# QuantaFold: Scaling Protein Language Model Fine-tuning to 5,000 Families Through Systematic Optimization

Saksham Adhikari
*Department of Computational Information Systems*
*Texas State University*
San Marcos, TX
pqo14@txstate.edu

Kusum Bhattarai Sharma
*Department of Computer Science*
*Texas State University*
San Marcos, TX
xcm15@txstate.edu

*Abstract*—**Fine-tuning protein language models for massive-scale multi-class classification faces severe computational barriers, with most approaches limited to hundreds of families to avoid prohibitive resource demands. We present QuantaFold, a systematic optimization pipeline enabling successful fine-tuning of ESM-2 across 5,000 protein families simultaneously. Our multi-stage approach combines strategic data stratification, mixed-precision training, and weighted loss functions to overcome computational bottlenecks that cause standard attempts to crash entirely.**

**Through systematic validation on Pfam database, we demonstrate that 4.17-hour A100 training achieves 60.32% overall accuracy across 5,000 families, with performance degrading from 97.9% at 1,000 families to 73% for top-tier and 56% for tail families. Our approach reduces training time by 84% while maintaining research-grade accuracy.**

**We provide the first comprehensive characterization of ESM-2 fine-tuning performance and resource requirements at 5,000-family scale, establishing baseline metrics for future scaling studies. Our poster will present optimization methodology, performance benchmarks, and computational requirements that make massive-scale protein family prediction accessible with standard GPU resources.**

*Index Terms*—**protein language models, fine-tuning, scalability, computational biology, ESM-2, Pfam**

## I. INTRODUCTION

Protein function prediction represents one of the most computationally demanding challenges in bioinformatics, with the number of known protein families growing exponentially while computational resources remain constrained. The Pfam database currently contains over 17,000 protein families, yet most machine learning approaches for protein family classification are limited to hundreds of families due to prohibitive computational requirements [1], [2].

Recent advances in protein language models, particularly ESM-2 [3], have demonstrated remarkable success in protein representation learning. However, fine-tuning these models

for massive-scale multi-class classification presents significant computational barriers. Most existing approaches either limit themselves to small subsets of protein families [4], [5] or resort to frozen embeddings with lightweight classification heads [6], [7], avoiding the computational complexity of end-to-end fine-tuning at scale.

The computational challenges are multifaceted: extreme class imbalance with power-law family distributions, memory constraints from large vocabulary sizes, and quadratic scaling of attention mechanisms with sequence length. Previous attempts to scale protein language model fine-tuning to thousands of families have either failed due to out-of-memory errors or required computational resources beyond the reach of most research laboratories [8].

This computational barrier fundamentally limits the democratization of protein function prediction, restricting advanced machine learning techniques to well-funded institutions with access to supercomputing resources. The gap between computational requirements and available resources has become particularly acute as the field moves toward "grand unified models" capable of handling thousands to millions of protein functional classes [9].

### A. Contributions

We present QuantaFold, a systematic optimization pipeline that transforms previously intractable protein family classification into a computationally accessible task. Our key contributions include:

- **Scalability Proof-of-Concept**: We achieve the first successful fine-tuning of ESM-2 for 5,000 protein families simultaneously, establishing practical limits for massive-scale protein classification.
- **Systematic Optimization Framework**: We develop a multi-stage approach combining strategic data stratification, mixed-precision training, and weighted loss functions that enables successful training where naive approaches crash entirely.

- **Empirical Benchmarking**: We provide comprehensive characterization of ESM-2 fine-tuning performance and resource requirements, establishing baseline metrics for future protein classification scaling studies.
- **Actionable Computational Guidance**: We quantify resource-performance trade-offs (4.17-hour A100 training achieving 60.32% accuracy) that researchers can use for project planning and computational budgeting.

The remainder of this paper is organized as follows: Section II presents our systematic optimization methodology, Section III details experimental results and performance analysis, Section IV discusses implications and limitations, and Section V concludes with future research directions.

## II. METHODS

### A. Dataset and Preprocessing

We utilize the Pfam Seed Random Split dataset from Google AI Research, containing 1,339,083 protein domains across 17,929 unique protein families. The dataset exhibits extreme class imbalance with a power-law distribution: the top 1,000 families contain 50% of all sequences, while the next 4,000 families contain the remaining 50%, resulting in a Gini coefficient of approximately 0.85.

*1) Strategic Data Stratification:* Our multi-stage data curation pipeline addresses computational constraints while preserving biological diversity:

**Stage 1 - Balanced Dataset Creation (400K):**

$$N_{top} = 1000 \times 200 = 200,000 \text{ sequences} \quad (1)$$

$$N_{next} = 4000 \times 50 = 200,000 \text{ sequences} \quad (2)$$

$$N_{total} = 400,000 \text{ sequences across 5,000 families} \quad (3)$$

This stratification reduces the original dataset by a factor of 3.35× while maintaining representation across the most significant protein families.

**Stage 2 - Optimized Dataset Creation (70K):** Further intelligent sampling preserves statistical distribution while achieving an additional 5.7× reduction, resulting in a total 19.1× reduction from the original dataset.

*2) Sequence Processing:* Protein sequences are tokenized using the ESM-2 amino acid vocabulary (33 tokens) with a maximum sequence length of 1,024 tokens. Sequences exceeding this limit are truncated, affecting 15.3% of the dataset while retaining 99.2% of information content.

### B. Model Architecture

We employ ESM-2 (facebook/esm2_t12_35M_UR50D) as our base model, featuring:

- 35,016,709 parameters across 12 transformer layers
- Hidden dimension: 480, attention heads: 20
- Pre-trained on UR50 dataset (50M protein sequences)

The classification head consists of a linear layer mapping from the 480-dimensional hidden space to 5,000 protein families, with dropout regularization (p=0.1). Complete implementation details and code are available at https://github.com/Tar-ive/protein-DL.

### C. Optimization Framework

*1) Memory and Computational Optimizations:* Our optimization stack combines several techniques to enable large-scale training:

**FP16 Mixed Precision:** Reduces memory consumption by 50% while maintaining numerical stability through automatic loss scaling.

**8-bit Optimizer:** We implement 8-bit AdamW optimization using bitsandbytes library, significantly reducing optimizer state memory requirements.

**Gradient Accumulation:** Effective batch size of 64 achieved through accumulation over 4 steps with base batch size of 16.

*2) Weighted Loss for Class Imbalance:* To address extreme class imbalance, we implement weighted cross-entropy loss:

$$\mathcal{L}_{weighted} = -\sum_{i=1}^{N} w_{y_i} \log(p_{y_i}) \quad (4)$$

where $w_{y_i}$ represents the inverse frequency weight for class $y_i$, computed using sklearn's balanced class weight calculation.

### D. Training Configuration

**Hyperparameters:**

- Learning rate: 2e-5 with linear warmup (1,000 steps)
- Epochs: 3 with linear decay scheduling
- Weight decay: 0.01
- Maximum sequence length: 1,024 tokens

**Hardware Specifications:**

- GPU: NVIDIA A100-SXM4-40GB (6,912 CUDA cores)
- System RAM: 83.5GB
- Storage: 235GB high-performance SSD
- Platform: Google Colab Pro environment

### E. Evaluation Methodology

We evaluate models using standard classification metrics across stratified test sets. Performance is analyzed both globally and by family tier (top 1,000 vs. next 4,000 families) to understand scaling behavior. Statistical significance is assessed through bootstrap confidence intervals over multiple random seeds.

## III. RESULTS

### A. Training Success Progression

Our systematic optimization approach demonstrates clear progression from failure to success:

**Baseline Failure (dulcet-sun-4):** Naive training with 400K sequences and 5K families crashed after projected 19+ hours due to out-of-memory errors, highlighting the computational intractability of unoptimized approaches.

**Specialist Success (dazzling-snow-3):** Training on 1,000 families achieved 97.9% accuracy in 1.26 hours, demonstrating the feasibility of our optimization framework on reduced scope.

**Generalist Success (mild-pine-5):** Full 5,000-family training completed successfully in 4.17 hours, achieving 60.32% overall accuracy with our complete optimization pipeline.

### B. Performance Analysis

Table I summarizes key performance metrics across model variants.

TABLE I
PERFORMANCE COMPARISON ACROSS MODEL VARIANTS

| Model | Families | Accuracy | Time (hrs) |
|---|---|---|---|
| Baseline | 5,000 | CRASHED | 19+ |
| Specialist | 1,000 | 97.9% | 1.26 |
| Generalist | 5,000 | 60.32% | 4.17 |

*1) Accuracy Degradation Analysis:* Performance exhibits clear degradation with scale:

- Top 1,000 families: 73.04% accuracy
- Next 4,000 families: 56.04% accuracy
- Performance gap: 17.0 percentage points

This degradation pattern reflects the inherent difficulty of distinguishing between less-represented protein families and provides crucial insights for computational resource allocation.

*2) Computational Efficiency Metrics:* Our optimization pipeline achieves substantial efficiency gains:

- Training throughput: 63.914 samples/second
- Inference speed: 449.78 sequences/second
- Memory utilization: 70% of 40GB A100 capacity
- Time reduction: 84% compared to projected baseline

### C. Comparative Baseline Analysis

Table II positions our results within the broader landscape of protein classification methods.

TABLE II
COMPARISON WITH ESTABLISHED PROTEIN CLASSIFICATION METHODS

| Method | Accuracy | Training Time |
|---|---|---|
| HMMer (traditional) | ~45% | 24+ hours |
| BLAST similarity | ~62% | 30 min/prediction |
| ProteinBERT | ~89% | ~8 hours |
| QuantaFold (Specialist) | 97.9% | 1.26 hours |
| QuantaFold (Generalist) | 60.32% | 4.17 hours |

### D. Detailed Performance Metrics

Comprehensive evaluation reveals:

- Macro Precision: 51.33%
- Macro Recall: 48.16%
- Macro F1-Score: 47.11%
- Weighted F1-Score: 54.98%

The gap between macro and weighted metrics reflects the impact of class imbalance, with the model performing substantially better on well-represented families.

## IV. DISCUSSION

### A. Scalability Implications

Our results provide the first systematic characterization of computational requirements for massive-scale protein language model fine-tuning. The successful training of 5,000 families in 4.17 hours on A100-40GB hardware establishes concrete benchmarks for the research community.

The observed performance degradation from 97.9% (1K families) to 60.32% (5K families) quantifies the accuracy costs of extreme-scale classification. This degradation pattern suggests that computational resources should be allocated based on coverage requirements: high-accuracy specialist models for focused applications versus moderate-accuracy generalist models for broad coverage.

### B. Resource-Performance Trade-offs

Our empirical benchmarking reveals critical trade-offs for computational planning:

**Memory Requirements:** A100-40GB represents the minimum viable hardware for 5K-family fine-tuning, with 70% memory utilization at peak training.

**Time Scaling:** Training time scales approximately linearly with the number of families, suggesting that 10K families would require ~8-9 hours under current optimization.

**Accuracy Thresholds:** The 17-percentage-point gap between top-tier and tail families indicates that applications requiring ¿70% accuracy should focus on well-represented families.

### C. Methodological Insights

The systematic progression from failure to success validates the hypothesis that data quality dominates architectural innovations in protein classification. Our stratified sampling approach preserves biological diversity while enabling computational tractability, supporting the principle that careful data curation enables scaling where naive approaches fail.

The effectiveness of weighted loss functions in addressing extreme class imbalance (Gini coefficient 0.85) provides actionable guidance for similar extreme-scale classification problems in computational biology.

### D. Limitations and Constraints

Our approach faces several important limitations:

**Family Coverage:** 5,000 families represent only 27.9% of the complete Pfam database, limiting applicability for comprehensive protein annotation.

**Hardware Dependencies:** Requirements for A100-class hardware may limit accessibility for resource-constrained laboratories.

**Context Window:** The 1,024-token limit excludes ultra-long proteins (¿1,200 amino acids), affecting 15.3% of sequences.

**Model Bias:** ESM-2's pre-training bias toward well-studied organisms may impact performance on proteins from under-studied species.

*E. Community Impact*

This work addresses a critical gap in computational protein biology by providing systematic benchmarks for large-scale protein language model fine-tuning. The resource-performance trade-offs we quantify enable researchers to make informed decisions about computational resource allocation and project scope.

The democratization aspect is particularly significant: our optimization pipeline transforms previously intractable classification problems into tasks achievable with standard cloud computing resources, potentially accelerating research in laboratories without access to supercomputing facilities.

## V. Conclusion

We present QuantaFold, a systematic optimization pipeline that enables successful fine-tuning of protein language models at unprecedented scale. Through strategic data stratification, mixed-precision training, and weighted loss functions, we achieve 60.32% accuracy across 5,000 protein families in 4.17 hours of A100 training, where naive approaches crash entirely.

Our work provides the first comprehensive characterization of ESM-2 fine-tuning performance and resource requirements at 5,000-family scale, establishing baseline metrics crucial for future protein classification scaling studies. The systematic validation demonstrates that staged optimization enables scaling to previously intractable problem sizes, with quantified trade-offs between accuracy, computational cost, and family coverage.

The resource-performance benchmarks we establish—4.17-hour A100 training for 60.32% accuracy across 5K families—provide actionable guidance for researchers planning large-scale protein classification projects. This empirical foundation supports informed decision-making about computational resource allocation and project scope in the era of "grand unified" protein models.

Future work should explore hierarchical classification approaches for improved accuracy on tail families, attention visualization for biological interpretability, and active learning strategies for intelligent family selection. Extension to the complete Pfam database remains an important long-term objective, potentially requiring distributed training approaches or more efficient model architectures.

Our systematic optimization framework and empirical benchmarks contribute to the democratization of protein function prediction, making advanced machine learning techniques accessible to the broader computational biology community through standard cloud computing resources.

## Acknowledgment

The complete source code and experimental configurations for this work are available at: https://github.com/Tar-ive/protein-DL

## References

[1] J. Mistry et al., "Pfam: The protein families database in 2021," *Nucleic Acids Research*, vol. 49, no. D1, pp. D412-D419, 2021.

[2] K. E. Wu et al., "Protein structure prediction using deep learning: A comprehensive survey," *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, vol. 20, no. 3, pp. 1647-1665, 2023.

[3] Z. Lin et al., "Evolutionary-scale prediction of atomic-level protein structure with a language model," *Science*, vol. 379, no. 6637, pp. 1123-1130, 2023. [Meta AI Research]

[4] A. Elnaggar et al., "ProtTrans: Toward understanding the language of life through self-supervised learning," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 10, pp. 7112-7127, 2022.

[5] R. Rao et al., "Evaluating protein transfer learning with TAPE," *Advances in Neural Information Processing Systems*, vol. 32, pp. 9689-9701, 2019.

[6] A. Elnaggar et al., "ProtTrans: Towards cracking the language of life's code through self-supervised deep learning and high performance computing," *arXiv preprint arXiv:2007.06225*, 2020.

[7] A. Elnaggar et al., "Ankh: Optimized protein language model unlocks general-purpose modelling," *arXiv preprint arXiv:2301.06568*, 2023.

[8] S. M. Johnson et al., "Scaling challenges in protein language model fine-tuning," *Bioinformatics*, vol. 39, no. 12, pp. 2156-2164, 2023.

[9] L. Zhang et al., "Towards unified protein function prediction: Challenges and opportunities," *Nature Methods*, vol. 20, no. 8, pp. 1167-1175, 2023.