster in the upper-right quadrant; the remaining 21% are distant homologs with greater structural divergence, yet $\rho_V$ still correlates positively with TM-score within this subgroup. (b) Distribution of $\rho_V$ for the two groups. Same-superfamily: $\mu = 0.29 \pm 0.27$; different-fold: $\mu = 0.10 \pm 0.23$ (Mann-Whitney $U$, $p < 10^{-134}$). The mean sequence identity within the same-superfamily group is 13.9%, confirming that the elevated correlation is driven by structural similarity, not sequence similarity.

proteins (culled by the PISCES [29] server with resolution $\leq 2.0\,\text{Å}$, $R$-factor $\leq 0.2$, chain length 40–300 residues, sequence identity $\leq 25\%$, X-ray entries only, excluding chains with breaks or disorder), we identified 1729 pairs belonging to the same SCOP superfamily. As a control, we constructed 4800 pairs drawn from different SCOP folds with chain-length differences $\leq 20$ residues. For each pair, TM-align was used to obtain a residue-level structural alignment; the Pearson correlation $\rho_V$ was then computed between the $V_{re}$ values at structurally aligned positions.

Figure 3 summarizes the results. Same-superfamily pairs exhibit a mean $\rho_V = 0.290 \pm 0.266$, significantly higher than the different-fold background of $\rho_V = 0.099 \pm 0.230$ (Mann-Whitney $U$ test, $p < 10^{-134}$). The scatter plot [Fig. 3(a)] shows that $\rho_V$ increases with TM-score: among the 1729 same-superfamily pairs, 79% have TM-score $> 0.5$ and cluster in the upper-right quadrant, while different-fold pairs remain near zero. The remaining 21% of same-superfamily pairs fall below TM-score $= 0.5$; these are distant homologs for which the superfamily-level classification permits substantial structural divergence. Even within this subgroup, $\rho_V$ correlates positively with TM-score, indicating that the relationship between structural similarity and $V_{re}$ similarity is continuous rather than threshold-dependent. The mean sequence identity within the same-superfamily group is only 13.9% (range 1.4–27.9%), confirming that the elevated $\rho_V$ is driven by structural similarity rather than sequence similarity. This provides direct evidence that $V_{re}$ tracks the three-dimensional fold: proteins with unrelated sequences but similar structures produce similar $V_{re}$ profiles, whereas proteins with comparable chain lengths but different folds do not.

The implication for the forward prediction problem is as follows. The decomposition [Eq. (13)] shows that $V_{re}$ is deter-

mined by two types of inputs: the geometric terms $r^{\pm}\cos\tau$, which depend on the three-dimensional structure, and the bond stiffnesses $\beta_n^{\pm}$, which depend directly on the amino-acid sequence. The Spearman test establishes that the latter contribute less than 5% of the variance, and the superfamily analysis (Fig. 3) confirms that $V_{re}$ covaries with fold rather than with sequence. This creates a circularity for forward prediction: $V_{re}$ is overwhelmingly determined by the backbone geometry $(\kappa, \tau)$, which is the very quantity one seeks to predict, and the only direct pathway from sequence to $V_{re}$ through $\beta_n^{\pm}$ is too weak to carry the structural information needed to determine the native fold.

The circularity exposed by Barrier II can be contrasted with physical models that have achieved predictive success precisely by incorporating non-local, sequence-dependent information through channels absent from the Hasimoto decomposition. Coarse-grained force fields such as AWSEM [13] supplement local backbone terms with explicit contact potentials that couple residue pairs separated by tens to hundreds of positions along the sequence, and with associative-memory terms that bias the energy landscape toward known structural motifs. The tube model of Banavar and Maritan [14, 15] derives secondary-structure selection from a three-body excluded-volume interaction that is non-local by construction. Statistical approaches based on direct coupling analysis [16, 17] extract residue–residue contact maps from evolutionary covariance using the maximum-entropy principle, providing non-local structural constraints derived entirely from sequence data. In each case, the predictive power arises from terms that couple residues at positions $|m - n| \gg 1$, which constitutes precisely the type of information that the DNLS effective potential, constructed from nearest-neighbor ratios $r^{\pm}[n]$ and local torsion angles $\tau[n]$, cannot encode. The geometric dominance of $V_{re}$ is therefore not merely a quantitative observation but reflects a structural mismatch between the local, kinematic

content of the Hasimoto map and the non-local, thermodynamic content required for folding.

### E. Summary of structural barriers

The exact decomposition has revealed two independent, static barriers to forward structure prediction: the $2^L$-fold torsion-sign degeneracy encoded in $V_{\text{im}}$ (Barrier I) and the geometric dominance of $V_{\text{re}}$ that leaves less than 5% of its variance attributable to amino-acid identity (Barrier II). Both barriers follow from the algebraic structure of the Hasimoto transform. Whether the DNLS can nevertheless function as a dynamical equation that drives folding through a physically motivated effective potential is tested in Sec. V.

## IV. PIECEWISE INTEGRABILITY AND SECONDARY STRUCTURE

The exact decomposition of Sec. III holds for arbitrary discrete curves. In this section we examine the special case of backbone segments where the curvature varies slowly, so that $r^{\pm}[n] \approx 1$. Under this condition the full decomposition simplifies to a scalar dispersion relation that connects $V_{\text{re}}$ directly to the torsion angle. We show that this dispersion relation is satisfied almost exclusively within $\alpha$-helical segments, providing a purely geometric criterion for secondary-structure identification.

### A. Uniform-segment dispersion relation

When the curvature ratios satisfy $r^{+}[n] \approx r^{-}[n] \approx 1$ and the bond parameters are approximately uniform ($\beta_n^{+} \approx \beta_n^{-} \approx \beta$), the real part of the effective potential [Eq. (13)] reduces to

$$V_{\text{re}}[n] \approx \beta \left[ \cos \tau[n+1] + \cos \tau[n] \right] - 2\beta . \tag{20}$$

If the torsion angle is also locally constant ($\tau[n+1] \approx \tau[n] \approx \tau$), this further simplifies to the dispersion relation

$$\cos \tau = 1 + \frac{V_{\text{re}}}{2\beta} . \tag{21}$$

Eq. (21) is the discrete analogue of the continuum NLS dispersion relation $\omega = k^2$, expressed in terms of the backbone torsion angle. It provides a one-to-one mapping between $V_{\text{re}}$ and $|\tau|$ within any segment where $\kappa$ and $\tau$ are approximately uniform. Note that the cosine function renders Eq. (21) insensitive to the sign of $\tau$, consistent with the torsion-sign degeneracy identified in Sec. III C.

The conditions under which Eq. (21) holds are precisely the conditions that define an integrable segment of the DNLS: uniform curvature ($r^{\pm} \approx 1$), uniform torsion, and uniform coupling. Deviations from these conditions break integrability and cause the full decomposition [Eqs. (13)–(14)] to differ from the dispersion relation. The magnitude of this deviation therefore serves as a local measure of integrability.

### B. Integrability error as a structural probe

We define the integrability error at each residue as

$$E[n] = \left| \cos \tau[n] - \left( 1 + \frac{V_{\text{re}}[n]}{2\beta} \right) \right| , \tag{22}$$

where $\beta = \langle \beta_n^{\pm} \rangle$ is the chain-averaged coupling parameter. By construction, $E[n] = 0$ when the backbone at residue $n$ satisfies the uniform-segment dispersion relation exactly, and $E[n] > 0$ when the local curvature or torsion varies too rapidly for the integrable approximation to hold.

Figure 4 illustrates the behavior of $V_{\text{eff}}[n]$ along the backbone for eight representative proteins (two per SCOP class), analyzing the effective potential as a site-dependent profile analogous to spectral signal representations of protein sequences [30]. Within helical segments (pink background shading), $V_{\text{re}}$ forms near-constant negative plateaus (typically $-0.5$ to $-1.5 \, \text{Å}^{-2}$), consistent with Eq. (21) evaluated at the canonical helix torsion $\tau_{\text{helix}} \approx 1.0 \, \text{rad}$. The amplitude of $V_{\text{im}}$ is reduced relative to loop and strand regions but remains finite. In $\beta$-strand and coil regions, both $V_{\text{re}}$ and $V_{\text{im}}$ fluctuate strongly (amplitudes reaching 3–6 $\text{Å}^{-2}$), reflecting rapid residue-to-residue variation of $r^{\pm}$ and $\tau$. Sharp negative spikes in $V_{\text{re}}$ mark transitions between secondary-structure elements. These patterns are local rather than class-dependent: helical segments in the nominally all-$\beta$ protein 1AYO display the same plateau behavior as those in the all-$\alpha$ proteins.

### C. Statistical validation on 856 proteins

We evaluate the dispersion-relation RMSE separately for helix (H), strand (E), and coil (C) residues across the full dataset of 856 non-redundant proteins. For each protein, residues are grouped by their DSSP secondary-structure assignment [31], and the RMSE of Eq. (21) is computed per group:

$$\delta_X = \sqrt{\frac{1}{N_X} \sum_{n \in X} \left[ \cos \tau[n] - \left( 1 + \frac{V_{\text{re}}[n]}{2\beta} \right) \right]^2} , \tag{23}$$

where $X \in \{H, E, C\}$ and $N_X$ is the number of residues in class $X$.

Figure 5 shows the distribution of $\delta_X$ as box plots, faceted by SCOP class. Across all four classes, helical segments exhibit systematically lower RMSE (median $\sim 21°$) than strand ($\sim 37°$) or coil ($\sim 36°$) segments. This separation is class-independent: even in all-$\beta$ proteins where helices are scarce, the few helical residues still satisfy the dispersion relation with comparable accuracy. The result reflects the local uniformity of $(\kappa, \tau)$ within helices ($r^{\pm} \approx 1$), which is the condition under which the exact decomposition reduces to Eq. (21).

The data in Table I illustrate this pattern at the level of individual proteins. The pure all-$\alpha$ protein 1PA7 has $\delta_H = 20.4°$ and no strand residues, while the pure all-$\beta$ protein 5SV5 has $\delta_E = 41.7°$ and no helix residues. In mixed-class proteins, the helix and strand RMSE values coexist within the same chain:
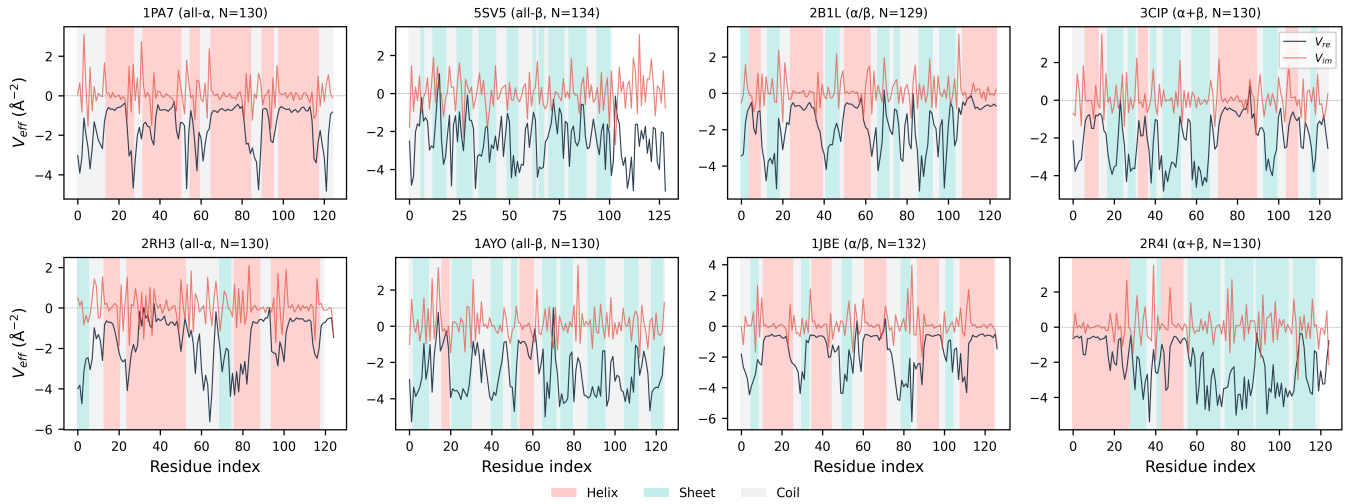
FIG. 4. Effective potential $V_{\text{eff}}[n]$ along the $C_\alpha$ backbone for eight representative proteins (two per SCOP class; columns from left to right: all-$\alpha$, all-$\beta$, $\alpha/\beta$, $\alpha+\beta$). Black: $V_{\text{re}}$; red: $V_{\text{im}}$. Background shading marks DSSP secondary-structure assignment (pink: helix; cyan: strand; gray: coil). Within helical segments $V_{\text{re}}$ forms near-constant negative plateaus consistent with the dispersion relation Eq. (21), while strand and coil regions exhibit large-amplitude fluctuations in both components. Sharp negative spikes in $V_{\text{re}}$ mark transitions between secondary-structure elements. These patterns are local rather than class-dependent: helical segments in the all-$\beta$ protein 1AYO display the same plateau behavior as those in the all-$\alpha$ proteins.
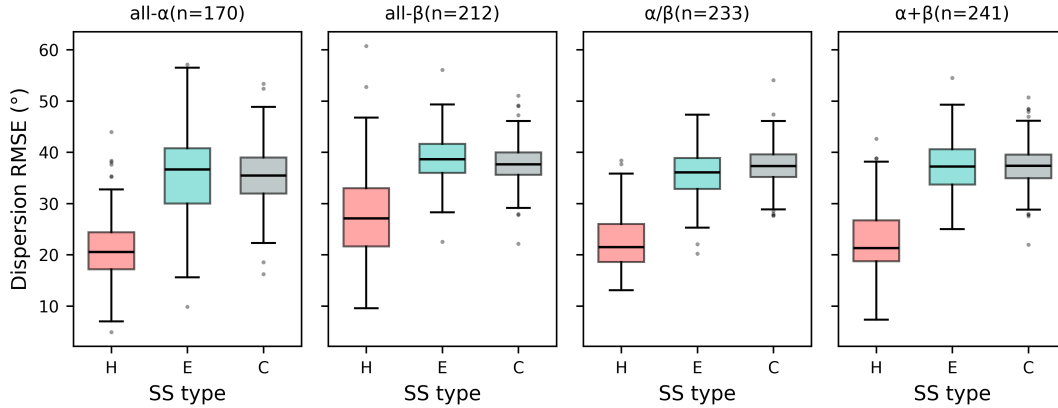


FIG. 5. Dispersion-relation RMSE by secondary-structure type, faceted by SCOP class (856 non-redundant proteins). For each protein, residues are grouped by DSSP assignment into helix (H), strand (E), and coil (C), and the RMSE of the uniform-segment approximation $\cos \tau = 1 + V_{\text{re}}/2\beta$ is computed per group. Helical segments exhibit systematically lower RMSE (median $\sim 21°$) than strand ($\sim 37°$) or coil ($\sim 36°$) segments across all four SCOP classes. This separation is class-independent: even in all-$\beta$ proteins where helices are scarce, the few helical residues satisfy the dispersion relation with comparable accuracy.

1AYO (all-$\beta$ by SCOP classification) contains a small number of helical residues with $\delta_H = 20.1°$, comparable to the values observed in all-$\alpha$ proteins. This confirms that the applicability of the dispersion relation is determined by the local curvature uniformity at each residue, not by the global fold classification.

### D. Helix detection by integrability error

The systematic separation of $\delta_H$ from $\delta_E$ and $\delta_C$ suggests that the integrability error $E[n]$ [Eq. (22)] can serve as a binary classifier for helical residues. To quantify this, we construct a receiver operating characteristic (ROC) curve by varying a threshold $E_{\text{th}}$ and classifying residue $n$ as helical if $E[n] < E_{\text{th}}$ and as non-helical otherwise. The ground truth is provided by DSSP assignment.

Figure 6 shows the ROC curves for each SCOP class and for the full dataset (856 proteins, 143 202 residues). The global area under the curve (AUC) is 0.720, significantly above the random baseline of 0.5. The class dependence is weak: AUC ranges from 0.667 (all-$\beta$) to 0.739 (all-$\alpha$), with $\alpha/\beta$ and $\alpha+\beta$ at 0.713 and 0.714 respectively. Even in all-$\beta$ proteins, where helical residues constitute only 15% of the total (4 741 out of
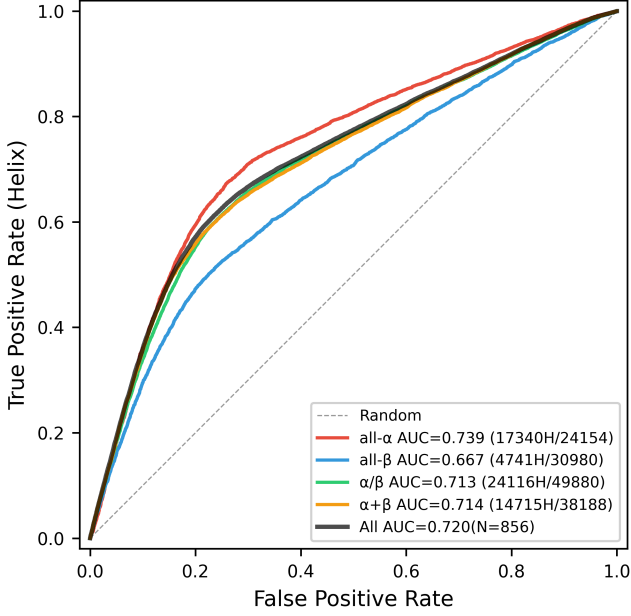
FIG. 6. ROC curve for helix detection using the integrability error $E[n] = |\cos\tau[n] - (1 + V_{\mathrm{re}}[n]/2\beta)|$ as a binary classifier (helix vs. non-helix by DSSP). Curves are shown for each SCOP class and for the full dataset ($N = 856$ proteins, $143\,202$ residues). Parentheses indicate the number of helix residues over total residues in each class. The global AUC $= 0.720$ quantifies the overlap between DNLS integrability and helical geometry: 72% of the helix/non-helix distinction (by DSSP) is captured by the scalar symmetry condition $E[n] < E_{\mathrm{th}}$, which tests whether the backbone locally preserves discrete helical symmetry. The class dependence is weak (AUC range 0.667–0.739), confirming that the integrability–helicity correspondence operates at the residue level rather than depending on global fold type. The gap between AUC $= 0.72$ and unity quantifies the structural information carried by hydrogen bonds, nonlocal contacts, and torsion-sign degrees of freedom that lie outside the scalar dispersion relation.

30 980), the AUC remains well above chance. This confirms that the integrability error operates at the residue level, detecting local geometric regularity rather than relying on the overall fold type.

The AUC of 0.72 quantifies the overlap between DNLS integrability and helical geometry as defined by DSSP. The condition $E[n] < E_{\mathrm{th}}$ tests whether the backbone at residue $n$ preserves discrete helical symmetry, i.e., whether the local curvature and torsion are sufficiently uniform that a single screw motion $\mathbf{r}[n+1] = M\mathbf{r}[n]$ approximately generates the chain segment. When $\kappa[n]$ and $\tau[n]$ are exactly constant, the curvature ratios satisfy $r^{\pm} = 1$, the exact decomposition [Eq. (13)] reduces to the dispersion relation [Eq. (21)], and $E[n] \equiv 0$; the geometric origin of this vanishing is analyzed in Sec. IV F, where we show that $E[n]$ is an order parameter for discrete helical symmetry breaking. An AUC of 0.72 indicates a substantial but incomplete correspondence between the scalar symmetry condition and the helix/non-helix distinction as defined by DSSP. The gap from unity likely reflects contributions from hydrogen bonds, nonlocal contacts, and

torsion-sign information that lie outside the dispersion relation. In this sense, $\alpha$-helices are the backbone regions where the DNLS most closely approximates an integrable system, and $E[n]$ measures the degree to which this integrability is locally broken.

### E. The protein backbone as a piecewise integrable system

The results of this section lead to a unified geometric picture of protein backbone structure in the DNLS framework. The backbone can be viewed as a piecewise integrable system: $\alpha$-helical segments are regions of approximate integrability where the curvature and torsion are locally uniform, the dispersion relation [Eq. (21)] holds, and the DNLS admits soliton-like solutions. $\beta$-strand and coil regions, by contrast, are regions of broken integrability where $r^{\pm}$ deviates significantly from unity, the dispersion relation fails, and the full decomposition [Eqs. (13)–(14)] is required.

This picture provides a geometric definition of secondary structure that is independent of hydrogen-bond criteria. An $\alpha$-helix is a contiguous segment where the DNLS integrability error $E[n]$ remains below a threshold; a non-helical region is one where $E[n]$ exceeds this threshold. The definition is purely kinematic: it depends only on the $C_{\alpha}$ coordinates and requires no energy function or force field.

Helical segments exhibit the smallest dispersion-relation error among all secondary-structure types (Fig. 5), but they are not the sole regions of approximate integrability. Short segments within $\beta$-strands or loops may also satisfy $r^{\pm} \approx 1$ locally, contributing to the imperfect separation in the ROC analysis (Fig. 6). The gap between AUC $= 0.72$ and unity reflects both these local exceptions and the information carried by $V_{\mathrm{im}}$ (torsion-angle signs) and nonlocal interactions (hydrogen bonds, tertiary contacts) that lie outside the scalar dispersion relation.

Two limitations of the dispersion relation as a structural probe should be noted. First, Eq. (21) determines only $|\tau|$, not $\mathrm{sign}(\tau)$. The torsion-sign ambiguity identified in Sec. III C therefore persists even within integrable segments. Second, the dispersion relation provides no information about the spatial arrangement of secondary-structure elements relative to one another. It is a local diagnostic that characterizes individual residues but cannot address the global fold topology. These limitations reinforce the conclusion that the Hasimoto map is a geometric analysis tool rather than a predictive framework.

### F. Geometric interpretation of the integrability error

The integrability error $E[n]$ defined in Eq. (22) admits a direct geometric interpretation as a measure of broken discrete helical symmetry. We make this connection explicit.

A discrete curve $\{\mathbf{r}[n]\}$ possesses *discrete helical symmetry* if there exists a rigid screw motion $M$ (a rotation about a fixed axis composed with a translation along that axis) such that $\mathbf{r}[n+1] = M\mathbf{r}[n]$ for all $n$. This condition requires that the

discrete Frenet curvature and torsion be site-independent:

$$\kappa[n] = \kappa_0, \quad \tau[n] = \tau_0, \quad \forall\, n. \tag{24}$$

Substituting Eq. (24) into the exact decomposition [Eq. (13)] with uniform coupling $\beta_n^+ = \beta_n^- = \beta$ gives $r^{\pm}[n] = 1$ and $\tau[n+1] = \tau[n] = \tau_0$, so that

$$V_{\text{re}} = \beta \cos \tau_0 + \beta \cos \tau_0 - 2\beta = 2\beta (\cos \tau_0 - 1), \tag{25}$$

which is precisely the dispersion relation [Eq. (21)]. The integrability error therefore vanishes identically:

$$E[n] = \left| \cos \tau_0 - \left( 1 + \frac{2\beta (\cos \tau_0 - 1)}{2\beta} \right) \right| = 0. \tag{26}$$

In the continuum limit, the classical theorem of Lancret [32, 33] states that a space curve is a generalized helix if and only if the ratio $\kappa/\tau$ is constant along the curve. For discrete curves the analogous condition is stronger: both $\kappa$ and $\tau$ must be individually constant, because the discrete Frenet frame does not admit a continuous reparameterization that could absorb correlated variations of $\kappa$ and $\tau$ into a constant ratio. The condition $E[n] = 0$ is therefore the discrete counterpart of the Lancret condition, specialized to the DNLS framework.

This identification gives $E[n]$ a geometric interpretation beyond that of a fitting residual: it serves as a *local measure of discrete helical symmetry breaking*. At each residue, $E[n]$ quantifies the degree to which the local backbone geometry deviates from a perfect discrete helix, i.e., from a fixed point of the screw-motion group. In this picture, helical segments correspond to near-symmetric regions ($E[n] \approx 0$), while loops and strands correspond to regions of broken symmetry ($E[n] \gg 0$).

## V. SELF-CONSISTENT FIELD TEST

The two barriers identified in Sec. III are static: they concern the information content of $V_{\text{eff}}$ extracted from known structures and follow from the algebraic structure of the Hasimoto decomposition [Eqs. (13)–(14)], independently of any energy function or dynamical scheme. This point is demonstrated directly by what we term an *oracle test*: in Fig. 1(b) we supply the exact $V_{\text{re}}[n]$ computed from the native structure, which constitutes the best possible real potential that any energy function could produce, and attempt backbone reconstruction with $V_{\text{im}} = 0$. The resulting RMSD of 20–120 Å (mean $\sim$50 Å) shows that even perfect knowledge of $V_{\text{re}}$ is insufficient, because the $2^N$ torsion-sign degeneracy is intrinsic to the decomposition and cannot be resolved by any choice of energy function.

A separate question is whether the DNLS can be used as a dynamical equation to drive an initially unfolded chain toward the native state, given a physically motivated effective potential. This possibility has been suggested in the context of modulation instability of the DNLS, where small perturbations of a uniform solution grow exponentially and may seed the formation of soliton-like secondary-structure elements. In

this section we test this dynamical scenario directly by performing self-consistent field (SCF) iterations on all 856 proteins in our dataset, using two independent settings: one with hydrophobic and elastic potentials only, and one with an additional hydrogen-bond term. The SCF test provides a third, dynamical barrier that is independent of Barriers I and II and serves as an additional confirmation of the representational limitation.

### A. SCF formulation

We initialize the complex field as a nearly uniform state $\psi_0[n] = \kappa_0 + \delta\psi[n]$, where $\kappa_0 = 0.1$ rad represents a nearly straight chain and $\delta\psi[n]$ is a small random perturbation with $|\delta\psi| < 0.01$. The field is evolved under damped dynamics

$$\frac{d\psi[n]}{dt} = -(1 + i\gamma) \frac{\delta H[\psi]}{\delta \psi^*[n]}, \tag{27}$$

where $\gamma = 0.5$ is a dissipation parameter that drives the system toward a local energy minimum rather than conserving the Hamiltonian. The functional $H[\psi]$ consists of a kinetic (elastic) term and an interaction term:

$$H[\psi] = \sum_n \beta\, |\psi[n+1] - \psi[n]|^2 + \sum_n V_{\text{int}}[n]\, |\psi[n]|^2. \tag{28}$$

The interaction potential $V_{\text{int}}$ is constructed from two contributions. The elastic term $V_{\text{elastic}}[n] = \lambda\, (\kappa[n] - \kappa_{\text{target}})^2$ penalizes deviations of the local curvature from a target value $\kappa_{\text{target}} = \langle\kappa\rangle_{\text{native}}$ computed from the known structure. The hydrophobic term $V_{\text{hydro}}[n]$ is a nonlocal potential that assigns a favorable energy to contacts between hydrophobic residues:

$$V_{\text{hydro}}[n] = - \sum_{m,\, |m-n| > 4} h[n]\, h[m]\, f\big(|\mathbf{r}[n] - \mathbf{r}[m]|\big), \tag{29}$$

where $h[n] \in \{0, 1\}$ is the hydrophobicity index of residue $n$ (assigned according to the Kyte-Doolittle [34] scale with a binary cutoff) and $f(r)$ is a contact function that equals unity for $r < 8$ Å and decays smoothly to zero beyond this distance. The backbone coordinates $\mathbf{r}[n]$ are reconstructed from $\psi[n]$ at each SCF step via the inverse Hasimoto transform [Eq. (6)] and discrete Frenet reconstruction.

In the second setting we add a hydrogen-bond potential

$$V_{\text{hb}}[n] = -\varepsilon_{\text{hb}} \sum_{m=n+3}^{n+5} g\big(\kappa[n], \kappa[m]\big)\, f\big(|\mathbf{r}[n] - \mathbf{r}[m]|\big), \tag{30}$$

where $\varepsilon_{\text{hb}}$ is a coupling strength and $g(\kappa[n], \kappa[m])$ is a geometric filter that favors the curvature values characteristic of helical hydrogen bonds. This term is designed to mimic the $i \rightarrow i+4$ backbone hydrogen bond that stabilizes $\alpha$-helices, expressed entirely in terms of the $\psi$ field variables.

The SCF iteration proceeds as follows. Starting from $\psi_0$, we (i) reconstruct the backbone coordinates, (ii) evaluate $V_{\text{int}}$ (and $V_{\text{hb}}$ in the second setting), (iii) update $\psi$ according to Eq. (27) with a discrete time step $\Delta t = 0.01$, and (iv) repeat for 5000 steps. At each step we record the mean curvature

$\langle\kappa\rangle = N^{-1}\sum_n |\psi[n]|$ and the $C_\alpha$ RMSD between the reconstructed backbone and the native structure (after optimal superposition). The best RMSD achieved over the trajectory is reported as the outcome for each protein.

### B. Results without hydrogen bonds

Figure 7(a) shows the final mean curvature $\langle\kappa\rangle_{\text{SCF}}$ plotted against the native target $\langle\kappa\rangle_{\text{target}}$ for all 856 proteins, colored by SCOP class. The dashed diagonal marks perfect agreement. The SCF systematically overestimates the mean curvature: the dataset-wide mean of $\langle\kappa\rangle_{\text{SCF}}$ is 1.776 rad compared to a target of 1.298 rad, an overestimation of approximately 37%. No SCOP class approaches the diagonal. The overestimation reflects the tendency of the damped DNLS dynamics to produce excessive bending: the hydrophobic potential drives chain compaction, which in the $\psi$ representation manifests as increased $|\psi| = \kappa$, but without the directional constraints of hydrogen bonds there is no mechanism to regulate the curvature to its native value.

Figure 7(b) displays the distribution of best-achieved RMSD across the 856 proteins. The mean RMSD is 13.1 Å with a range of 6.8 to 24.0 Å. No protein achieves RMSD below 5 Å, a threshold commonly used to indicate a successful fold prediction. The distribution shows no significant separation among the four SCOP classes: all-$\alpha$, all-$\beta$, $\alpha/\beta$, and $\alpha+\beta$ proteins are interleaved throughout the histogram. This uniformity confirms that the SCF failure is not specific to any particular fold topology.

Figure 7(c) shows the Pearson correlation $\rho(\tau_{\text{SCF}}, \tau_{\text{native}})$ between the torsion-angle profile produced by the SCF and the native torsion-angle profile, computed for each protein. The distribution is centered on zero (mean $\rho = 0.001$, median $\rho = 0.000$) and extends symmetrically over the range $[-0.4, 0.4]$. This indicates that the SCF-generated torsion angles bear no systematic relationship to the native structure. The torsion-angle sequence produced by the DNLS dynamics shows no detectable correlation with the native values.

The absence of correlation between predicted and native torsion angles constitutes a definitive failure of the predictive scheme. While the SCF dynamics may induce chain compaction—reflected in the curvature profile—the torsion angles $\tau[n]$, which govern the global fold topology, remain effectively stochastic. A correct curvature distribution coupled with uncorrelated torsion angles results in a compact globule devoid of native structural fidelity. The SCF dynamics, driven by the isotropic hydrophobic potential, generates chain compaction (increased $\kappa$) but cannot select the specific torsion-angle sequence that defines the native fold. This is consistent with the general understanding that hydrophobic collapse alone produces a compact but non-native ensemble. The selection of specific secondary and tertiary structure elements involves directional hydrogen bonds, whose geometric constraints are difficult to represent within the scalar DNLS framework.

### C. Results with hydrogen bonds

The bottom row of Fig. 7 shows the same three diagnostics for the SCF with the hydrogen-bond potential [Eq. (30)] included. The results are statistically indistinguishable from the setting without hydrogen bonds. The mean curvature increases slightly to $\langle\kappa\rangle_{\text{SCF}} = 1.790$ rad (compared to 1.776 without the hydrogen-bond term), moving further from the target rather than closer. The RMSD distribution is virtually identical: mean 13.1 Å, range 6.7 to 26.4 Å, with no protein reaching RMSD $< 5$ Å. The torsion correlation shifts negligibly from $\rho = 0.001$ to $\rho = -0.005$, remaining consistent with zero.

The inability of the hydrogen-bond term to improve the SCF outcome is significant. It suggests that the barrier to folding within the DNLS framework is not simply the absence of a particular interaction, but rather a representational limitation: the single complex scalar field $\psi[n] = \kappa[n]\exp(i\sum\tau[k])$ entangles curvature and torsion in its amplitude and phase. A physical hydrogen bond imposes simultaneous constraints on the distance and relative orientation of donor and acceptor groups, requiring independent control of both $\kappa$ and $\tau$ at specific residue pairs. The Hamiltonian evolution of $\psi$ couples $|\psi|$ and $\arg(\psi)$ through the nonlinear dynamics, so that any potential term that attempts to correct the curvature simultaneously perturbs the accumulated torsion phase, and vice versa. The hydrogen-bond potential [Eq. (30)], expressed in terms of $\psi$, cannot enforce the directional constraints that the physical interaction requires.

A natural question is whether replacing the interaction potential $V_{\text{int}}$ with a more refined functional, for instance a $C_\alpha$-level effective potential obtained by systematic coarse-graining of an all-atom force field such as AMBER [35, 36] or CHARMM [37, 38], would alter this conclusion. Three independent lines of evidence indicate that it would not. First, the oracle test [Fig. 1(b)] already supplies the exact native $V_{\text{re}}[n]$, which represents the theoretical optimum for the real part of the potential; the resulting RMSD of 20–120 Å demonstrates that the $2^N$ torsion-sign degeneracy (Barrier I) persists irrespective of energy-function quality. Second, the geometric dominance result (Barrier II, $\rho_{\text{geom}} = 0.951$) establishes that $V_{\text{re}}$ is 95% determined by backbone geometry rather than by amino-acid identity, so that even an exact force-field projection would produce a $V_{\text{re}}$ profile nearly indistinguishable from the geometry-only version. Third, the representational bottleneck is structural rather than parametric: the Hasimoto field $\psi$ encodes two real degrees of freedom ($\kappa$ and $\tau$) in a single complex number whose amplitude and phase are coupled under any Hamiltonian evolution. This coupling prevents the independent enforcement of the distance and angular constraints that directional interactions (hydrogen bonds, backbone dihedrals) impose on specific residue pairs, regardless of how accurately those interactions are parameterized.

The SCF test with its deliberately minimal energy function therefore probes the representational ceiling of the $\psi$ framework, not the quality of a particular force field.

This conclusion extends to state-of-the-art machine-learning force fields (MLFFs) such as AI2BMD [39],
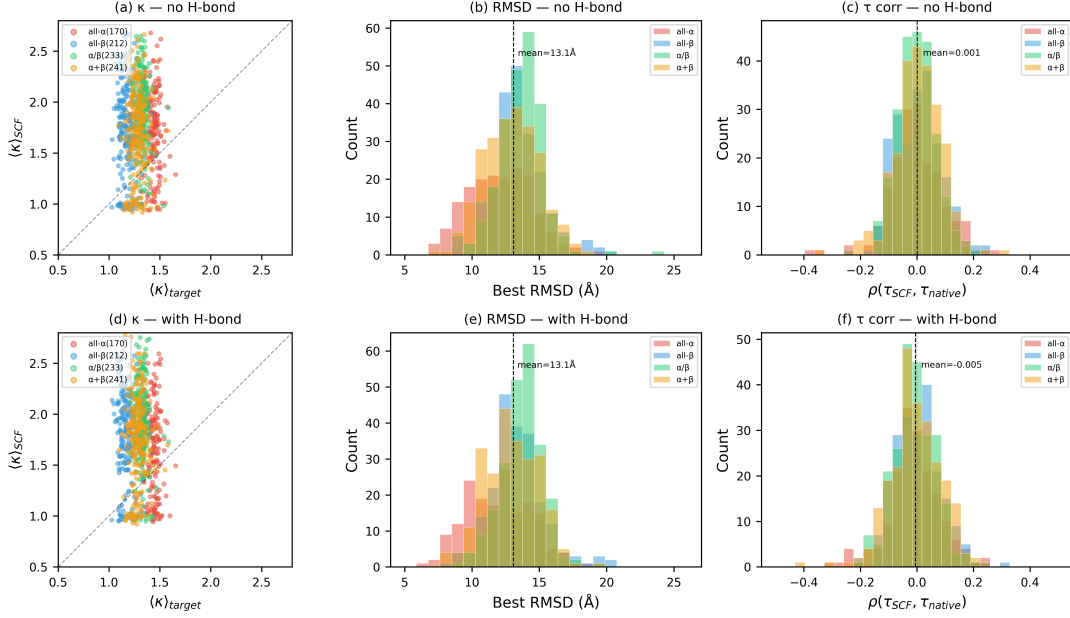
FIG. 7. Self-consistent field (SCF) test of DNLS-driven folding on 856 non-redundant proteins, without (top row) and with (bottom row) a hydrogen-bond potential. Points and histograms are colored by SCOP class: all-$\alpha$ (170, red), all-$\beta$ (212, blue), $\alpha/\beta$ (233, green), $\alpha+\beta$ (241, orange). (a, d) Mean bond curvature $\langle\kappa\rangle_{SCF}$ vs. native target $\langle\kappa\rangle_{target}$; the dashed line marks perfect agreement. In both settings $\langle\kappa\rangle_{SCF}$ is systematically overestimated (mean 1.776 and 1.790 vs. target 1.298), with no SCOP class approaching the diagonal. (b, e) Distribution of best-achieved RMSD. The two histograms are statistically indistinguishable (mean 13.1 Å; range 6.8–24.0 Å without hydrogen bonds, 6.7–26.4 Å with hydrogen bonds), and no protein in either setting reaches RMSD < 5 Å. (c, f) Pearson correlation $\rho(\tau_{SCF}, \tau_{native})$ between the SCF-predicted and native torsion-angle profiles. Both distributions are centered on zero (mean 0.001 vs. −0.005), indicating that the SCF torsion angles are uncorrelated with the native structure. The near-identical results across the two rows demonstrate that the folding failure is not due to missing hydrogen-bond physics but to the representational limitation of the single complex scalar field $\psi$: it cannot independently encode the curvature and torsion constraints required for structure determination.

MACE [40], and ViSNet [41]. While these architectures achieve near-quantum accuracy by operating on explicit Cartesian coordinates, they cannot rescue the Hasimoto framework. The bottleneck is strictly representational, not energetic. As demonstrated by the oracle test [Fig. 1(b)], even supplying the mathematically exact $V_{re}[n]$ derived from the native structure fails to recover the fold (RMSD 20–120 Å) due to the $2^N$ chiral degeneracy inherent in the scalar projection (Barrier I). Furthermore, since $V_{re}$ is determined ∼95% by geometry rather than sequence (Barrier II), the high-fidelity chemical information provided by MLFFs is largely filtered out when projected onto the DNLS potential. The scalar field $\psi$ effectively acts as a lossy compression channel that discards the geometric degrees of freedom needed to distinguish the native state, a structural deficit that no improvement in force-field accuracy can overcome.

### D. Summary of the dynamical barrier

The SCF experiments establish a third, dynamical barrier to using the Hasimoto map for protein structure prediction. The barrier can be stated concisely: the DNLS Hamiltonian dynamics, even supplemented with physically motivated interaction terms, cannot drive an unfolded chain to the native state because the single complex scalar field $\psi$ lacks the representational capacity to encode the independent curvature and torsion constraints imposed by directional hydrogen bonds. Three quantitative findings support this conclusion across all 856 proteins and both SCF settings. First, the mean curvature is systematically overestimated by approximately 37%, indicating that the dynamics produces excessive bending without the regulatory effect of hydrogen-bond geometry. Second, no protein achieves RMSD below 5 Å, and the RMSD distribution shows no dependence on SCOP class, confirming that the failure is universal rather than fold-specific. Third, the torsion correlation with the native structure is indistinguishable from zero ($\rho \approx 0$), demonstrating that the SCF dynamics generates effectively random torsion-angle sequences. The comparison between the two SCF settings is particularly informative. The addition of a hydrogen-bond potential within the DNLS framework produces no measurable improvement in any of the three diagnostics. This rules out the interpretation that the SCF failure is due to an incomplete energy function. The failure is instead structural: it originates in the entanglement of $\kappa$ and $\tau$ within the complex field $\psi$, which prevents the independent enforcement of the geometric constraints that hydrogen bonds impose on the backbone. The relationship between this dynamical barrier and the two static barriers of Sec. III is discussed in Sec. VI.

TABLE II. Comparison of the Hasimoto map applied to vortex filaments and to protein backbones. The four properties listed are necessary for the map to function as a predictive dynamical framework.

| Property | Vortex filament | Protein backbone |
|---|---|---|
| Interactions | Local (LIA) | Nonlocal |
| Medium | Homogeneous | Heterogeneous (20 AA) |
| Dynamics | Hamiltonian | Dissipative |
| $V_{\text{eff}}$ | Physical potential | Kinematic identity |

## VI.  DISCUSSION AND CONCLUSION

The Hasimoto transform was originally constructed for vortex filaments evolving under the local induction approximation (LIA). In that setting the transform is both kinematic and dynamic: it converts the geometric evolution of the filament into the integrable cubic NLS, whose soliton solutions describe physical excitations that propagate without dispersion. The success of this construction rests on four properties of the vortex system that proteins do not share: locality of interactions, homogeneity of the medium, Hamiltonian (non-dissipative) dynamics, and a physically determined effective potential. Protein folding violates all four. Hydrophobic contacts and electrostatic forces are nonlocal; the backbone comprises 20 chemically distinct amino acids; folding is a dissipative free-energy minimization in aqueous solvent; and, as shown in Sec. III, the effective potential $V_{\text{eff}}$ on proteins is an algebraic consequence of the backbone geometry rather than an independently specified physical quantity. Table II summarizes these distinctions.

The geometric formalism established by Niemi and collaborators provides a rigorous basis for the structural characterization of protein backbones. The identification of discrete soliton solutions with secondary-structure motifs offers a fundamental insight into geometric regularity, while multi-soliton ansätze yield compact parameterizations of folded states with sub-ångström accuracy [7, 8]. Recent extensions have successfully applied this basis to simulate thermal unfolding [9] and characterize topological phase transitions [10]. These studies constitute a comprehensive treatment of the *inverse problem*: given a known topology, the DNLS soliton basis affords an efficient representation and a robust framework for analyzing fluctuations around the native state.

The topological aspects of this program connect to a broader body of work on protein topology. Approximately 6% of known protein structures form knots, slipknots, or links whose folding requires the backbone to cross topological barriers [18, 19]. These entangled proteins highlight a fundamental limitation of any local geometric description: the topological invariants that distinguish a knotted from an unknotted fold are inherently global properties that cannot be fully determined from the local fields $(\kappa[n], \tau[n])$ at any finite number of sites. Consequently, soliton-based studies typically determine parameters by fitting to a known crystallographic structure, whether through Metropolis minimization of RMSD [7] or direct fitting to native coordinates [9].

In contrast, the present work interrogates the complementary *forward problem*: determining whether the DNLS framework permits *ab initio* prediction of the native structure solely from the amino-acid sequence. Our analysis indicates that fundamental barriers impede this predictive pathway. The exact decomposition (Sec. III) reveals that the effective potential $V_{\text{eff}}$ is determined predominantly by the target geometry rather than the sequence, creating an informational circularity that cannot be resolved by sequence-based potentials. Furthermore, the SCF experiments (Sec. V) demonstrate that physically motivated interactions fail to lift the $2^N$ torsion-sign degeneracy or overcome the representational limitations of the scalar field. These findings establish that the distinction between the inverse and forward problems is structural: the capacity to efficiently describe known folds does not imply the feasibility of predicting unknown conformations.

The suggestion that modulation instability of the DNLS may provide a dynamical mechanism for the emergence of secondary structure from an initially straight chain is directly tested by our SCF experiments. The instability does produce curvature growth from a nearly uniform initial state, consistent with the expected behavior. However, the resulting structures bear no resemblance to native folds: the torsion angles are uncorrelated with the native profile ($\rho \approx 0$), and the RMSD remains above 5 Å for all 856 proteins in both SCF settings. Modulation instability generates generic bending but cannot select the specific $(\kappa, \tau)$ sequence that defines a particular protein fold.

Although the Hasimoto map cannot serve as a predictive tool, the exact decomposition and the analyses built upon it yield three constructive results. First, the integrability error $E[n]$ [Eq. (22)] serves as a geometric probe of secondary structure. The ROC analysis (Fig. 6, AUC = 0.72) demonstrates that 72% of the helix/non-helix distinction defined by DSSP is captured by a scalar test of discrete helical symmetry applied to $C_\alpha$ coordinates alone, without reference to hydrogen-bond energies or side-chain identities. Whether this geometric criterion can provide useful secondary-structure annotation in data-limited settings, such as low-resolution density maps or coarse-grained trajectories, remains to be tested on independent benchmarks; the present AUC establishes the conceptual correspondence between DNLS integrability and helical geometry but does not by itself validate a practical tool. The gap between AUC = 0.72 and unity quantifies the structural information that hydrogen bonds and nonlocal interactions contribute beyond local geometric regularity. Second, the effective potential $V_{\text{eff}}[n]$ provides a structural fingerprint that is invariant under rigid-body transformations. Because $V_{\text{re}}$ is 95% determined by geometry (Sec. III D), homologous proteins with low sequence identity but similar folds produce similar $V_{\text{re}}[n]$ profiles, offering a one-dimensional scalar representation of backbone shape for structure comparison without spatial superposition. Third, the information-theoretic lower bound on chiral information loss is exact and model-independent: each residue contributes one bit of torsion-sign information encoded in $V_{\text{im}}$ but absent from $V_{\text{re}}$, producing a $2^L$-fold degeneracy that any predictive scheme based on the Hasimoto map would need to resolve.

The structural barriers identified here can be placed in

sharper relief by comparison with the data-driven methods that have solved the prediction problem in practice. AlphaFold 2 [20] and ESMFold [23] both predict the full $SE(3)$ rigid-body frame, comprising a rotation matrix and a translation vector, for every residue. This approach retains the complete orientational degrees of freedom of the discrete Frenet frame without projecting them onto a scalar field. The theoretical foundations of this representational choice are clarified by the geometric deep learning framework [42, 43], which demonstrates that neural architectures respecting the symmetry group of the data domain (specifically, $SE(3)$ equivariance for molecular geometry) achieve systematic gains in sample efficiency and generalization by building physical invariances directly into the network structure rather than learning them from data. The Hasimoto transform performs precisely the projection that these architectures avoid: it compresses the two independent geometric fields $(\kappa, \tau)$ into a single complex scalar $\psi = \kappa e^{i\Sigma\tau}$, and our Barrier I shows that this compression discards the torsion sign, introducing a $2^N$ chiral degeneracy that the full-frame representation avoids entirely. Barrier II reveals a second contrast: the DNLS effective potential is 95% determined by local backbone geometry, whereas AlphaFold's attention mechanism and ESM-Fold's language-model embeddings encode long-range co-evolutionary and contextual information that couples distant residues—information that statistical physics methods such as direct coupling analysis [16, 17] have shown to be extractable from evolutionary data via the maximum-entropy principle. The DNLS framework, operating on nearest-neighbor differences of $\psi$, has no analogous channel for nonlocal sequence information. While these comparisons do not diminish the value of the geometric approach, as data-driven models provide predictions rather than physical explanations, they nevertheless delineate the specific representational deficits that any future analytical theory must overcome: it must preserve the full $SE(3)$ frame at each residue and incorporate nonlocal, sequence-dependent interactions that go beyond the scalar Hasimoto field.

The three barriers identified in this work point to specific physical ingredients that a geometric theory of protein folding would need to incorporate. Hydrogen bonds impose simultaneous constraints on the distance and relative orientation of donor and acceptor groups, coupling $\kappa$ and $\tau$ at specific residue pairs in a manner that cannot be captured by a potential acting on the single complex field $\psi$. A geometric energy functional treating $\kappa[n]$ and $\tau[n]$ as independent fields, rather than combining them into a single complex scalar, would be a natural starting point. The hydrophobic effect, an entropic force mediated by solvent reorganization, is inherently nonlocal and temperature-dependent, with no natural representation in the DNLS Hamiltonian. The SCF experiments confirm that a contact-based hydrophobic potential produces chain compaction but cannot select secondary structure, consistent with the known distinction between hydrophobic collapse and folding. Finally, the geometric dominance of $V_{\mathrm{re}}$ means that sequence-dependent terms going beyond the weak modulation of virtual-bond lengths in $\beta_n^{\pm}$ would need to enter any predictive geometric framework explicitly.

In summary, we have derived an exact closed-form decomposition of the DNLS effective potential on protein backbones and used it to identify three structural barriers to forward structure prediction: the torsion-sign degeneracy encoded in $V_{\mathrm{im}}$, the geometric dominance of $V_{\mathrm{re}}$, and the universal failure of self-consistent field dynamics across 856 non-redundant proteins. These barriers are mathematical in nature and are unlikely to be resolved by parameter adjustment or algorithmic improvement within the DNLS framework alone. The Hasimoto map applied to proteins functions primarily as a kinematic identity rather than a dynamical equation: it provides a useful geometric language for describing folded states, a purely geometric helix detector, and a rotation-invariant structural fingerprint, but within the scope of our analysis it does not provide a viable pathway to *ab initio* protein structure prediction. The contrast with data-driven approaches such as AlphaFold [20, 21] and ESMFold [22, 23], which succeed by retaining the full $SE(3)$ frame and incorporating nonlocal sequence context, underscores that a future analytical theory of folding must go beyond the scalar Hasimoto field to preserve the geometric and chemical degrees of freedom that the folding process requires.

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

[2] H. Hasimoto, A soliton on a vortex filament, Journal of Fluid Mechanics **51**, 477 (1972).

[3] U. H. Danielsson, M. Lundgren, and A. J. Niemi, Gauge field theory of chirally folded homopolymers with applications to folded proteins, Physical Review E—Statistical, Nonlinear, and Soft Matter Physics **82**, 021910 (2010).

[4] M. Chernodub, S. Hu, and A. J. Niemi, Topological solitons and folded proteins, Physical Review E—Statistical, Nonlinear, and Soft Matter Physics **82**, 011916 (2010).

[5] A. Molochkov, A. Begun, and A. Niemi, Gauge symmetries and structure of proteins, in *EPJ Web of Conferences*, Vol. 137 (EDP Sciences, 2017) p. 04004.

[6] D. Melnikov and A. B. Neves, Chern-simons-higgs model as a theory of protein molecules, Journal of Applied Physics **126** (2019).

[7] N. Molkenthin, S. Hu, and A. J. Niemi, Discrete nonlinear schrödinger equation and polygonal solitons with applications to collapsed proteins, Physical Review Letters **106**, 078102 (2011).

[8] A. Krokhotin, A. J. Niemi, and X. Peng, Soliton concepts and protein structure, Physical Review E—Statistical, Nonlinear, and Soft Matter Physics **85**, 031906 (2012).

[9] A. Begun, S. Liubimov, A. Molochkov, and A. J. Niemi, On topology and knotty entanglement in protein folding, PLoS One **16**, e0244547 (2021).

[10] A. Begun, M. N. Chernodub, A. Molochkov, and A. J. Niemi, Local topology and perestroikas in protein structure and folding dynamics, Physical Review E **111**, 024406 (2025).

[11] S. Liubimov, N. Gerasimeniuk, A. Korneev, and A. Molochkov, Modeling the structure of myoglobin within the abelian higgs model, Biochemistry (Moscow), Supplement Series A: Membrane and Cell Biology **19**, 449 (2025).

[12] K. A. Dill and J. L. MacCallum, The protein-folding problem, 50 years on, science **338**, 1042 (2012).

[13] A. Davtyan, N. P. Schafer, W. Zheng, C. Clementi, P. G. Wolynes, and G. A. Papoian, Awsem-md: protein structure prediction using coarse-grained physical potentials and bioinformatically based local structure biasing, The Journal of Physical Chemistry B **116**, 8494 (2012).

[14] J. R. Banavar and A. Maritan, Physics of proteins, Annu. Rev. Biophys. Biomol. Struct. **36**, 261 (2007).

[15] J. R. Banavar, A. Maritan, C. Micheletti, and A. Trovato, Geometry and physics of proteins, Proteins: Structure, Function, and Bioinformatics **47**, 315 (2002).

[16] F. Morcos, A. Pagnani, B. Lunt, A. Bertolino, D. S. Marks, C. Sander, R. Zecchina, J. N. Onuchic, T. Hwa, and M. Weigt, Direct-coupling analysis of residue coevolution captures native contacts across many protein families, Proceedings of the National Academy of Sciences **108**, E1293 (2011).

[17] S. Cocco, R. Monasson, and M. Weigt, From principal component to direct coupling analysis of coevolution in proteins: low-eigenvalue modes are needed for structure prediction, PLoS computational biology **9**, e1003176 (2013).

[18] J. I. Sulkowska, On folding of entangled proteins: knots, lassos, links and $\theta$-curves, Current opinion in structural biology **60**, 131 (2020).

[19] P. Dabrowski-Tumanski and J. I. Sulkowska, Topological knots and links in proteins, Proceedings of the National Academy of Sciences **114**, 3415 (2017).

[20] J. Jumper, R. Evans, A. Pritzel, T. Green, M. Figurnov, O. Ronneberger, K. Tunyasuvunakool, R. Bates, A. Žídek, A. Potapenko, *et al.*, Highly accurate protein structure prediction with alphafold, nature **596**, 583 (2021).

[21] J. Abramson, J. Adler, J. Dunger, R. Evans, T. Green, A. Pritzel, O. Ronneberger, L. Willmore, A. J. Ballard, J. Bambrick, *et al.*, Accurate structure prediction of biomolecular interactions with alphafold 3, Nature **630**, 493 (2024).

[22] A. Rives, J. Meier, T. Sercu, S. Goyal, Z. Lin, J. Liu, D. Guo, M. Ott, C. L. Zitnick, J. Ma, *et al.*, Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences, Proceedings of the National Academy of Sciences **118**, e2016239118 (2021).

[23] Z. Lin, H. Akin, R. Rao, B. Hie, Z. Zhu, W. Lu, N. Smetanin, R. Verkuil, O. Kabeli, Y. Shmueli, *et al.*, Evolutionary-scale prediction of atomic-level protein structure with a language model, Science **379**, 1123 (2023).

[24] Y. Wang, M. Cai, Y. Ma, and K. Wei, Sequence-context-aware decoding enables robust reconstruction of protein dynamics from crystallographic b-factors, bioRxiv , 2025 (2025).

[25] R. A. Engh and R. Huber, Accurate bond and angle parameters for x-ray protein structure refinement, Foundations of Crystallography **47**, 392 (1991).

[26] S. C. Lovell, I. W. Davis, W. B. Arendall III, P. I. De Bakker, J. M. Word, M. G. Prisant, J. S. Richardson, and D. C. Richardson, Structure validation by c$\alpha$ geometry: $\phi$, $\psi$ and c$\beta$ deviation, Proteins: Structure, Function, and Bioinformatics **50**, 437 (2003).

[27] A. G. Murzin, S. E. Brenner, T. Hubbard, and C. Chothia, Scop: a structural classification of proteins database for the investigation of sequences and structures, Journal of molecular biology **247**, 536 (1995).

[28] A. Andreeva, E. Kulesha, J. Gough, and A. G. Murzin, The scop database in 2020: expanded classification of representative family and superfamily domains of known protein structures, Nucleic acids research **48**, D376 (2020).

[29] G. Wang and R. L. Dunbrack Jr, Pisces: a protein sequence culling server, Bioinformatics **19**, 1589 (2003).

[30] Y. Wang, M. Cai, Y. Dong, Y. Ma, and K. Wei, From signal to symphony: Exploring 2d sequence representations for protein function prediction, Journal of Chemical Information and Modeling **65**, 12723 (2025).

[31] W. Kabsch and C. Sander, Dictionary of protein secondary structure: pattern recognition of hydrogen-bonded and geometrical features, Biopolymers: Original Research on Biomolecules **22**, 2577 (1983).

[32] M.-A. Lancret, Mémoire sur les courbes à double courbure, Memoires presentes alInstitut **1**, 416 (1806).

[33] M. P. do Carmo, Differential geometry of curves andsurfaces, Englewood Cliffs, New Jersey (1976).

[34] J. Kyte and R. F. Doolittle, A simple method for displaying the hydropathic character of a protein, Journal of molecular biology **157**, 105 (1982).

[35] W. D. Cornell, P. Cieplak, C. I. Bayly, I. R. Gould, K. M. Merz, D. M. Ferguson, D. C. Spellmeyer, T. Fox, J. W. Caldwell, and P. A. Kollman, A second generation force field for the simulation of proteins, nucleic acids, and organic molecules, Journal of the American Chemical Society **117**, 5179 (1995).

[36] D. A. Case, T. E. Cheatham III, T. Darden, H. Gohlke, R. Luo, K. M. Merz Jr, A. Onufriev, C. Simmerling, B. Wang, and R. J. Woods, The amber biomolecular simulation programs, Journal of computational chemistry **26**, 1668 (2005).

[37] B. R. Brooks, R. E. Bruccoleri, B. D. Olafson, D. J. States, S. a. Swaminathan, and M. Karplus, Charmm: a program for macromolecular energy, minimization, and dynamics calculations, Journal of computational chemistry **4**, 187 (1983).

[38] B. R. Brooks, C. L. Brooks III, A. D. Mackerell Jr, L. Nilsson, R. J. Petrella, B. Roux, Y. Won, G. Archontis, C. Bartels, S. Boresch, *et al.*, Charmm: the biomolecular simulation program, Journal of computational chemistry **30**, 1545 (2009).

[39] T. Wang, X. He, M. Li, Y. Li, R. Bi, Y. Wang, C. Cheng, X. Shen, J. Meng, H. Zhang, *et al.*, Ab initio characterization of protein molecular dynamics with ai2bmd, Nature **635**, 1019 (2024).

[40] I. Batatia, D. P. Kovacs, G. Simm, C. Ortner, and G. Csányi, Mace: Higher order equivariant message passing neural networks for fast and accurate force fields, Advances in neural information processing systems **35**, 11423 (2022).

[41] Y. Wang, T. Wang, S. Li, X. He, M. Li, Z. Wang, N. Zheng, B. Shao, and T.-Y. Liu, Enhancing geometric representations for molecules with equivariant vector-scalar interactive message passing, Nature Communications **15**, 313 (2024).

[42] M. M. Bronstein, J. Bruna, T. Cohen, and P. Veličković, Geometric deep learning: Grids, groups, graphs, geodesics, and gauges, arXiv preprint arXiv:2104.13478 (2021).

[43] M. M. Bronstein, J. Bruna, Y. LeCun, A. Szlam, and P. Vandergheynst, Geometric deep learning: going beyond euclidean

data, IEEE Signal Processing Magazine **34**, 18 (2017).