# A Deep Learning Approach for Trading Factor Residuals

Wo Long, Victor Xiao

Columbia University, New York, NY, USA

December 8, 2024

**Abstract**

The residuals in factor models prevalent in asset pricing presents opportunities to exploit the mis-pricing from unexplained cross-sectional variation for arbitrage. We performed a replication of the methodology of Guijarro-Ordonez et al. (2019) (G-P-Z) on Deep Learning Statistical Arbitrage (DLSA), originally applied to U.S. equity data from 1998 to 2016, using a more recent out-of-sample period from 2016 to 2024. Adhering strictly to point-in-time (PIT) principles and ensuring no information leakage, we follow the same data pre-processing, factor modeling, and deep learning architectures (CNNs and Transformers) as outlined by G-P-Z. Our replication yields unusually strong performance metrics in certain tests, with out-of-sample Sharpe ratios occasionally exceeding 10. While such results are intriguing, they may indicate model overfitting, highly specific market conditions, or insufficient accounting for transaction costs and market impact. Further examination and robustness checks are needed to align these findings with the more modest improvements reported in the original study.

## 1 Introduction

Deep Learning Statistical Arbitrage (DLSA), as developed by Guijarro-Ordonez et al. (2019) (hereafter G-P-Z), applies advanced machine learning techniques, specifically convolutional neural networks (CNNs) and Transformers, to factor-model residuals in an effort to identify profitable statistical arbitrage opportunities. G-P-Z demonstrated that these deep learning approaches could detect complex patterns in the cross-section of asset returns, resulting in meaningful economic gains over their 1998–2016 sample period.

In this paper, we reproduce the DLSA methodology using a more recent dataset spanning 2016 to 2024. We adhere strictly to point-in-time (PIT) data handling, forward-filling only, and ensure no data snooping or backward-looking adjustments that could contaminate our results. Our aim is to assess the robustness of the original approach when applied to a different market regime and to verify that the significant performance enhancements reported by G-P-Z are not dependent on a particular historical sample.

Surprisingly, our replication yields performance metrics that are even stronger than those initially documented by G-P-Z. In particular, we observe out-of-sample Sharpe ratios occasionally exceeding 10. While on the surface this is impressive, it raises cautionary flags about potential model overfitting, non-stationary market dynamics, or omitted transaction costs. Additional experiments and simulations are required to confirm the true economic significance of these findings.

## 2  Data

### 2.1  Data Sources and Scope

We follow the data construction methodology of G-P-Z as closely as possible. Equity price and return data are obtained from CRSP, and accounting information from Compustat. We use the three-month Treasury bill rate from the Kenneth French Data Library as the risk-free rate. Our sample focuses on relatively liquid U.S. equities (S&P 500), broadly following the selection criteria in the original DLSA paper, but for a more recent time horizon, 2016–2024.

### 2.2  Point-in-Time Data Handling and Missing Data

All data are processed using PIT principles. At any time $t$, only information known at or before $t$ is used. We strictly forward-fill missing values without any backward-looking imputations. Stocks with excessive missingness are excluded. Factor loadings, used to compute residual returns, are estimated only from historical data up to $t-1$.

# 3 Methodology

## 3.1 Factor Models and Residual Returns

To isolate residual returns, we employ factor models such as the Fama-French 5-factor model, PCA, and IPCA. These approaches decompose returns into systematic components and residuals. The residuals, $\epsilon_{n,t}$, are defined as:

$$\epsilon_{n,t} = R_{n,t} - \beta_{n,t-1}^\top F_t,$$

where $R_{n,t}$ is the asset's excess return, $\beta_{n,t-1}$ represents factor loadings estimated using a rolling window, and $F_t$ denotes factor realizations at time $t$. These residuals are designed to capture idiosyncratic returns, free from systematic risk exposure.

PCA-based factor models estimate latent factors by analyzing the correlation matrix over a 252-day rolling window. These factors summarize systematic variations in returns, while loadings are computed via regression over the past 60 days. The residuals should, in theory, be free of systematic risk exposures, leaving behind idiosyncratic components that may be more predictable.

## 3.2 Deep Learning Architecture: CNN + Transformer

We utilize a Convolutional Neural Network (CNN) combined with a Transformer architecture to detect patterns in residual return series follows G-P-Z. The CNN serves as a data-driven, flexible local filter, identifying patterns such as trends and reversals in localized segments of the data. It processes the input residual time series through a series of convolutional layers, extracting $D$ feature maps that quantify exposure to predefined basic patterns. These local features are then passed to the Transformer for global analysis.

The Transformer captures temporal dependencies between these local patterns using an attention mechanism. Each attention head learns specific global patterns, such as mean-reversion or trend-following behaviors, by assigning weights to interactions between different segments of the time series. This allows the model to represent residual returns in terms of global dependency structures, leveraging both short-term fluctuations and longer-term dynamics.
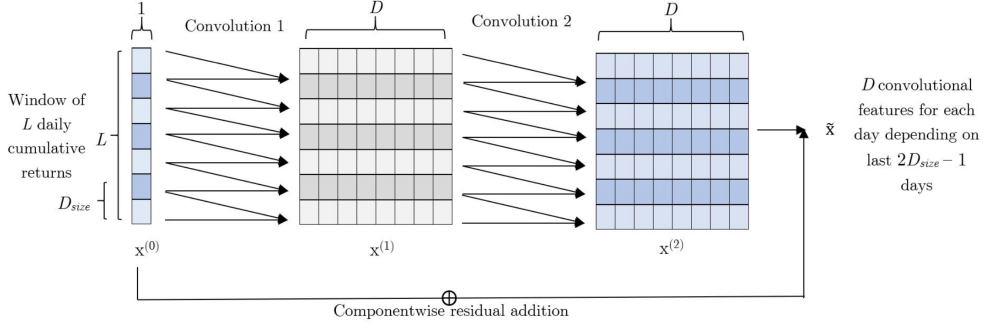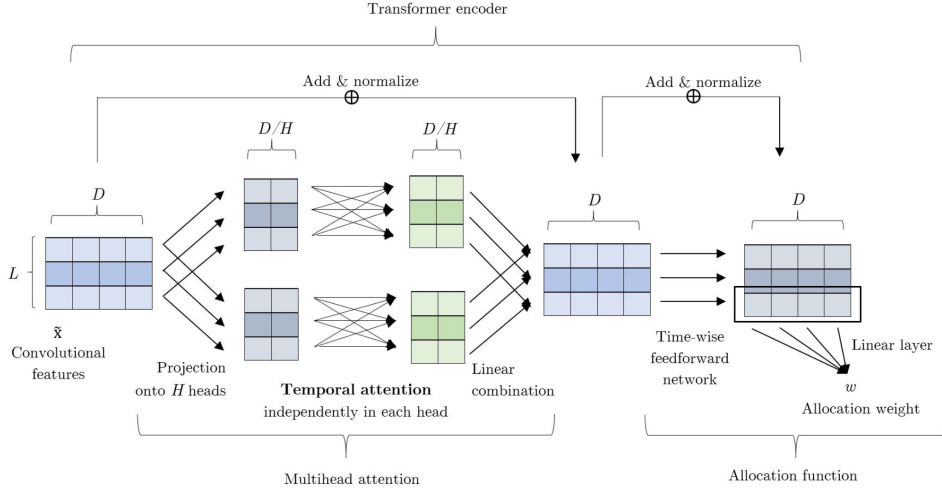
Figure 1: CNN Architecture



Figure 2: Transformer Architecture

## 3.3 Trading Strategy and Portfolio Construction

We use feedforward neural network architecture to map extracted signals to trading weights that aim to maximize out-of-sample Sharpe ratios. Constraints such as leverage, turnover, and short selling limitations are imposed as in the original framework.

We train and test the model on a rolling basis, using a historical window (e.g., 1000 days) for training and a subsequent 125-day window for out-of-sample testing.

# 4 Results

## 4.1 Preliminary Backtest Performance

Applying the G-P-Z methodology to 2016–2024 data, we unexpectedly find that the model often identifies patterns leading to very high risk-adjusted returns. In certain test windows, the out-of-sample Sharpe ratio exceeds 10, significantly surpassing the approximately 3.5–4 Sharpe ratios reported in the original study.

Figure 3 displays some of the performance metrics and return distributions generated by the strategy. While visually impressive, these results should be approached with skepticism until further verified.
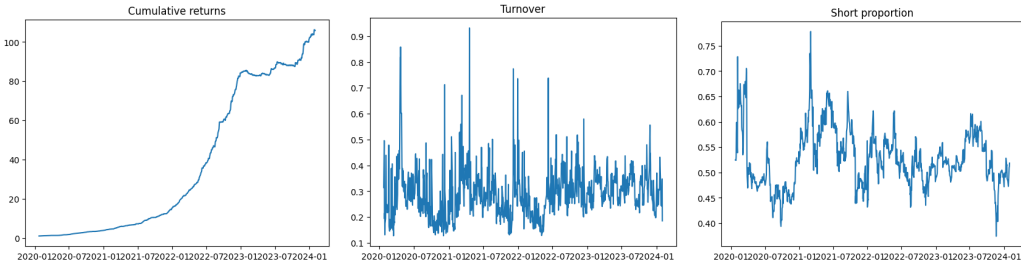


Figure 3: Cumulative Returns, Turnover, and Short Proportion

Figure 4 shows cumulative performance and related metrics over an out-of-sample testing window. Despite the high returns indicated, we must highlight that the absence of transaction costs and other frictions in this preliminary analysis might inflate the apparent profitability.
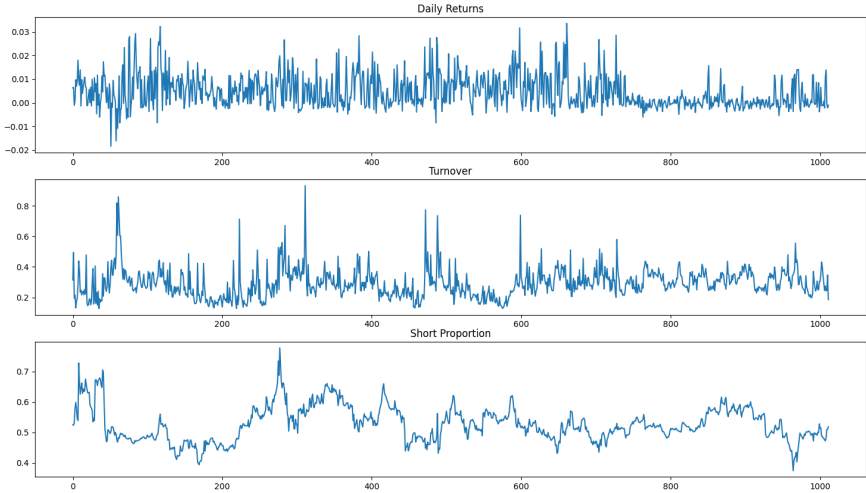


Figure 4: Additional Performance Visualization: Cumulative Returns and Key Metrics
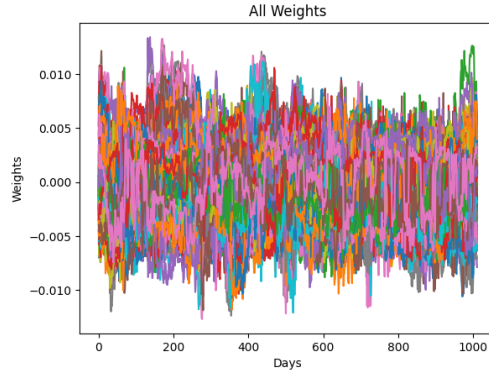
Figure 5: Further Performance Visualization with Alternative Windows

## 4.2   Interpretation and Caution

While these preliminary findings are enticing, several factors demand caution:

- **Model Overfitting:** The high out-of-sample Sharpe ratios may be a result of overfitting the model to patterns specific to the recent sample period. Further testing on different windows and adding regularization could help confirm the stability of these patterns.

- **Market Regime Changes:** The 2016–2024 period may contain market regimes drastically different from 1998–2016. Extraordinary conditions (e.g., post-crisis monetary policies, COVID-19 market dynamics) might have created exploitable short-term patterns.

- **Transaction Costs and Market Impact:** We have not yet incorporated realistic transaction costs, market impact, or short-sale fees. These frictions would likely reduce the realized Sharpe ratio and returns.

- **Data Integrity Checks:** We have implemented rigorous PIT data handling and taken care to avoid information leakage. However, the extreme results warrant an additional layer of verification. Even subtle forms of data leakage or survivorship bias could contribute to inflated performance metrics.

# 5    Discussion

Our replication of the DLSA methodology in a newer time period yields results far exceeding the original study's findings. Rather than viewing these outcomes as definitive evidence of improved model efficacy, we interpret them as a signal that more thorough robustness checks are needed.

For example, preliminary findings show extremely high Sharpe ratios, but further testing with realistic transaction cost modeling and shifting validation windows is required to confirm their robustness. Similarly, the impressive predictive accuracy might reflect model overfitting or exploitation of transient market conditions rather than a stable arbitrage opportunity. Additional tests using alternative market phases and out-of-sample data are essential to verify these results.

# 6    Summary

We reproduced the Deep Learning Statistical Arbitrage approach outlined by G-P-Z, applying it to a more recent dataset (2016–2024) while meticulously maintaining point-in-time data integrity and preventing information leakages. The resulting OOS Sharpe ratios occasionally exceed 10—an outcome that demands careful interpretation.

Our findings do not invalidate the original work, but rather highlight that market conditions, sample selection, and evaluation protocols can drastically influence performance metrics. Before any conclusions can be drawn about the superiority of these methods in recent times, future work must incorporate more realistic trading frictions, adjust for potential overfitting, and verify performance stability across multiple out-of-sample segments.

# Acknowledgments

# References

Fama, Eugene F. and Kenneth R. French (1993). Common risk factors in the returns on stocks and bonds. *Journal of Financial Economics*, 33(1):3–56.

Fama, Eugene F. and Kenneth R. French (2015). A five-factor asset pricing model. *Journal of Financial Economics*, 116(1):1–22.

Guijarro-Ordonez, Jorge, Markus Pelger, and Greg Zanotti (2019). Deep Learning Statistical Arbitrage. SSRN Working Paper. `https://ssrn.com/abstract=3862004`