

Language-Guided Tuning: Enhancing Numeric Optimization with Textual Feedback

Yuxing Lu^{1,*,}, Yucheng Hu^{2,#}, Nan Sun³, Xukai Zhao⁴

¹Peking University

²Southeast University

³Huazhong University of Science and Technology

⁴Tsinghua University

[#]These authors contribute equally.

Abstract

Configuration optimization remains a critical bottleneck in machine learning, requiring coordinated tuning across model architecture, training strategy, feature engineering, and hyperparameters. Traditional approaches treat these dimensions independently and lack interpretability, while recent automated methods struggle with dynamic adaptability and semantic reasoning about optimization decisions. We introduce Language-Guided Tuning (LGT), a novel framework that employs multi-agent Large Language Models to intelligently optimize configurations through natural language reasoning. We apply textual gradients - qualitative feedback signals that complement numerical optimization by providing semantic understanding of training dynamics and configuration interdependencies. LGT coordinates three specialized agents: an Advisor that proposes configuration changes, an Evaluator that assesses progress, and an Optimizer that refines the decision-making process, creating a self-improving feedback loop. Through comprehensive evaluation on six diverse datasets, LGT demonstrates substantial improvements over traditional optimization methods, achieving performance gains while maintaining high interpretability.

Code — <https://github.com/YuxingLu613/Language-Guided-Tuning>

Introduction

The success of deep learning has been accompanied by an increasingly complex challenge: the optimization of neural network performance through systematic configuration tuning across multiple interdependent components (Huang, Li, and Yao 2019; Alom et al. 2019). Modern deep learning systems require optimization across four critical dimensions: model architecture, feature engineering, training strategy, and hyperparameters (Deng 2014; Pouyanfar et al. 2018; Goodfellow, Bengio, and Courville 2016). The process of identifying optimal configurations across this multi-dimensional space has emerged as a critical bottleneck in the machine learning development pipeline, often requiring extensive computational resources and domain expertise to navigate the complex interactions between these components (Yu and Zhu 2020). For instance, the choice of data augmentation strategy may require corresponding adjustments to learning rate and model architecture (Shorten

and Khoshgoftaar 2019), while the selection of loss functions can impact the effectiveness of different optimization algorithms (Wang et al. 2022).

Current approaches to configuration optimization fall into several categories, each with significant limitations that hinder their effectiveness in addressing the multi-dimensional nature of the problem (Huang, Li, and Yao 2019). Manual tuning is inherently unscalable and heavily dependent on practitioner expertise, especially when coordinating across multiple configuration dimensions. Grid search and random search methods, though systematic, suffer from the curse of dimensionality and fail to leverage information from previous evaluations, requiring exhaustive exploration of the configuration space (Liashchynskiy and Liashchynskiy 2019). More sophisticated approaches such as Bayesian optimization and evolutionary algorithms offer improvements in efficiency but remain computationally expensive and struggle with the dynamic nature of training processes (Victoria and Maragatham 2021). Furthermore, these methods typically operate in a static framework, unable to adapt configurations during training based on real-time performance feedback, and often treat configuration dimensions as independent variables rather than recognizing their inherent interdependencies. Additionally, the decision-making process in these approaches lacks interpretability and reasoning capabilities, making it difficult to understand and trust the optimization choices made by the system.

Recent advances in large language models (LLMs) have demonstrated remarkable capabilities for optimization tasks, with applications in hyperparameter tuning and configuration selection (Liu, Gao, and Li 2025; Chen et al. 2022; Yang et al. 2023). However, existing LLM-based optimization approaches face critical limitations: they typically operate as single-agent systems without coordinated multi-dimensional reasoning, lack systematic feedback mechanisms for iterative improvement, and fail to provide interpretable optimization trajectories that combine semantic understanding with numerical precision (Yuksekgonul et al. 2024).

To address these limitations, we propose Language-Guided Tuning (LGT), a multi-agent framework that leverages the strength of LLMs and their textual feedback signals that complement numerical loss functions through natural language reasoning. Our system employs three specialized LLM agents: an **Advisor** that proposes configuration

changes, an **Evaluator** that assesses optimization progress, and an **Optimizer** that refines prompts for continuous improvement, creating a self-improving feedback loop for coordinated multi-dimensional optimization. Experimental validation on six diverse datasets shows that LGT achieves up to 23.3% absolute accuracy improvement and 49.3% error reduction while providing high interpretability through natural language reasoning, establishing a new paradigm for intelligent and transparent configuration optimization.

Our key contributions include: (1) Textual gradients as qualitative optimization signals that provide interpretable guidance alongside numerical gradients, (2) Multi-agent system that enables coordinated optimization across all four configuration dimensions while recognizing their interdependencies, (3) Self-improving feedback mechanism that enhances optimization effectiveness over time through prompt refinement, and (4) Comprehensive evaluation demonstrating superior performance and high interpretability compared to established methods across diverse tasks.

Related Works

Our work builds upon two main research areas: configuration tuning and emerging LLM-based optimization. We review these areas to position our contribution within the broader landscape of automated machine learning.

Configuration Tuning Methods

Traditional configuration optimization approaches can be categorized into static and dynamic methods. Static approaches include grid search (Pontes et al. 2016; LaValle, Branicky, and Lindemann 2004) and random search (Bergstra and Bengio 2012; Gharehchopogh et al. 2023), which explore predefined configuration spaces but suffer from the curse of dimensionality and fail to leverage information from previous evaluations. Bayesian optimization emerged as a more sophisticated alternative (Wang et al. 2023; Frazier 2018), using probabilistic surrogate models to guide the search process. However, these methods remain computationally expensive and struggle with the dynamic nature of training processes, often treating configuration dimensions as independent variables.

Dynamic optimization methods, including Neural Architecture Search (NAS) (Liashchynskyi and Liashchynskyi 2019; Ren et al. 2021) and Automated Machine Learning (AutoML) (Karmaker et al. 2021; Zöller and Huber 2021) frameworks, aim to adapt configurations during training. While NAS has produced state-of-the-art architectures, it requires much computational resources and focuses primarily on architectural choices rather than comprehensive configuration optimization. AutoML systems automate the entire machine learning workflow but often lack interpretability and struggle with the complex interdependencies between different configuration dimensions. Furthermore, these approaches typically operate in a static framework, unable to leverage real-time training feedback for dynamic adaptation.

LLM-Based Optimization and Textual Gradients

Recent advances in Large Language Models have opened new possibilities for intelligent optimization through natu-

ral language reasoning (Zhao et al. 2023; Wang et al. 2025; Yu and Liu 2024; Zhang et al. 2023; Lu and Wang 2025). LLMs have demonstrated remarkable capabilities in understanding complex problems, reasoning about multi-step processes, and providing actionable advice across diverse domains (Yuksekgonul et al. 2024; Lu et al. 2025). This capability extends naturally to configuration optimization, where LLMs can leverage their extensive knowledge of machine learning principles and training dynamics (Liu, Gao, and Li 2025; Yang et al. 2023; Ye et al. 2024).

The concept of using language models for optimization has been explored in various contexts, including code generation, mathematical reasoning, and decision-making processes (Liu, Gao, and Li 2025; Yang et al. 2023; Huang et al. 2025; Nie et al. 2024; Ma et al. 2024). However, the application of LLMs to configuration optimization introduces a novel paradigm where semantic reasoning guides numerical optimization. The interaction between different LLM agents can create what terms "textual gradients" (Yuksekgonul et al. 2024)—qualitative feedback signals that complement traditional numerical loss functions.

Our work differs from existing LLM-based approaches by introducing a multi-agent system specifically designed for configuration optimization, where the textual gradient concept enables a self-improving feedback loop. LGT combines the interpretability of natural language reasoning with the precision of mathematical optimization, creating a new paradigm for intelligent and adaptive configuration tuning.

Methods

Problem Formulation and Framework Overview

Figure 1 presents the overall architecture of our Language-Guided Tuning (LGT) framework, illustrating the integration of textual and numerical gradient flows for comprehensive configuration optimization.

Configuration optimization. We formulate the configuration optimization problem as follows. Let $\mathcal{C} = \mathcal{C}_A \times \mathcal{C}_F \times \mathcal{C}_T \times \mathcal{C}_H$ be the multi-dimensional configuration space, where \mathcal{C}_A , \mathcal{C}_F , \mathcal{C}_T , and \mathcal{C}_H represent the spaces of model architecture, feature engineering, training strategy, and hyperparameters respectively. Each configuration $c \in \mathcal{C}$ parameterizes a model f_c that maps inputs $x \in \mathcal{X}$ to outputs $y \in \mathcal{Y}$.

The objective is to find the optimal configuration $c^* \in \mathcal{C}$ that minimizes the expected loss over the data distribution:

$$c^* = \arg \min_{c \in \mathcal{C}} \mathcal{L}(c) = \arg \min_{c \in \mathcal{C}} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(f_c(x), y)] \quad (1)$$

where ℓ is the task-specific loss function and \mathcal{D} is the data distribution. The key challenge lies in the complexity and interdependencies between configuration dimensions, making \mathcal{C} non-separable and requiring coordinated optimization across all dimensions.

Multi-Agent system. Our framework employs a multi-agent system $\mathcal{A} = \{A_{\text{adv}}, A_{\text{eval}}, A_{\text{opt}}\}$ consisting of three specialized agents: an Advisor (A_{adv}) that generates configuration modifications, an Evaluator (A_{eval}) that assesses optimization effectiveness, and a Prompt Optimizer (A_{opt}) that refines agent decision-making capabilities.

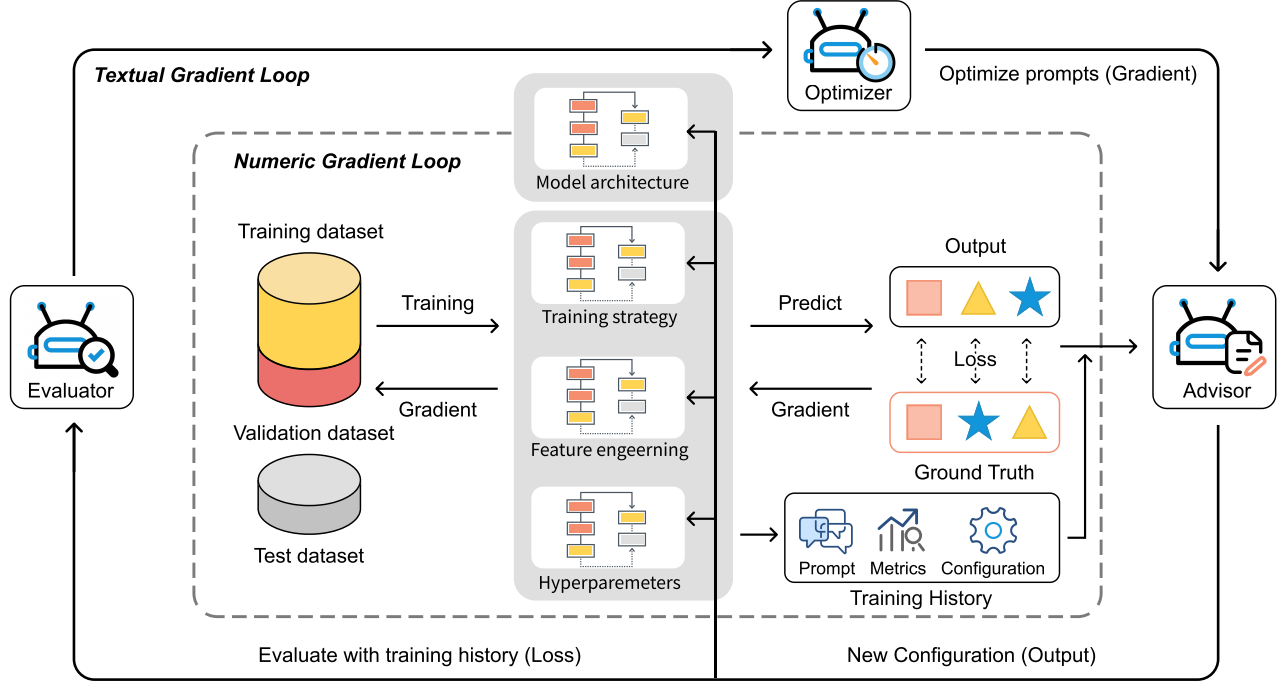


Figure 1: Overall framework of Language-Guided Tuning (LGT). The system operates through dual loops: an inner numeric gradient loop for traditional training, and an outer textual gradient loop where three LLM agents (Advisor, Evaluator, Optimizer) coordinate to optimize configurations across four dimensions (model architecture, training strategy, feature engineering, hyperparameters). Training history feeds back to agents for dynamic adaptation and self-improving optimization.

Each agent A_i operates on a state space \mathcal{S}_i with an action space \mathcal{A}_i , following an interaction protocol defined by the feedback function $\mathcal{F} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}'$ that governs state transitions. This creates a closed-loop system where agents communicate through structured pipelines, enabling coordinated optimization while maintaining interpretability through natural language reasoning.

Configuration Components

Model Architecture Optimization The architecture space \mathcal{C}_A encompasses all possible neural network configurations, including layer types, connectivity patterns, and structural parameters. The architecture optimization function $A_{\text{arch}} : \mathcal{M}_{\text{hist}} \rightarrow \mathcal{C}_A$ analyzes complete training history $\mathcal{M}_{\text{hist}}$ of a whole training iteration to generate improved architectures.

This optimization occurs independently before training cycles of other configurations. The generated architectures must satisfy structural constraints $\theta_a \in \mathcal{F}_a$, where \mathcal{F}_a defines the space of feasible architectures based on computational and architectural requirements.

Feature Engineering and Data Augmentation The feature engineering space $\mathcal{C}_F = \{\tau : \mathcal{X} \rightarrow \mathcal{X}' | \tau \in \mathcal{T}\}$ consists of all possible data transformation functions, where \mathcal{T} represents the set of available augmentation methods including duplication, rotation, shift, flip, scale, noise, contrast, and domain-specific transformations. The dynamic selection function $A_{\text{aug}} : \mathcal{M}_t \times \mathcal{C}_F \rightarrow \mathcal{C}_F$ chooses augmentation

strategies based on current training metrics \mathcal{M}_t , implementing an adaptive strategy:

$$\tau_{t+1} = \arg \max_{\tau \in \mathcal{T}} \mathbb{E}[\mathcal{L}(\tau(x), y)] \quad (2)$$

This enables the system to respond to overfitting/underfitting patterns by selecting appropriate augmentation methods that maximize expected performance improvement.

Training Strategy Optimization The training strategy space $\mathcal{C}_T = \mathcal{L}_{\text{loss}} \times \mathcal{O}_{\text{opt}} \times \mathcal{S}_{\text{sched}}$ encompasses loss functions, optimizers, and schedulers. The dynamic selection function $A_{\text{strat}} : \mathcal{M}_t \times \mathcal{C}_T \rightarrow \mathcal{C}_T$ adapts training strategies based on current performance metrics \mathcal{M}_t . The optimal strategy selection follows:

$$(l^*, o^*, s^*) = \arg \min_{(l, o, s) \in \mathcal{C}_T} \mathcal{L}_{\text{val}}(l, o, s) \quad (3)$$

where l , o , and s represent loss function, optimizer, and scheduler respectively, and \mathcal{L}_{val} is the validation loss.

Hyperparameter tuning. The hyperparameter space $\mathcal{C}_H = \{\lambda \in \mathbb{R}^d | \lambda_i \in [\lambda_i^{\min}, \lambda_i^{\max}]\}$ includes learning rates, class weights, weight decay, and other fine-grained configuration parameters. The epoch-level update rule follows:

$$\lambda_{t+1} = \lambda_t + \Delta \lambda_t \quad (4)$$

where $\Delta \lambda_t = A_{\text{adv}}(\mathcal{M}_t, \lambda_t)$ represents the advisor's suggested parameter modifications. All updates must satisfy constraint $\lambda_{t+1} \in \mathcal{C}_H$ to ensure parameter validity.

Multi-Agent System and Prompt Optimization

Advisor (Output): configuration generator. The Advisor processes current training state through the input processing function $\phi_{\text{adv}} : (\mathcal{M}_t, c_t, \mathcal{C}) \rightarrow \mathcal{S}_{\text{adv}}$, which transforms training metrics, current configuration, and configuration space into an internal state representation. And the decision function $A_{\text{adv}} : \mathcal{S}_{\text{adv}} \rightarrow \Delta c_t$ generates configuration modifications Δc_t based on the processed state. This process involves natural language reasoning:

$$r_t^{\text{adv}} = \mathcal{G}_{\text{adv}}(\text{prompt}_{\text{adv}}(\mathcal{S}_{\text{adv}})) \quad (5)$$

where \mathcal{G}_{adv} is the LLM generation function and $\text{prompt}_{\text{adv}}$ constructs structured prompts from the agent state. The final configuration changes are extracted through parsing: $\Delta c_t = \text{parse}(r_t^{\text{adv}})$.

Evaluator (Loss): performance assessor. The Evaluator processes current training state through the input processing function $\phi_{\text{eval}} : (\mathcal{M}_t, c_t, \mathcal{M}_{\text{base}}, c_{\text{base}}) \rightarrow \mathcal{S}_{\text{eval}}$, which transforms current and baseline metrics and configurations into an internal state representation. The decision function $A_{\text{eval}} : \mathcal{S}_{\text{eval}} \rightarrow \{0, 1\}$ generates a binary success indicator. This process involves natural language reasoning:

$$r_t^{\text{eval}} = \mathcal{G}_{\text{eval}}(\text{prompt}_{\text{eval}}(\mathcal{S}_{\text{eval}})) \quad (6)$$

where $\mathcal{G}_{\text{eval}}$ is the LLM generation function and $\text{prompt}_{\text{eval}}$ constructs structured prompts from the agent state. The final success assessment is extracted through parsing: $\text{success}_t = \text{parse}(r_t^{\text{eval}})$.

Optimizer (Gradient): prompt updater. The Prompt Optimizer processes optimization history through the input processing function $\phi_{\text{opt}} : (\mathcal{H}_t, \mathcal{M}_t) \rightarrow \mathcal{S}_{\text{opt}}$, which transforms historical data and current metrics into prompt improvement signals. The prompt update function $A_{\text{opt}} : \mathcal{S}_{\text{opt}} \rightarrow \Delta \text{prompt}_t$ generates prompt modifications based on the processed state. This process involves natural language reasoning:

$$r_t^{\text{opt}} = \mathcal{G}_{\text{opt}}(\text{prompt}_{\text{opt}}(\mathcal{S}_{\text{opt}})) \quad (7)$$

where \mathcal{G}_{opt} is the LLM generation function and $\text{prompt}_{\text{opt}}$ constructs structured prompts for prompt optimization. The final prompt updates are extracted through parsing: $\text{prompt}_{\text{adv}}^{t+1} = \text{prompt}_{\text{adv}}^t + \text{parse}(r_t^{\text{opt}})$.

Self-Improving Feedback Loop The three-agent collaboration creates a self-improving, meta-level feedback loop that operates across all configuration dimensions. The system maintains an optimization history \mathcal{H}_t that accumulates comprehensive information from each epoch, providing the necessary context for leveraging interdependencies between configuration components that traditional methods treat independently:

$$\mathcal{H}_t = \mathcal{H}_{t-1} \cup \{(c_t, \mathcal{M}_t, \text{prompt}_{\text{adv}}^t)\} \quad (8)$$

The feedback loop among Advisor, Evaluator and Optimizer ensures the training process becomes more adaptive and targeted by enabling agents to refer to both training results and configurations for making appropriate adjustments.

Epoch-Level Training Algorithm The complete training algorithm operates at epoch boundaries:

Algorithm 1: Epoch-Level Language-Guided Tuning

```

1: for  $t = 1, 2, \dots, T$  do
2:    $\mathcal{M}_t = \text{train\_epoch}(f_{c_t}, \mathcal{D})$  {Standard training}
3:    $\Delta c_t = A_{\text{adv}}(\mathcal{M}_t, c_t)$  {Configuration changes}
4:    $c_{t+1} = c_t + \Delta c_t$  {Apply updates}
5:    $\text{success}_t = A_{\text{eval}}(\mathcal{M}_t, c_t)$  {Evaluate changes}
6:    $\text{prompt}_{\text{adv}}^{t+1} = A_{\text{opt}}(\mathcal{H}_t, \mathcal{M}_t)$  {Update prompt}
7:    $\mathcal{H}_t = \mathcal{H}_{t-1} \cup \{(c_t, \mathcal{M}_t, \text{prompt}_{\text{adv}}^t)\}$  {History}
8: end for

```

Theoretical Analysis

Convergence Analysis The convergence of the LGT framework depends on the stability of the three-agent feedback loop and the boundedness of configuration updates. Under the assumption that the configuration space \mathcal{C} is compact and the loss function \mathcal{L} is Lipschitz continuous, the system converges to a local optimum:

$$\lim_{t \rightarrow \infty} \mathbb{E}[\mathcal{L}(c_t)] = \mathcal{L}^* + \delta \quad (9)$$

where \mathcal{L}^* is the optimal loss and $\delta \geq 0$ represents the approximation gap due to the discrete nature of LLM-based optimization. The stability of the feedback loop is ensured by bounded configuration changes:

$$\|\Delta c_t\| \leq \epsilon \quad \forall t > T_0 \quad (10)$$

and the boundedness of prompt updates:

$$\|\Delta \text{prompt}_t\| \leq \epsilon_p \quad \forall t > T_0 \quad (11)$$

The convergence rate depends on the learning dynamics of the prompt optimization process:

$$\mathcal{L}(c_t) - \mathcal{L}^* \leq O(1/\sqrt{t}) + O(\epsilon_p \cdot t) \quad (12)$$

Multi-Agent System Analysis The three-agent system exhibits emergent properties through their interaction. The stability of the feedback loop requires:

$$\|A_{\text{adv}} \circ A_{\text{eval}} \circ A_{\text{opt}}\| \leq 1 \quad (13)$$

where the composition represents the complete feedback cycle. The meta-level learning process converges when:

$$\lim_{t \rightarrow \infty} \|\text{prompt}_{\text{adv}}^{t+1} - \text{prompt}_{\text{adv}}^t\| = 0 \quad (14)$$

Computational Complexity The time complexity of the complete system is:

$$O(T \cdot (|\mathcal{A}| \cdot \text{cost}_{\mathcal{G}} + \text{cost}_{\text{train}} + \text{cost}_{\text{prompt}})) \quad (15)$$

where $\text{cost}_{\mathcal{G}}$ is the language generation cost, $\text{cost}_{\text{train}}$ is the training cost per epoch, and $\text{cost}_{\text{prompt}}$ is the prompt optimization overhead. The space complexity is:

$$O(|\mathcal{H}_T| + |\mathcal{C}| + |\mathcal{P}|) \quad (16)$$

where $|\mathcal{P}|$ represents the space required for storing agent prompts. Compared to other methods, LGT provides:

- **Adaptive complexity:** $O(1)$ configuration updates per epoch vs. $O(|\mathcal{C}|)$ for grid search
- **Contextual optimization:** Leverages interdependencies vs. independent dimension optimization
- **Interpretable decisions:** Natural language reasoning vs. black-box optimization

Experiments

Datasets

We evaluate LGT on six diverse datasets. The datasets range from small-scale classic problems (Iris (Goodfellow, Bengio, and Courville 2016), Water Potability (Simeonov et al. 2001), House Price (Harrison Jr and Rubinfeld 1978), Wine Quality (Cortez et al. 2009)) to large-scale deep learning benchmarks (MNIST (Deng 2012), CIFAR-10 (Krizhevsky, Hinton et al. 2009)), and include both classification and regression tasks to demonstrate the framework’s versatility across different domains and data scales.

Table 1: Dataset characteristics

Dataset	Samples	Features	Task Type	Classes
MNIST	70,000	784	Classification	10
CIFAR-10	60,000	3,072	Classification	10
Iris	150	4	Classification	3
Water Potability	3,276	9	Classification	2
House Price	50000	5	Regression	-
Wine Quality	1599	11	Regression	-

Baselines

We compare LGT against five baseline methods: **No Tuning** uses default configurations without optimization, employing standard architectures and hyperparameters. **Random Search** randomly samples configurations from predefined spaces (Bergstra and Bengio 2012). **Grid Search** systematically explores discretized configuration spaces through exhaustive enumeration (Liashchynskiy and Liashchynskiy 2019). **Neural Architecture Search (NAS)** optimizes model structure using DARTS (Liu, Simonyan, and Yang 2018). **Bayesian Optimization** uses Gaussian process surrogate models with expected improvement acquisition (Frazier 2018).

Experimental Setup

LLM Backbone: LLM agents interact via DeepSeek API (Liu et al. 2024) with temperature 0.2. All datasets use 80%/20% train/test splits with standard preprocessing.

Search Spaces: Architecture optimization covers 2-5 layers with 32-512 neurons. Feature engineering includes rotation, scaling, noise injection, and task-specific augmentations. Training strategies span 3 optimizers (SGD, Adam, AdamW). Hyperparameters include class weights (0.1-10.0) and learning rates from 10^{-4} to 10^{-1} .

Evaluation: Each method evaluates maximum 50 configurations with 10 epochs per evaluation. All experiments conducted over 10 independent runs with different random seeds (42-51) for statistical rigor. Metrics include accuracy, F1-score, AUC for classification; MAE, MSE, R^2 for regression.

Results

Overall Performance Comparison

Table 2 presents the comprehensive performance comparison between LGT and established optimization methods

Table 2: Classification performance comparison (mean \pm std over 10 independent runs). LGT achieves consistent improvements across diverse datasets. Best results in **bold**.

Dataset	Method	Acc.	F1	AUC
MNIST	No Tuning	78.41 _{2.1}	77.88 _{2.3}	97.55 _{0.4}
	Random Search	89.23 _{1.8}	88.67 _{1.9}	97.89 _{0.3}
	Grid Search	87.15 _{2.0}	86.92 _{2.1}	98.02 _{0.2}
	NAS	96.34 _{1.1}	96.12 _{1.2}	98.67 _{0.2}
	Bayesian Opt.	93.67 _{1.6}	93.45 _{1.7}	98.45 _{0.3}
	LGT	98.99_{0.8}	98.99_{0.8}	99.99_{0.1}
CIFAR-10	No Tuning	49.01 _{2.4}	49.01 _{2.5}	88.60 _{1.2}
	Random Search	51.78 _{3.1}	51.45 _{3.2}	89.12 _{1.4}
	Grid Search	58.34 _{2.8}	58.12 _{2.9}	91.23 _{1.1}
	NAS	55.89 _{3.4}	55.67 _{3.5}	90.45 _{1.3}
	Bayesian Opt.	62.45 _{2.1}	62.23 _{2.2}	92.95 _{0.9}
	LGT	69.64_{2.0}	69.61_{2.1}	95.42_{0.7}
Water Potability	No Tuning	65.17 _{2.8}	60.94 _{3.1}	66.33 _{1.9}
	Random Search	66.34 _{3.2}	62.15 _{3.4}	67.12 _{2.1}
	Grid Search	65.89 _{2.9}	61.67 _{3.2}	66.78 _{2.0}
	NAS	66.12 _{3.1}	61.89 _{3.3}	66.89 _{2.2}
	Bayesian Opt.	66.67 _{2.7}	62.45 _{2.9}	67.23 _{1.8}
	LGT	66.92_{2.5}	63.52_{2.7}	67.89_{1.6}
Iris	No Tuning	73.33 _{4.2}	80.95 _{3.8}	86.60 _{2.1}
	Random Search	81.33 _{3.7}	86.45 _{3.2}	90.78 _{1.8}
	Grid Search	78.67 _{4.1}	84.23 _{3.9}	89.45 _{2.0}
	NAS	89.33 _{2.8}	92.67 _{2.1}	95.12 _{1.2}
	Bayesian Opt.	86.67 _{3.2}	90.45 _{2.7}	93.78 _{1.5}
	LGT	96.67_{1.8}	96.67_{1.7}	98.74_{0.8}

across all datasets. Results show mean \pm standard deviation over 10 independent runs, demonstrating LGT’s consistent superiority with low variance. LGT achieves substantial improvements across diverse tasks, consistently outperforming all baselines with statistically significant margins, validating the effectiveness of our textual gradient approach.

The experimental evaluation reveals LGT’s superior optimization capabilities across diverse datasets. While traditional methods like Random Search and Grid Search struggle with consistency, and even sophisticated approaches like Neural Architecture Search and Bayesian Optimization show domain-specific limitations, LGT achieves remarkable improvements of up to 23.3% absolute accuracy on classification tasks and 49.3% error reduction on regression problems. The framework’s strength becomes particularly evident on challenging vision datasets, where MNIST accuracy jumps from 78.41% to 98.99% and CIFAR-10 improves from 49.01% to 69.64%—gains that surpass what conventional optimization methods typically achieve even with extensive computational resources.

Beyond raw performance gains, LGT demonstrates an important advantage over black-box optimization approaches: it maintains high interpretability throughout the optimization process. While methods like Bayesian Optimization and NAS optimize effectively but provide little insight into their decision-making, LGT’s textual outputs offer transparent reasoning for every configuration change, enabling to understand and trust the optimization trajectory.

Table 3: Regression performance comparison (mean \pm std over 10 independent runs). LGT achieves consistent improvements across diverse datasets. Best results in **bold**.

Dataset	Method	MAE	MSE	R ²
House Price	No Tuning	58.77 _{2.3}	51.37 _{3.1}	0.757 _{0.02}
	Random Search	53.21 _{3.8}	45.82 _{4.2}	0.781 _{0.03}
	Grid Search	55.34 _{2.1}	47.91 _{2.9}	0.773 _{0.02}
	NAS	49.35 _{4.1}	38.67 _{3.8}	0.821 _{0.03}
	Bayesian Opt.	49.23 _{2.9}	41.45 _{3.4}	0.804 _{0.02}
	LGT	40.08_{2.1}	25.30_{2.8}	0.864_{0.02}
Wine Quality	No Tuning	51.84 _{2.7}	33.62 _{2.9}	0.673 _{0.03}
	Random Search	48.67 _{3.1}	31.24 _{3.4}	0.698 _{0.03}
	Grid Search	49.23 _{2.4}	32.15 _{2.8}	0.685 _{0.02}
	NAS	45.78 _{3.9}	28.91 _{3.7}	0.721 _{0.04}
	Bayesian Opt.	47.12 _{2.8}	30.45 _{3.2}	0.706 _{0.03}
	LGT	37.80_{2.2}	18.12_{2.4}	0.823_{0.02}

These remarkable improvements highlight a fundamental shift from traditional optimization approaches. While conventional methods like Grid Search exhaustively explore predefined spaces and Bayesian Optimization relies on probabilistic models, LGT leverages semantic understanding to guide dynamic optimization through natural language reasoning. This enables coordinated multi-dimensional adaptation that traditional approaches cannot achieve, as they typically treat configuration dimensions as independent variables. The textual gradient framework allows LGT to understand and leverage the complex interdependencies between architecture, training strategy, feature engineering, and hyperparameters—delivering superior performance across diverse domains while maintaining high interpretability throughout the optimization process.

Dynamic Convergence Analysis

Figure 2 illustrates the training dynamics of LGT compared to No Tuning baseline across six datasets. LGT consistently achieves faster convergence and superior final performance compared to the baseline, with particularly dramatic improvements visible across all datasets. The performance gaps are especially substantial on MNIST (final loss: 0.0469 vs 1.5827 baseline) and Iris (final loss: 0.0337 vs 0.7041 baseline), demonstrating LGT’s ability to achieve significantly better optimization outcomes.

The convergence analysis reveals why LGT outperforms traditional optimization approaches in both speed and stability. Unlike conventional methods that follow predetermined schedules or rely on fixed heuristics, LGT’s adaptive multi-agent system continuously evaluates training dynamics and implements intelligent configuration adjustments throughout the training process without requiring additional computational overhead. This dynamic adaptation enables dramatically faster convergence compared to static baseline approaches, with the red curves consistently descending more steeply than orange baseline curves across all datasets.

The stability advantage becomes particularly evident when comparing LGT’s smooth trajectories to the oscillatory behavior typical of grid search or random search meth-

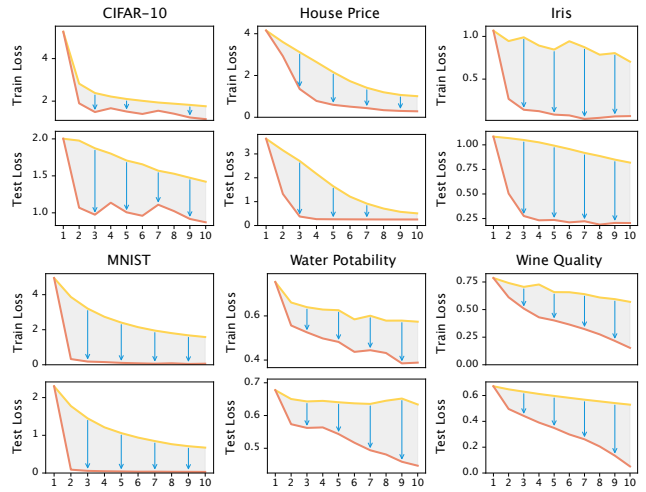


Figure 2: Train and test loss curves across six datasets comparing LGT (red) with baseline (orange). LGT demonstrates consistently faster convergence and better performance.

ods. While traditional approaches often struggle with premature convergence or unstable training dynamics, LGT maintains steady progress toward optimal configurations through its semantic understanding of training feedback. The performance gaps are especially pronounced across different dataset complexities: MNIST demonstrates LGT’s ability to escape poor initial configurations and achieve near-zero loss, while Iris showcases improvements in both train and test phases that baseline methods cannot match, highlighting the framework’s versatility across diverse learning scenarios.

Iterations	i = 1	...	i = 2	...	i = 3	...
Model Architecture	filter1: 8 \rightarrow 16 filter2: 16 \rightarrow 32 Add kernel3 dropout: 0.2 \rightarrow 0.4	—	/	—	filter1: 16 \rightarrow 32 filter2: 32 \rightarrow 64 filter3: 64 \rightarrow 128 dropout: 0.4 \rightarrow 0.3	—
Epochs	t = 1	...	t = 3	...	t = 7	...
Training Strategy	loss: cross entropy optimizer: adam	—	loss: focal loss optimizer: adamw	—	optimizer: adam	—
Feature Engineering	/	—	Flip	—	Rotation	—
Hyper Parameter	lr: 0.01 weight decay: 0.0001	—	lr: 0.02 weight_decay: 0.0005 class_3_weight=1.2	—	lr: 0.02 weight_decay: 0.0005 class_8_weight=0.95 class_9_weight=0.95	—

Figure 3: Configuration evolution showing coordinated optimization across four dimensions.

Figure 3 demonstrates the systematic evolution of configuration dimensions across both strategic iterations and tactical epoch-level adjustments. The dual timeline reveals LGT’s intelligent coordination strategy: iteration-level planning establishes major architectural foundations, while epoch-level adaptation fine-tunes training dynamics in real-time. At the iteration level, LGT follows a structured progression: Iteration 1 focuses on architectural modifications (filter1: 8 \rightarrow 16, filter2: 16 \rightarrow 32, adding kernel3, dropout: 0.2 \rightarrow 0.4), while Iteration 3 scales up the architecture further with refined regularization (dropout: 0.4 \rightarrow 0.3). Simultane-

ously, the epoch-level timeline shows tactical adjustments within training cycles: at epoch $t=1$, the system begins with cross-entropy loss and Adam optimizer, transitions to focal loss and AdamW at epoch $t=3$ to address emerging class imbalance issues, then optimizes back to Adam at epoch $t=7$ for final convergence acceleration. The feature engineering progression from basic augmentation (flip) to sophisticated transformations (rotation) and the hyperparameter refinements (learning rates from $0.01 \rightarrow 0.02 \rightarrow 0.02$, weight decay adjustments, and dynamic class weighting) demonstrate how LGT coordinates multiple configuration dimensions without creating optimization conflicts.

Textual Evolution and Interpretability

Figure 4 demonstrates the interpretable optimization process of our multi-agent system, showing how qualitative reasoning (textual gradients) enables transparent decision-making through natural language. This figure provides concrete evidence of how LGT achieves high interpretability while maintaining optimization effectiveness.

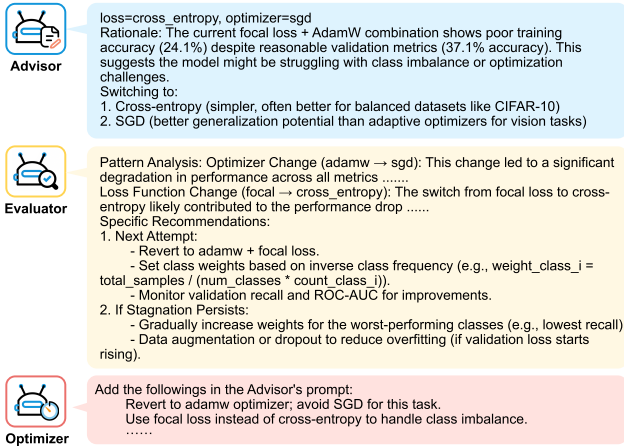


Figure 4: Example agent outputs showing multi-agent interaction. The Advisor provides specific configuration recommendations with detailed reasoning, while the Evaluator assesses optimization progress and provides meta-level guidance. The Optimizer refines prompts based on optimization history, creating a self-improving feedback loop.

The textual gradient analysis in Figure 4 showcases LGT’s unique transparency advantage. Unlike black-box methods like Bayesian Optimization or NAS, LGT’s multi-agent system provides detailed justifications for every modification: the Advisor articulates specific rationales (e.g., switching to focal loss for class imbalance), while the Evaluator provides meta-level analysis of performance patterns. This interpretability enables iterative learning and trust-building, with the Optimizer synthesizing lessons learned to create increasingly sophisticated guidance based on accumulated experience.

Ablation Study and Component Analysis

Table 4 presents ablation study demonstrating the individual effectiveness of each tuning dimension in our LGT frame-

work. We evaluate each of the four core optimization components—Model Architecture, Feature Engineering, Training Strategy, and Hyperparameters—in isolation to understand their relative contributions on classification (MNIST) and regression (Housing) tasks.

Table 4: Ablation study showing the effectiveness of individual optimization components on representative datasets (mean \pm std over 10 independent runs).

Configuration	MNIST AUC	Housing MSE
Baseline (No LLM)	97.55 _{0.4}	51.37 _{3.1}
Model Architecture Only	99.51 _{0.3}	25.69 _{2.9}
Feature Engineering Only	98.18 _{0.5}	47.51 _{3.4}
Training Strategy Only	99.96 _{0.2}	32.42 _{2.7}
Hyperparameters Only	99.89 _{0.3}	28.60 _{2.5}
Full LGT System	99.99_{0.1}	25.30_{2.8}

While conventional methods typically focus on isolated aspects, LGT’s strength emerges from coordinating across all dimensions. Individual components show context-dependent effectiveness. For MNIST, Training Strategy optimization (AUC: 99.96) proves most critical, suggesting the initial architecture suffices for this task but training dynamics require refinement. For Housing regression, Model Architecture optimization (MSE: 25.30) becomes paramount, indicating the default structure inadequately captures complex tabular relationships and requires architectural modifications. Rather than broad task-type generalizations, these results reflect LGT’s ability to identify where original configurations are most deficient and allocate optimization effort accordingly.

This coordinated approach addresses a fundamental limitation of existing methods: the assumption that configuration dimensions can be optimized independently. The ablation results demonstrate that while individual LGT components outperform traditional single-dimension approaches, the complete LGT achieves superior results by revealing and addressing the most limiting aspects of initial configurations, validating the multi-agent coordination paradigm.

Conclusion

This paper introduced Language-Guided Tuning (LGT), a novel framework that leverages Large Language Models for intelligent configuration optimization through **textual gradients**—qualitative feedback signals that complement numerical gradients. Our multi-agent system achieves up to 23.3% absolute accuracy improvement and 49.3% error reduction while consistently outperforming Neural Architecture Search and Bayesian Optimization across six diverse datasets. The framework demonstrates dynamic convergence with strategic epoch-level adaptations, provides high interpretability through natural language reasoning, and enables coordinated optimization across interdependent configuration dimensions. By bridging semantic reasoning and numerical optimization, LGT represents a paradigm shift toward more interpretable and intelligent machine learning systems.

References

- Alom, M. Z.; Taha, T. M.; Yakopcic, C.; Westberg, S.; Sidike, P.; Nasrin, M. S.; Hasan, M.; Van Essen, B. C.; Awwal, A. A.; and Asari, V. K. 2019. A state-of-the-art survey on deep learning theory and architectures. *electronics*, 8(3): 292.
- Bergstra, J.; and Bengio, Y. 2012. Random search for hyperparameter optimization. *The journal of machine learning research*, 13(1): 281–305.
- Chen, Y.; Song, X.; Lee, C.; Wang, Z.; Zhang, R.; Dohan, D.; Kawakami, K.; Kochanski, G.; Doucet, A.; Ranzato, M.; et al. 2022. Towards learning universal hyperparameter optimizers with transformers. *Advances in Neural Information Processing Systems*, 35: 32053–32068.
- Cortez, P.; Cerdeira, A.; Almeida, F.; Matos, T.; and Reis, J. 2009. Modeling wine preferences by data mining from physicochemical properties. *Decision support systems*, 47(4): 547–553.
- Deng, L. 2012. The mnist database of handwritten digit images for machine learning research [best of the web]. *IEEE signal processing magazine*, 29(6): 141–142.
- Deng, L. 2014. A tutorial survey of architectures, algorithms, and applications for deep learning. *APSIPA transactions on Signal and Information Processing*, 3: e2.
- Frazier, P. I. 2018. A tutorial on Bayesian optimization. *arXiv preprint arXiv:1807.02811*.
- Gharehchopogh, F. S.; Namazi, M.; Ebrahimi, L.; and Abdollahzadeh, B. 2023. Advances in sparrow search algorithm: a comprehensive survey. *Archives of Computational Methods in Engineering*, 30(1): 427–455.
- Goodfellow, I.; Bengio, Y.; and Courville, A. 2016. Deep Learning.
- Harrison Jr, D.; and Rubinfeld, D. L. 1978. Hedonic housing prices and the demand for clean air. *Journal of environmental economics and management*, 5(1): 81–102.
- Huang, C.; Li, Y.; and Yao, X. 2019. A survey of automatic parameter tuning methods for metaheuristics. *IEEE transactions on evolutionary computation*, 24(2): 201–216.
- Huang, Z.; Wu, W.; Wu, K.; Wang, J.; and Lee, W.-B. 2025. Calm: Co-evolution of algorithms and language model for automatic heuristic design. *arXiv preprint arXiv:2505.12285*.
- Karmaker, S. K.; Hassan, M. M.; Smith, M. J.; Xu, L.; Zhai, C.; and Veeramachaneni, K. 2021. Automl to date and beyond: Challenges and opportunities. *Acm computing surveys (csur)*, 54(8): 1–36.
- Krizhevsky, A.; Hinton, G.; et al. 2009. Learning multiple layers of features from tiny images.
- LaValle, S. M.; Branicky, M. S.; and Lindemann, S. R. 2004. On the relationship between classical grid search and probabilistic roadmaps. *The International Journal of Robotics Research*, 23(7-8): 673–692.
- Liashchynskyi, P.; and Liashchynskyi, P. 2019. Grid search, random search, genetic algorithm: a big comparison for NAS. *arXiv preprint arXiv:1912.06059*.
- Liu, A.; Feng, B.; Xue, B.; Wang, B.; Wu, B.; Lu, C.; Zhao, C.; Deng, C.; Zhang, C.; Ruan, C.; et al. 2024. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*.
- Liu, H.; Simonyan, K.; and Yang, Y. 2018. Darts: Differentiable architecture search. *arXiv preprint arXiv:1806.09055*.
- Liu, S.; Gao, C.; and Li, Y. 2025. AgentHPO: Large language model agent for hyper-parameter optimization. In *The Second Conference on Parsimony and Learning (Proceedings Track)*.
- Lu, Y.; Fu, G.; Wu, W.; Zhao, X.; Goi, S. Y.; and Wang, J. 2025. DoctorRAG: Medical RAG Fusing Knowledge with Patient Analogy through Textual Gradients. *arXiv preprint arXiv:2505.19538*.
- Lu, Y.; and Wang, J. 2025. KARMA: Leveraging Multi-Agent LLMs for Automated Knowledge Graph Enrichment. *arXiv:2502.06472*.
- Ma, Z.; Guo, H.; Chen, J.; Peng, G.; Cao, Z.; Ma, Y.; and Gong, Y.-J. 2024. Llamoco: Instruction tuning of large language models for optimization code generation. *arXiv preprint arXiv:2403.01131*.
- Nie, A.; Cheng, C.-A.; Kolobov, A.; and Swaminathan, A. 2024. The importance of directional feedback for llm-based optimizers. *arXiv preprint arXiv:2405.16434*.
- Pontes, F. J.; Amorim, G.; Balestrassi, P. P.; Paiva, A.; and Ferreira, J. R. 2016. Design of experiments and focused grid search for neural network parameter optimization. *Neuro-computing*, 186: 22–34.
- Pouyanfar, S.; Sadiq, S.; Yan, Y.; Tian, H.; Tao, Y.; Reyes, M. P.; Shyu, M.-L.; Chen, S.-C.; and Iyengar, S. S. 2018. A survey on deep learning: Algorithms, techniques, and applications. *ACM computing surveys (CSUR)*, 51(5): 1–36.
- Ren, P.; Xiao, Y.; Chang, X.; Huang, P.-Y.; Li, Z.; Chen, X.; and Wang, X. 2021. A comprehensive survey of neural architecture search: Challenges and solutions. *ACM Computing Surveys (CSUR)*, 54(4): 1–34.
- Shorten, C.; and Khoshgoftaar, T. M. 2019. A survey on image data augmentation for deep learning. *Journal of big data*, 6(1): 1–48.
- Simeonov, V.; Barbieri, P.; Walczak, B.; Massart, D.; and Tsakovski, S. 2001. Environmetric modeling of a potable water data set. *Toxicological & Environmental Chemistry*, 79(1-2): 55–72.
- Victoria, A. H.; and Maragatham, G. 2021. Automatic tuning of hyperparameters using Bayesian optimization. *Evolutionary Systems*, 12(1): 217–223.
- Wang, Q.; Ma, Y.; Zhao, K.; and Tian, Y. 2022. A comprehensive survey of loss functions in machine learning. *Annals of Data Science*, 9(2): 187–212.
- Wang, X.; Jin, Y.; Schmitt, S.; and Olhofer, M. 2023. Recent advances in Bayesian optimization. *ACM Computing Surveys*, 55(13s): 1–36.
- Wang, Z.; Chu, Z.; Doan, T. V.; Ni, S.; Yang, M.; and Zhang, W. 2025. History, development, and principles of large language models: an introductory survey. *AI and Ethics*, 5(3): 1955–1971.

- Yang, C.; Wang, X.; Lu, Y.; Liu, H.; Le, Q. V.; Zhou, D.; and Chen, X. 2023. Large language models as optimizers. In *The Twelfth International Conference on Learning Representations*.
- Ye, H.; Wang, J.; Cao, Z.; Berto, F.; Hua, C.; Kim, H.; Park, J.; and Song, G. 2024. Reevo: Large language models as hyper-heuristics with reflective evolution. *Advances in neural information processing systems*, 37: 43571–43608.
- Yu, H.; and Liu, J. 2024. Deep insights into automated optimization with large language models and evolutionary algorithms. *arXiv preprint arXiv:2410.20848*.
- Yu, T.; and Zhu, H. 2020. Hyper-parameter optimization: A review of algorithms and applications. *arXiv preprint arXiv:2003.05689*.
- Yuksekgonul, M.; Bianchi, F.; Boen, J.; Liu, S.; Huang, Z.; Guestrin, C.; and Zou, J. 2024. Textgrad: Automatic” differentiation” via text. *arXiv preprint arXiv:2406.07496*.
- Zhang, M. R.; Desai, N.; Bae, J.; Lorraine, J.; and Ba, J. 2023. Using large language models for hyperparameter optimization. *arXiv preprint arXiv:2312.04528*.
- Zhao, W. X.; Zhou, K.; Li, J.; Tang, T.; Wang, X.; Hou, Y.; Min, Y.; Zhang, B.; Zhang, J.; Dong, Z.; et al. 2023. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 1(2).
- Zöllner, M.-A.; and Huber, M. F. 2021. Benchmark and survey of automated machine learning frameworks. *Journal of artificial intelligence research*, 70: 409–472.