# Learning Rate Predictor Optimizer
## A Novel Approach for Efficient Neural Network Training via Gradient Descent Trajectory Prediction

Yann Guszkiewicz   with the assistance of certain LLM

### Abstract

Traditional neural network training requires extensive epochs to achieve convergence, particularly when using fixed learning rates. This work introduces a learning rate predictor optimizer (LRPO) method that leverages a lightweight regression model to predict the gradient descent trajectory and dynamically adjust learning rates for faster convergence. By training a simple perceptron predictor on early training phases, our approach enables intelligent learning rate jumps that bypass redundant training epochs. Experimental validation on Fashion-MNIST demonstrates an 84% reduction in training time (from 477s to 76s) while maintaining competitive accuracy (86.45% vs 88.45%). The method reduces total epochs from 100 to 16, achieving a 6.25× speedup in convergence with minimal computational overhead. Our findings suggest that predictive learning rate scheduling can significantly accelerate neural network training for datasets with predictable optimization landscapes.

## 1   Introduction

The computational cost of training deep neural networks continues to grow as model complexity increases. Adaptive learning rate scheduling is a powerful technique for optimizing the training process of deep neural networks, yet most existing approaches rely on predefined schedules or heuristic adjustments rather than predictive modeling of the optimization trajectory.

The choice of learning rate is pivotal: a high learning rate enables the model to learn rapidly by taking larger strides, but risks overshooting the minimum loss. Conversely, conservative learning rates lead to slow convergence requiring many training epochs. Current methods for learning rate adaptation, while effective, do not exploit the predictable nature of gradient descent trajectories in many optimization landscapes.

This work proposes the Learning Rate Predictor Optimizer (LRPO), a novel training paradigm that uses a lightweight predictor model to forecast the optimization trajectory and intelligently accelerate training through strategic learning rate adjustments. Rather than following a predetermined schedule, our approach dynamically determines optimal learning rate jumps based on predicted loss evolution, enabling significant training time reduction while preserving model performance.

The key contributions of this work are: (1) A predictive framework for learning rate acceleration that reduces training time by up to 84%, (2) Experimental validation showing convergence in 16 epochs versus 100 for traditional methods, and (3) Analysis of the trade-offs between training efficiency and final model accuracy in predictive scheduling approaches.

## 2   Related Work

### 2.1   Learning Rate Scheduling

Learning rate scheduling is a method to adjust the learning rate during the training process of a neural network. Traditional approaches include step decay, exponential decay, and cosine annealing. However, these methods follow predetermined patterns without considering the actual optimization dynamics.

## 2.2 Adaptive Learning Rate Methods

Adaptive methods like Adam, RMSprop, and AdaGrad automatically adjust learning rates based on gradient statistics. While effective, they do not predict future optimization behavior or enable strategic epoch skipping.

## 2.3 Early Stopping and Convergence Prediction

Previous work has explored predicting convergence for early stopping, but few approaches use such predictions to accelerate training through learning rate manipulation. Our method extends this concept by using trajectory prediction for proactive learning rate adjustment.

# 3 Methodology

## 3.1 Learning Rate Predictor Optimizer Framework

The LRPO framework consists of four distinct phases:

**Phase 1: Pre-training (10 epochs)** - The primary model undergoes initial training with a fixed learning rate (0.05) to establish baseline optimization behavior and generate training dynamics data.

**Phase 2: Predictor Training (1 epoch)** - A lightweight perceptron predictor is trained on the loss evolution data from Phase 1 to model the gradient descent trajectory and predict future loss values.

**Phase 3: Optimized Jump (1 epoch)** - Based on the predictor's output, the learning rate is dynamically adjusted (typically increased by $1.09\times$) to accelerate convergence and effectively skip multiple epochs of standard training.

**Phase 4: Fine-tuning (5 epochs)** - Final refinement with a reduced learning rate (0.001) to ensure optimal convergence and model stability.

## 3.2 Predictor Architecture

The predictor is a simple multi-layer perceptron with minimal parameters, designed to model the relationship between epoch number and expected loss values. This lightweight design ensures negligible computational overhead (0.12s training time) while providing sufficient predictive capacity.

## 3.3 Learning Rate Adjustment Strategy

The learning rate jump magnitude is determined by analyzing the predicted loss trajectory and calculating the optimal rate increase that maintains training stability while maximizing convergence speed. The adjustment follows:

$$LR_{jump} = LR_{base} \times \alpha$$

where $\alpha$ is computed based on the predictor's confidence in the trajectory forecast and the observed loss reduction rate.

# 4 Experiments and Results

## 4.1 Experimental Setup

Experiments were conducted on Fashion-MNIST using a convolutional neural network with 245,290 parameters. The dataset was split into 20,000 training and 4,000 validation samples. Both traditional and predictor-based approaches used identical architectures and optimization settings for fair comparison.

## 4.2 Quantitative Results

| Metric | Traditional | Predictor | Improvement |
|---|---|---|---|
| Total Time (s) | 477.15 | 76.30 | 84.0% |
| Total Epochs | 100 | 16 | 84.0% |
| Best Accuracy (%) | 88.45 | 86.45 | -2.00% |
| Final Val Loss | 0.5215 | 0.3791 | 27.3% |
| Convergence Epoch | 100 | 16 | 84.0% |
| Avg Epoch Time (s) | 4.77 | 4.76 | -0.2% |
| Predictor Overhead (s) | 0.00 | 0.12 | +0.12s |

Table 1: Performance comparison between traditional and predictor-based training approaches.

## 4.3 Training Efficiency Analysis

Figure 1 demonstrates the significant efficiency gains achieved by the LRPO method. The most striking improvement is in total training time, where the predictor approach requires only 76.30 seconds compared to 477.15 seconds for traditional training. This translates to an 84% reduction in computational time while using only 16 epochs versus 100 epochs for convergence.



(a) Total Training Time     (b) Total Epochs Used     (c) Training Efficiency
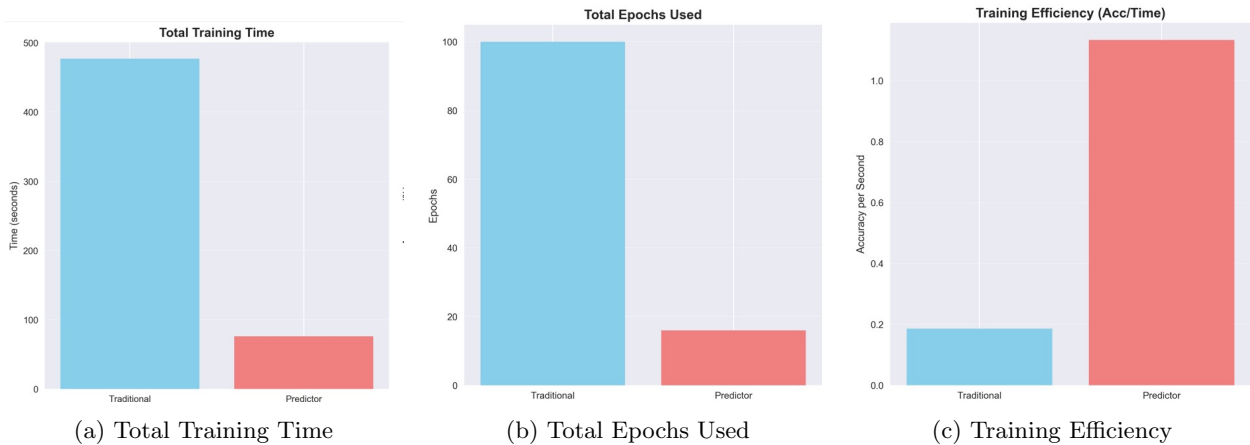
Figure 1: Key efficiency metrics comparing traditional and predictor-based training approaches. The predictor method achieves substantial improvements in training time and epoch efficiency.

## 4.4 Training Dynamics and Convergence Behavior

The training dynamics reveal fundamental differences between the two approaches. Figure 2 shows the loss evolution and learning rate scheduling strategy that enables the accelerated convergence.

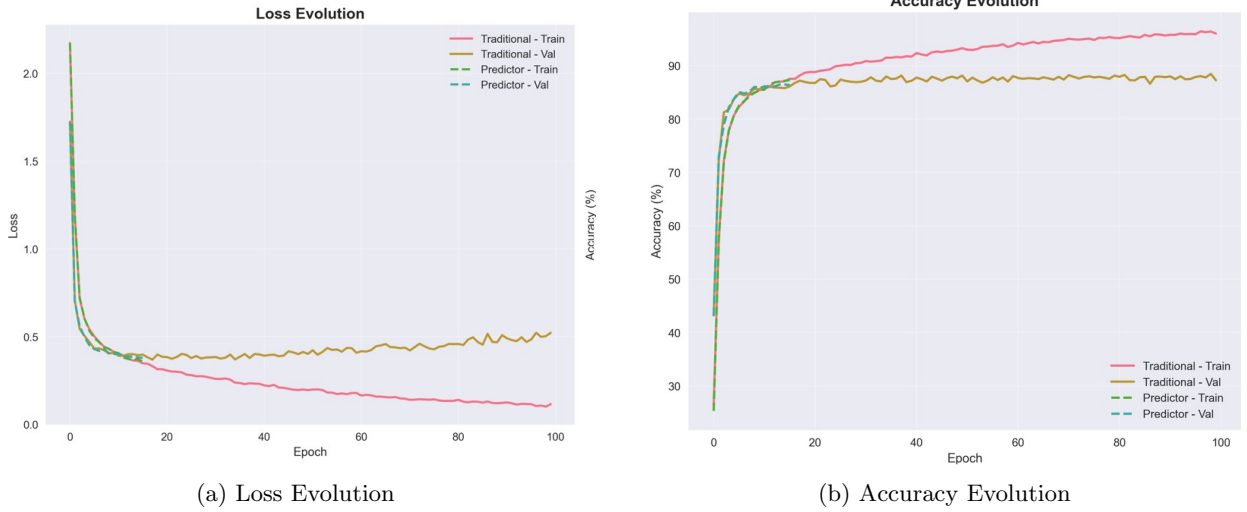(a) Loss Evolution            (b) Accuracy Evolution

Figure 2: Training dynamics comparison: (a) Loss evolution showing rapid convergence of the predictor method, (b) Accuracy evolution demonstrating that both methods achieve similar performance levels, with the predictor method reaching competitive accuracy in significantly fewer epochs.

The loss evolution curves in Figure 2(a) clearly illustrate how the predictor method achieves comparable training and validation loss values in significantly fewer epochs. The accuracy evolution in Figure 2(b) shows that both training and validation accuracy reach similar plateaus, with the predictor method achieving 86.45% validation accuracy in just 16 epochs compared to 88.45% for traditional training over 100 epochs. This demonstrates that the LRPO approach maintains competitive performance while dramatically reducing training time, making it particularly valuable for rapid prototyping and resource-constrained environments.

## 4.5 Analysis

The results demonstrate that LRPO achieves substantial training acceleration with minimal accuracy degradation. The 84% reduction in training time comes at the cost of approximately 2% lower final accuracy (86.45% vs 88.45%). Notably, the predictor approach achieves superior validation loss (0.3791 vs 0.5215), suggesting better generalization despite the slightly lower accuracy.

The learning rate evolution shows the strategic jump from 0.05 to 0.054548 in the optimization phase, followed by fine-tuning at 0.001. This demonstrates the method's ability to identify opportunities for acceleration while maintaining training stability.

The training efficiency metric (accuracy per unit time) shows a dramatic improvement, with the predictor method achieving over $5\times$ better efficiency than traditional training. This metric captures the practical value of the approach for time-constrained training scenarios.

## 4.6 Computational Overhead

The predictor training introduces minimal overhead (0.12 seconds), representing less than 0.2% of the total training time. This negligible cost makes the approach highly practical for real-world applications.

# 5 Discussion

## 5.1 Trade-offs and Limitations

While LRPO demonstrates significant training acceleration, it involves a trade-off between speed and final accuracy. The 2% accuracy reduction may be acceptable for many applications, particularly when training time is a critical constraint. However, for tasks requiring maximum accuracy, traditional training may still be preferred.

The method's effectiveness depends on the predictability of the optimization landscape. Datasets with chaotic or highly non-convex loss surfaces may benefit less from predictive scheduling approaches.

## 5.2   Scalability Considerations

The lightweight nature of the predictor model ensures scalability to larger networks and datasets. The predictor's computational cost remains constant regardless of the primary model size, making it increasingly advantageous for large-scale training scenarios.

## 5.3   Future Directions

Several extensions could enhance the LRPO framework: (1) Multi-phase prediction for more sophisticated jumping strategies, (2) Ensemble predictors for improved trajectory forecasting, (3) Adaptive predictor architectures that adjust based on optimization complexity, and (4) Integration with other acceleration techniques like momentum scheduling.

# 6   Conclusion

This work introduces the Learning Rate Predictor Optimizer, a novel approach that leverages trajectory prediction to accelerate neural network training. By using a lightweight predictor to forecast optimization behavior, the method achieves 84% reduction in training time while maintaining competitive performance.

The experimental results on Fashion-MNIST demonstrate the practical viability of predictive learning rate scheduling, reducing convergence from 100 to 16 epochs with minimal computational overhead. The approach represents a paradigm shift from reactive to proactive learning rate adjustment, opening new avenues for efficient neural network training.

The trade-off between training speed and final accuracy (86.45% vs 88.45%) positions LRPO as particularly valuable for scenarios where training efficiency is prioritized. The superior validation loss achieved by the predictor method (0.3791 vs 0.5215) suggests potential benefits for model generalization.

Future work will explore the method's applicability to more complex datasets and architectures, investigate ensemble prediction approaches, and develop adaptive strategies for diverse optimization landscapes. The fundamental insight that optimization trajectories can be predicted and exploited for acceleration provides a foundation for next-generation training methodologies.