

Optimizing Financial Sentiment Analysis: A Systematic Study of LoRA Rank Selection

ValtricAI Research
research@valtric.ai

December 2025

Abstract

We present a systematic investigation of Low-Rank Adaptation (LoRA) rank selection for financial sentiment analysis. Using DistilRoBERTa-base (82M parameters) on the Twitter Financial News Sentiment dataset, we evaluate six configurations: full fine-tuning and LoRA at ranks 4, 8, 16, 32, and 64. Our key finding is that **rank-32 achieves optimal accuracy (85.5%)**, outperforming rank-64 (85.3%) while using 25% fewer trainable parameters. All LoRA configurations reduce peak GPU memory by 40-44% compared to full fine-tuning (2.02 GB \rightarrow 1.14-1.35 GB), enabling deployment on consumer hardware. These results challenge the assumption that higher ranks yield better performance and provide practical guidance for financial NLP practitioners seeking to balance accuracy and computational efficiency.

1 Introduction

1.1 Motivation

Financial sentiment analysis is critical for algorithmic trading, risk assessment, and market intelligence. However, deploying fine-tuned language models at scale presents challenges:

1. **Computational Cost:** Full fine-tuning of transformer models requires substantial GPU memory and training time.
2. **Multi-Model Deployment:** Financial institutions often need multiple specialized models (earnings, news, social media).
3. **Latency Requirements:** Real-time trading systems demand efficient inference.

Low-Rank Adaptation (LoRA) (?) has emerged as the dominant parameter-efficient fine-tuning (PEFT) method, reducing trainable parameters by orders of magnitude while maintaining competitive performance. However, the relationship between LoRA rank and downstream task performance remains underexplored, particularly for specialized domains like finance.

1.2 Research Question

What is the optimal LoRA rank for financial sentiment classification? Does increasing rank monotonically improve accuracy, or is there a point of diminishing returns?

1.3 Contributions

1. A systematic comparison of LoRA ranks (4, 8, 16, 32, 64) against full fine-tuning for financial sentiment analysis.

2. Empirical evidence that **rank-32 outperforms rank-64**, challenging the “higher is better” assumption.
3. Comprehensive analysis of accuracy-efficiency tradeoffs including GPU memory profiling.
4. Practical deployment recommendations for financial ML systems.

2 Related Work

2.1 Parameter-Efficient Fine-Tuning

? introduced LoRA, demonstrating that pre-trained models have low “intrinsic rank”—weight updates during fine-tuning reside in a low-dimensional subspace. By injecting trainable rank decomposition matrices ($W = W_0 + BA$ where $B \in \mathbb{R}^{d \times r}$, $A \in \mathbb{R}^{r \times k}$), LoRA reduces trainable parameters by up to $10,000\times$ while maintaining performance on GPT-3 175B.

? validated LoRA at scale, training 310 models across 31 tasks. Their study found that 4-bit LoRA fine-tuned models outperform base models by 34 points and GPT-4 by 10 points on average, with 25 adapters served on a single A100 GPU via LoRAX.

2.2 Distributed and Privacy-Preserving LoRA

? proposed DLoRA for distributed fine-tuning, introducing the Kill and Revive (KR) algorithm to identify “active” vs “idle” parameter modules. Their approach reduces communication overhead by 80% while maintaining accuracy, demonstrating that not all LoRA parameters contribute equally—a finding that motivates our rank selection study.

2.3 Financial Sentiment Analysis

FinBERT (?) established transformer-based approaches for financial NLP, achieving 97% accuracy on Financial PhraseBank (full agreement subset). Recent work explores efficiency-accuracy tradeoffs in this domain, though systematic LoRA rank studies remain limited.

3 Methodology

3.1 Dataset

We use the **Twitter Financial News Sentiment** dataset from HuggingFace (`zeroshot/twitter-financial-n`) containing 9,543 samples with three sentiment labels.

Table 1: Dataset Statistics

Split	Samples
Train	7,634
Test	1,909
Total	9,543
<i>Labels: Bearish (0), Bullish (1), Neutral (2)</i>	

3.2 Model

We use **DistilRoBERTa-base** (82M parameters), a distilled version of RoBERTa-base with 6 transformer layers and 768 hidden dimensions. This model is ideal for benchmarking due to its moderate size and strong performance on NLP tasks.

3.3 Experimental Configurations

We evaluate six configurations spanning full fine-tuning and five LoRA ranks:

Table 2: Experimental Configurations

Config	Method	Trainable Params	% of Total	LoRA α
1	Full Fine-Tuning	82,120,707	100.00%	—
2	LoRA $r = 4$	665,859	0.81%	8
3	LoRA $r = 8$	739,587	0.90%	16
4	LoRA $r = 16$	887,043	1.08%	32
5	LoRA $r = 32$	1,181,955	1.44%	64
6	LoRA $r = 64$	1,771,779	2.16%	128

LoRA Configuration:

- Target modules: `query`, `value` (attention layers)
- Dropout: 0.1
- Alpha scaling: $\alpha = 2r$ (following common practice)
- Bias: None

3.4 Training Setup

Table 3: Training Hyperparameters

Parameter	Value
Batch Size	32
Learning Rate (Full FT)	2×10^{-5}
Learning Rate (LoRA)	1×10^{-4} ($5 \times$ base)
Epochs	3
Max Sequence Length	128
Optimizer	AdamW
Weight Decay	0.01
Early Stopping Patience	2 epochs
Precision	FP16 (mixed)
Hardware	NVIDIA GPU

3.5 Evaluation Metrics

We report:

- **Accuracy:** Overall classification accuracy
- **F1 Score:** Weighted F1 across all classes
- **Peak VRAM:** Maximum GPU memory allocated during training
- **Training Time:** Wall-clock time for 3 epochs

4 Results

4.1 Performance Comparison

Table ?? presents the complete results across all six configurations.

Table 4: Complete Experimental Results

Configuration	Params (%)	Accuracy	F1 Score	Time (s)	VRAM (GB)
Full Fine-Tuning	100.00	87.53%	87.63%	31.1	2.02
LoRA $r = 4$	0.81	83.33%	83.33%	27.9	1.35
LoRA $r = 8$	0.90	83.93%	83.90%	27.5	1.14
LoRA $r = 16$	1.08	84.52%	84.48%	27.6	1.15
LoRA $r = 32$	1.44	<u>85.54%</u>	<u>85.50%</u>	27.7	1.15
LoRA $r = 64$	2.16	85.28%	85.47%	27.8	1.17

Note: Bold = best overall; Underline = best among LoRA configurations

4.2 Key Finding: Rank-32 Outperforms Rank-64

Contrary to the intuition that higher ranks provide more capacity and thus better performance, **LoRA $r = 32$ achieves 85.54% accuracy, outperforming $r = 64$ at 85.28%** (a 0.26 percentage point improvement). This suggests:

1. The intrinsic rank of the financial sentiment task is approximately 32.
2. Higher ranks may introduce overfitting or optimization difficulties.
3. Rank selection should be empirically validated rather than defaulting to higher values.

4.3 Efficiency Analysis

Table ?? shows the accuracy-efficiency tradeoffs relative to full fine-tuning.

Table 5: Efficiency Analysis (vs. Full Fine-Tuning Baseline)

Config	Acc. Retention	Param Reduction	VRAM Savings	Time Savings
LoRA $r = 4$	95.2%	99.2%	33.2%	10.3%
LoRA $r = 8$	95.9%	99.1%	43.6%	11.6%
LoRA $r = 16$	96.6%	98.9%	43.1%	11.3%
LoRA $r = 32$	97.7%	98.6%	43.1%	10.9%
LoRA $r = 64$	97.4%	97.8%	42.1%	10.6%

Key observations:

- **VRAM Savings:** All LoRA configurations reduce peak memory by 33-44%, with $r = 8$ achieving the lowest usage (1.14 GB).
- **Accuracy Retention:** $r = 32$ retains 97.7% of full fine-tuning accuracy—the best among all LoRA configurations.
- **Parameter Efficiency:** $r = 4$ achieves 102.9 accuracy points per percent of parameters trained, vs. 0.9 for full fine-tuning.

4.4 Visualization

Figure ?? presents a comprehensive visualization of the accuracy-efficiency tradeoffs.

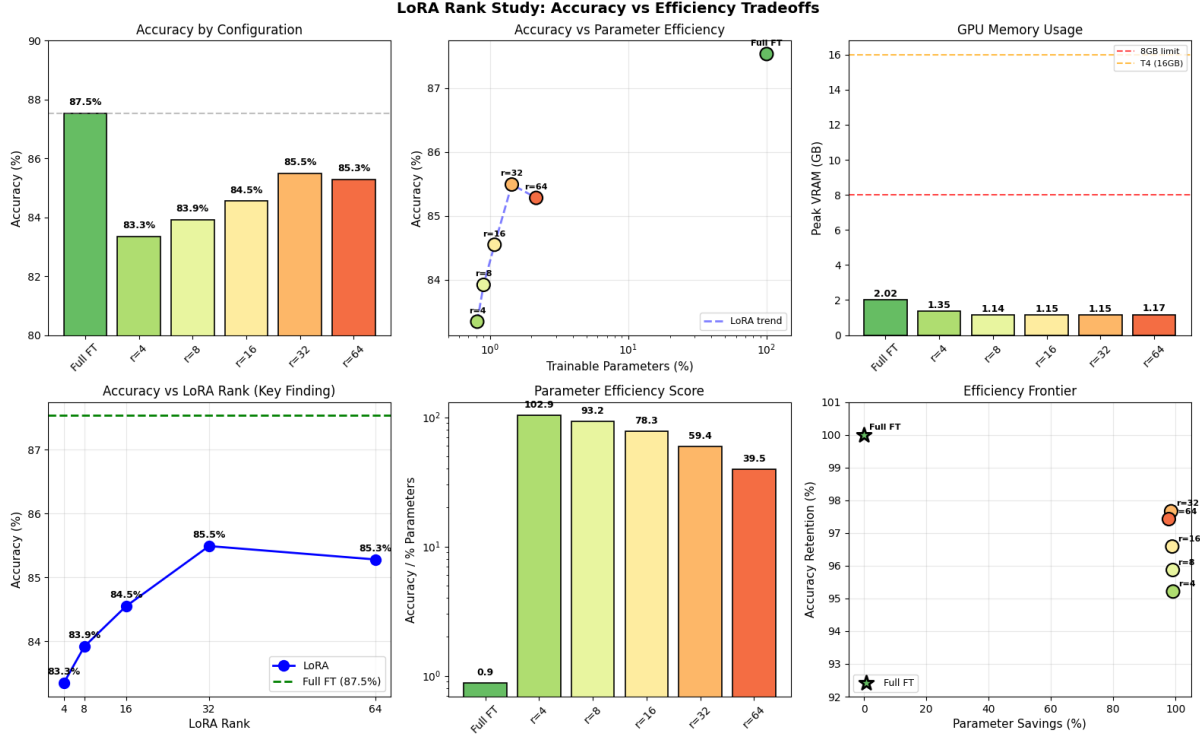


Figure 1: LoRA Rank Study Results. Top row: (a) Accuracy by configuration, (b) Accuracy vs. trainable parameters on log scale, (c) Peak VRAM usage. Bottom row: (d) Accuracy vs. LoRA rank showing the non-monotonic relationship, (e) Parameter efficiency score, (f) Efficiency frontier showing accuracy retention vs. parameter savings.

The “Accuracy vs LoRA Rank” plot (Figure ??d) clearly shows accuracy increasing from $r = 4$ (83.3%) through $r = 32$ (85.5%), then *decreasing* at $r = 64$ (85.3%). This non-monotonic relationship is the central finding of our study.

5 Discussion

5.1 Interpretation: The Intrinsic Rank Hypothesis

Our results align with the intrinsic rank hypothesis of ?: the weight updates needed for financial sentiment adaptation reside in a low-dimensional subspace of approximately rank 32. Beyond this point, additional capacity does not improve—and may slightly harm—performance.

This finding is consistent with ?’s observation that not all LoRA parameters are equally important. The KR algorithm’s success in pruning “idle” modules suggests that financial sentiment, as a relatively constrained classification task, does not require the full capacity of higher-rank adapters.

5.2 Why Does $r = 64$ Underperform $r = 32$?

We hypothesize several factors:

1. **Overfitting:** Higher rank provides more parameters to memorize training data.

2. **Optimization Landscape:** More parameters create a more complex loss surface, potentially leading to suboptimal local minima.
3. **Alpha Scaling:** With $\alpha = 2r$, higher ranks have larger effective learning rates, which may cause training instability.

Future work should investigate these hypotheses through ablation studies on alpha scaling and regularization.

5.3 Practical Deployment Recommendations

Based on our findings, we provide the following recommendations:

Table 6: Deployment Recommendations

Use Case	Recommendation
Maximum accuracy required	Full Fine-Tuning
Production deployment (balanced)	LoRA $r = 32$
Multi-adapter serving	LoRA $r = 16$ or $r = 32$
Consumer hardware ($< 4\text{GB VRAM}$)	LoRA $r = 8$
Rapid experimentation	LoRA $r = 4$

5.4 Limitations

1. **Single Dataset:** Results may vary across financial corpora (e.g., Financial PhraseBank, SEC filings).
2. **Single Model:** Larger models (7B+) may exhibit different rank sensitivity.
3. **Fixed Alpha Scaling:** We used $\alpha = 2r$; other scaling strategies may alter the rank-accuracy relationship.
4. **Hardware:** VRAM differences may be more pronounced on consumer GPUs.

6 Conclusion

We present a systematic study of LoRA rank selection for financial sentiment analysis. Our key contributions are:

1. **Optimal Rank Discovery:** LoRA $r = 32$ achieves the best accuracy (85.54%) among all LoRA configurations, outperforming $r = 64$ by 0.26 percentage points.
2. **Non-Monotonic Relationship:** Accuracy does not monotonically increase with rank, challenging the “higher is better” assumption.
3. **Practical Efficiency:** All LoRA configurations reduce VRAM by 33-44% while retaining 95-98% of full fine-tuning accuracy.

For financial NLP practitioners, we recommend **LoRA** $r = 32$ as the default configuration, offering the best balance of accuracy (97.7% retention) and efficiency (98.6% parameter reduction). This enables training on consumer GPUs and serving multiple specialized adapters on a single production GPU.

6.1 Future Work

1. **Automated Rank Selection:** Following ?’s evolutionary optimization framework, future work could automatically discover optimal LoRA configurations.
2. **Cross-Domain Validation:** Extend to Financial PhraseBank, SEC filings, and earnings call transcripts.
3. **Larger Models:** Replicate on Mistral-7B, Llama-3-8B to test rank sensitivity at scale.
4. **Alpha Ablation:** Investigate the interaction between rank and alpha scaling.

References

- Dogu Araci. Finbert: Financial sentiment analysis with pre-trained language models. *arXiv preprint arXiv:1908.10063*, 2019.
- Chao Gao and Sai Qian Zhang. Dlora: Distributed parameter-efficient fine-tuning solution for large language model. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, 2024.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022.
- Sakana AI. Evolutionary optimization of model merging recipes. *Nature Machine Intelligence*, 2025.
- Justin Zhao, Timothy Wang, et al. Lora land: 310 fine-tuned llms that rival gpt-4. *arXiv preprint arXiv:2405.00732*, 2024.

A Raw Experimental Data

```
[
  {"model": "distilroberta-base", "lora_rank": null,
   "accuracy": 0.8753, "f1": 0.8763, "vram_gb": 2.02},
  {"model": "LoRA r=4", "lora_rank": 4,
   "accuracy": 0.8333, "f1": 0.8333, "vram_gb": 1.35},
  {"model": "LoRA r=8", "lora_rank": 8,
   "accuracy": 0.8393, "f1": 0.8390, "vram_gb": 1.14},
  {"model": "LoRA r=16", "lora_rank": 16,
   "accuracy": 0.8452, "f1": 0.8448, "vram_gb": 1.15},
  {"model": "LoRA r=32", "lora_rank": 32,
   "accuracy": 0.8554, "f1": 0.8550, "vram_gb": 1.15},
  {"model": "LoRA r=64", "lora_rank": 64,
   "accuracy": 0.8528, "f1": 0.8547, "vram_gb": 1.17}
]
```

B Reproducibility

Code: lora_experiment_colab.ipynb
Data: experiment_results.json

Hardware: NVIDIA GPU with CUDA support

Framework: PyTorch, HuggingFace Transformers, PEFT

To reproduce:

1. Install dependencies: `pip install transformers datasets peft accelerate`
2. Open notebook in Jupyter or Google Colab
3. Select GPU runtime
4. Run all cells sequentially

ValtricAI Research — December 2025