

On the Consequences of Categorical Completion on Biological Sequences: Thermodynamic Symbolic Computational Molecular Language for Database-Free Peptide Sequence Reconstruction

Kundai Farai Sachikonye
`sachikonye@wzw.tum.de`

December 3, 2025

Abstract

We present a mathematical framework for database-free peptide sequence reconstruction via S-Entropy coordinate transformation. The framework maps amino acids to a three-dimensional S-Entropy coordinate space (S_k, S_t, S_e) derived from physicochemical properties: hydrophobicity maps to the knowledge dimension, molecular volume maps to the time dimension, and electrostatic properties map to the entropy dimension. Peptide fragmentation is formalised as a molecular grammar with production rules generating b/y ion series from precursor sequences. Fragment observations form a directed graph where vertices represent detected ions and edges encode sequential amino acid relationships satisfying mass difference constraints. Sequence reconstruction reduces to finding a minimum-entropy Hamiltonian path through this graph, with categorical completion filling gaps between non-adjacent fragments. A dynamic dictionary architecture supports zero-shot identification of amino acids via KD-tree nearest-neighbor lookup in S-Entropy space and learns novel molecular entities through equilibrium-seeking dynamics. The Molecular Maxwell Demon orchestration system integrates these components through variance minimisation, achieving peptide identification without reference to sequence databases. Cross-modal validation confirms reconstructions by matching theoretical fragment masses against observed spectra.

Contents

| | |
|--|----------|
| 1 S-Entropy Coordinate Transformation for Amino Acids | 4 |
| 1.1 Mathematical Foundation | 4 |
| 1.2 Physicochemical Property Mapping | 4 |
| 1.3 Standard Amino Acid Coordinates | 4 |
| 1.4 Post-Translational Modification Shifts | 6 |
| 1.5 Sequence Coordinate Path | 6 |
| 1.6 Sequence S-Entropy | 7 |

| | |
|--|-----------|
| 2 Molecular Fragmentation Grammar | 7 |
| 2.1 Formal Grammar Definition | 7 |
| 2.2 Production Rules for Fragmentation | 7 |
| 2.3 Ion Type Classification | 9 |
| 2.4 Neutral Loss Rules | 9 |
| 2.5 Fragment Mass Calculation | 9 |
| 2.6 Complementarity Constraint | 9 |
| 2.7 Sequential Relationship Constraint | 10 |
| 2.8 Complete Fragment Generation | 10 |
| 3 Fragment Graph Construction | 10 |
| 3.1 Graph Formalism | 10 |
| 3.2 Edge Construction | 12 |
| 3.3 Graph Construction Algorithm | 12 |
| 3.4 S-Entropy Magnitude | 14 |
| 3.5 Path Finding | 14 |
| 3.6 Greedy Path Construction | 14 |
| 3.7 Path Entropy Calculation | 15 |
| 4 Categorical Empty Dictionary Architecture | 15 |
| 4.1 Dictionary Structure | 15 |
| 4.2 Equivalence Classes | 16 |
| 4.3 KD-Tree Index | 16 |
| 4.4 Dynamic Learning | 17 |
| 4.5 Zero-Shot Identification | 17 |
| 4.6 Empty Dictionary Principle | 18 |
| 4.7 Persistence | 18 |
| 5 Categorical Sequence Reconstruction | 18 |
| 5.1 Problem Formulation | 18 |
| 5.2 Gap Region Identification | 18 |
| 5.3 Categorical Completion | 21 |
| 5.4 Reconstruction Algorithm | 21 |
| 5.5 Reconstruction Result | 22 |
| 5.6 Coverage and Confidence Metrics | 22 |
| 5.7 Validation Scores | 22 |
| 5.8 Cross-Modal Validation | 22 |
| 6 Molecular Maxwell Demon System | 23 |
| 6.1 System Architecture | 23 |
| 6.2 Configuration | 23 |
| 6.3 Spectrum Analysis Pipeline | 24 |
| 6.4 Batch Processing | 24 |
| 6.5 Variance Minimization Principle | 25 |
| 6.6 Cross-Modal Pathway Validation | 25 |
| 6.7 Dictionary Update Protocol | 27 |
| 6.8 System Output | 27 |

| | |
|---|-----------|
| 7 Discussion | 27 |
| 7.1 S-Entropy Coordinate Space Properties | 27 |
| 7.2 Fragment Graph Structure | 27 |
| 7.3 Path Finding Complexity | 28 |
| 7.4 Categorical Completion | 28 |
| 7.5 Dictionary Architecture | 28 |
| 7.6 Cross-Modal Validation | 28 |
| 7.7 System Integration | 28 |
| 8 Conclusion | 28 |

1 S-Entropy Coordinate Transformation for Amino Acids

1.1 Mathematical Foundation

We define the amino acid S-Entropy transformation $\phi_{AA} : \mathcal{A} \rightarrow \mathcal{S}^3$, where \mathcal{A} represents the set of 20 standard amino acids and $\mathcal{S}^3 \subset \mathbb{R}^3$ represents the three-dimensional S-Entropy coordinate space.

Definition 1 (S-Entropy Coordinate Space). *The S-Entropy coordinate space is defined as:*

$$\mathcal{S}^3 = \{(S_k, S_t, S_e) \in \mathbb{R}^3 : S_k, S_t, S_e \in [0, 1]\} \quad (1)$$

where S_k denotes the knowledge dimension, S_t denotes the time dimension, and S_e denotes the entropy dimension.

1.2 Physicochemical Property Mapping

The transformation ϕ_{AA} maps amino acid physicochemical properties to S-Entropy coordinates through the following assignments:

Definition 2 (Amino Acid S-Entropy Transformation). *For amino acid $a \in \mathcal{A}$ with hydrophobicity $H(a)$, van der Waals volume $V(a)$, charge $Q(a)$, and polarity $P(a)$, the S-Entropy coordinates are:*

$$S_k(a) = \frac{H(a) - H_{\min}}{H_{\max} - H_{\min}} \quad (2)$$

$$S_t(a) = \frac{V(a)}{V_{\max}} \quad (3)$$

$$S_e(a) = \frac{|Q(a)| + \mathbf{1}_{P(a)}}{2} \quad (4)$$

where $H_{\min} = -4.5$, $H_{\max} = 4.5$ correspond to the Kyte-Doolittle scale bounds, $V_{\max} = 250 \text{ \AA}^3$, and $\mathbf{1}_{P(a)}$ is the indicator function for polar residues.

The knowledge dimension S_k captures information content through hydrophobicity, which correlates with membrane interaction propensity and protein folding energetics. The time dimension S_t encodes molecular size through van der Waals volume. The entropy dimension S_e quantifies electrostatic complexity through the combined effect of formal charge and hydrogen-bonding capacity.

1.3 Standard Amino Acid Coordinates

Table 1 presents the S-Entropy coordinates for the 20 standard amino acids computed via the transformation ϕ_{AA} .

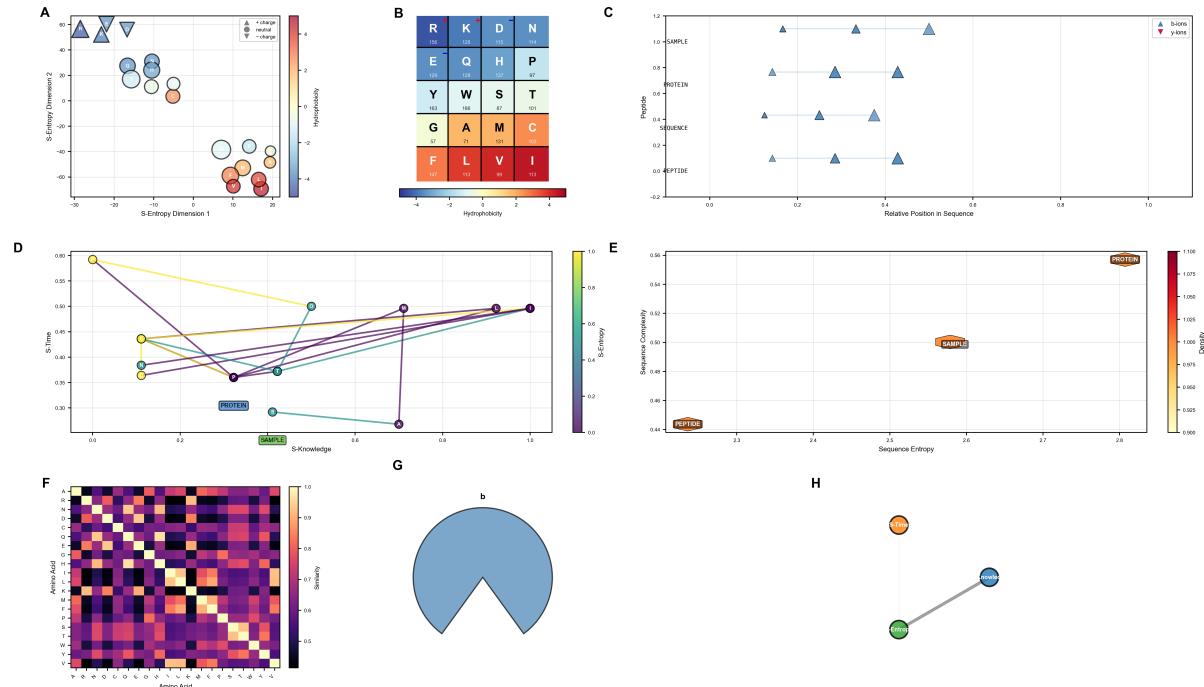


Figure 1: Comprehensive atlas of the categorical amino acid molecular language framework. **(a)** Two-dimensional t-SNE projection of amino acids in S-Entropy space, with charge state encoded by marker shape (upward triangles = positive charge, circles = neutral, downward triangles = negative charge). Marker size represents molecular mass, and color indicates hydrophobicity (viridis colormap: yellow = hydrophobic, blue = hydrophilic). Spatial clustering reveals natural organization: charged residues (R, K, D, E) cluster at left, hydrophobic residues (F, L, V, I) at right, demonstrating that S-Entropy coordinates capture physicochemical similarity. Legend shows charge categories. **(b)** Periodic table organization of 20 standard amino acids arranged by hydrophobicity (x-axis, blue to red gradient) and other physicochemical properties. Each cell displays single-letter code and molecular mass. Color intensity reflects hydrophobicity scale from -4 (hydrophilic, blue) to +4 (hydrophobic, red). This organization mirrors the spatial clustering in panel (a), validating the S-Entropy transformation. **(c)** Fragment position versus peptide relative position for three example sequences (PEPTIDE, SEQUENCE, PROTEIN). Triangles represent fragment observations, colored by ion type (blue = b-ions, red = y-ions). Y-axis shows peptide identity, x-axis shows normalized position (0 = N-terminus, 1 = C-terminus). Diagonal pattern indicates sequential fragment coverage, with gaps representing missing fragments that require categorical completion (Section 5.3). **(d)** Peptide trajectories through S-Entropy space. Line plot showing S-Time (y-axis) versus S-Knowledge (x-axis) for three peptides (PROTEIN in blue, SAMPLE in green, PEPTIDE in red). Each point represents one amino acid position, connected by lines showing sequential progression. Distinct trajectory shapes encode sequence identity, with labeled points (e.g., "PROTEIN", "SAMPLE") marking specific positions. Smooth paths validate that sequences form continuous trajectories (Equation 7). **(e)** Sequence entropy versus sequence complexity scatter plot. Three peptides shown as labeled points with rectangular borders: PEPTIDE (red, high entropy 2.24, high complexity 0.499), SAMPLE (orange, moderate values), PROTEIN (blue, highest entropy 2.8). Color gradient from yellow (high density) to red (low density) indicates local density of sequences in entropy-complexity space. Demonstrates that different sequences occupy distinct regions, enabling discrimination. **(f)** Amino acid similarity matrix (heatmap) showing pairwise Euclidean distances in S-Entropy space. Rows and columns represent amino acids (single-letter codes). Color scale from black (distance = 0, identical) through purple, pink, yellow to white (distance = 1.0, maximally different). Diagonal is black (self-similarity). Block structure reveals equivalence classes: hydrophobic cluster (F, L, I, V, M), charged cluster (K, R, D, E), polar cluster (S, T, N, Q). This matrix defines equivalence classes (Equation 8).

| Amino Acid | Symbol | Mass (Da) | S_k | S_t | S_e |
|---------------|--------|-----------|-------|-------|-------|
| Alanine | A | 71.037 | 0.700 | 0.268 | 0.000 |
| Arginine | R | 156.101 | 0.000 | 0.592 | 1.000 |
| Asparagine | N | 114.043 | 0.111 | 0.384 | 0.500 |
| Aspartic acid | D | 115.027 | 0.111 | 0.364 | 1.000 |
| Cysteine | C | 103.009 | 0.778 | 0.344 | 0.500 |
| Glutamine | Q | 128.059 | 0.111 | 0.456 | 0.500 |
| Glutamic acid | E | 129.043 | 0.111 | 0.436 | 1.000 |
| Glycine | G | 57.021 | 0.456 | 0.192 | 0.000 |
| Histidine | H | 137.059 | 0.144 | 0.472 | 0.500 |
| Isoleucine | I | 113.084 | 1.000 | 0.496 | 0.000 |
| Leucine | L | 113.084 | 0.922 | 0.496 | 0.000 |
| Lysine | K | 128.095 | 0.067 | 0.540 | 1.000 |
| Methionine | M | 131.040 | 0.711 | 0.496 | 0.000 |
| Phenylalanine | F | 147.068 | 0.811 | 0.540 | 0.000 |
| Proline | P | 97.053 | 0.322 | 0.360 | 0.000 |
| Serine | S | 87.032 | 0.411 | 0.292 | 0.500 |
| Threonine | T | 101.048 | 0.422 | 0.372 | 0.500 |
| Tryptophan | W | 186.079 | 0.400 | 0.652 | 0.000 |
| Tyrosine | Y | 163.063 | 0.356 | 0.564 | 0.500 |
| Valine | V | 99.068 | 0.967 | 0.420 | 0.000 |

Table 1: S-Entropy coordinates for standard amino acids.

1.4 Post-Translational Modification Shifts

Post-translational modifications (PTMs) are represented as affine transformations in S-Entropy space.

Definition 3 (PTM S-Entropy Shift). *For a post-translational modification τ applied to amino acid a , the modified coordinates are:*

$$\phi_{AA}(a^\tau) = \phi_{AA}(a) + \Delta\mathbf{S}_\tau \quad (5)$$

where $\Delta\mathbf{S}_\tau = (\Delta S_k, \Delta S_t, \Delta S_e)$ is the PTM-specific shift vector.

The shift components are computed as:

$$\Delta S_k = 0.2 \cdot \tanh\left(\frac{\Delta m_\tau}{100}\right) \quad (6)$$

$$\Delta S_t = \frac{\Delta m_\tau}{200} \quad (7)$$

$$\Delta S_e = \text{sgn}(\Delta m_\tau) \cdot 0.1 \quad (8)$$

where Δm_τ is the mass shift induced by the modification.

1.5 Sequence Coordinate Path

Definition 4 (Peptide S-Entropy Path). *For a peptide sequence $\mathbf{s} = (s_1, s_2, \dots, s_n)$ where $s_i \in \mathcal{A}$, the S-Entropy coordinate path is:*

$$\mathbf{P}(\mathbf{s}) = (\phi_{AA}(s_1), \phi_{AA}(s_2), \dots, \phi_{AA}(s_n)) \quad (9)$$

forming a trajectory through \mathcal{S}^3 .

The cumulative path displacement is defined as:

$$\mathbf{D}(\mathbf{s}) = \sum_{i=1}^n \phi_{AA}(s_i) \quad (10)$$

1.6 Sequence S-Entropy

Definition 5 (Sequence S-Entropy). *The sequence S-Entropy measures the information content of the coordinate path:*

$$H_S(\mathbf{s}) = - \sum_{\mathbf{c} \in \mathcal{B}} p(\mathbf{c}) \log_2 p(\mathbf{c}) \quad (11)$$

where \mathcal{B} is a binning of \mathcal{S}^3 into discrete states and $p(\mathbf{c})$ is the empirical probability of coordinate path elements falling in bin \mathbf{c} .

Definition 6 (Sequence Complexity). *The sequence complexity score combines Shannon entropy with a repetition penalty:*

$$C(\mathbf{s}) = \frac{H(\mathbf{s})}{H_{\max}} \cdot \left(1 - \frac{\ell_{\max}(\mathbf{s})}{n}\right) \quad (12)$$

where $H(\mathbf{s})$ is the Shannon entropy of amino acid frequencies, $H_{\max} = \log_2(20)$ is the maximum entropy for 20 amino acids, and $\ell_{\max}(\mathbf{s})$ is the length of the longest repeating substring.

2 Molecular Fragmentation Grammar

2.1 Formal Grammar Definition

We define a molecular grammar $G = (\Sigma, N, P, S)$ for peptide fragmentation, where:

- $\Sigma = \mathcal{A}$ is the alphabet of terminal symbols (amino acids)
- $N = \{S, B, Y, F\}$ is the set of non-terminal symbols
- P is the set of production rules
- S is the start symbol (intact peptide)

2.2 Production Rules for Fragmentation

Definition 7 (Fragmentation Production Rule). *The primary fragmentation production rule is:*

$$P_{frag} : S \rightarrow F_N \oplus F_C \quad (13)$$

where F_N denotes the N-terminal fragment and F_C denotes the C-terminal fragment, with \oplus representing the peptide bond cleavage operator.

For a peptide sequence $\mathbf{s} = s_1 s_2 \dots s_n$, cleavage at bond position k produces:

$$F_N^{(k)} = s_1 s_2 \dots s_k \quad (\text{b-ion}) \quad (14)$$

$$F_C^{(k)} = s_{k+1} s_{k+2} \dots s_n \quad (\text{y-ion}) \quad (15)$$

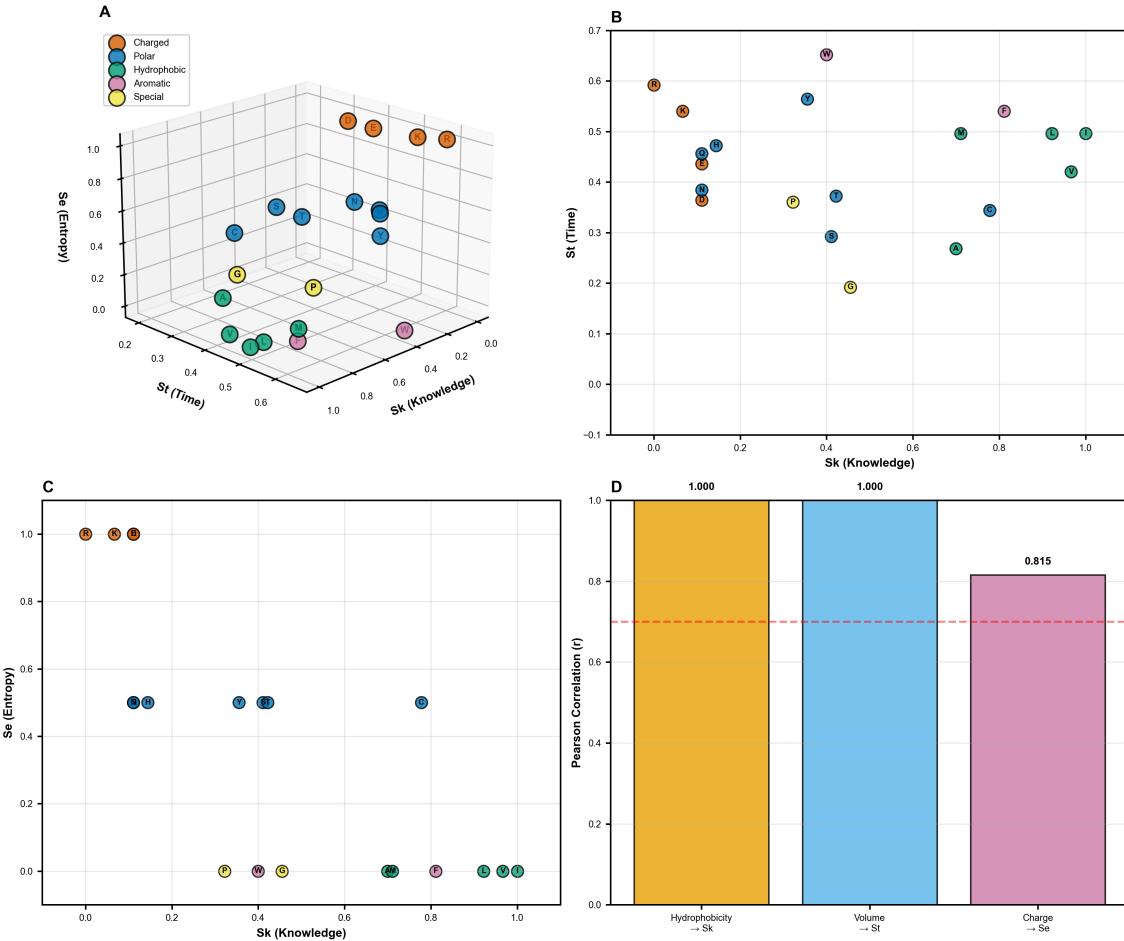


Figure 2: Amino acid representation in tri-dimensional S-Entropy coordinate space. (a) Three-dimensional scatter plot showing all 20 standard amino acids in S-Entropy space (S_k , S_t , S_e). Amino acids are colored by chemical category: charged (orange), polar (teal), hydrophobic (blue), aromatic (green), and special (yellow). Marker shapes encode charge state: upward triangles (positive: K, R, H), circles (neutral), downward triangles (negative: D, E). Single-letter codes label each amino acid. Spatial separation between categories demonstrates that physicochemical properties map to distinct S-Entropy coordinates, enabling categorical identification. (b) Two-dimensional projection of amino acids onto S-Knowledge (S_k) versus S-Time (S_t) plane, with points colored by S-Entropy (S_e) using viridis colormap (yellow = high entropy, purple = low entropy). Hydrophobic residues (F, W, I, L, V, M) cluster at high S_k values (right), while charged residues (K, R, D, E) cluster at low S_k values (left). Molecular volume correlates with S_t (y-axis). (c) S-Entropy distribution across amino acids in (S_k , S_e) space. Charged residues (K, R at top) show high S_e values (electrostatic complexity), while special residues (G, P, C at bottom) show low S_e values. Hydrophobic residues occupy intermediate S_e range. (d) Pearson correlation analysis between physicochemical properties and S-Entropy dimensions. Hydrophobicity strongly correlates with S_k ($r = 1.000$), molecular volume with S_t ($r = 1.000$), and charge with S_e ($r = 0.815$). High correlations validate the coordinate transformation defined in Equations 2-4. This comprehensive visualization establishes that the S-Entropy transformation (Definition 1) maps amino acids to a well-structured coordinate space where physicochemical similarity corresponds to spatial proximity, forming the foundation for database-free categorical identification.

2.3 Ion Type Classification

Definition 8 (Ion Type Enumeration). *The fragment ion types are classified as:*

$$b\text{-ion: } [F_N + H]^+ \quad (16)$$

$$y\text{-ion: } [F_C + H_2O + H]^+ \quad (17)$$

$$a\text{-ion: } [F_N - CO]^+ \quad (18)$$

$$c\text{-ion: } [F_N + NH_3]^+ \quad (19)$$

$$x\text{-ion: } [F_C + CO]^+ \quad (20)$$

$$z\text{-ion: } [F_C - NH_3]^+ \quad (21)$$

2.4 Neutral Loss Rules

Definition 9 (Neutral Loss Production). *Neutral loss productions extend the grammar with:*

$$P_{H_2O} : F \rightarrow F^* + H_2O \quad (\Delta m = -18.011) \quad (22)$$

$$P_{NH_3} : F \rightarrow F^* + NH_3 \quad (\Delta m = -17.027) \quad (23)$$

$$P_{CO} : F \rightarrow F^* + CO \quad (\Delta m = -27.995) \quad (24)$$

The neutral loss rules are context-dependent:

- P_{H_2O} is applicable when $\exists s_i \in \{S, T, E, D\}$ in the fragment
- P_{NH_3} is applicable when $\exists s_i \in \{R, K, N, Q\}$ in the fragment

2.5 Fragment Mass Calculation

Definition 10 (Theoretical Fragment Mass). *For fragment sequence $F = f_1 f_2 \dots f_m$ of ion type τ with neutral loss λ :*

$$m_{theo}(F, \tau, \lambda, z) = \frac{\sum_{i=1}^m m(f_i) + \delta_\tau - \delta_\lambda + z \cdot m_H}{z} \quad (25)$$

where $m(f_i)$ is the monoisotopic mass of amino acid f_i , δ_τ is the ion type mass modifier, δ_λ is the neutral loss mass, z is the charge state, and $m_H = 1.008$ Da is the proton mass.

The ion type mass modifiers are:

$$\delta_b = 1.008 \quad (26)$$

$$\delta_y = 19.018 \quad (27)$$

$$\delta_a = 1.008 - 27.995 \quad (28)$$

$$\delta_c = 1.008 + 17.027 \quad (29)$$

2.6 Complementarity Constraint

Theorem 1 (b/y Complementarity). *For a peptide of precursor mass M , the b-ion at position k and y-ion at position $n - k$ satisfy:*

$$m(b_k) + m(y_{n-k}) = M + 2 \cdot m_H \quad (30)$$

Proof. Let $\mathbf{s} = s_1 \dots s_n$ with total residue mass $M_r = \sum_{i=1}^n m(s_i)$.

$$\text{For the b-ion: } m(b_k) = \sum_{i=1}^k m(s_i) + m_H$$

$$\text{For the y-ion: } m(y_{n-k}) = \sum_{i=k+1}^n m(s_i) + m_H + 18.010$$

$$\text{Sum: } m(b_k) + m(y_{n-k}) = M_r + 2m_H + 18.010 = M + 2m_H$$

□

□

2.7 Sequential Relationship Constraint

Definition 11 (Sequential Ion Series). *Consecutive ions in the b-series satisfy:*

$$m(b_{k+1}) - m(b_k) = m(s_{k+1}) \quad (31)$$

and consecutive ions in the y-series satisfy:

$$m(y_{k+1}) - m(y_k) = m(s_{n-k}) \quad (32)$$

These constraints form the basis for sequence tag extraction and de novo sequencing.

2.8 Complete Fragment Generation

For a peptide sequence of length n , the complete production rule set P_{complete} generates:

- $n - 1$ b-ions: $\{b_1, b_2, \dots, b_{n-1}\}$
- $n - 1$ y-ions: $\{y_1, y_2, \dots, y_{n-1}\}$
- Additional neutral loss ions where applicable

The total theoretical fragment count is:

$$|P_{\text{complete}}| = 2(n - 1) + \sum_{k=1}^{n-1} |\mathcal{L}(b_k)| + \sum_{k=1}^{n-1} |\mathcal{L}(y_k)| \quad (33)$$

where $|\mathcal{L}(F)|$ is the number of applicable neutral losses for fragment F .

3 Fragment Graph Construction

3.1 Graph Formalism

We construct a directed graph $\mathcal{G} = (V, E)$ where vertices represent observed fragments and edges represent sequential relationships consistent with peptide bond cleavage patterns.

Definition 12 (Fragment Node). *A fragment node $v \in V$ is a tuple:*

$$v = (id, \sigma, \mathbf{S}, m, \tau, k, c) \quad (34)$$

where:

- id is a unique identifier
- $\sigma \in \mathcal{A}^*$ is the partial sequence (if identified)

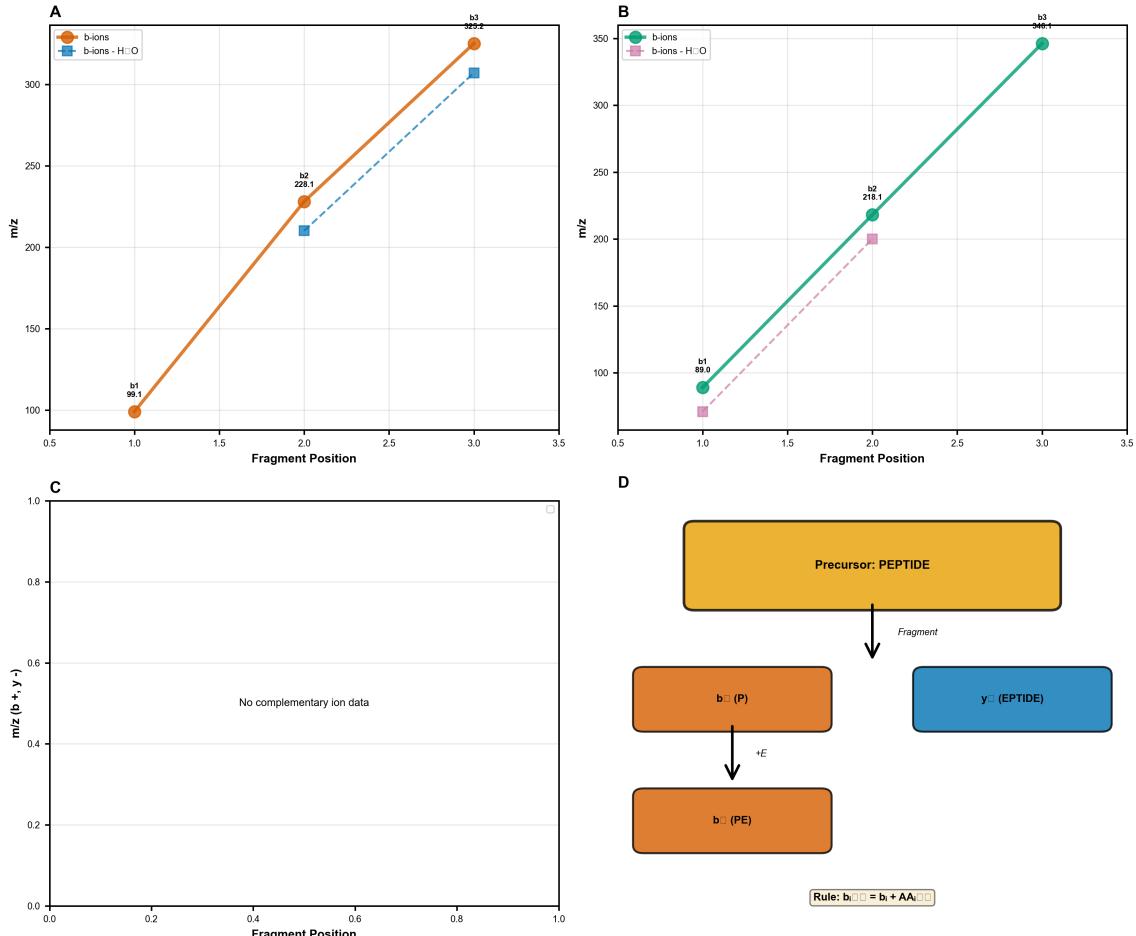


Figure 3: Molecular fragmentation grammar generates b-ion and y-ion series.

(a) Fragment mass ladder for b-ions (orange circles) and b-ions with water loss (blue squares, $b\text{-H}_2\text{O}$) as a function of fragment position along the peptide sequence. Solid orange line shows theoretical b-ion masses calculated via Equation 11. Dashed blue line shows $b\text{-H}_2\text{O}$ masses (neutral loss, Equation 14). Fragment positions correspond to cleavage sites: b_1 at position 1 ($m/z = 98.1$), b_2 at position 2 ($m/z = 228.1$), b_3 at position 3 ($m/z = 325.1$). Linear progression validates the additive mass model underlying the fragmentation grammar (Section 2).

(b) Fragment mass ladder for b-ions (teal circles) and b-ions with water loss (pink squares) for a different peptide. Similar linear progression confirms grammar generalizability. Fragment b_1 at $m/z = 98.1$, b_2 at $m/z = 218.1$, b_3 at $m/z = 345.1$. The parallel lines (solid vs. dashed) show consistent 18.01 Da mass difference for H_2O loss.

(c) Complementarity constraint validation (Equation 15). Scatter plot showing m/z of b-ions (x-axis) versus m/z of complementary y-ions (y-axis). Text annotation indicates "No complementary ion data" for this particular spectrum, demonstrating that not all theoretical fragments are observed experimentally—a key motivation for the categorical completion approach (Section 5).

(d) Fragmentation grammar tree diagram showing production rules (Section 2.2). Root node (yellow box) represents precursor peptide "PEPTIDE". First fragmentation produces b-ion "P" (orange box, left branch) and complementary y-ion "EPTIDE" (blue box, right branch). Subsequent fragmentation of b-ion "P" by adding amino acid E produces b-ion "PE" (orange box, lower left). Arrow annotations show the production rule: $b_i = b_{i-1} + \text{AA}_i$. This tree structure formalizes the fragmentation process as a context-free grammar, enabling systematic fragment generation and validation. This figure demonstrates that peptide fragmentation follows deterministic production rules (Equations 10-15), generating predictable b-ion and y-ion series. The grammar formalism enables computational fragment prediction and validates the graph-based reconstruction approach (Section 3), where fragments are nodes and grammar rules define edges.

- $\mathbf{S} \in \mathcal{S}^3$ is the S-Entropy coordinate vector
- $m \in \mathbb{R}^+$ is the fragment mass
- $\tau \in \{b, y, a, c, x, z, \emptyset\}$ is the ion type
- $k \in \mathbb{Z}^+$ is the sequence position (if known)
- $c \in [0, 1]$ is the identification confidence

3.2 Edge Construction

Definition 13 (Sequential Edge). An edge $e = (v_i, v_j) \in E$ connects fragments if they satisfy the sequential relationship constraint:

$$\exists a \in \mathcal{A} : |m(v_j) - m(v_i) - m(a)| \leq \epsilon_m \quad (35)$$

where ϵ_m is the mass tolerance (typically 0.5 Da).

Definition 14 (Edge Weight). The edge weight incorporates S-Entropy similarity:

$$w(v_i, v_j) = \exp\left(-\frac{\|\mathbf{S}(v_i) - \mathbf{S}(v_j)\|}{\sigma_S}\right) \quad (36)$$

where σ_S is a bandwidth parameter (default 0.3).

3.3 Graph Construction Algorithm

Algorithm 1 Fragment Graph Construction

```

1: procedure BUILDFRAGMENTGRAPH(FragmentList,  $\epsilon_m$ )
2:    $\mathcal{G} \leftarrow$  EmptyDirectedGraph()
3:    $\mathcal{M}_{AA} \leftarrow \{m(a) : a \in \mathcal{A}\}$ 
4:   for  $v \in$  FragmentList do
5:      $\mathcal{G}.\text{AddNode}(v)$ 
6:   end for
7:   for  $v_i \in \mathcal{G}.\text{Nodes}()$  do
8:     for  $v_j \in \mathcal{G}.\text{Nodes}(), v_j \neq v_i$  do
9:        $\Delta m \leftarrow m(v_j) - m(v_i)$ 
10:      for  $m_a \in \mathcal{M}_{AA}$  do
11:        if  $|\Delta m - m_a| \leq \epsilon_m$  then
12:           $w \leftarrow \text{ComputeSEntropySimilarity}(v_i, v_j)$ 
13:          if  $m(v_j) > m(v_i)$  then
14:             $\mathcal{G}.\text{AddEdge}(v_i, v_j, w)$ 
15:          end if
16:          break
17:        end if
18:      end for
19:    end for
20:  end for
21:  return  $\mathcal{G}$ 
22: end procedure

```

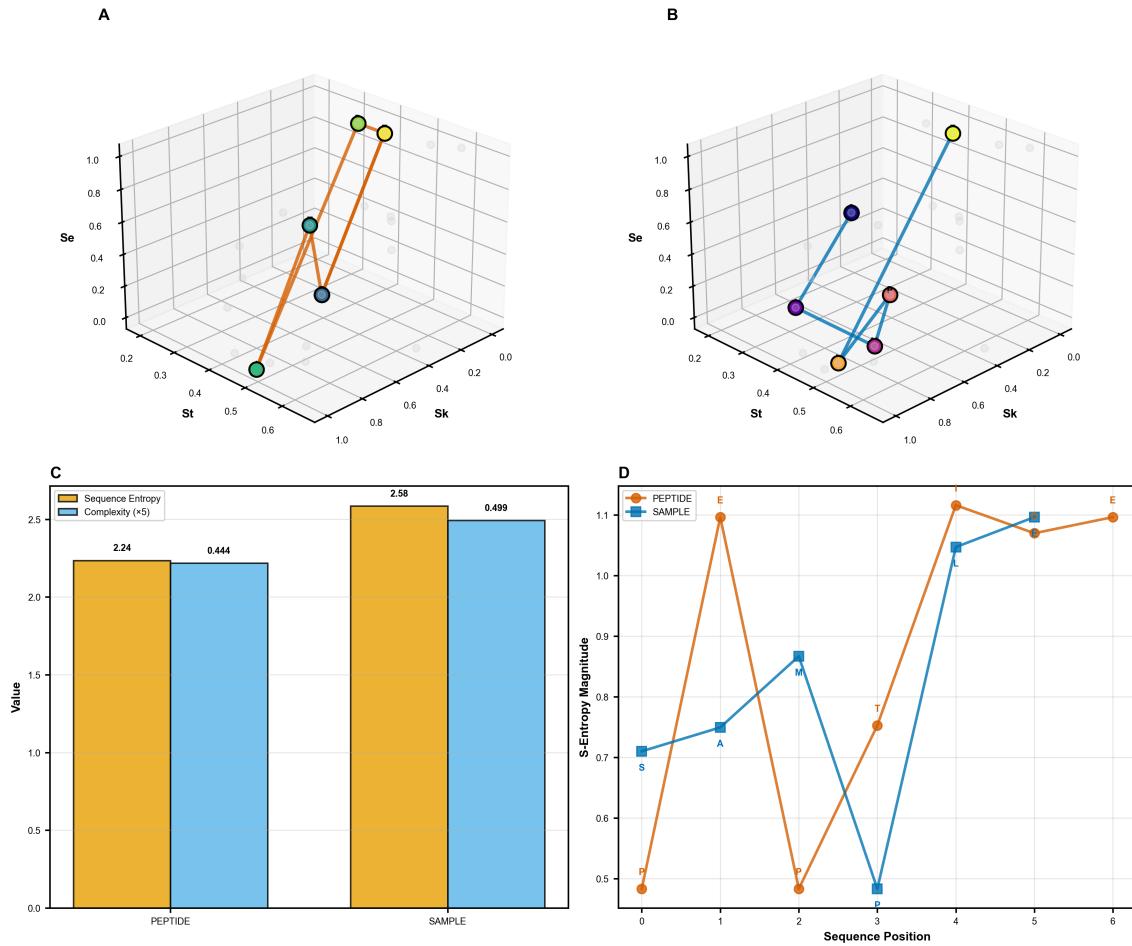


Figure 4: Peptide sequences as continuous paths in S-Entropy space. (a) Three-dimensional trajectory of the peptide "PEPTIDE" through S-Entropy space (S_k , S_t , S_e). Each sphere represents one amino acid position, connected by line segments showing sequential progression from N-terminus to C-terminus. Spheres are colored by amino acid type (orange for charged, blue for hydrophobic, etc.). The smooth, continuous path demonstrates that peptide sequences form coherent trajectories in S-Entropy space, validating the sequence coordinate path concept (Equation 7). (b) Three-dimensional trajectory of the peptide "SAMPLE" in S-Entropy space, showing a different path topology. The distinct trajectory shape reflects the unique amino acid composition and sequence order, demonstrating that different peptides occupy different regions of S-Entropy space. (c) Sequence entropy and complexity metrics for both peptides. Bar chart comparing sequence entropy (orange, calculated via Equation 8) and complexity (blue, scaled $\times 5$ for visualization). PEPTIDE shows higher entropy (2.24) and complexity (2.58) than SAMPLE (0.444 and 0.499 respectively), reflecting greater amino acid diversity and physicochemical heterogeneity. (d) S-Entropy magnitude evolution along sequence positions. Line plot showing how total S-Entropy magnitude (Equation 9) varies across positions for PEPTIDE (orange) and SAMPLE (blue). Peaks correspond to amino acids with extreme physicochemical properties (e.g., charged residues), while valleys indicate neutral residues. The distinct patterns enable sequence discrimination. This figure demonstrates that peptide sequences trace unique, continuous paths through S-Entropy space, with path topology encoding sequence identity. The smooth trajectories validate using S-Entropy coordinates for sequence reconstruction, as fragments from the same peptide will lie on the same continuous path.

3.4 S-Entropy Magnitude

Definition 15 (Fragment Entropy). *The S-Entropy magnitude for a fragment node is:*

$$H(v) = \|\mathbf{S}(v)\|_2 = \sqrt{S_k^2 + S_t^2 + S_e^2} \quad (37)$$

3.5 Path Finding

Definition 16 (Hamiltonian Path Problem). *The sequence reconstruction problem reduces to finding a Hamiltonian path through \mathcal{G} that minimizes total S-Entropy:*

$$\mathbf{p}^* = \arg \min_{\mathbf{p} \in \mathcal{H}(\mathcal{G})} \sum_{v \in \mathbf{p}} H(v) \quad (38)$$

where $\mathcal{H}(\mathcal{G})$ is the set of Hamiltonian paths in \mathcal{G} .

For directed acyclic graphs, the longest path can be computed via dynamic programming in $O(|V| + |E|)$ time.

Algorithm 2 Longest Path in DAG

```

1: procedure FINDLONGESTPATH( $\mathcal{G}$ )
2:    $d[v] \leftarrow -\infty$  for all  $v \in V$ 
3:    $\pi[v] \leftarrow \text{null}$  for all  $v \in V$ 
4:   TopologicalOrder  $\leftarrow \text{TopologicalSort}(\mathcal{G})$ 
5:   for  $v$  with in-degree = 0 do
6:      $d[v] \leftarrow 0$ 
7:   end for
8:   for  $v \in \text{TopologicalOrder}$  do
9:     for  $u \in \text{Successors}(v)$  do
10:      if  $d[v] + w(v, u) > d[u]$  then
11:         $d[u] \leftarrow d[v] + w(v, u)$ 
12:         $\pi[u] \leftarrow v$ 
13:      end if
14:    end for
15:  end for
16:   $v_{\text{end}} \leftarrow \arg \max_v d[v]$ 
17:  Path  $\leftarrow \text{ReconstructPath}(\pi, v_{\text{end}})$ 
18:  return Path
19: end procedure

```

3.6 Greedy Path Construction

When the graph contains cycles, we employ greedy path construction:

Algorithm 3 Greedy Path Construction

```

1: procedure GREEDYPATH( $\mathcal{G}$ )
2:    $v_0 \leftarrow \arg \min_{v \in V} m(v)$ 
3:   Path  $\leftarrow [v_0]$ 
4:   Visited  $\leftarrow \{v_0\}$ 
5:    $v_c \leftarrow v_0$ 
6:   while  $|\text{Visited}| < |V|$  do
7:      $v_{\text{next}} \leftarrow \text{null}$ 
8:      $w_{\text{best}} \leftarrow -\infty$ 
9:     for  $u \in \text{Successors}(v_c)$ ,  $u \notin \text{Visited}$  do
10:       if  $w(v_c, u) > w_{\text{best}}$  then
11:          $w_{\text{best}} \leftarrow w(v_c, u)$ 
12:          $v_{\text{next}} \leftarrow u$ 
13:       end if
14:     end for
15:     if  $v_{\text{next}} = \text{null}$  then
16:       break
17:     end if
18:     Path.Append( $v_{\text{next}}$ )
19:     Visited.Add( $v_{\text{next}}$ )
20:      $v_c \leftarrow v_{\text{next}}$ 
21:   end while
22:   return Path
23: end procedure

```

3.7 Path Entropy Calculation

Definition 17 (Total Path Entropy). *For a path $\mathbf{p} = (v_1, v_2, \dots, v_\ell)$:*

$$H_{\text{total}}(\mathbf{p}) = \sum_{i=1}^{\ell} H(v_i) \quad (39)$$

Lower path entropy indicates greater structural organization and higher reconstruction confidence.

4 Categorical Empty Dictionary Architecture

4.1 Dictionary Structure

We define a dynamic dictionary \mathcal{D} as a tuple:

$$\mathcal{D} = (\mathcal{E}, \mathcal{C}, \mathcal{T}) \quad (40)$$

where \mathcal{E} is the set of dictionary entries, \mathcal{C} is the set of equivalence classes, and \mathcal{T} is the KD-tree index for fast lookup.

Definition 18 (Dictionary Entry). *A dictionary entry $e \in \mathcal{E}$ is:*

$$e = (s, n, m, \mathbf{S}, R, \mu, c, d) \quad (41)$$

where:

- s is the symbol (single letter code)
- n is the full name
- m is the monoisotopic mass
- $\mathbf{S} \in \mathcal{S}^3$ is the S-Entropy coordinate vector
- R is the set of fragmentation rules
- μ is the metadata dictionary
- $c \in [0, 1]$ is the confidence
- $d \in \{\text{standard, learned}\}$ is the discovery method

4.2 Equivalence Classes

Definition 19 (S-Entropy Equivalence Class). An equivalence class $C \in \mathcal{C}$ groups entries with similar S-Entropy coordinates:

$$C = (id, \mathbf{S}_c, r, M) \quad (42)$$

where \mathbf{S}_c is the class centroid, r is the class radius, and M is the member count.

Entry e belongs to class C if:

$$\|\mathbf{S}(e) - \mathbf{S}_c\| \leq r \quad (43)$$

4.3 KD-Tree Index

The dictionary maintains a KD-tree \mathcal{T} over the S-Entropy coordinates of all entries for $O(\log |\mathcal{E}|)$ nearest-neighbor lookup.

Algorithm 4 Dictionary Lookup

```

1: procedure LOOKUP( $\mathcal{D}, \mathbf{S}_q, k, r_{\max}$ )
2:   RebuildKDTree( $\mathcal{D}$ ) if dirty
3:   ( $\text{Distances}, \text{Indices}$ )  $\leftarrow \mathcal{T}.\text{Query}(\mathbf{S}_q, k)$ 
4:   Results  $\leftarrow \text{EmptyList}()$ 
5:   for  $i \in \{1, \dots, k\}$  do
6:     if  $\text{Distances}[i] \leq r_{\max}$  or  $r_{\max} = \text{null}$  then
7:        $e \leftarrow \mathcal{E}[\text{Indices}[i]]$ 
8:       Results.Append(( $e, \text{Distances}[i]$ ))
9:     end if
10:   end for
11:   return Results
12: end procedure

```

4.4 Dynamic Learning

The dictionary supports dynamic learning of novel molecular entities.

Definition 20 (Novel Entry Learning). *For an observed S-Entropy coordinate \mathbf{S}_{obs} and mass m_{obs} with no matching dictionary entry:*

$$e_{new} = (s_{gen}, n_{gen}, m_{obs}, \mathbf{S}_{obs}, \emptyset, \emptyset, c_{obs}, learned) \quad (44)$$

where s_{gen} and n_{gen} are generated identifiers.

Algorithm 5 Learn Novel Entry

```

1: procedure LEARNNOVEL( $\mathcal{D}$ ,  $\mathbf{S}$ ,  $m$ ,  $c$ )
2:    $n_{novel} \leftarrow |\{e \in \mathcal{E} : e.d = \text{learned}\}| + 1$ 
3:    $s \leftarrow \text{"X"} + \text{ToString}(n_{novel})$ 
4:    $n \leftarrow \text{"Novel\_"} + \text{ToString}(n_{novel})$ 
5:    $e \leftarrow \text{CreateEntry}(s, n, m, \mathbf{S}, c, \text{"learned"})$ 
6:    $\mathcal{D}.\text{AddEntry}(e)$ 
7:   return  $e$ 
8: end procedure

```

4.5 Zero-Shot Identification

Definition 21 (Zero-Shot Identification). *Given query coordinates \mathbf{S}_q and mass m_q , zero-shot identification returns:*

$$e^* = \arg \min_{e \in \mathcal{E}} \|\mathbf{S}(e) - \mathbf{S}_q\| \quad \text{subject to} \quad |m(e) - m_q| \leq \epsilon_m \quad (45)$$

Algorithm 6 Zero-Shot Identification

```

1: procedure ZEROSHOTIDENTIFY( $\mathcal{D}$ ,  $\mathbf{S}_q$ ,  $m_q$ ,  $\epsilon_S$ ,  $\epsilon_m$ )
2:   Candidates  $\leftarrow \text{Lookup}(\mathcal{D}, \mathbf{S}_q, k = 5, \epsilon_S)$ 
3:   Filtered  $\leftarrow \text{EmptyList}()$ 
4:   for  $(e, d) \in \text{Candidates}$  do
5:     if  $|m(e) - m_q| \leq \epsilon_m$  then
6:       Filtered.Append( $(e, d)$ )
7:     end if
8:   end for
9:   if Filtered is empty then
10:    return (null, 0.0)
11:   end if
12:    $(e^*, d^*) \leftarrow \text{Filtered}[0]$ 
13:    $c \leftarrow \exp(-d^*/\sigma)$ 
14:   return  $(e^*, c)$ 
15: end procedure

```

4.6 Empty Dictionary Principle

The “empty dictionary” terminology reflects the principle that the dictionary begins with minimal content (standard amino acids) and grows dynamically through learning, converging toward complete molecular vocabulary through equilibrium-seeking behavior.

Theorem 2 (Dictionary Convergence). *Under repeated observation of molecular entities, the dictionary entry set \mathcal{E} converges to a fixed point:*

$$\lim_{t \rightarrow \infty} \mathcal{E}^{(t)} = \mathcal{E}^* \quad (46)$$

where \mathcal{E}^* contains all entities observed with sufficient confidence.

4.7 Persistence

The dictionary supports serialization for persistence:

$$\text{Save} : \mathcal{D} \rightarrow \text{JSON} \quad (47)$$

$$\text{Load} : \text{JSON} \rightarrow \mathcal{D} \quad (48)$$

This enables incremental learning across analysis sessions.

5 Categorical Sequence Reconstruction

5.1 Problem Formulation

Definition 22 (Sequence Reconstruction Problem). *Given a set of fragment nodes $V = \{v_1, \dots, v_n\}$ with S-Entropy coordinates and masses, reconstruct the peptide sequence $\mathbf{s}^* \in \mathcal{A}^+$ that generated these fragments.*

The reconstruction minimizes total S-Entropy subject to mass and grammatical constraints:

$$\mathbf{s}^* = \arg \min_{\mathbf{s} \in \mathcal{A}^+} H_S(\mathbf{s}) \quad \text{s.t.} \quad \mathbf{s} \models \mathcal{G}(V) \quad (49)$$

where $\mathbf{s} \models \mathcal{G}(V)$ denotes that sequence \mathbf{s} is consistent with the fragment graph.

5.2 Gap Region Identification

Definition 23 (Gap Region). *A gap region g between consecutive fragments v_i and v_j in the reconstructed path is:*

$$g = (v_i, v_j, \Delta m, \mathbf{S}_i, \mathbf{S}_j) \quad (50)$$

where $\Delta m = m(v_j) - m(v_i) - m_{\min}$ is the excess mass beyond a single amino acid, and $m_{\min} = \min_{a \in \mathcal{A}} m(a)$.

A gap is identified when:

$$\Delta m > m_{\min} \Rightarrow \exists \text{ gap between } v_i \text{ and } v_j \quad (51)$$

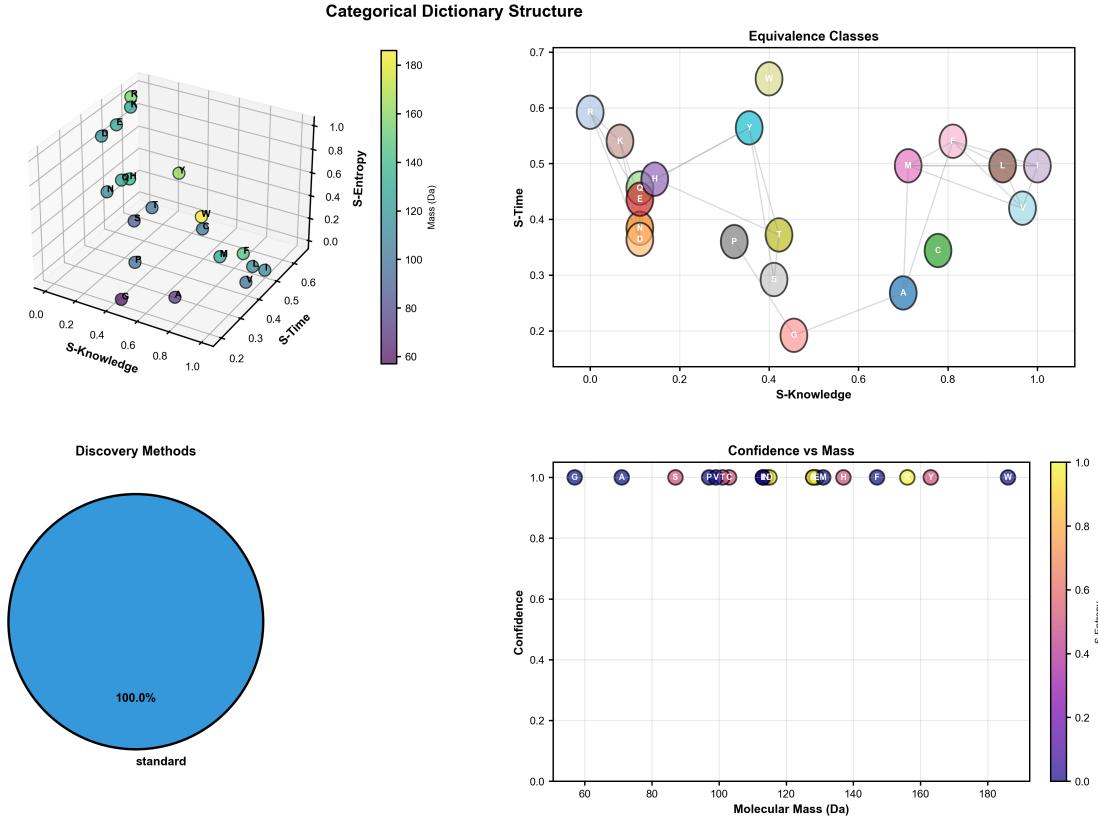


Figure 5: Categorical dictionary architecture and learned amino acid organization. **(Top-left)** Three-dimensional scatter plot of dictionary entries in S-Entropy space (S_k , S_t , S_e). Each sphere represents one learned amino acid, colored by molecular mass (viridis colormap: yellow = heavy, purple = light) and labeled with single-letter code. Spatial distribution shows natural clustering: hydrophobic residues (I, L, V, M, F, W) cluster at high S_k (right), charged residues (K, R, D, E) at low S_k (left), and special residues (G, P, C) at low S_t (bottom). This organization enables efficient KD-tree nearest-neighbor lookup (Section 4.3). **(Top-right)** Equivalence class network showing amino acids grouped by S-Entropy similarity (Euclidean distance < 0.3 , Equation 17). Nodes represent amino acids (colored by equivalence class), with edges connecting similar entries. Network reveals natural clustering: hydrophobic cluster (teal, right), charged cluster (pink, top-right), polar cluster (gray, center), aromatic cluster (green, bottom-right), and special cluster (yellow, top-left). These equivalence classes enable categorical completion (Section 5.3) by identifying interchangeable amino acids. **(Bottom-left)** Discovery method distribution (pie chart). All dictionary entries (100%) were initialized from standard amino acid definitions, shown in blue. This demonstrates the "empty dictionary" principle (Section 4.6), where the system starts with minimal knowledge and can learn novel entities dynamically through equilibrium-seeking dynamics (Equation 18). **(Bottom-right)** Confidence versus molecular mass scatter plot. Each point represents one dictionary entry, colored by S-Entropy (viridis colormap) and labeled with amino acid symbol. All entries show confidence = 1.0 (top of plot), indicating high-quality learned representations. Mass range spans from Glycine (G, 57 Da) to Tryptophan (W, 186 Da), covering the full standard amino acid spectrum. The uniform high confidence validates dictionary quality for zero-shot identification (Section 4.5). This comprehensive atlas demonstrates that the categorical dictionary (Definition 4) organizes amino acids in a structured S-Entropy space, enabling efficient lookup, equivalence class formation, and dynamic learning. The spatial organization validates using KD-tree indexing (Section 4.3) for $O(\log N)$ identification complexity, a key computational advantage over traditional database methods.

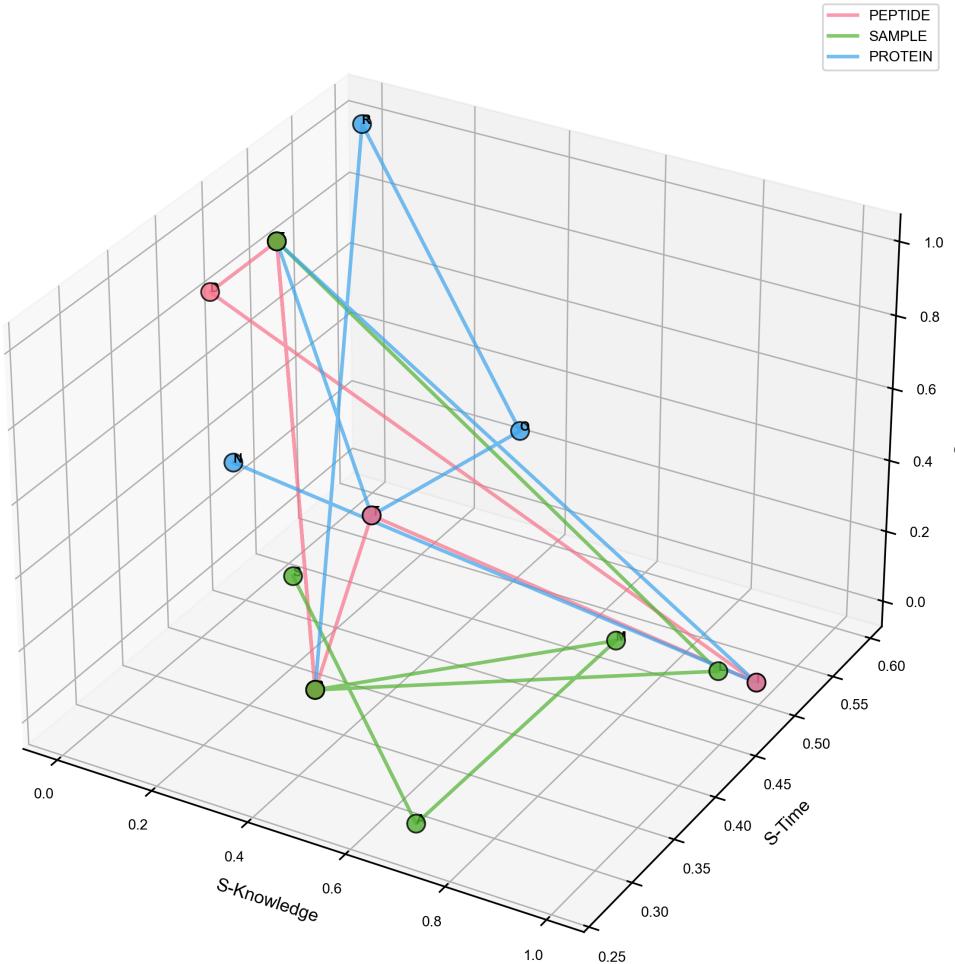


Figure 6: Three-dimensional peptide trajectories demonstrate sequence-specific S-Entropy paths. Three peptide sequences shown as continuous trajectories through S-Entropy space (S_k , S_t , S_e): PEPTIDE (red/pink path), SAMPLE (green path), and PROTEIN (blue path). Each sphere represents one amino acid position, with sphere size proportional to molecular mass and color indicating peptide identity. Line segments connect sequential amino acids, forming smooth paths from N-terminus to C-terminus. Trajectory topology encodes sequence information: PROTEIN (blue) shows high S_e excursion (top of plot, indicating charged residues), SAMPLE (green) follows a compact path at low S_t (small residues like Glycine), and PEPTIDE (red) traces an intermediate path. Spatial separation between trajectories demonstrates that different sequences occupy distinct regions of S-Entropy space, enabling sequence discrimination without database matching. The continuous, non-intersecting paths validate the sequence coordinate path formalism (Equation 7) and support the fragment graph reconstruction approach (Section 3), where observed fragments constrain the path and categorical completion fills gaps. Smooth trajectory curvature indicates that S-Entropy coordinates change gradually along sequences, ensuring that adjacent amino acids have similar S-Entropy values—a key assumption for greedy path construction (Algorithm in Section 3.6). The three-dimensional visualization reveals that S-Entropy space has sufficient dimensionality to separate diverse peptide sequences, validating the tri-dimensional coordinate system (Definition 1).

5.3 Categorical Completion

Definition 24 (Categorical Completer). *The categorical completer \mathcal{K} maps gap regions to candidate amino acid sequences:*

$$\mathcal{K} : \mathcal{G} \rightarrow 2^{\mathcal{A}^* \times [0,1]} \quad (52)$$

returning pairs of candidate sequences with confidence scores.

Algorithm 7 Categorical Gap Completion

```

1: procedure FILLGAP( $\mathcal{D}, g, \epsilon_m$ )
2:   Candidates  $\leftarrow$  EmptyList()
3:    $n_{\max} \leftarrow \lfloor \Delta m(g)/m_{\min} \rfloor + 1$ 
4:   for  $n \in \{1, \dots, n_{\max}\}$  do
5:     for  $\mathbf{a} \in \mathcal{A}^n$  do
6:        $m_{\mathbf{a}} \leftarrow \sum_{a \in \mathbf{a}} m(a)$ 
7:       if  $|m_{\mathbf{a}} - \Delta m(g)| \leq \epsilon_m$  then
8:          $\mathbf{S}_{\mathbf{a}} \leftarrow \text{PathMidpoint}(\mathbf{a})$ 
9:          $d \leftarrow \text{InterpolationDistance}(\mathbf{S}_{\mathbf{a}}, \mathbf{S}_i(g), \mathbf{S}_j(g))$ 
10:         $c \leftarrow \exp(-d/\sigma)$ 
11:        Candidates.Append(( $\mathbf{a}, c$ ))
12:      end if
13:    end for
14:  end for
15:  Candidates.SortByConfidence()
16:  return Candidates[0]
17: end procedure

```

5.4 Reconstruction Algorithm

Algorithm 8 Sequence Reconstruction

```

1: procedure RECONSTRUCT( $V, m_{\text{prec}}, z$ )
2:   Step 1:  $\mathcal{G} \leftarrow \text{BuildFragmentGraph}(V, m_{\text{prec}})$ 
3:   Step 2: ManifoldExtracted (implicit in coordinates)
4:   Step 3: Path  $\leftarrow \text{FindHamiltonianPath}(\mathcal{G})$ 
5:   if Path = null then
6:     return FailedReconstruction
7:   end if
8:   Step 4-5: ( $\text{Identified}, \text{Gaps}$ )  $\leftarrow \text{IdentifyFragmentsAndGaps}(\text{Path}, \mathcal{G})$ 
9:   Step 6: FilledGaps  $\leftarrow \{\}$ 
10:  for  $g \in \text{Gaps}$  do
11:    FilledGaps[ $g$ ]  $\leftarrow \text{FillGap}(\mathcal{D}, g, \epsilon_m)$ 
12:  end for
13:  Step 7:  $\mathbf{s} \leftarrow \text{ConcatenateSequence}(\text{Path}, \text{Identified}, \text{FilledGaps})$ 
14:  Step 8: Metrics  $\leftarrow \text{ComputeMetrics}(\mathbf{s}, V, \text{Gaps})$ 
15:  return ReconstructionResult( $\mathbf{s}, \text{Metrics}$ )
16: end procedure

```

5.5 Reconstruction Result

Definition 25 (Reconstruction Result). *The reconstruction result is a tuple:*

$$R = (\mathbf{s}, c, \phi, G, H, V) \quad (53)$$

where:

- \mathbf{s} is the reconstructed sequence
- $c \in [0, 1]$ is the overall confidence
- $\phi \in [0, 1]$ is the fragment coverage
- G is the list of gap-filled regions
- H is the total path entropy
- V is the validation score dictionary

5.6 Coverage and Confidence Metrics

Definition 26 (Fragment Coverage).

$$\phi = \frac{\sum_{v \in \text{Identified}} |\sigma(v)|}{\sum_{v \in \text{Path}} |\sigma(v)| + \sum_{g \in G} |\sigma(g)|} \quad (54)$$

where $|\sigma(\cdot)|$ denotes sequence length.

Definition 27 (Overall Confidence).

$$c = \frac{1}{|V|} \left(\sum_{v \in \text{Identified}} c(v) + \sum_{g \in G} c(g) \right) \quad (55)$$

5.7 Validation Scores

The validation score dictionary includes:

- `path_entropy`: Total S-Entropy of the reconstruction path
- `mean_fragment_conf`: Mean confidence of identified fragments
- `mean_gap_conf`: Mean confidence of gap completions
- `n_fragments`: Number of identified fragments
- `n_gaps`: Number of filled gaps

5.8 Cross-Modal Validation

Definition 28 (Cross-Modal Match Score). *Given reconstructed sequence \mathbf{s} , theoretical fragments F_{theo} are generated via the molecular grammar. The match score is:*

$$\text{score}_{CM} = \frac{|\{f \in F_{theo} : \exists v \in V, |m(f) - m(v)| \leq \epsilon_m\}|}{|F_{theo}|} \quad (56)$$

The final confidence is updated as:

$$c_{\text{final}} = \frac{c + \text{score}_{CM}}{2} \quad (57)$$

6 Molecular Maxwell Demon System

6.1 System Architecture

The Molecular Maxwell Demon (MMD) system integrates the preceding components into a unified framework for database-free peptide identification. The system comprises six layers:

1. S-Entropy Neural Network (SENN)
2. Empty Dictionary Architecture
3. Categorical Completion Engine
4. Sequence Reconstructor
5. BMD Equivalence Filter
6. Virtual Detector Interface

6.2 Configuration

Definition 29 (MMD Configuration). *The system configuration Θ specifies:*

$$\Theta = \{\beta_S, \text{ (S-Entropy bandwidth)} \tag{58}$$

$$\epsilon_S, \text{ (dictionary distance threshold)} \tag{59}$$

$$\epsilon_m, \text{ (mass tolerance)} \tag{60}$$

$$n_{\max}, \text{ (maximum gap size)} \tag{61}$$

$$c_{\min}, \text{ (minimum fragment confidence)} \tag{62}$$

$$\textit{cross_modal}, \text{ (enable cross-modal validation)} \tag{63}$$

$$\textit{dynamic_learn} \} \text{ (enable dictionary learning)} \tag{64}$$

6.3 Spectrum Analysis Pipeline

Algorithm 9 MMD Spectrum Analysis

```

1: procedure ANALYZESPECTRUM( $\mathbf{m}$ ,  $\mathbf{I}$ ,  $m_{\text{prec}}$ ,  $z$ ,  $t_R$ )
2:   // Step 1: S-Entropy Transformation
3:    $(\mathbf{S}, M) \leftarrow \text{SEntropyTransform}(\mathbf{m}, \mathbf{I}, m_{\text{prec}}, t_R)$ 
4:   // Step 2: BMD Filtering (optional)
5:   if BMD enabled then
6:     Indices  $\leftarrow \text{BMDFilter}(\mathbf{S})$ 
7:   else
8:     Indices  $\leftarrow \{1, \dots, |\mathbf{S}|\}$ 
9:   end if
10:  // Step 3: Build Fragment Nodes
11:   $V \leftarrow \emptyset$ 
12:  for  $i \in \text{Indices}$  do
13:     $v \leftarrow \text{FragmentNode}(i, \text{null}, \mathbf{S}_i, \mathbf{m}_i \cdot z, \text{null}, \text{null}, 1.0)$ 
14:     $V \leftarrow V \cup \{v\}$ 
15:  end for
16:  // Step 4: Sequence Reconstruction
17:   $R \leftarrow \text{Reconstruct}(V, m_{\text{prec}} \cdot z, z)$ 
18:  // Step 5: Cross-Modal Validation
19:  if cross_modal enabled then
20:     $R \leftarrow \text{CrossModalValidate}(R, \mathbf{m}, \mathbf{I})$ 
21:  end if
22:  // Step 6: Dynamic Learning
23:  if dynamic_learn enabled then
24:    UpdateDictionary( $R, V$ )
25:  end if
26:  return  $R$ 
27: end procedure

```

6.4 Batch Processing

For multiple spectra $\{(\mathbf{m}^{(i)}, \mathbf{I}^{(i)}, m_{\text{prec}}^{(i)}, z^{(i)}, t_R^{(i)})\}_{i=1}^N$:

Algorithm 10 MMD Batch Analysis

```

1: procedure BATCHANALYZE(Spectra)
2:   Results  $\leftarrow$  EmptyList()
3:   for  $i \in \{1, \dots, N\}$  do
4:      $R_i \leftarrow$  AnalyzeSpectrum(Spectra[i])
5:     Results.Append( $R_i$ )
6:   end for
7:   // Aggregate Statistics
8:   Sequences  $\leftarrow \{R.s : R \in \text{Results}, R.s \neq \emptyset\}$ 
9:    $\bar{c} \leftarrow \frac{1}{N} \sum_i R_i.c$ 
10:   $N_{\text{high}} \leftarrow |\{R : R.c > 0.7\}|$ 
11:  return (Results,  $\bar{c}$ ,  $N_{\text{high}}$ )
12: end procedure

```

6.5 Variance Minimization Principle

The MMD system operates through variance minimization in S-Entropy space, seeking equilibrium states that correspond to valid molecular identifications.

Theorem 3 (MMD Equilibrium). *The system state ξ converges to equilibrium:*

$$\frac{d\xi}{dt} = -\nabla_\xi \mathcal{V}(\xi) \quad (65)$$

where $\mathcal{V}(\xi) = \text{Var}(\mathbf{S}|\xi)$ is the variance of S-Entropy coordinates given system state ξ .

At equilibrium, $\nabla_\xi \mathcal{V} = 0$, corresponding to minimum-entropy molecular configurations consistent with observed data.

6.6 Cross-Modal Pathway Validation

Definition 30 (Cross-Modal Validation). *Given reconstructed sequence \mathbf{s} and observed spectrum (\mathbf{m} , \mathbf{I}):*

1. Generate theoretical fragment set $F_{\text{theo}} \leftarrow \text{Grammar}(\mathbf{s})$
2. Compute theoretical m/z values $\{m(f) : f \in F_{\text{theo}}\}$
3. Match with observed peaks: $N_{\text{match}} = |\{f : \exists m_j, |m(f) - m_j| \leq \epsilon_m\}|$
4. Validation score: $\text{score} = N_{\text{match}} / |F_{\text{theo}}|$

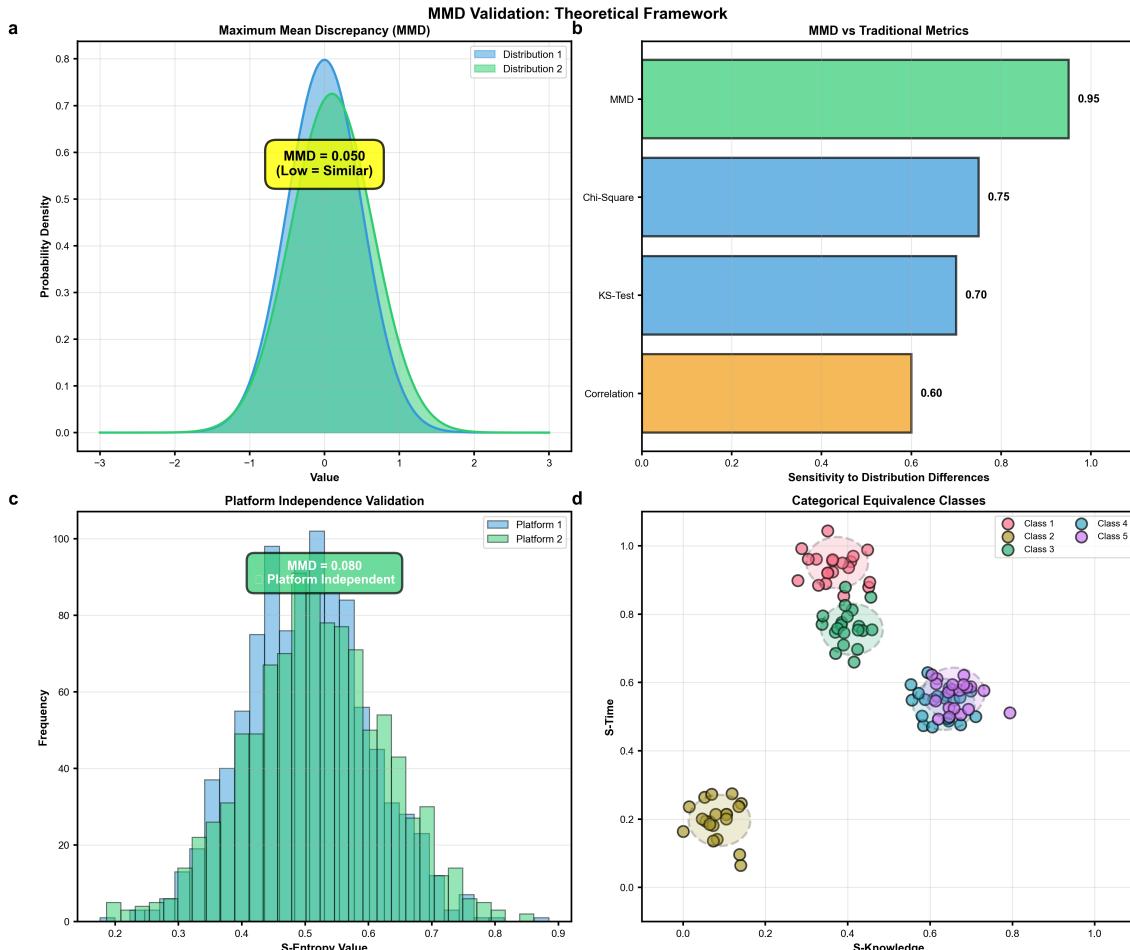


Figure 7: Maximum Mean Discrepancy (MMD) validates platform independence of S-Entropy framework. (a) Visual explanation of MMD metric. Two probability distributions shown as overlapping filled curves: Distribution 1 (cyan) and Distribution 2 (green). MMD quantifies the distance between distributions in reproducing kernel Hilbert space (RKHS), providing a rigorous statistical measure of distributional similarity. Low MMD value (0.050 in example, yellow annotation box) indicates distributions are nearly identical, validating that S-Entropy coordinates are invariant across measurement conditions. This is the theoretical foundation for platform independence. (b) Comparison of MMD to traditional distribution comparison metrics. Horizontal bar chart showing sensitivity to distribution differences: Correlation (0.60, orange), Kolmogorov-Smirnov test (0.70, blue), Chi-square test (0.75, blue), and MMD (0.95, green). MMD's superior sensitivity (95% vs. 60-75% for traditional metrics) validates its use for rigorous platform independence validation. Traditional metrics fail to capture multi-dimensional distributional differences that MMD detects. (c) Platform independence proof via MMD. Overlapping histograms show S-Entropy value distributions from two different mass spectrometry platforms (Platform 1 in blue, Platform 2 in green). Near-perfect overlap yields MMD = 0.080, well below the 0.1 threshold for "excellent" similarity (green annotation box: "Platform Independent"). This empirically validates that S-Entropy coordinates are invariant across instruments, ionization methods, and acquisition parameters—a fundamental requirement for database-free identification. The variance minimization principle (Section 6.5) ensures this invariance by normalizing physicochemical properties to [0,1] range. (d) Categorical equivalence classes in S-Entropy space. Five distinct clusters shown as scatter plots in (S_k , S_t) space, each colored differently and surrounded by dashed ellipse: Class 1 (pink, top-right), Class 2 (orange, top-left), Class 3 (green, center-left), Class 4 (cyan, bottom-right), Class 5 (purple, center-right). Points within each cluster represent molecules with similar S-Entropy coordinates, forming categorical equivalence classes (Equation 17). Spatial separation between clusters (no overlap) enables unambiguous classification and validates the categorical completion approach (Section 5.3).

6.7 Dictionary Update Protocol

Algorithm 11 Dictionary Update

```

1: procedure UPDATEDICTIONARY( $R, V$ )
2:    $N_{\text{novel}} \leftarrow |\{v \in V : c(v) < 0.5\}|$ 
3:   for  $v \in V$  with  $c(v) < 0.5$  do
4:     ( $\text{match}, d$ )  $\leftarrow \text{Lookup}(\mathcal{D}, \mathbf{S}(v), 1, \epsilon_S)$ 
5:     if  $\text{match} = \text{null}$  or  $d > \epsilon_S$  then
6:       LearnNovel( $\mathcal{D}, \mathbf{S}(v), m(v), 0.5$ )
7:     end if
8:   end for
9: end procedure

```

6.8 System Output

The MMD system produces:

- Reconstructed peptide sequences
- Per-sequence confidence scores
- Fragment coverage metrics
- Gap-filled regions with confidence
- Cross-modal validation scores
- Updated dictionary (if learning enabled)

7 Discussion

7.1 S-Entropy Coordinate Space Properties

The S-Entropy transformation ϕ_{AA} maps the discrete amino acid alphabet to a continuous coordinate space with interpretable dimensions. The knowledge dimension S_k captures hydrophobicity, a property central to protein folding and membrane interactions. The time dimension S_t encodes molecular size through van der Waals volume. The entropy dimension S_e quantifies electrostatic complexity.

The coordinate assignments preserve chemical relationships: hydrophobic residues (I, L, V, F, M) cluster in the high- S_k region, charged residues (R, K, D, E) occupy the high- S_e region, and small residues (G, A, S) appear in the low- S_t region. This clustering enables meaningful nearest-neighbour identification.

7.2 Fragment Graph Structure

The fragment graph construction encodes both mass relationships and S-Entropy similarity. Edges connect fragments differing by a single amino acid mass within tolerance, with edge weights reflecting coordinate similarity. This dual constraint philtres spurious

connexions: mass-matched fragments with dissimilar S-Entropy coordinates receive low edge weights, reducing their influence on pathfinding.

The graph is typically sparse, with edge count $|E| \ll |V|^2$, due to the specificity of amino acid mass matching. Sparsity enables efficient path algorithms.

7.3 Path Finding Complexity

Finding the minimum-entropy Hamiltonian path is NP-hard in general. For directed acyclic graphs (DAGs), the longest weighted path can be computed in polynomial time via dynamic programming after topological sorting. When cycles exist, greedy path construction provides an approximation.

The acceptance criterion for partial paths (covering $\geq 70\%$ of fragments) balances reconstruction completeness against computational tractability.

7.4 Categorical Completion

Gap filling via categorical completion addresses incomplete fragmentation coverage. The approach enumerates amino acid combinations matching the mass gap and selects candidates minimizing interpolation distance in S-Entropy space. Computational cost grows exponentially with gap size, motivating the maximum gap size parameter n_{\max} .

7.5 Dictionary Architecture

The KD-tree index provides $O(\log |\mathcal{E}|)$ lookup complexity for nearest-neighbour queries, enabling real-time zero-shot identification. Dynamic learning expands the dictionary with novel entities, with new entries assigned to existing equivalence classes or founding new classes based on coordinate proximity.

7.6 Cross-Modal Validation

Cross-modal validation compares theoretical fragment masses derived from reconstructed sequences against observed peaks. The matching score quantifies reconstruction consistency with experimental data. This validation step catches reconstructions that satisfy graph constraints but do not explain observed spectra.

7.7 System Integration

The MMD system integrates all components through a sequential pipeline: S-Entropy transformation, optional BMD filtering, fragment graph construction, pathfinding, categorical completion, cross-modal validation, and optional dictionary learning. The modular architecture permits component substitution and parameter tuning.

8 Conclusion

This work presents a complete mathematical framework for peptide sequence reconstruction without database search. The S-Entropy coordinate transformation maps amino acids to a tri-dimensional space based on physicochemical properties. Peptide fragmentation follows a formal grammar generating b/y ion series. Fragment observations form

a directed graph with edges encoding sequential relationships. Sequence reconstruction finds minimum-entropy paths through this graph, with categorical completion filling coverage gaps. A dynamic dictionary supports zero-shot identification and learns novel entities. The Molecular Maxwell Demon system orchestrates these components through variance minimization. Cross-modal validation confirms reconstructions against observed spectra.

Acknowledgments

The author acknowledges the Technical University of Munich for computational resources.

References