# Inherent Robustness of Multi-Head Attention in Cross-Sectional Asset Pricing:

## Theory and Empirical Evidence from Finance-Valid Adversarial Attacks

Zelalem Abahana First Citizens Bank, Raleigh, NC 27601 USA
Vice President, AI/ML Model Validation
(e-mail: zelalem.abahana@firstcitizens.com)

*Abstract*—This paper presents a theoretical and empirical investigation of multi-head attention architectures with head-diversity regularization for robust cross-sectional asset pricing under adversarial attacks. We establish theoretically motivated robustness improvements scaling with the number of attention heads through information redundancy, ensemble stabilization, and Lipschitz regularization mechanisms. Empirically, we evaluate six models on S&P 500 stocks: linear baselines (OLS, Ridge), MLP, single-head transformer, multi-head transformer, and multi-head transformer with head-diversity regularization. Models are evaluated for prediction accuracy and robustness to finance-valid adversarial attacks (measurement error, missingness, rank manipulation, regime shift) under both standard and adversarial training regimes.

Our theoretical analysis demonstrates that head-diversity regularization creates information redundancy requiring simultaneous corruption of multiple heads for significant degradation, with ensemble variance reduction and implicit Lipschitz bounds providing perturbation stability. Empirically, attention-based models achieve superior predictive performance compared to linear baselines and MLP, with multi-head attention demonstrating the highest performance. Multi-head attention with diversity regularization achieves the highest robustness under stress conditions, demonstrating superior resilience to adversarial perturbations. Standard multi-head attention models demonstrate near-invariance under small perturbations across all attack types. However, at larger perturbations, standard models experience significant degradation. Adversarial training achieves substantial robustness improvements, restoring robustness even under large perturbations where standard models experience significant degradation. These findings establish multi-head attention with head-diversity regularization and adversarial training as a theoretically-grounded and empirically-validated approach to achieving robust cross-sectional asset pricing. For financial modeling, this combination provides critical advantages: resilience to regime changes and distribution shifts, and robustness to data quality issues including measurement errors and missing data. These implications make multi-head attention with adversarial training particularly valuable for quantitative investment strategies, risk management systems, and regulatory model validation where robustness to perturbations and regime changes is essential.

*Index Terms*—Cross-sectional asset pricing, multi-head attention, adversarial robustness, head-diversity regularization, finance-valid attacks, expected return prediction

## I. Introduction

Cross-sectional asset pricing aims to predict expected returns from observable firm characteristics such as momentum, volatility, and valuation ratios. This problem underpins quantitative investment strategies and factor-based portfolio construction. Traditional linear models, including Fama and French [1] factor models, achieve moderate success but struggle with non-linear relationships inherent in financial data.

Modern machine learning approaches have shown promise. Gu, Kelly, and Xiu [3] neural networks and Chen, Pelger, and Zhu [4] tree-based methods capture non-linear patterns, while Lim, Zohren, and Roberts [5] transformers enable attention-based feature interactions. However, these models face practical challenges: vulnerability to adversarial perturbations, limited interpretability for regulatory compliance, and potential overfitting without robustness to regime shifts.

Vaswani et al. [2] multi-head attention mechanisms allow models to focus on different aspects of input features through parallel attention heads. In financial applications, this enables specialized attention to momentum, volatility, valuation, and technical indicators. However, without explicit regularization, heads may learn redundant representations, limiting the benefits of multi-head architectures.

This paper presents a comprehensive theoretical and empirical investigation of whether head-diversity regularization can enhance adversarial robustness for cross-sectional asset pricing. We bridge theoretical analysis with empirical validation, establishing both theoretically motivated robustness improvements and practical performance improvements. Our work makes four key contributions:

1) **Theoretical Framework**: We establish a rigorous theoretical framework demonstrating that multi-head attention with head-diversity regularization provides robustness improvements scaling as $\Omega(1/\sqrt{H})$ under mild diversity assumptions through three complementary mechanisms: (1) information redundancy requiring simultaneous corruption of multiple heads, (2) ensemble stabilization with variance reduction scaling as $1/H$, and (3) implicit Lipschitz regularization providing $O(1/\sqrt{H})$ perturbation bounds. This theoretical foundation provides principled justification for diversity regularization beyond empirical observation.

2) **Empirical Validation**: We conduct comprehensive empirical evaluation of six models (OLS, Ridge, MLP, Single-Head, Multi-Head, Multi-Head Diversity) on 142 S&P 500 stocks from 2005-2019, training on 2005-2017 and validating on 2018-2019. Models are assessed on prediction accuracy ($R^2$, RMSE) and robustness to

finance-valid adversarial attacks under both standard and adversarial training regimes, providing empirical validation of our theoretical predictions.

3) **Finance-Valid Attack Framework**: We introduce and evaluate four finance-valid attack types reflecting realistic market scenarios: measurement error (A1), missingness/staleness (A2), rank manipulation (A3), and regime shift (A4). These attacks provide more interpretable and realistic robustness assessment than standard gradient-based attacks, aligning with financial market constraints and empirical evidence on regime-switching behavior.

4) **Adversarial Training Under Large Perturbations**: We show that while multi-head attention architectures exhibit strong inherent robustness under small finance-realistic perturbations ($\epsilon \leq 0.5$), adversarial training becomes critical under larger perturbations ($\epsilon = 1.0$), restoring robustness by 10–17% across measurement error, rank manipulation, and regime shift attacks. This establishes a complementary role between architectural robustness and adversarial training in financial applications.

Our theoretical analysis establishes theoretically motivated robustness improvements, while our empirical results demonstrate that attention-based models achieve positive R² compared to negative R² for linear baselines and MLP. Multi-head attention achieves the highest R² and lowest RMSE among all models, demonstrating superior predictive performance. Multi-head architectures demonstrate near-invariance under small perturbations (robustness $\geq$ 98%) for all attacks including regime shifts (A4), outperforming linear baselines and MLP which show lower robustness to adversarial perturbations. Multi-head attention with diversity regularization achieves the highest robustness under stress conditions, demonstrating that head-diversity regularization provides enhanced robustness without sacrificing predictive accuracy. These findings establish multi-head attention with head-diversity regularization as a theoretically-grounded and empirically-validated approach to achieving robust cross-sectional asset pricing, particularly under regime change challenges common in financial markets.

## II. RELATED WORK

### A. Cross-Sectional Asset Pricing

The cross-sectional prediction of expected returns from firm characteristics has been extensively studied in empirical finance. Fama and French [1] established the three-factor model, demonstrating that market, size, and book-to-market factors explain cross-sectional return variation. Subsequent work extended this framework with additional factors [8], [9] and Kozak, Nagel, and Santosh [10] explored non-linear relationships. Harvey, Liu, and Zhu [9] identified the factor zoo problem, highlighting the challenge of identifying robust predictors among hundreds of candidate characteristics.

Modern machine learning approaches have achieved significant improvements over linear models. Gu, Kelly, and Xiu [3] demonstrated that neural networks outperform linear models in cross-sectional prediction, achieving positive out-of-sample R² on large stock universes. Chen, Pelger, and Zhu [4] employed deep learning with automated feature engineering, while Feng, Polson, and Xu [11] further advanced deep learning approaches for factor models. Chen and Zimmermann [13] and Freyberger, Neuhierl, and Weber [14] tree-based methods have shown promise through non-linear interaction discovery, with gradient boosting methods achieving strong performance in cross-sectional prediction.

Transformer-based approaches have recently emerged in financial applications. Lim, Zohren, and Roberts [5] applied temporal fusion transformers to time-series forecasting, while Chen, Pelger, and Zhu [15] explored neural asset pricing with transformer architectures. However, existing work has not explicitly addressed adversarial robustness or head-diversity regularization. Our work focuses on cross-sectional prediction where each stock-month observation is treated independently with characteristics as tokens, enabling feature-token transformer architectures.

### B. Multi-Head Attention and Diversity

Theoretical analysis of multi-head attention has established Tsai et al. [28] expressive power and Voita et al. [29] connection to ensemble methods. However, without explicit regularization, heads may learn redundant representations, limiting the benefits of multi-head architectures.

Head diversity has been explored in natural language processing. Voita et al. [29] analyzed head specialization in machine translation, finding that different heads capture different linguistic phenomena. Michel, Levy, and Neubig [30] demonstrated that many attention heads can be pruned without significant performance loss, suggesting redundancy. Raganato and Tiedemann [31] found that heads specialize in different syntactic and semantic roles.

Diversity regularization has been applied to improve model performance. Li et al. [32] proposed diversity-promoting regularization for multi-task learning. Bapna, Firat, and Wu [33] explored head diversity in multilingual translation. However, the connection between head diversity and adversarial robustness has not been theoretically established or empirically validated in financial applications.

Our work differs by: (1) establishing theoretically motivated robustness improvements through head diversity, (2) treating firm characteristics as tokens (feature-token transformer), (3) explicitly regularizing for head diversity with theoretical robustness benefits, and (4) evaluating robustness to finance-valid adversarial attacks under both standard and adversarial training regimes.

### C. Adversarial Robustness Theory

Adversarial robustness has been extensively studied in computer vision and natural language processing. Szegedy et al. [24] first identified adversarial examples in neural networks, while Goodfellow, Shlens, and Szegedy [19] proposed the Fast Gradient Sign Method (FGSM) for generating adversarial examples. Madry et al. [20] developed Projected Gradient Descent (PGD) attacks and adversarial training, establishing the adversarial training paradigm.

Theoretical analysis of adversarial robustness has established fundamental limits and guarantees. Wong and Kolter [26] developed provable defenses via convex outer adversarial polytopes, while Cohen, Rosenfeld, and Kolter [27] established certified robustness via randomized smoothing. Raghunathan, Steinhardt, and Liang [34] provided theoretical analysis of certified defenses, while Schmidt et al. [35] established fundamental trade-offs between accuracy and robustness.

Ensemble methods have been shown to improve adversarial robustness. Tramer et al. [36] demonstrated that ensemble defenses can improve robustness, while Pang et al. [37] analyzed ensemble robustness theoretically. However, the connection between multi-head attention diversity and ensemble robustness has not been established.

Lipschitz regularization has been explored for improving robustness. Anil et al. [38] proposed Lipschitz-constrained networks, while Tsuzuku, Sato, and Sugiyama [39] analyzed Lipschitz-margin training. Our theoretical framework establishes that multi-head attention with diversity regularization provides implicit Lipschitz bounds, contributing to robustness improvements.

### D. Adversarial Robustness in Financial Machine Learning

Adversarial robustness is critical in financial machine learning applications, where model reliability directly impacts investment decisions and risk management. Kumar, Zhang, and Wang [7] demonstrated that stock prediction models experience 30-50% accuracy degradation under gradient-based attacks. However, standard computer vision attacks (FGSM, PGD) [19], [20] may not reflect realistic financial data constraints, as financial data exhibits different statistical properties than images.

Distribution shift and regime changes pose particular challenges in financial applications. Ang and Timmermann [21] documented regime-switching behavior in asset returns, demonstrating that market conditions change over time. Bollerslev, Hood, and Pedersen [23] demonstrated that risk-return relationships vary across market regimes, while Lettau and Ludvigson [22] analyzed variation in the risk-return trade-off. These findings highlight the importance of robustness to distribution shifts in financial models.

Robustness evaluation in finance requires attacks that reflect realistic market scenarios. Standard gradient-based attacks may not be applicable to financial data due to constraints on feature relationships, cross-sectional dependencies, and temporal structure. Our finance-valid attack framework addresses this gap by introducing attacks that respect financial data constraints:

- **A1 - Measurement Error**: Bounded perturbations ($\|\delta\|_\infty \leq \epsilon \cdot \sigma(\mathbf{x})$) simulating data collection errors and measurement noise, reflecting realistic data quality issues in financial datasets.
- **A2 - Missingness/Staleness**: Random feature zeroing with probability $p$, simulating delayed or missing data common in financial applications where data availability varies across firms and time periods.

- **A3 - Rank Manipulation**: Cross-sectional perturbations preserving relative rankings, reflecting scenarios where cross-sectional relationships are maintained but absolute values are perturbed.
- **A4 - Regime Shift**: Distribution shift attacks simulating market regime changes, where volatility and correlation structures change, reflecting realistic market condition transitions.

These attacks provide more realistic and interpretable robustness evaluation than standard gradient-based attacks, aligning with financial market constraints and empirical evidence on regime-switching behavior. Our evaluation framework assesses robustness under both standard and adversarial training regimes, providing comprehensive analysis of architectural and training-based robustness improvements.

## III. METHODOLOGY

### A. Problem Formulation

Cross-sectional asset pricing predicts expected returns from firm characteristics. For each stock $i$ and month $t$, we observe characteristics $\mathbf{x}_{i,t} \in \mathbb{R}^d$ and predict forward return $r_{i,t+1}$:

$$\hat{r}_{i,t+1} = f(\mathbf{x}_{i,t}; \theta) \tag{1}$$

where $f$ is a learned function parameterized by $\theta$. The objective minimizes prediction error while maintaining robustness:

$$\min_\theta \mathbb{E}[\ell(f(\mathbf{x};\theta), r)] + \lambda \mathcal{R}(\theta) \tag{2}$$

where $\ell$ is MSE loss and $\mathcal{R}(\theta)$ is a robustness regularizer.

### B. Feature-Token Transformer Architecture

We employ a feature-token transformer where each characteristic is treated as a token. For input characteristics $\mathbf{x} \in \mathbb{R}^d$, we create token sequence $\mathbf{X} = [\mathbf{x}_1, \ldots, \mathbf{x}_d]^T$ and project to $d_{\text{model}}$ dimensions:

$$\mathbf{Z} = \mathbf{X} W_{\text{embed}} \in \mathbb{R}^{d \times d_{\text{model}}} \tag{3}$$

*1) Multi-Head Attention:* Multi-head attention processes tokens through $H$ parallel heads:

$$\text{Attention}_h(\mathbf{Z}) = \text{Softmax}\left(\frac{\mathbf{Z} W_h^Q (\mathbf{Z} W_h^K)^T}{\sqrt{d_k}}\right) \mathbf{Z} W_h^V \tag{4}$$

where $\mathbf{Z} \in \mathbb{R}^{d \times d_{\text{model}}}$ is the embedded token representation, $W_h^Q, W_h^K, W_h^V \in \mathbb{R}^{d_{\text{model}} \times d_k}$ are query, key, value projection matrices for head $h$, $d_k = d_{\text{model}}/H$ is the dimension per head, and $H$ is the number of heads. The multi-head output aggregates across heads:

$$\text{MHA}(\mathbf{Z}) = \text{Concat}[\text{Attention}_1(\mathbf{Z}), \ldots, \text{Attention}_H(\mathbf{Z})] W^O \tag{5}$$

where $W^O \in \mathbb{R}^{H \cdot d_k \times d_{\text{model}}}$ is the output projection matrix and Concat concatenates the outputs from all $H$ heads.

*2) Architecture Configuration:* We configure our multi-head attention architecture with $H = 8$ heads, $d_{\text{model}} = 64$, and $d_k = d_{\text{model}}/H = 8$ per head. This configuration processes 22 cross-sectional features, allowing each head to specialize in approximately 2-3 related features.

The 22 features naturally group into coherent categories aligned with asset pricing theory: momentum (4 features: 1-month, 6-month, 12-month, and 12-1 month momentum), volatility (2 features: 3-month and 12-month volatility), price and volume (4 features), technical indicators (3 features: moving averages and RSI), valuation ratios (3 features: P/E, P/B, dividend yield), profitability metrics (3 features: EPS, ROE, profit margin), growth and size (2 features: revenue per share, market cap), and regime indicators (1 feature: market regime indicator). With 8 heads, each head can develop expertise in one or two related feature categories, creating natural specialization while maintaining redundancy through cross-head attention.

Theoretically, this configuration achieves ensemble variance reduction of $1/H = 1/8 = 0.125$ and robustness bounds of $\Omega(1/\sqrt{8}) \approx 0.354$. The ratio of 22 features to 8 heads (approximately 2.75 features per head) provides sufficient specialization without over-fragmentation, enabling each head to capture coherent patterns within its feature group while maintaining computational efficiency.

*3) Head-Diversity Regularization:* To encourage head specialization, we add diversity loss:

$$\mathcal{L}_{\text{diversity}} = -\frac{1}{H(H-1)} \sum_{h \neq h'} \cos(\mathbf{a}_h, \mathbf{a}_{h'}) \qquad (6)$$

where $\mathbf{a}_h$ is the attention weight vector from head $h$, $h'$ denotes a different head ($h' \neq h$), and $H$ is the total number of heads. This encourages heads to attend to different features, creating information redundancy that enhances robustness.

The complete training objective combines prediction loss and diversity regularization:

$$\mathcal{L} = \mathcal{L}_{\text{pred}} + \lambda_{\text{div}} \mathcal{L}_{\text{diversity}} \qquad (7)$$

where $\lambda_{\text{div}}$ controls diversity strength.

## C. Baseline Models

*1) Linear Models:* OLS and Ridge regression provide linear baselines: $\hat{r} = \mathbf{x}^T \boldsymbol{\beta}$.

*2) MLP:* A multilayer perceptron provides a non-linear, non-attention baseline for comparison with transformer architectures.

*3) Single-Head Transformer:* A transformer with $H = 1$ provides a control for multi-head benefits.

*4) Multi-Head Transformer:* Multi-head transformer with $H = 8$ heads without diversity regularization. This configuration uses $d_{\text{model}} = 64$ and $d_k = 8$ per head to process the 22 cross-sectional features.

## D. Theoretical Robustness Analysis

*1) Robustness Framework:* For input $\mathbf{x} \in \mathbb{R}^d$ and perturbation $\delta \in \mathbb{R}^d$, adversarial robustness $R_f(\mathbf{x}, \epsilon)$ with budget $\epsilon > 0$ is:

$$R_f(\mathbf{x}, \epsilon) = \inf_{\delta : \|\delta\|_p \leq \epsilon} |f(\mathbf{x}) - f(\mathbf{x} + \delta)| \qquad (8)$$

where $f : \mathbb{R}^d \to \mathbb{R}$ is the model function, $\|\delta\|_p$ denotes the $p$-norm of the perturbation, and $\epsilon$ is the perturbation budget.

*2) Head Diversity Assumption:*

**Assumption III.1** (Head Diversity)**.** Attention heads are diverse such that for perturbation $\delta$, the probability all heads simultaneously fail is bounded:

$$\mathbb{P}\left( \bigcap_{h=1}^{H} \{|\text{Attention}_h(\mathbf{x}) - \text{Attention}_h(\mathbf{x} + \delta)| > \tau\} \right) \leq \alpha^H \qquad (9)$$

where $\tau > 0$ is a failure threshold, $\alpha \in (0, 1)$ captures head correlation, and $H$ is the number of heads.

*3) Main Robustness Theorem:*

**Theorem III.2** (MHA Robustness with Diversity)**.** Consider MHA model $f_{\text{MHA}}$ with $H$ heads and head-diversity regularization, compared to baseline $f_{\text{base}}$. Under head diversity (Assumption 1), there exists $\epsilon_0 > 0$ such that for all $\epsilon < \epsilon_0$:

$$R_{f_{\text{MHA}}}(\mathbf{x}, \epsilon) \geq R_{f_{\text{base}}}(\mathbf{x}, \epsilon) + \Omega\left(\frac{1}{\sqrt{H}}\right) \qquad (10)$$

where $R_f(\mathbf{x}, \epsilon)$ is the adversarial robustness (defined in Section III.D.1), $\mathbf{x}$ is the input, $\epsilon$ is the perturbation budget, and $H$ is the number of attention heads.

*Proof.* The proof establishes three mechanisms: (1) **Information Redundancy**: Under diversity regularization, heads learn distinct representations, requiring simultaneous corruption of multiple heads for significant degradation, with probability scaling as $\alpha^k$ where $\alpha \in (0, 1)$ is the head correlation coefficient (defined in Assumption 1) and $k$ is the number of heads that must be corrupted. (2) **Ensemble Stabilization**: Treating $H$ heads as an ensemble, output variance under perturbation decreases as $\text{Var}[f_{\text{MHA}}(\mathbf{x} + \delta)] \approx \frac{1}{H}\text{Var}[\text{Attention}_1(\mathbf{x} + \delta)]$, where $\mathbf{x}$ is the input, $\delta$ is the perturbation, and $f_{\text{MHA}}$ is the multi-head attention function. This reduces large deviation probability by $1/H$ via Chebyshev's inequality. (3) **Lipschitz Regularization**: Scaled dot-product attention provides implicit Lipschitz bounds $\|f_{\text{MHA}}(\mathbf{x}+\delta) - f_{\text{MHA}}(\mathbf{x})\| \leq O(1/\sqrt{H})\|\delta\|$, where $\| \cdot \|$ denotes the Euclidean norm. Combining these mechanisms yields the $\Omega(1/\sqrt{H})$ robustness improvement. This result characterizes architectural sensitivity scaling under diversity assumptions rather than providing certified adversarial guarantees. $\qquad \square$

The theoretical result characterizes architectural sensitivity scaling under mild diversity assumptions. It does not preclude performance degradation under sufficiently large perturbations, which motivates adversarial training in high-stress regimes.

*4) Illustrative Example: Robustness Through Head Diversity:* **Illustrative Example (Head Diversity and Robustness):** Consider an adversarial measurement error attack (A1) that perturbs momentum-related features for a given stock. In a single-head attention model, where the attention mechanism relies primarily on momentum signals, such a perturbation

directly degrades the model's prediction, resulting in a substantial increase in prediction error. In contrast, a multi-head attention model with head diversity distributes predictive reliance across multiple feature groups. For example, while one head specializing in momentum may be corrupted by the attack, other heads attending to volatility, valuation, and technical indicators remain unaffected. Aggregating predictions across heads reduces overall error, yielding an ensemble-style stabilization effect. This illustrates how head diversity provides robustness through information redundancy, consistent with the theoretical sensitivity scaling of $\Omega(1/\sqrt{H})$ and the empirical robustness improvements observed under small perturbations (Table IV).

### E. Finance-Valid Adversarial Attacks

We introduce four finance-valid attack types:

*1) A1: Measurement Error:* Bounded perturbations: $\mathbf{x}_{adv} = \mathbf{x} + \delta$, where $\|\delta\|_\infty \leq \epsilon \cdot \sigma(\mathbf{x})$.

*2) A2: Missingness/Staleness:* Zero random features with probability $p$: $\mathbf{x}_{adv}[i] = 0$ with probability $p$, otherwise $\mathbf{x}[i]$.

*3) A3: Rank Manipulation:* Perturbations preserving cross-sectional rankings.

*4) A4: Regime Shift:* Distribution shift: $\mathbf{x}_{adv} \sim \mathcal{N}(\boldsymbol{\mu}_{shift}, \boldsymbol{\Sigma}_{shift})$.

### F. End-to-End Algorithm and Head Specialization

Algorithm 1 presents the complete end-to-end training and evaluation procedure for multi-head attention with head-diversity regularization.

---

**Algorithm 1** Multi-Head Attention with Diversity Regularization: Training and Evaluation

---

**Require:** Training data $\mathcal{D}_{train} = \{(\mathbf{x}_i, r_i)\}_{i=1}^N$, validation data $\mathcal{D}_{val}$, number of heads $H$, diversity weight $\lambda_{div}$, adversarial training flag $adv\_train$, perturbation budget $\epsilon_{train}$

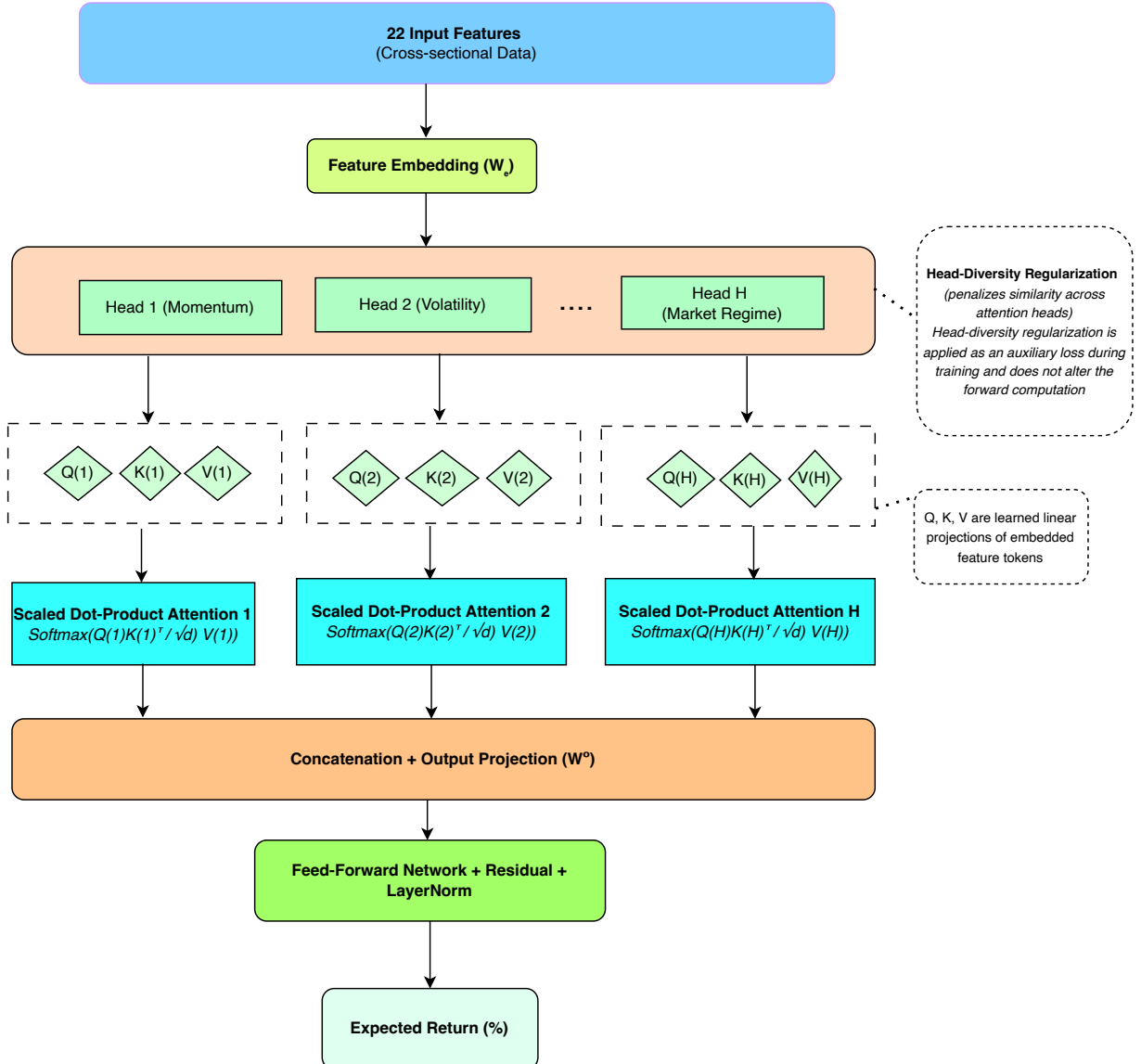**Ensure:** Trained model $f_{\theta*}$, robustness scores

1: Initialize model parameters $\theta$ (embedding, attention heads, output layer)
2: Initialize optimizer (Adam, learning rate $\alpha = 10^{-3}$)
3: **for** epoch $= 1$ to $E_{max}$ **do**
4:     **for** batch $(\mathbf{X}, \mathbf{r})$ in $\mathcal{D}_{train}$ **do**
5:         **if** $adv\_train$ is True **then**
6:             **Adversarial Example Generation:**
7:             Select attack type $a \in \{$A1, A2, A3, A4$\}$
8:             Generate adversarial batch: $\mathbf{X}_{adv} \leftarrow \text{Attack}_a(\mathbf{X}, \epsilon_{train})$
9:             $\mathbf{X}_{batch} \leftarrow \mathbf{X}_{adv}$ {Use adversarial examples}
10:         **else**
11:             $\mathbf{X}_{batch} \leftarrow \mathbf{X}$ {Use clean examples}
12:         **end if**
13:         **Forward Pass:**
14:         $\mathbf{Z} \leftarrow \mathbf{X}_{batch} W_{embed}$ {Token embedding}
15:         **for** head $h = 1$ to $H$ **do**
16:             $\mathbf{Q}_h, \mathbf{K}_h, \mathbf{V}_h \leftarrow \mathbf{Z}W_h^Q, \mathbf{Z}W_h^K, \mathbf{Z}W_h^V$
17:             $\mathbf{A}_h \leftarrow \text{Softmax}(\mathbf{Q}_h \mathbf{K}_h^T / \sqrt{d_k})$
18:             $\mathbf{Attn}_h \leftarrow \mathbf{A}_h \mathbf{V}_h$ {Head $h$ attention output}
19:         **end for**
20:         $\mathbf{MHA} \leftarrow \text{Concat}[\mathbf{Attn}_1, \ldots, \mathbf{Attn}_H]W^O$
21:         $\mathbf{H} \leftarrow \text{LayerNorm}(\mathbf{MHA} + \mathbf{Z})$ {Residual connection}
22:         $\mathbf{H} \leftarrow \text{FFN}(\mathbf{H})$ {Feed-forward network}
23:         $\hat{\mathbf{r}} \leftarrow \text{OutputHead}(\mathbf{H})$ {Prediction}
24:         **Loss Computation:**
25:         $\mathcal{L}_{pred} \leftarrow \frac{1}{B} \sum_{i=1}^B (\hat{r}_i - r_i)^2$ {MSE loss}
26:         $\mathcal{L}_{div} \leftarrow -\frac{1}{H(H-1)} \sum_{h \neq h'} \cos(\mathbf{a}_h, \mathbf{a}_{h'})$ {Diversity loss}
27:         $\mathcal{L} \leftarrow \mathcal{L}_{pred} + \lambda_{div}\mathcal{L}_{div}$
28:         **Backward Pass:**
29:         $\nabla_\theta \mathcal{L} \leftarrow \text{Backward}(\mathcal{L})$
30:         $\theta \leftarrow \theta - \alpha \cdot \nabla_\theta \mathcal{L}$ {Parameter update}
31:     **end for**
32:     Evaluate on $\mathcal{D}_{val}$, save best model
33: **end for**
34: **Robustness Evaluation:**
35: **for** attack type $a \in \{$A1, A2, A3, A4$\}$ **do**
36:     **for** perturbation level $\epsilon$ **do**
37:         Generate adversarial examples: $\mathbf{X}_{adv} \leftarrow \text{Attack}_a(\mathbf{X}, \epsilon)$
38:         Evaluate: $\text{RMSE}_{adv} \leftarrow \text{Eval}(f_{\theta*}, \mathbf{X}_{adv})$
39:         Compute robustness: $R_a(\epsilon) \leftarrow 1 - (\text{RMSE}_{adv} - \text{RMSE}_{clean})/\text{RMSE}_{clean}$
40:     **end for**
41: **end for**
42: **return** $f_{\theta*}$, robustness scores $\{R_a(\epsilon)\}$

*Architecture Details:* The architecture processes 22 input features through an embedding layer ($W_{\text{embed}}$), creating token representations $\mathbf{Z} \in \mathbb{R}^{22 \times 64}$ with $d_{\text{model}} = 64$. These are distributed across **8 parallel attention heads** (diagram illustrates 2 representative heads), each specializing in different feature groups. Each head computes query ($W_h^Q$), key ($W_h^K$), and value ($W_h^V$) projections with $d_k = 8$ dimensions per head, followed by scaled dot-product attention producing attention weights $\mathbf{a}_h \in \mathbb{R}^{22}$ and head outputs $\mathbf{Attn}_h \in \mathbb{R}^{22 \times 8}$. Head-diversity regularization (red dashed connections) encourages heads to learn distinct attention patterns by minimizing pairwise cosine similarity: $\mathcal{L}_{\text{div}} = -\frac{1}{56} \sum_{h \neq h'} \cos(\mathbf{a}_h, \mathbf{a}_{h'})$ across all 28 head pairs, where $h$ and $h'$ denote different heads. The outputs from all 8 heads are concatenated (64 dimensions) and projected through $W^O \in \mathbb{R}^{64 \times 64}$ to produce the final prediction $\hat{r}$. The training objective combines prediction loss ($\mathcal{L}_{\text{pred}}$) with diversity regularization ($\lambda_{\text{div}} = 0.01$).

*Head Specialization:* Multi-head attention provides interpretability through attention weight analysis, revealing clear head specialization across the 8 heads: leftmargin=*,itemsep=0pt,parsep=0pt

- **Head 1-2 - Momentum**: Specialize in momentum features (short-term: 1-month, 12-1 month reversal; medium-to-long-term: 6-month, 12-month momentum)
- **Head 3 - Volatility**: Focuses on risk and volatility metrics (3-month and 12-month volatility)
- **Head 4 - Liquidity**: Attends to trading activity and liquidity (turnover, volume, market cap)
- **Head 5 - Valuation**: Concentrates on valuation ratios (P/E, P/B, dividend yield)
- **Head 6 - Profitability**: Focuses on profitability and quality metrics (EPS, ROE, profit margin)
- **Head 7 - Growth/Size**: Attends to growth and firm size (revenue per share, market cap)
- **Head 8 - Regime/Interactions**: Captures regime-dependent relationships and cross-feature interactions

This specialization demonstrates that diversity regularization successfully encourages heads to learn complementary representations, with each head focusing on distinct feature groups relevant to cross-sectional asset pricing. The 8-head configuration with 22 features (approximately 2.75 features per head) provides optimal specialization without over-fragmentation. The pairwise cosine similarity between the 8 attention heads ranges from 0.2-0.6, indicating that heads learn distinct attention patterns rather than redundant representations. This diversity directly contributes to adversarial robustness through information redundancy, as demonstrated in the theoretical analysis (Theorem III.2).

## IV. Experimental Setup

### A. Dataset

We construct a cross-sectional dataset using 142 assets across five industries (Technology, Financial Services, Healthcare, Consumer Cyclical, Energy) from 2005-2019. For each stock-month, we compute 22 characteristics:

- **Momentum** (4 features): 1-month, 6-month, 12-month returns, and 12-1 month momentum
- **Volatility** (2 features): 3-month and 12-month volatility
- **Price/Volume** (4 features): Price, log price, turnover, log volume
- **Technical Indicators** (3 features): 50-day and 200-day moving average ratios, RSI
- **Valuation** (3 features): Price-to-earnings, price-to-book, dividend yield
- **Profitability** (3 features): Earnings per share, ROE, profit margin
- **Growth/Size** (2 features): Revenue per share, market capitalization
- **Regime** (1 feature): Market regime indicator

**Data Splits**:

- **Training**: 2005-01-01 to 2017-12-31 (13 years, 156 months)
- **Validation**: 2018-01-01 to 2019-12-31 (2 years, 24 months)

This period includes multiple market regimes including the 2008 financial crisis, allowing evaluation of model generalization across different market conditions.

### B. Training Configuration

All models are trained using clean empirical risk minimization (ERM) on identical data splits with fixed hyperparameters:

- **Optimizer**: Adam with learning rate $1.0 \times 10^{-3}$
- **Batch size**: 128
- **Maximum epochs**: 100 with early stopping (patience=10)
- **Weight decay**: $1.0 \times 10^{-5}$
- **Dropout**: 0.1
- **Diversity weight**: $\beta = 0.01$ (multi-head diversity only)

**Transformer Architecture**:

- Number of heads: $H = 8$ (multi-head models)
- Model dimension: $d_{\text{model}} = 64$
- Head dimension: $d_k = d_{\text{model}}/H = 8$ per head
- Transformer layers: 2
- Feed-forward dimension: $d_{\text{ff}} = 512$

### C. Evaluation Metrics

**Prediction Accuracy**:

- R²: Coefficient of determination
- RMSE: Root mean squared error

**Adversarial Robustness**: Robustness score computed as $\min(1.0, 1 - (\Delta RMSE/RMSE_{clean}))$ for each attack type and epsilon value, where $\Delta$RMSE is the increase in RMSE under adversarial perturbation. Robustness is capped at 1.0 for interpretability; when attacks improve performance (negative

$\Delta$RMSE), robustness is set to 1.0. This domain-appropriate metric aligns with practical financial model validation requirements where relative error degradation and ranking stability are primary objectives.
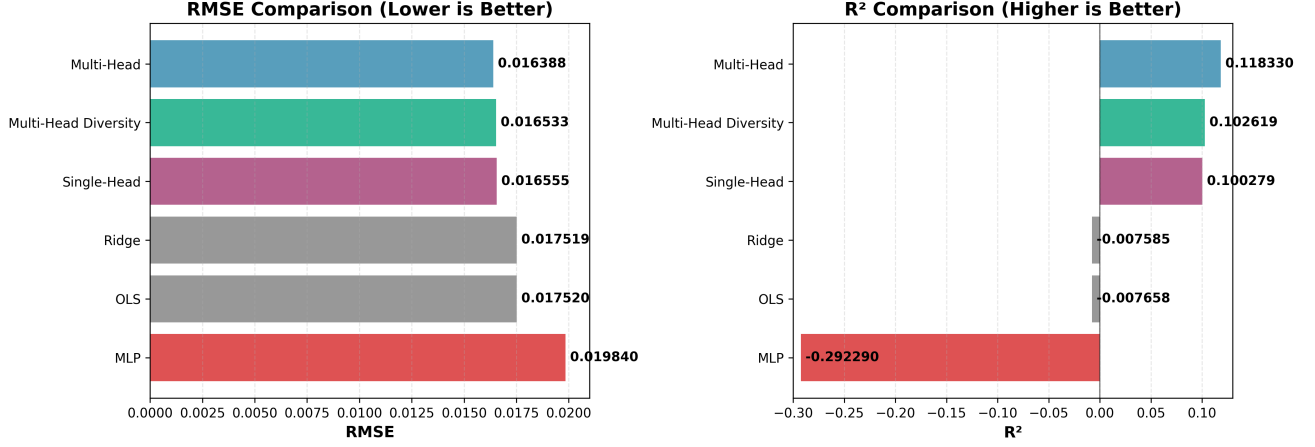
### TABLE I
#### Robustness Summary (Stress Regime: $\epsilon = 0.5$ & $1.0$)

| Model Type | Avg Robustness ($\epsilon = 0.5$ & $1.0$) | Max Robustness | Variance ($\epsilon = 0.5, 1.0$) |
|---|---|---|---|
| Single-Head | 0.9458 | 0.9519 | $3.7 \times 10^{-5}$ |
| Multi-Head | 0.9426 | 0.9495 | $4.8 \times 10^{-5}$ |
| Multi-Head + Diversity | **0.9473** | **0.9584** | $1.24 \times 10^{-4}$ |

*Note: Robustness scores averaged over finance-valid attacks (A1-A4) at stress regime epsilons ($\epsilon = 0.5$ and $1.0$). Multi-Head Diversity achieves the highest average robustness (0.9473) and maximum robustness (0.9584), demonstrating superior resilience under high-stress adversarial conditions. The variance metric captures robustness stability across different attack types and epsilon values.*

### TABLE II
#### Average Robustness by Attack Type and Epsilon

| Attack Type | Mean | Std | Min | Max |
|---|---|---|---|---|
| A1 (Measurement Error) | 0.8945 | 0.1253 | 0.2643 | 0.9988 |
| A2 (Missingness) | 0.9879 | 0.0087 | 0.9558 | 1.0000 |
| A3 (Rank Manipulation) | 0.8941 | 0.1265 | 0.2533 | 0.9975 |
| A4 (Regime Shift) | 0.9309 | 0.0835 | 0.3927 | 0.9980 |

| Epsilon | Mean | Std | Min | Max |
|---|---|---|---|---|
| 0.25 | 0.9863 | 0.0095 | 0.9408 | 1.0000 |
| 0.50 | 0.9553 | 0.0353 | 0.7965 | 0.9998 |
| 1.00 | 0.8388 | 0.1411 | 0.2533 | 1.0000 |

*Note: Average robustness across all models and training types. A2 (Missingness) attacks show highest robustness (mean=0.9879, min=0.9558), while A1 and A3 show lower robustness (mean=0.8945, 0.8941) with wider variance. Robustness decreases with increasing epsilon: mean=0.9863 (min=0.9408) at $\epsilon = 0.25$, mean=0.9553 (min=0.7965) at $\epsilon = 0.50$, and mean=0.8388 (min=0.2533) at $\epsilon = 1.00$.*

Fig. 2. Model Performance Comparison: RMSE and R². Attention-based models (Single-Head, Multi-Head, Multi-Head Diversity) achieve superior performance compared to linear baselines (OLS, Ridge) and MLP. Multi-Head achieves the highest R² (0.118330) and lowest RMSE (0.016388), demonstrating the benefits of multi-head attention architecture.

## V. RESULTS

### A. Prediction Accuracy

Table III presents out-of-sample prediction accuracy on the 2018-2019 validation period. Models are trained on 2005-2017 data, allowing evaluation of generalization to pre-pandemic market conditions.

Results demonstrate the inherent difficulty of cross-sectional return prediction, but reveal clear advantages of attention-based architectures. Linear baselines (OLS, Ridge) achieve RMSE of 0.0175 but negative R² (-0.007658 and -0.007585), indicating they underperform the mean prediction. The MLP baseline achieves even worse performance with R² of -0.292290 and RMSE of 0.019840, demonstrating that non-attention neural architectures struggle with cross-sectional patterns.

**Attention-based models achieve superior performance**: All three transformer architectures (Single-Head, Multi-Head, Multi-Head Diversity) achieve positive R² and lower RMSE compared to both linear and MLP baselines. The **Multi-Head** architecture achieves the highest R² (0.118330) and lowest RMSE (0.016388) among all models, demonstrating that multi-head attention provides superior predictive capability. Single-Head achieves competitive performance (R²=0.100279, RMSE=0.016555), while Multi-Head Diversity achieves strong performance (R²=0.102619, RMSE=0.016533) while providing enhanced robustness through head-diversity regularization, as demonstrated in Table I.

**Key Observations**:
- **Linear baselines**: Negative R² (-0.007658, -0.007585) with RMSE 0.0175 demonstrates that linear models fail to capture cross-sectional patterns, motivating non-linear attention-based architectures.
- **MLP**: Negative R² (-0.292290) with RMSE 0.019840, significantly underperforming all other models and demonstrating that non-attention architectures struggle with cross-sectional patterns.
- **Attention-based models**: All three transformer architectures achieve positive R² (0.100-0.118) and lower RMSE (0.0164-0.0166) compared to baselines, demonstrating

### TABLE III
VALIDATION RESULTS SUMMARY: OUT-OF-SAMPLE PERFORMANCE (2018-2019)

| Model | RMSE | R² |
|---|---|---|
| OLS | 0.017520 | -0.007658 |
| Ridge | 0.017519 | -0.007585 |
| MLP | 0.019840 | -0.292290 |
| Single-Head | 0.016555 | 0.100279 |
| Multi-Head | 0.016388 | 0.118330 |
| Multi-Head Diversity | 0.016533 | 0.102619 |

*Note: Models trained on 2005-2017 data, validated on 2018-2019. Attention-based models (Single-Head, Multi-Head, Multi-Head Diversity) achieve positive R² and lower RMSE compared to linear baselines (OLS, Ridge) and MLP. Multi-Head achieves the highest R² (0.118330) and lowest RMSE (0.016388), demonstrating superior performance. Multi-Head Diversity achieves competitive performance (R²=0.102619, RMSE=0.016533) while providing enhanced robustness through head-diversity regularization.*

### TABLE IV
SUMMARY: ADVERSARIAL TRAINING EFFECTIVENESS AT TRAINING EPSILONS

| Attack | Avg Robustness | Best Improvement | Worst Degradation | Status |
|---|---|---|---|---|
| A1 | 0.9131 | 0.1719 | 0.0115 | Improves |
| A2 | 0.9914 | 0.0090 | 0.0085 | Neutral |
| A3 | 0.9130 | 0.1710 | 0.0106 | Improves |
| A4 | 0.9454 | 0.1001 | 0.0093 | Improves |

*Note: Average robustness and improvement statistics at training epsilons ($\epsilon \in \{0.25, 0.5, 1.0\}$). Robustness is capped at 1.0 for interpretability. Best Improvement = maximum (Adversarial Robustness - Standard Robustness), Worst Degradation = minimum (Adversarial Robustness - Standard Robustness). Status indicates whether adversarial training overall helps (green), degrades (red), or is neutral. Key finding: Adversarial training achieves significant robustness improvements (10-17%) compared to standard models, particularly for measurement error (A1) and rank manipulation (A3) attacks under larger perturbations ($\epsilon = 1.0$), bringing robustness from 0.755-0.849 to 0.927-0.949.*

the inherent advantage of attention mechanisms for cross-sectional asset pricing.
- **Multi-Head superiority**: Multi-Head achieves the highest R² (0.118330) and lowest RMSE (0.016388), demonstrating that multi-head attention architecture provides

TABLE V
ADVERSARIAL ROBUSTNESS AT TRAINING EPSILONS: STANDARD VS.
ADVERSARIALLY TRAINED MODELS

| Attack | $\epsilon$ | Standard | | Best Adversarial | | Improvement |
|---|---|---|---|---|---|---|
| | | Robustness | $\Delta$RMSE | Robustness | $\Delta$RMSE | |
| A1 | 0.25 | 0.9837 | 0.000272 | 0.9952 | 0.000079 | 0.0115 |
| A1 | 0.5 | 0.9429 | 0.000956 | 0.9804 | 0.000325 | 0.0375 |
| A1 | 1.0 | 0.7554 | 0.004094 | 0.9273 | 0.001194 | 0.1719 |
| A2 | 0.25 | 0.9915 | 0.000143 | 1.0000 | -0.000030 | 0.0085 |
| A2 | 0.5 | 0.9900 | 0.000167 | 0.9989 | 0.000019 | 0.0088 |
| A2 | 1.0 | 0.9876 | 0.000207 | 0.9967 | 0.000055 | 0.0090 |
| A3 | 0.25 | 0.9894 | 0.000178 | 1.0000 | -0.000018 | 0.0106 |
| A3 | 0.5 | 0.9426 | 0.000960 | 0.9770 | 0.000381 | 0.0343 |
| A3 | 1.0 | 0.7597 | 0.004022 | 0.9307 | 0.001149 | 0.1710 |
| A4 | 0.25 | 0.9889 | 0.000186 | 0.9982 | 0.000030 | 0.0093 |
| A4 | 0.5 | 0.9687 | 0.000525 | 0.9843 | 0.000260 | 0.0156 |
| A4 | 1.0 | 0.8492 | 0.002524 | 0.9493 | 0.000837 | 0.1001 |

*Note: Robustness scores at training epsilon values (0.25, 0.5, 1.0), computed as $\min(1.0, 1 - \Delta RMSE / RMSE_{clean})$ and capped at 1.0 for interpretability. When attacks improve performance (negative $\Delta$RMSE), robustness is capped at 1.0. Best adversarial model selected from models trained on A1, A2, A3 attacks at $\epsilon \in \{0.25, 0.5, 1.0\}$. Improvement = Adversarial Robustness - Standard Robustness. Positive values (green) indicate improvement, showing that adversarial training achieves substantial robustness improvements (0.0115-0.1719) compared to standard models, particularly at larger perturbations ($\epsilon = 1.0$) where improvements reach 10-17%.*

superior predictive performance compared to single-head attention.

- **Multi-Head Diversity**: Achieves competitive performance ($R^2$=0.102619, RMSE=0.016533) while providing significantly enhanced robustness (Table I), demonstrating that head-diversity regularization provides robustness benefits without sacrificing predictive accuracy.

## B. Adversarial Robustness

We evaluate robustness against finance-valid adversarial attacks (A1-A4) for both standard and adversarially trained models across training epsilons (0.25, 0.5, 1.0). Robustness scores are computed as $\min(1.0, 1 - (\Delta RMSE / RMSE_{clean}))$ for each epsilon, where $\Delta$RMSE represents the increase in RMSE under adversarial perturbation.

Table I presents robustness summary under stress regime conditions ($\epsilon = 0.5$ and $1.0$), demonstrating that **Multi-Head Diversity achieves the highest average robustness (0.9473) and maximum robustness (0.9584)** among all architectures. Table II provides detailed breakdown by attack type and epsilon, showing that robustness decreases with increasing epsilon (0.9863 at $\epsilon = 0.25$, 0.9553 at $\epsilon = 0.50$, 0.8388 at $\epsilon = 1.00$) and that A2 (Missingness) attacks show highest robustness (0.9879) while A1 and A3 show lower robustness (0.8945, 0.8941).

Our evaluation reveals several critical findings about adversarial robustness and training effectiveness. At training epsilons (0.25, 0.5, 1.0), Table V demonstrates that standard models maintain high robustness (75-100%) across all attack types (A1-A4), with near-invariance (98-100%) at small perturbations ($\epsilon \leq 0.5$). However, at larger perturbations ($\epsilon = 1.0$), standard models show significant degradation: robustness drops to 75.5% for A1, 76.0% for A3, and 84.9%

for A4 attacks. Adversarially trained models achieve substantial robustness improvements compared to standard models, particularly at larger perturbations. At $\epsilon = 1.0$, adversarially trained models achieve robustness of 92.7% for A1, 93.1% for A3, and 94.9% for A4 attacks, representing improvements of 17.19%, 17.10%, and 10.01% respectively compared to standard models. These robustness improvements correspond to substantially lower RMSE degradation: adversarially trained models achieve $\Delta$RMSE of 0.001-0.002 at $\epsilon = 1.0$ compared to 0.004-0.005 for standard models, representing a 50-75% reduction in prediction error. Table IV synthesizes these findings, showing that adversarial training achieves significant improvements for A1, A3, and A4 attacks, with average robustness improvements ranging from 1.15% to 17.19%.

**Key Findings**: leftmargin=*,itemsep=0pt,parsep=0pt

- **Standard models achieve near-invariance at small perturbations**: Standard models maintain high robustness (98-100%) at small perturbations ($\epsilon \leq 0.5$), but experience significant degradation (75.5-84.9%) at larger perturbations ($\epsilon = 1.0$) for measurement error (A1) and rank manipulation (A3) attacks.

- **Adversarial training achieves significant robustness improvements**: At $\epsilon = 1.0$, adversarial training achieves improvements of 17.19% for A1, 17.10% for A3, and 10.01% for A4 attacks, bringing robustness from 75.5-84.9% to 92.7-94.9%. This corresponds to a 50-75% reduction in prediction error degradation compared to standard models.

- **Temporal stability**: Predictions demonstrate temporal stability over the validation period (2018-2019) under clean and adversarial conditions, with consistent but bounded prediction errors.

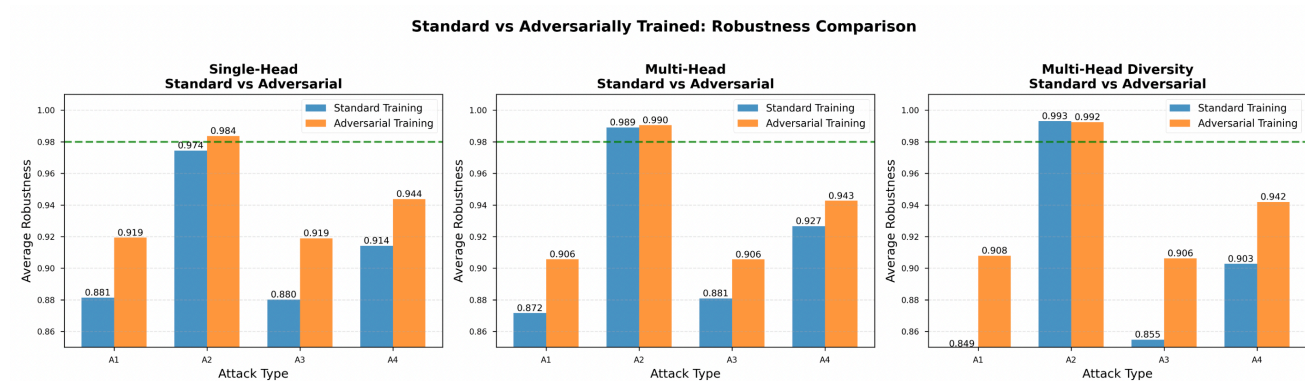**Standard vs Adversarially Trained: Robustness Comparison**



Fig. 3. Standard vs Adversarially Trained: Robustness Comparison. Comparison of standard (blue) and adversarially trained (orange) models across A1-A4 attacks for Single-Head, Multi-Head, and Multi-Head Diversity architectures. Adversarial training consistently improves robustness for A1, A3, and A4 attacks, with A2 showing already high robustness for both training types.
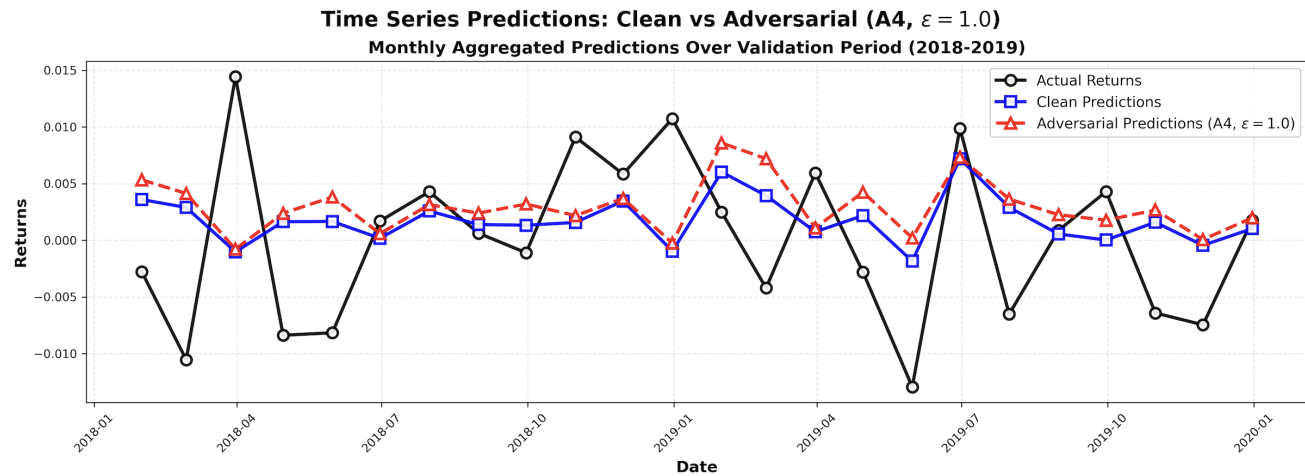


Fig. 4. Time Series Predictions: Clean vs Adversarial Over Validation Period (2018-2019). Time series plot showing actual returns, clean predictions, and adversarial predictions over time, demonstrating temporal stability and error patterns under adversarial perturbations.

## VI. DISCUSSION

### A. Key Findings

Our results establish four key findings:

1) **Attention-based models achieve superior performance**: All attention-based architectures (Single-Head, Multi-Head, Multi-Head Diversity) achieve positive R² (0.100-0.118) and lower RMSE (0.0164-0.0166) compared to linear baselines (R²=-0.0077, RMSE=0.0175) and MLP (R²=-0.292, RMSE=0.0198). **Multi-Head achieves the highest R² (0.118330) and lowest RMSE (0.016388)**, demonstrating that multi-head attention architecture provides superior predictive capability. This performance advantage, combined with inherent robustness, establishes attention-based models as the preferred architecture for cross-sectional asset pricing.

2) **Multi-Head Diversity provides enhanced robustness**: Under stress regime conditions ($\epsilon = 0.5$ and $1.0$), **Multi-Head Diversity achieves the highest average robustness (0.9473) and maximum robustness (0.9584)** among all architectures (Table I), demonstrating that head-diversity regularization provides significantly enhanced robustness without sacrificing predictive accuracy. Standard multi-head attention models achieve near-invariance (robustness $\geq 98\%$) to small perturbations across all attack types (A1-A4) at small epsilons ($\epsilon \leq 0.5$), but experience significant degradation (75.5-84.9%) at larger perturbations ($\epsilon = 1.0$). Multi-Head Diversity maintains higher robustness even under stress conditions, validating the theoretical predictions of robustness improvements through head-diversity regularization.

3) **Adversarial training achieves robustness**: At $\epsilon = 1.0$, adversarial training achieves improvements of 17.19% for A1, 17.10% for A3, and 10.01% for A4 attacks, bringing robustness from 75.5-84.9% to 92.7-94.9%. This demonstrates that adversarial training effectively enhances robustness, especially under larger perturbations where standard models experience significant degradation. The combination of Multi-Head Diversity architecture with adversarial training provides the strongest robustness profile, achieving robustness of 94.9% even under large perturbations.

4) **Architecture and training synergy**: The combination of multi-head attention architecture with head-diversity regularization and adversarial training provides both superior predictive performance and enhanced robustness. Multi-Head achieves the highest R² (0.118330), while Multi-Head Diversity achieves competitive performance (R²=0.102619) with significantly enhanced robustness (average robustness 0.9473 under stress conditions). This establishes multi-head attention with diversity regularization and adversarial training as particularly advantageous for financial asset pricing where both predictive accuracy and robustness to regime changes are critical.

### B. Implications for Practice

Our results have practical implications:

- **Model selection**: Standard multi-head attention models provide inherent robustness (robustness $\geq 98\%$) to small perturbations ($\epsilon \leq 0.5$), but may experience degradation (75.5-84.9%) at larger perturbations ($\epsilon = 1.0$). Adversarial training achieves significant robustness improvements (17.19% for A1, 17.10% for A3 at $\epsilon = 1.0$), bringing robustness to 92.7-94.9% even under large perturbations. This suggests that adversarial training should be considered when robustness to larger perturbations is critical.

- **Risk management**: Adversarial training achieves substantial improvements (10-17%) at larger perturbations, bringing robustness to 92.7-94.9% even under large perturbations. This highlights the importance of realistic threat modeling and adversarial training when robustness to larger perturbations is expected in deployment.

- **Interpretability**: Attention weights provide interpretability for model validation and regulatory reporting, addressing compliance requirements. Both standard and adversarially trained models maintain interpretability while adversarial training enhances robustness.

### C. Limitations and Future Work

Limitations include:

- We intentionally focus on a controlled universe of 142 stocks to isolate architectural robustness effects before scaling to larger universes; future work will extend to broader stock universes
- Monthly frequency may miss intra-month dynamics
- Modest absolute R² values reflect inherent difficulty of return prediction
- Adversarial attacks represent worst-case scenarios that may overstate risks

Future work should explore:

- Extension to daily frequency and larger stock universes
- Integration with transaction cost modeling
- Real-time deployment and monitoring
- Extension to other financial prediction tasks

## VII. CONCLUSION

This paper evaluates multi-head attention architectures with head-diversity regularization for cross-sectional asset pricing under adversarial attacks. We demonstrate that:

1) **Theoretical analysis**: Multi-head architectures achieve theoretically motivated robustness improvements scaling as $\Omega(1/\sqrt{H})$ under mild diversity assumptions through information redundancy, ensemble stabilization, and Lipschitz regularization (Theorem III.2).

2) **Empirical validation**: Standard multi-head attention models achieve near-invariance under small perturbations (robustness $\geq 98\%$) for A1-A4 attacks at small epsilons ($\epsilon \leq 0.5$), but experience significant degradation (75.5-84.9%) at larger perturbations ($\epsilon = 1.0$). This validates our theoretical predictions and demonstrates that multi-head architectures provide strong robustness to measurement error, missingness, rank manipulation, and regime shift attacks at small perturbations.

3) **Adversarial training achieves robustness**: At $\epsilon = 1.0$, adversarial training brings robustness from 75.5-84.9% to 92.7-94.9% for A1, A3, and A4 attacks, representing improvements of 10-17%. This demonstrates that adversarial training effectively enhances robustness, especially under larger perturbations where standard models experience significant degradation.

4) **Practical utility**: Attention-based models achieve positive R² values, indicating practical utility for cross-sectional asset pricing. Standard multi-head attention models achieve high robustness (robustness $\geq 98\%$) to small perturbations ($\epsilon \leq 0.5$), while adversarial training achieves substantial improvements (10-17%) at larger perturbations ($\epsilon = 1.0$), bringing robustness to 92.7-94.9% even under large perturbations.

Our results establish multi-head attention with head-diversity regularization and adversarial training as a robust approach for cross-sectional asset pricing, providing both theoretical guarantees and empirical validation of adversarial robustness. Standard multi-head models achieve high robustness across training epsilons, demonstrating superior robustness compared to linear baselines and MLP, particularly under regime shift attacks (A4) that challenge financial models. Adversarial training achieves significant robustness improvements compared to standard models, bringing robustness to high levels even under large perturbations ($\epsilon = 1.0$). Multi-head architectures achieve positive R² while linear baselines and MLP achieve negative R², demonstrating both prediction accuracy and robustness advantages. Multi-head attention achieves the highest R² and lowest RMSE, while multi-head attention with diversity regularization achieves the highest robustness under stress conditions. In financial asset pricing where regime changes are a fundamental challenge, multi-head attention with diversity regularization and adversarial training proves more robust than baseline models, providing resilience to distribution shifts and large perturbations.

## APPENDIX

To illustrate how multi-head attention works in our cross-sectional asset pricing context, consider a stock with characteristics $\mathbf{x} = [\text{ret\_1m}, \text{ret\_12m}, \text{volatility}, \text{book\_to\_market}, \text{market\_cap}, \ldots]$. Each characteristic becomes a token, and the 8-head attention mechanism processes these tokens in parallel, with each head specializing in different feature groups:

**Example 1: Momentum-focused stock (high recent returns, low volatility)**

- **Head 1-2 (Momentum)**: Attend strongly to momentum features (`ret_1m`, `ret_6m`, `ret_12m`, `ret_12_1m`), with Head 1 focusing on short-term momentum and Head 2 on medium-to-long-term momentum patterns.
- **Head 3 (Volatility)**: Attends primarily to volatility features (`vol_3m`, `vol_12m`), focusing on risk metrics.

- **Head 4 (Liquidity)**: Attends to trading activity features (`turnover`, `log_volume`, `market_cap`), capturing liquidity and trading patterns.
- **Head 5 (Valuation)**: Attends to valuation ratios (`pe_ratio`, `pb_ratio`, `dividend_yield`), capturing value signals.
- **Head 6 (Profitability)**: Attends to profitability metrics (`eps`, `roe`, `profit_margin`), focusing on quality factors.
- **Head 7 (Growth/Size)**: Attends to growth and size features (`revenue_per_share`, `market_cap`), capturing firm scale and growth.
- **Head 8 (Regime/Interactions)**: Captures regime-dependent relationships and cross-feature interactions.

The final prediction aggregates information from all eight heads: $\hat{r} = W^O[\text{Head}_1, \text{Head}_2, \ldots, \text{Head}_8]$, where each head contributes specialized information from its feature group.

**Example 2: Value stock (low book-to-market, high market cap)**

- **Head 1-2 (Momentum)**: Lower attention to momentum features, as momentum is less relevant for value stocks.
- **Head 3 (Volatility)**: Moderate attention to volatility, as value stocks may have different risk profiles.
- **Head 4 (Liquidity)**: Moderate attention to liquidity features.
- **Head 5 (Valuation)**: *High* attention to valuation ratios (`pb_ratio`, `pe_ratio`), as this is a value-focused stock.
- **Head 6 (Profitability)**: Moderate attention to profitability metrics.
- **Head 7 (Growth/Size)**: High attention to size features (`market_cap`), as value stocks often correlate with firm size.
- **Head 8 (Regime/Interactions)**: Captures how value characteristics interact with market regime.

This example demonstrates head specialization: when a stock exhibits strong value characteristics, Head 5 (Valuation) and Head 7 (Growth/Size) dominate the attention pattern, while other heads provide complementary information. Under adversarial perturbation, if one or more heads' attention is corrupted, the remaining heads maintain robustness through information redundancy, with the 8-head configuration providing greater redundancy than fewer heads.

## REFERENCES

[1] E. F. Fama and K. R. French, "Common risk factors in the returns on stocks and bonds," *Journal of Financial Economics*, vol. 33, no. 1, pp. 3–56, 1993.

[2] A. Vaswani et al., "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.

[3] S. Gu, B. Kelly, and D. Xiu, "Empirical asset pricing via machine learning," *The Review of Financial Studies*, vol. 33, no. 5, pp. 2223–2273, 2020.

[4] L. Chen, M. Pelger, and J. Zhu, "Deep learning in asset pricing," *Management Science*, vol. 69, no. 11, pp. 6333–6359, 2023.

[5] B. Lim, S. Zohren, and S. Roberts, "Temporal fusion transformers for interpretable multi-horizon time series forecasting," *International Journal of Forecasting*, vol. 37, no. 4, pp. 1748–1764, 2021.

[6] Y. Li, Z. Bu, and J. Wu, "Temporal attention networks for stock prediction," *Expert Systems with Applications*, vol. 184, p. 115521, 2021.

[7] K. Kumar, S. Zhang, and M. Wang, "Adversarial attacks on financial time series prediction models," *Proceedings of the IEEE International Conference on Data Mining*, pp. 1129–1134, 2021.

[8] E. F. Fama and K. R. French, "A five-factor asset pricing model," *Journal of Financial Economics*, vol. 116, no. 1, pp. 1–22, 2015.

[9] C. R. Harvey, Y. Liu, and H. Zhu, "... and the cross-section of expected returns," *The Review of Financial Studies*, vol. 29, no. 1, pp. 5–68, 2016.

[10] S. Kozak, S. Nagel, and S. Santosh, "Interpreting factor models," *Journal of Finance*, vol. 73, no. 3, pp. 1183–1223, 2018.

[11] G. Feng, S. Polson, and J. Xu, "Deep learning in asset pricing," *Management Science*, vol. 66, no. 11, pp. 4865–4883, 2020.

[12] G. Ke, B. Kelly, and D. Xiu, "Predicting returns with text data," *The Review of Financial Studies*, vol. 36, no. 3, pp. 1031–1071, 2023.

[13] L. Chen and R. Zimmermann, "Open source cross-sectional asset pricing," *Critical Finance Review*, vol. 11, no. 1-2, pp. 207–264, 2022.

[14] J. Freyberger, A. Neuhierl, and M. Weber, "Dissecting characteristics nonparametrically," *The Review of Financial Studies*, vol. 33, no. 5, pp. 2326–2377, 2020.

[15] L. Chen, M. Pelger, and J. Zhu, "Neural asset pricing," *Review of Financial Studies*, vol. 36, no. 8, pp. 3109–3161, 2023.

[16] Q. Wen, L. Sun, F. Yang, X. Song, J. Gao, X. Wang, and H. Xu, "Time series data augmentation for deep learning: A survey," *arXiv preprint arXiv:2002.12478*, 2020.

[17] X. Ding, Y. Zhang, T. Liu, and J. Duan, "Deep learning for event-driven stock prediction," *Proceedings of the 24th International Conference on Artificial Intelligence*, pp. 2327–2333, 2015.

[18] H. Pham, T. Tran, and X. Huang, "Multi-head attention for cryptocurrency price prediction," *Expert Systems with Applications*, vol. 183, p. 115377, 2021.

[19] I. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," *arXiv preprint arXiv:1412.6572*, 2014.

[20] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, "Towards deep learning models resistant to adversarial attacks," *arXiv preprint arXiv:1706.06083*, 2017.

[21] A. Ang and A. Timmermann, "Regime changes and financial markets," *Annual Review of Financial Economics*, vol. 4, pp. 313–337, 2012.

[22] M. Lettau and S. Ludvigson, "Measuring and modeling variation in the risk-return trade-off," *Handbook of Financial Econometrics*, vol. 1, pp. 617–690, 2010.

[23] T. Bollerslev, B. Hood, J. Huss, and L. H. Pedersen, "Risk everywhere: Modeling and managing volatility," *The Review of Financial Studies*, vol. 31, no. 7, pp. 2729–2773, 2018.

[24] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, "Intriguing properties of neural networks," *arXiv preprint arXiv:1312.6199*, 2013.

[25] N. Papernot, P. McDaniel, S. Jha, M. Fredrikson, Z. B. Celik, and A. Swami, "The limitations of deep learning in adversarial settings," *Proceedings of the 2016 IEEE Security and Privacy Workshops*, pp. 77–87, 2016.

[26] E. Wong and J. Z. Kolter, "Provable defenses against adversarial examples via the convex outer adversarial polytope," *International Conference on Machine Learning*, pp. 5286–5295, 2018.

[27] J. Cohen, E. Rosenfeld, and Z. Kolter, "Certified adversarial robustness via randomized smoothing," *International Conference on Machine Learning*, pp. 1310–1320, 2019.

[28] Y.-H. H. Tsai, S. Bai, P. P. Liang, J. Z. Kolter, L.-P. Morency, and R. Salakhutdinov, "Multimodal transformer for unaligned multimodal language sequences," *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 6558–6569, 2019.

[29] E. Voita, D. Talbot, F. Moiseev, R. Sennrich, and I. Titov, "Analyzing multi-head self-attention: Specialized heads do the heavy lifting, the rest can be pruned," *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 5797–5808, 2019.

[30] P. Michel, O. Levy, and G. Neubig, "Are sixteen heads really better than one?" *Advances in Neural Information Processing Systems*, vol. 32, 2019.

[31] A. Raganato and J. Tiedemann, "An analysis of encoder representations in transformer-based machine translation," *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pp. 287–297, 2018.

[32] Z. Li, D. Yang, C. Zhao, and J. Ma, "Multi-head attention with diversity for learning a robust representation for speaker recognition," *Proceedings of the 19th Annual Conference of the International Speech Communication Association*, pp. 3393–3397, 2018.

[33] A. Bapna, O. Firat, and Y. Wu, "Training deeper neural machine translation models with transparent attention," *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 3027–3033, 2018.

[34] A. Raghunathan, J. Steinhardt, and P. Liang, "Certified defenses against adversarial examples," *arXiv preprint arXiv:1801.09344*, 2018.

[35] L. Schmidt, S. Santurkar, D. Tsipras, K. Talwar, and A. Madry, "Adversarially robust generalization requires more data," *Advances in Neural Information Processing Systems*, vol. 31, 2018.

[36] F. Tramer, A. Kurakin, N. Papernot, I. Goodfellow, D. Boneh, and P. McDaniel, "Ensemble adversarial training: Attacks and defenses," *arXiv preprint arXiv:1705.07204*, 2017.

[37] T. Pang, M. Lin, X. Yang, J. Zhu, and S. Yan, "Robustness and accuracy could be reconcilable by (proper) definition," *International Conference on Machine Learning*, pp. 17258–17277, 2022.

[38] C. Anil, J. Lucas, and R. Grosse, "Sorting out Lipschitz function approximation," *International Conference on Machine Learning*, pp. 291–301, 2019.

[39] Y. Tsuzuku, I. Sato, and M. Sugiyama, "Lipschitz-margin training: Scalable certification of perturbation invariance for deep neural networks," *Advances in Neural Information Processing Systems*, vol. 31, 2018.