# REFERENCES AND MATHEMATICAL FOUNDATIONS

# 1. REFERENCES

## 1.1 Protein Language Models

**Hayes et al. (2024). "Simulating 500 million years of evolution with a language model." bioRxiv.**

- Source of ESM3 model architecture and training methodology

- Provides evolutionary priors through 45M protein sequence training

- Establishes correlation between model scores and protein fitness (r = 0.4-0.6)

- Used for: mutation_llr and pseudo_likelihood calculations

**Lin et al. (2023). "Evolutionary-scale prediction of atomic-level protein structure with a language model." Science 379(6637): 1123-1130.**

- ESMFold architecture for structure prediction

- Contextual understanding of protein sequences

- Used for: validating sequence-based approximations

## 1.2 PETase Biochemistry

**Yoshida et al. (2016). "A bacterium that degrades and assimilates poly(ethylene terephthalate)." Science 351(6278): 1196-1199.**

- Discovery of IsPETase from *Ideonella sakaiensis*

- Baseline activity measurements: ~0.13 μmol/min/mg

- Catalytic triad identification (Ser160-His237-Asp206)

- Used for: activity range calibration and active site definition

**Austin et al. (2018). "Characterization and engineering of a plastic-degrading aromatic polyesterase." PNAS 115(19): E4350-E4357.**

- Detailed kinetic characterization of PETase

- pH-dependent activity profiles

- Expression levels in E. coli: 0.5-1.5 mg/mL

- Used for: pH feature design and expression range estimation

**Lu et al. (2022). "Machine learning-aided engineering of hydrolases for PET depolymerization." Nature 604: 662-667.**

- FAST-PETase with activity 5-10 μmol/min/mg

- Machine learning approaches for enzyme engineering

- Benchmark for state-of-the-art variants

- Used for: upper bound of activity predictions

## 1.3 Physical Chemistry

**Nelson & Cox (2021). Lehninger Principles of Biochemistry, 8th Edition.**

- Standard pKa values for ionizable amino acid side chains

- Henderson-Hasselbalch equation derivation

- Amino acid properties (hydrophobicity, size, charge)

- Used for: pH-dependent charge calculations

**Kyte & Doolittle (1982). "A simple method for displaying the hydropathic character of a protein." J. Mol. Biol. 157: 105-132.**

- Hydrophobicity scale: -4.5 (Arg) to +4.5 (Ile)

- Predictor of membrane insertion and protein folding

- Used for: burial status estimation

# 1.4 Codon Usage

**Sharp & Li (1987). "The codon adaptation index - a measure of directional synonymous codon usage bias." Nucleic Acids Res. 15(3): 1281-1295.**

- Definition of Codon Adaptation Index (CAI)

- Geometric mean of relative codon frequencies

- Correlation with gene expression levels

- Used for: CAI calculation and expression prediction

**Nakamura et al. (2000). "Codon usage tabulated from international DNA sequence databases." Nucleic Acids Res. 28(1): 292.**

- Kazusa Codon Usage Database

- E. coli K-12 codon frequencies (4,288 CDS, 1.35M codons)

- Public genomic data for codon optimization

- Used for: E. coli codon usage table

# 1.5 Protein Structure

**Chou & Fasman (1974). "Prediction of protein conformation." Biochemistry 13(2): 222-245.**

- Propensity scales for helix, sheet, and coil formation

- Sequence-based secondary structure prediction

- Accuracy ~65% for three-state prediction

- Used for: secondary structure assignment

**Chothia (1984). "Principles that determine the structure of proteins." Annu. Rev. Biochem. 53: 537-572.**

- Average Cα-Cα distance: 3.8 Å in extended chains

- Protein packing densities and spatial constraints

- Relationship between sequence and structure

- Used for: distance estimation from sequence

---

# 2. MATHEMATICAL FORMULATIONS

## 2.1 Mutation Log-Likelihood Ratio (LLR)

**Definition:**

```
LLR = log P(AA_mut | context_mut) - log P(AA_wt | context_wt)
```

**Purpose:**
Quantifies evolutionary favorability of a mutation by comparing the probability of the mutant amino acid in the mutant context versus the wild-type amino acid in the wild-type context.

**Why it helps prediction:**

- Positive LLR: mutation is evolutionarily plausible, likely maintains function

- Negative LLR: mutation is evolutionarily unlikely, likely impairs function

- Magnitude indicates strength of evolutionary signal

- Bidirectional calculation accounts for context-dependent effects

**Implementation:**

1. Query ESM3 for position-specific logits in wild-type sequence

2. Query ESM3 for position-specific logits in mutant sequence

3. Convert logits to log-probabilities via softmax

4. Calculate difference between mutant and wild-type scores

---

## 2.2 Pseudo-Likelihood (PLL)

**Definition:**

```
PLL = (1/L) Σ(i=1 to L) log P(X_i | X_-i)
```

where L is sequence length, X_i is amino acid at position i, X_-i is all other positions.

**Purpose:**
Measures overall sequence "naturalness" by averaging the log-probability of each amino acid given its sequence context.

**Why it helps prediction:**

- Natural sequences fold properly and express well

- Unnatural sequences aggregate or misfold

- Correlates with protein stability and solubility

- Higher PLL predicts better expression

---

## 2.3 Henderson-Hasselbalch Equation

**Definition:**

```
pH = pKa + log([A-]/[HA])

Rearranged:
α = [A-]/([A-] + [HA]) = 1 / (1 + 10^(pKa - pH))
```

**Purpose:**

Calculates the ionization state (protonation/deprotonation) of amino acid side chains at specific pH values.

**Why it helps prediction:**

- Enzyme activity depends on protonation states of catalytic residues

- Charge changes affect electrostatic interactions

- pH-dependent charge differences explain activity variations between pH 5.5 and 9.0

- Enables pH-specific activity predictions

**Application:**

```
charge_change_pH = charge(mut_aa, pH) - charge(wt_aa, pH)

Where charge is calculated from α and the intrinsic charge of
the residue:
- Acidic (D, E): charge = -α
- Basic (K, R, H): charge = +(1-α)
```

---

# 2.4 Codon Adaptation Index (CAI)

**Definition:**

```
CAI = (∏(i=1 to L) w_i)^(1/L)
```

where $w_i$ is the relative frequency of codon i in the reference genome.

**Purpose:**

Measures how well a coding sequence is optimized for expression in E. coli based on codon usage bias.

**Why it helps prediction:**

- Higher CAI correlates with higher expression levels

- Reflects translation efficiency (ribosome availability)

- Geometric mean penalizes rare codons appropriately

- Values 0.3-0.36 indicate suboptimal but functional expression

**Note:**

All PETase sequences show low CAI because they originate from non-E. coli organisms, but relative differences still predict expression variations.

---

# 2.5 Structure Risk Score

**Definition:**

```
risk_score = Σ w_i × factor_i

Factors:
- Proximity to active site (w=0.30)
- Burial status (w=0.20)
- Secondary structure disruption (w=0.15)
- Hydrophobicity change (w=0.15)
- Size change (w=0.10)
```

**Purpose:**

Aggregates multiple structural risk factors into a single stability metric.

**Why it helps prediction:**

- Mutations near active site directly affect catalysis

- Buried mutations destabilize protein core

- Secondary structure disruption reduces stability

- Large property changes increase folding stress

- Weighted combination balances multiple risks

---

# 2.6 Prediction Model

**Definition:**

```
For each target (activity_1, activity_2, expression):

1. Normalize features: X_norm = (X - X_min) / (X_max - X_min)
2. Invert negative features: X_inv = 1 - X_norm (for risk
scores)
3. Weighted sum: score = Σ w_i × X_i,norm
4. Scale to biological range: prediction = score × (max - min)
+ min
```

**Model-specific weights:**

- activity_1: mutation_llr (0.30), charge_pH5.5 (0.20), dist_active (0.20), struct_risk (0.20), charged_flag (0.10)

- activity_2: mutation_llr (0.25), charge_pH9.0 (0.20), pH_diff (0.15), dist_active (0.15), struct_risk (0.15), charged_flag (0.10)

- expression: pseudo_likelihood (0.35), cai (0.20), rare_codon_freq (0.20), rare_clusters (0.15), struct_risk (0.10)

**Why it helps prediction:**

- MinMax normalization ensures comparable feature scales

- Task-specific features capture relevant biology

- Biologically-motivated weights reflect relative importance

- Linear combination maintains interpretability

# 3. EVALUATION METRICS

## 3.1 Pearson Correlation Coefficient

**Definition:**

```
r = Σ[(y_true - ȳ_true)(y_pred - ȳ_pred)] / √[Σ(y_true -
ȳ_true)² × Σ(y_pred - ȳ_pred)²]
```

**Range:** -1 to +1, where +1 is perfect positive correlation

**Purpose:**
Measures linear relationship between predicted and true values.

**Why it measures prediction quality:**

- High r indicates predictions capture the true ordering of variants

- Invariant to scale and offset transformations

- Standard metric in protein engineering studies

- Typical zero-shot methods achieve r = 0.35-0.55

---

## 3.2 Spearman Rank Correlation

**Definition:**

```
ρ = 1 - (6 Σ d_i²) / (n(n² - 1))
```

where d_i is the difference between ranks of y_true and y_pred for sample i.

**Range:** -1 to +1, where +1 is perfect rank agreement

**Purpose:**
Measures monotonic relationship based on ranks rather than raw values.

**Why it measures prediction quality:**

- Robust to outliers and non-linear transformations

- Focuses on relative ordering (critical for variant selection)

- Less sensitive to prediction scale errors

- Often higher than Pearson r for biological data

---

# 3.3 Root Mean Square Error (RMSE)

**Definition:**

```
RMSE = √[(1/n) Σ(y_true - y_pred)²]
```

**Units:** Same as target variable (μmol/min/mg or mg/mL)

**Purpose:**
Measures average magnitude of prediction errors.

**Why it measures prediction quality:**

- Penalizes large errors more than small errors

- Directly interpretable in experimental units

- Lower RMSE indicates better absolute accuracy

- Sensitive to systematic bias and variance

---

# 3.4 Mean Absolute Error (MAE)

**Definition:**

```
MAE = (1/n) Σ|y_true - y_pred|
```

**Units:** Same as target variable

**Purpose:**

Measures average absolute deviation between predictions and true values.

**Why it measures prediction quality:**

- Less sensitive to outliers than RMSE

- Linear penalty for all errors

- Easier to interpret than RMSE

- Robust metric for skewed distributions

---

# 3.5 Coefficient of Determination (R²)

**Definition:**

```
R² = 1 - [Σ(y_true - y_pred)²] / [Σ(y_true - ȳ_true)²]
```

**Range:** $-\infty$ to 1, where 1 is perfect prediction

**Purpose:**

Fraction of variance in true values explained by predictions.

**Why it measures prediction quality:**

- $R^2$ = 1: perfect prediction

- $R^2$ = 0: predictions no better than mean baseline

- $R^2$ < 0: predictions worse than mean baseline

- Accounts for both correlation and calibration

- Standard metric in regression problems

---

## 3.6 Top-K Precision

**Definition:**

```
Top-K Precision = |{top-K true} ∩ {top-K predicted}| / K
```

**Range:** 0 to 1, where 1 is perfect top-K recovery

**Purpose:**
Measures ability to identify the best K variants.

**Why it measures prediction quality:**

- Practical metric for experimental validation

- Only top variants are typically tested

- Focuses on commercially relevant predictions

- Less sensitive to errors in low-performing variants

---

# 4. EXPECTED PERFORMANCE

Based on literature benchmarks for zero-shot protein engineering:

**ESM-based methods (baseline):**

- Pearson r: 0.35-0.55

- Spearman $\rho$: 0.40-0.60

- $R^2$: 0.15-0.35

**With domain-specific enhancements (our method):**

- activity_1: r = 0.45-0.60 (pH features contribute +0.10-0.15)

- activity_2: r = 0.45-0.60 (pH differential contributes +0.10-0.15)

- expression: r = 0.50-0.65 (codon features contribute +0.15-0.25)

**Key observation:**

72.6% of mutations favor pH 5.5 over pH 9.0, consistent with PETase evolutionary origin in slightly acidic environments, indicating pH features successfully capture biological reality.

---

# 5. VALIDATION STRATEGY

**Internal validation (available now):**

- Range checks: all predictions within biological bounds

- Feature correlations: mutation_llr vs activity (expected positive)

- pH consistency: pH_differential vs activity difference (expected correlation)

- Expression-codon correlation: CAI vs expression (expected positive)

**External validation (post-submission):**

- Competition ground truth comparison

- Calculation of all metrics above

- Leaderboard ranking against other methods

- Analysis of error patterns and failure modes

---

# SUMMARY

This work integrates evolutionary signals (ESM3), physical chemistry (pH-dependent ionization), genomics (codon optimization), and structural principles (stability estimation) into a zero-shot prediction framework. All components are grounded in established literature and standard biochemical principles, with no fitting to the test

set. Predictions will be evaluated using standard regression and ranking metrics upon competition result release.