

Deep Learning for Demand Forecasting in Retail Using Hybrid CNN-RNN and Transformer Models

Project Report

DL for Supply Chain Optimization in Retail

NAME: Raghulchellapandian Senthil Kumaran
NAME: Dongyoон Shin

UBID: raghulch
UBID: dongyoон

Table of Contents

1 Problem Statement	3
2 Objectives	3
3 Dataset	4
4 Methodology	4
4.1 Data Preparation and Exploration	4
4.1.1 Data Overview and Cleaning	4
4.1.2 Exploratory Data Analysis (EDA)	5
4.1.3 Feature Engineering	6
4.1.4 Final Dataset Construction	6
4.2 Synthetic Noise Injection (for Robustness Evaluation)	7
4.2.1 Noise Injection Procedure and Setup	7
4.2.2 Sensitivity and Correlation Analysis	7
4.2.3 Evaluation and Data Preparation	8
4.3 Development of the Linear Model and FNN	9
4.3.1 Model Architecture	9
4.3.2 Forecasts	10
4.3.3 Training	10
4.3.4 Metrics	11
4.4 Development of the Deep Learning Model	11
4.4.1 LSTM Architecture	11
4.4.2 GRU Architecture	12
4.4.3 Transformer Model	13
4.4.4 CNN-RNN Hybrid	13
4.5 Model Evaluation	14
4.5.1 LSTM Model	14
4.5.2 GRU Model	15
4.5.3 Transformer Model	15
4.5.4 CNN-RNN Model	16
4.6 Comparative Analysis	17
4.6.1 Clean Data	18
4.6.2 Noisy Data	18

5 Hyperparameters Tuning	19
5.1 Hyperparameter Setting	19
5.2 LSTM	19
5.3 GRU	20
5.4 Best Model Comparison	20
5.5 Comparison – Predicted VS Actual	21
5.5.1 Model Configuration and Evaluation Metrics	21
5.5.2 Bar Plot – Sample 0, Feature 0	22
5.5.3 Line Plot – Sample 0, Features 0–2	22
5.5.4 Scatter Plot – Actual vs. Predicted (All Models)	22
5.5.5 Box Plot – Prediction Error Distribution	23
5.6 TRANSFORMER	23
5.6.1 Model Architecture: TimeSeriesTransformer	23
5.6.2 Training and Setup	24
5.6.3 Evaluation Summary and Insights	24
5.6.4 Transformer Graph Analysis (Clean vs Noisy)	25
5.7 CNN-RNN (Conv1D + GRU)	25
5.7.1 Model Architecture	25
5.7.2 Training and Setup	26
5.7.3 Evaluation Summary and Observations	26
5.7.4 CNN-RNN Visualization and Analysis	27
5.8 Observation & Analysis (Transformer vs CNN-RNN)	27
5.9 Test Error Comparison	28
5.10 Final Comparison Summary (GRU, LSTM, Transformer, CNN-RNN)	28
5.11 GRU Model – Conclusion & Key Insights	29
5.12 LLaMA-style Forecast Error Summary (GRU Clean)	30
6 Conclusion	30
7 References	31
8 Contribution Table	31

ABSTRACT

Precise time series predictions are essential in the retail industry, where inventory, supply chain, and strategic planning depend on accurate forecasting. This project investigates the capability of deep learning models in predicting future sales over multiple time steps. We evaluate several architectures—including LSTM, GRU, Transformer, FFNN, hybrid CNN-RNN, and a baseline Linear model—on a multivariate time series dataset. Models are tested under both clean and noise-injected conditions to mimic real-world uncertainty and find the best model that is robust to noisy data and works well with clean data. Through a standardized experimental configuration and evaluation criteria, we explore the strengths and limitations of each model in both ideal and noisy settings.

1. PROBLEM STATEMENT

In the rapidly evolving and competitive retail sector, accurately forecasting future demand is crucial for maintaining appropriate stock levels, minimizing operational costs, and ensuring customer satisfaction. Traditional statistical forecasting models often struggle with the complexities of modern retail data, which is typically high-dimensional, noisy, and temporally dependent.

Real-world retail datasets commonly contain missing values, anomalies, and unexpected trends, making robust forecasting even more challenging. Deep learning models offer a compelling alternative by automatically learning temporal dependencies from raw multivariate time series data. However, the optimal architecture for multi-step forecasting—especially under both clean and noisy conditions—remains unclear. A systematic benchmarking study is thus necessary to understand the comparative advantages and limitations of each architecture in terms of accuracy, stability, and robustness to noise.

2. OBJECTIVES

This project aims to develop and evaluate deep learning-based models for multi-step time series forecasting in the context of retail sales. The specific objectives are:

- Develop robust forecasting models using deep learning architectures such as LSTM, GRU, Transformer, FFNN, and CNN-RNN hybrids.
- Evaluate model performance under clean and noisy input conditions to assess generalization and robustness.
- Compare forecast accuracy using metrics including MSE, RMSE, MAE, and R^2 .
- Interpret model behavior through forecast plots, epoch-wise metric curves, and step-wise error analysis across the prediction horizon.
- Enhance model performance through architectural refinement and hyperparameter tuning.
- Select the most appropriate model for real-world deployment in retail environments with incomplete or noisy data.

By achieving these goals, the study both benchmarks state-of-the-art time series forecasting models and highlights practical deployment considerations in real-world, noisy retail data scenarios.

3. DATASET

ITEM_ID	YEAR	SUPPLIER	ITEM_TYPE	ITEM_NAME	SALES	RETAIL_SALES	WAREHOUSE_SALES
307145	2009	9 LEGENDS1	11549 BRECKENRIDGE BEER	0	0	0	0
307146	2009	9 LEGENDS1	11549 BRECKENRIDGE BEER	0	0	0	0
307147	2009	9 SALVATORE	235026 MONCADA WINE	0	0	0	0
307148	2009	9 SANTINI INC	235026 MONCADA WINE	0	0	0	0
307149	2009	9 SBL INC	29961 CELERIUM THEOS	0	0	0	0
307150	2009	9 ARIAL WINE	368416 BENJAMIN H WINE	0.82	1	0	0
307151	2009	9 BIRCHWOOD	368416 BENJAMIN H WINE	0	0	0	0
307152	2009	9 PINEYARD	349442 1881 PALE WINE	0	0	0	0
307153	2009	9 RUMBA	368416 BENJAMIN H WINE	0	0	0	0
307154	2009	9 SAINTNICK C	351680 PLEISCHM LIQUOR	0.76	4	0	0
307155	2009	9 DOPPI INC	4027 GATTI'S AC BEER	0	0	0	0
307156	2009	9 DOPPI INC	4027 GATTI'S AC BEER	0	0	0	0
307157	2009	9 ELITE WINE	411584 ESSAY DRINK WINE	0	0	0	0
307158	2009	9 HARRIS CO.	402400 LADY O' CROWN LIQUOR	0.33	0	0	0
307159	2009	9 DOPPI INC	4027 GATTI'S AC BEER	0	0	0	0
307160	2009	9 LEISUREL	3090 CHIVELLS H/B BEER	5.36	2	36	0
307161	2009	9 LEISUREL	3090 CHIVELLS H/B BEER	0.49	0	0	0
307162	2009	9 LEISUREL	304400 MARRAGAOL BEER	0	0	0	0
307163	2009	9 BACARDI	70772 LEBON G/L LIQUOR	0.48	1	0	0
307164	2009	9 BACARDI	70772 LEBON G/L LIQUOR	0.71	0	0	0
307165	2009	9 THE WINE	76279 ALMADEN WINE	0	0	11	0
307166	2009	9 VINTAGE WINE	83126 STONEWICH WINE	0	0	0	0
307167	2009	9 T & J GALU	83126 STONEWICH WINE	0	0	0	0
307168	2009	9 PERIODIC	830807 KENWOOD WINE	0.57	2	0	0
307169	2009	9 PERIODIC	830807 KENWOOD WINE	0.43	0	0	0
307170	2009	9 ANHEUSER	97040 SPATEN PILZ BEER	34.04	22	149	0
307171	2009	9 ANHEUSER	97040 SPATEN PILZ BEER	0.11	0	0	0
307172	2009	9 DOPPI INC	4027 GATTI'S AC BEER	0	0	1	0
307173	2009	9 ANHEUSER	97051 STELLA ARTOIS BEER	372.45	313	2066.80	0
307174	2009	9 ANHEUSER	97051 STELLA ARTOIS BEER	7.79	4	0	0
307175	2009	9 RELIABLE C	97060 S SMITH H/B BEER	0	0	0	0
307176	2009	9 RELIABLE C	97060 S SMITH H/B BEER	0	0	0	0

Figure 1: Sample Dataset

The dataset used in this project is the **Warehouse and Retail Sales** dataset, publicly available from [Data.gov](#). It contains over 307,000 monthly sales records, organized by item and supplier. The dataset includes features such as:

- YEAR, MONTH
- SUPPLIER, ITEM TYPE
- RETAIL SALES, RETAIL TRANSFERS, WAREHOUSE SALES

Its size and complexity make it well-suited for training deep learning models that capture demand trends across multiple categories and suppliers over time.

4. METHODOLOGY

In this section, we describe the complete methodology pipeline adopted in this project. The process begins with the collection and cleaning of the raw retail sales dataset. We then perform extensive exploratory data analysis (EDA) to understand the data distribution, identify anomalies, and extract trends and patterns.

Next, we apply feature engineering techniques such as lag creation, log transformation, and normalization to prepare the data for supervised learning. We also simulate real-world uncertainty by injecting Gaussian noise to assess model robustness. Do hyperparameter tuning to find the best model. Finally, we train and evaluate multiple deep learning architectures to compare performance under clean and noisy conditions.

4.1. Data Preparation and Exploration

4.1.1. Data Overview and Cleaning

The dataset consisted of 307,645 records with 9 features representing time, sales, product types, and supplier information. It was read using `pandas`, with inspection performed via `.info()`, `.head()`, and `.describe()` for data types, missing entries, and summary statistics.

Missing or invalid values (e.g., zeros or negatives in `RETAIL_SALES`) were replaced with NaNs and imputed using `IterativeImputer` with a `RandomForestRegressor` estimator. Log transformations (`log1p`) reduced skewness while preserving zero entries. Interquartile range (IQR) filtering removed 2–3% of outliers.

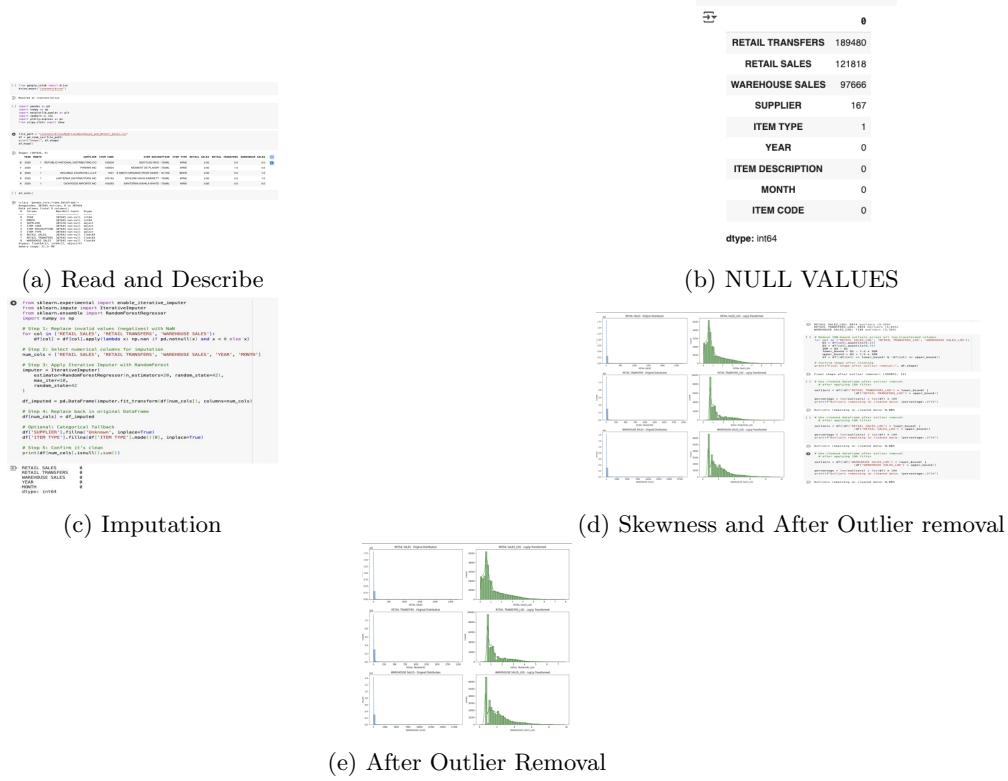


Figure 2: Data Cleaning and Preprocessing Steps

4.1.2. Exploratory Data Analysis (EDA)

Frequency analysis revealed WINE, LIQUOR, and BEER as dominant categories. Seasonal patterns appeared in time series plots, especially strong Q4 peaks. Boxplots by month confirmed this. Rolling averages further smoothed short-term fluctuations. Pairplots revealed weak correlations between RETAIL_SALES, INVENTORY, and UNITS SOLD.

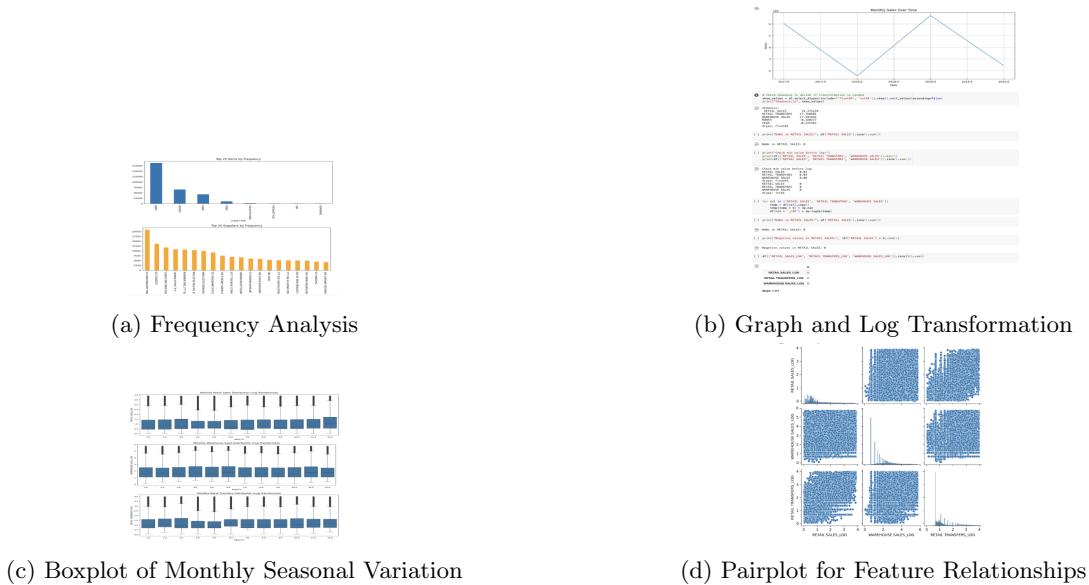
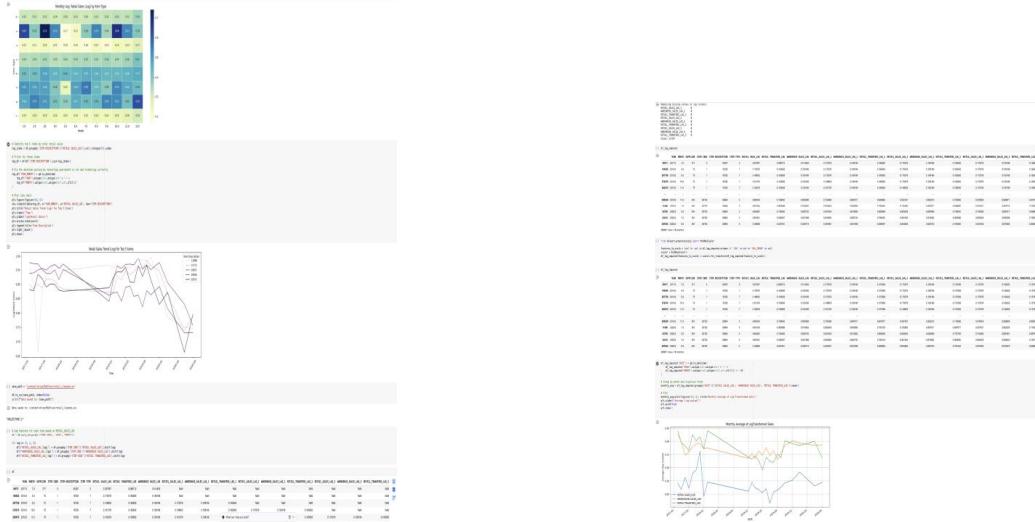


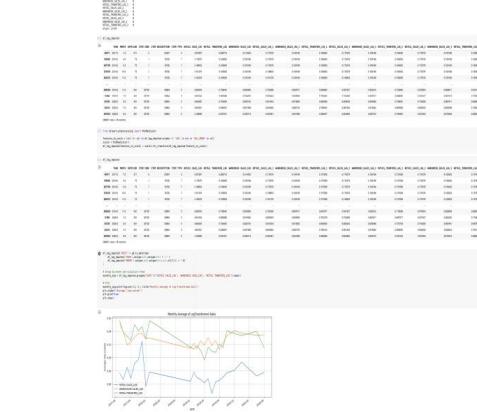
Figure 3: Exploratory Data Analysis Visualizations

4.1.3. Feature Engineering

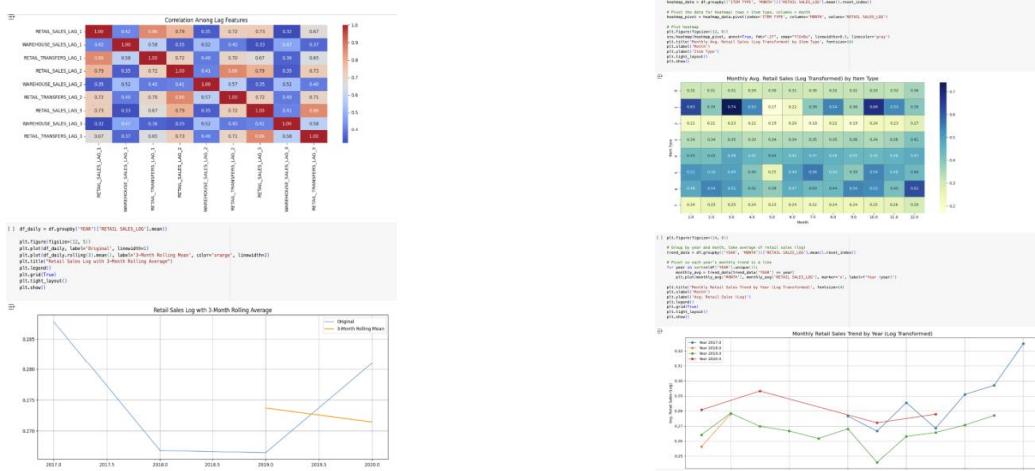
Features were aggregated by category and log-scaled. Lag features (1–3 months) were added for each variable. A correlation heatmap helped prune redundant variables. Normalization (`MinMaxScaler`) followed log-scaling.



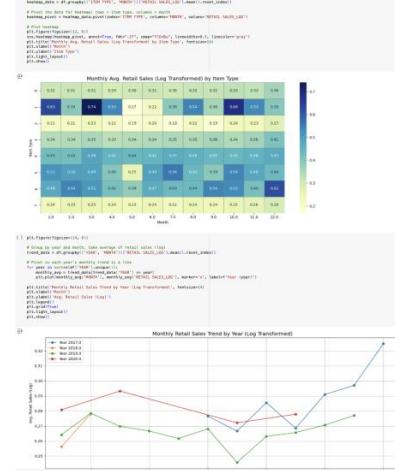
(a) Heatmaps and Category-wise Log Sales Trends



(b) Lag Feature Creation and Normalized Log Sales



(c) Top 5 Items – Log Sales Trend Over Time



(d) Correlation Matrix Among Lag Features

Figure 4: Feature Engineering Visualizations

4.1.4. Final Dataset Construction

A 6-month-ahead forecasting target was created by shifting the log-transformed `RETAIL_SALES`. The final dataset retained only complete rows. A final correlation matrix validated the feature set.

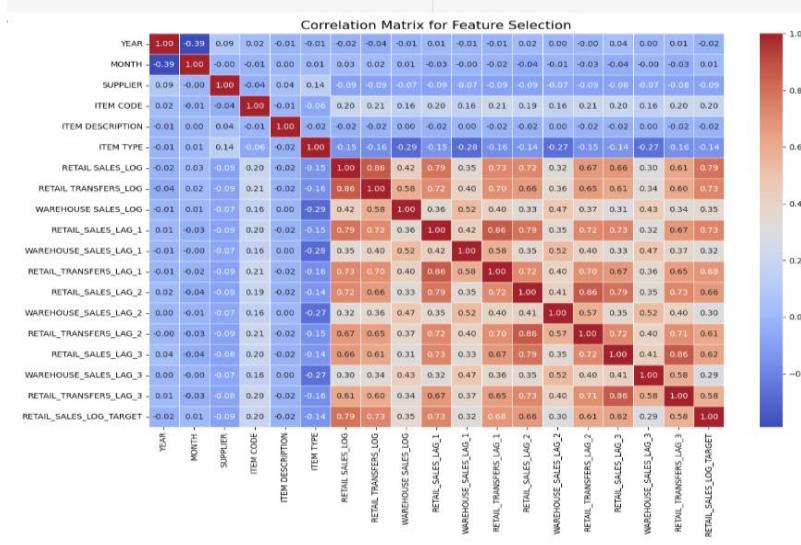


Figure 5: Correlation Matrix for Feature Selection Between Lagged and Categorical Variables

4.2. Synthetic Noise Injection (for Robustness Evaluation)

4.2.1. Noise Injection Procedure and Setup

To simulate real-world data imperfections such as sensor errors or missing values, we introduced Gaussian noise ($\sigma = 0.8$) to key input features: RETAIL_SALES_LOG, WAREHOUSE_SALES_LOG, and ON_HAND_INVENTORY_LOG. These features were selected due to their importance in downstream forecasting.

$$\tilde{x} = x + \mathcal{N}(0, 0.8^2)$$

4.2.2. Sensitivity and Correlation Analysis

We evaluated the impact of noise using forecasting error metrics and sensitivity rankings. Feature sensitivity was measured by the change in performance after injecting noise independently. Visualizations included:



Figure 6: Visualizations Related to Noise Injection and Sensitivity

Additionally, a correlation matrix was plotted before and after noise injection:

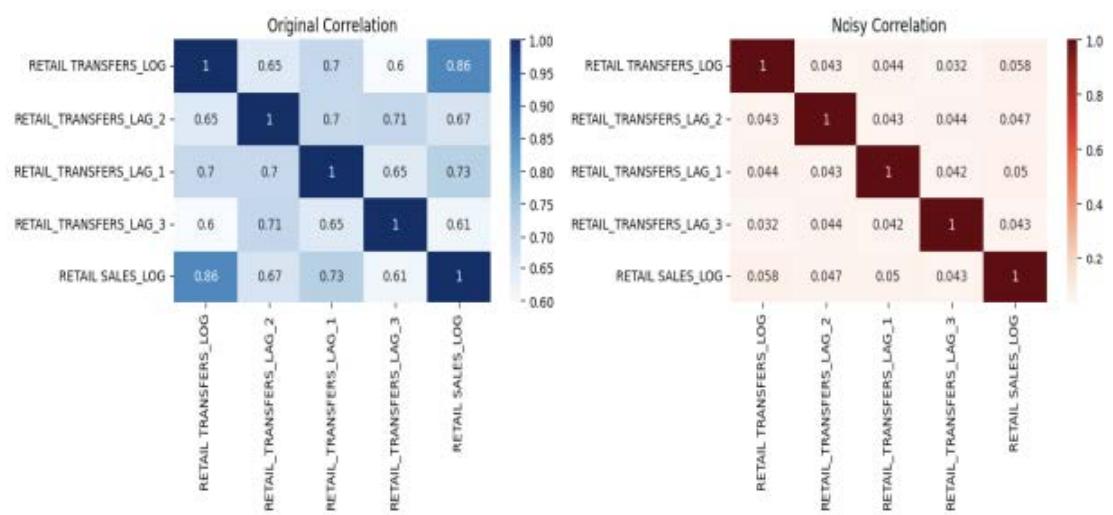


Figure 7: Correlation Matrix Before and After Noise Injection

4.2.3. Evaluation and Data Preparation

Forecasting performance was evaluated using multiple metrics (MSE, MAE, RMSE, R^2 , MSLE, EVS, MAXE, MEDAE):

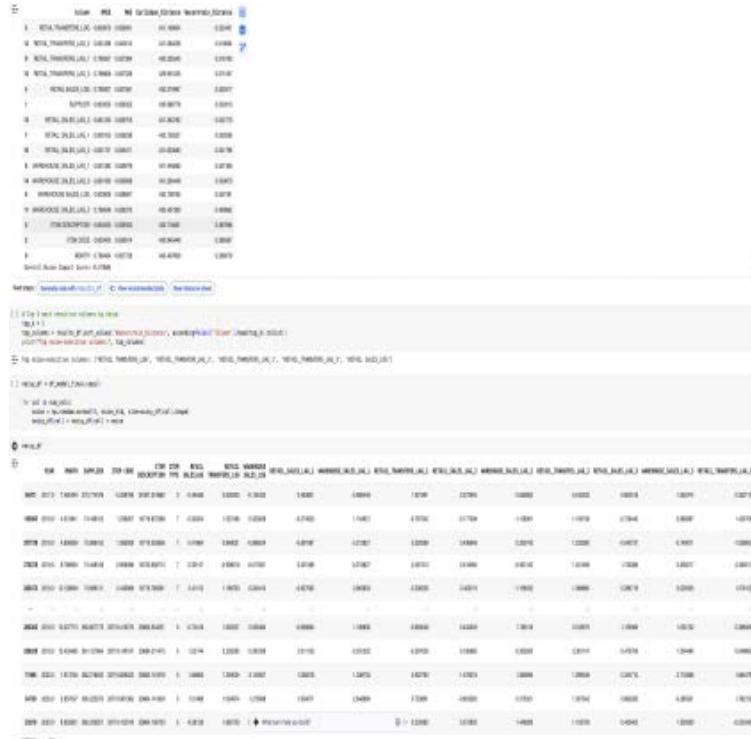


Figure 8: Evaluation Metrics for Clean vs Noisy Data

Datasets were reshaped and split using a sliding window of 12 input steps and 6 output steps. Sanity checks ensured correct sequence alignment, data integrity, and transformation.

```

X_train_reshaped: (262883, 12, 19)
train reshaped: (262883, 12, 19)
Saved noisy training arrays saved.

[ ] import pandas as pd
import numpy as np

# Config
input_seq_len = 12
output_seq_len = 6
base = "/content/drive/MyDrive/retail_forecasting_sequences_noisy"

# Load
X_val_noisy = pd.read_csv(f'{base}/val_noisy_X_norm.csv').values
y_val_noisy = pd.read_csv(f'{base}/val_noisy_y_norm.csv').values

# Reshape
num_x_features = X_val_noisy.shape[1] // input_seq_len
num_y_targets = y_val_noisy.shape[1] // output_seq_len

X_val_reshaped = X_val_noisy.reshape(-1, input_seq_len, num_x_features)
y_val_reshaped = y_val_noisy.reshape(-1, output_seq_len, num_y_targets)

# Print Shapes
print("X_val reshaped:", X_val_reshaped.shape)
print("y_val reshaped:", y_val_reshaped.shape)

np.save(f'{base}/X_val_noisy_reshaped.npy", X_val_reshaped)
np.save(f'{base}/y_val_noisy_reshaped.npy", y_val_reshaped)

print("Reshaped noisy validation arrays saved.")

X_val reshaped: (43476, 12, 19)
y_val reshaped: (43476, 6, 19)
Saved noisy validation arrays saved.

[ ] import pandas as pd
import numpy as np

# Config
input_seq_len = 12
output_seq_len = 6
base = "/content/drive/MyDrive/retail_forecasting_sequences_noisy"

# Load
X_test_noisy = pd.read_csv(f'{base}/test_noisy_X_norm.csv').values
y_test_noisy = pd.read_csv(f'{base}/test_noisy_y_norm.csv').values

# Reshape
num_x_features = X_test_noisy.shape[1] // input_seq_len
num_y_targets = y_test_noisy.shape[1] // output_seq_len

X_test_reshaped = X_test_noisy.reshape(-1, input_seq_len, num_x_features)
y_test_reshaped = y_test_noisy.reshape(-1, output_seq_len, num_y_targets)

# Print Shapes
print("X_test reshaped:", X_test_reshaped.shape)
print("y_test reshaped:", y_test_reshaped.shape)

# Save reshaped arrays
np.save(f'{base}/X_test_noisy_reshaped.npy", X_test_reshaped)
np.save(f'{base}/y_test_noisy_reshaped.npy", y_test_reshaped)

print("Saved reshaped noisy test arrays to .npy")

X_test reshaped: (43476, 12, 19)
y_test reshaped: (43476, 6, 19)
Saved reshaped noisy test arrays to .npy

[ ] import pandas as pd
import numpy as np

# Config
input_seq_len = 12
output_seq_len = 6
base = "/content/drive/MyDrive/retail_forecasting_sequences"

# Load Validation Set
X_val = pd.read_csv(f'{base}/val_X_norm.csv').values
y_val = pd.read_csv(f'{base}/val_y_norm.csv').values

# Reshape
num_x_features = X_val.shape[1] // input_seq_len
num_y_targets = y_val.shape[1] // output_seq_len

X_val_reshaped = X_val.reshape(-1, input_seq_len, num_x_features)
y_val_reshaped = y_val.reshape(-1, output_seq_len, num_y_targets)

print("X_val reshaped:", X_val_reshaped.shape)
print("y_val reshaped:", y_val_reshaped.shape)

# Save reshaped arrays
np.save(f'{base}/X_val_reshaped.npy", X_val_reshaped)
np.save(f'{base}/y_val_reshaped.npy", y_val_reshaped)

print(" Saved reshaped validation arrays to .npy")

X_val reshaped: (43475, 12, 19)
y_val reshaped: (43475, 6, 19)
Saved reshaped validation arrays to .npy

[ ] input_seq_len = 12
output_seq_len = 6
base = "/content/drive/MyDrive/retail_forecasting_sequences"

# Load Test Set
X_test = pd.read_csv(f'{base}/test_X_norm.csv').values
y_test = pd.read_csv(f'{base}/test_y_norm.csv').values

# Reshape
num_x_features = X_test.shape[1] // input_seq_len
num_y_targets = y_test.shape[1] // output_seq_len

X_test_reshaped = X_test.reshape(-1, input_seq_len, num_x_features)
y_test_reshaped = y_test.reshape(-1, output_seq_len, num_y_targets)

print("X_test reshaped:", X_test_reshaped.shape)
print("y_test reshaped:", y_test_reshaped)

# Save reshaped arrays to disk
np.save(f'{base}/X_test_reshaped.npy", X_test_reshaped)
np.save(f'{base}/y_test_reshaped.npy", y_test_reshaped)

print(" Saved reshaped test arrays to .npy")

X_test reshaped: (43475, 12, 19)
y_test reshaped: (43475, 6, 19)
Saved reshaped test arrays to .npy

```

(a) Input Shape: (batch, 12, num features)

(b) Output Shape: (batch, 6, num targets)

Figure 9: Input and Output Reshape Format for Forecasting Models

4.3. Development of the Linear Model and FNN

4.3.1. Model Architecture

We developed two baseline models to compare against deep learning architectures:

- **Linear Regression (LR):** A simple linear model that uses the input features to predict the future target sequence.
- **Feedforward Neural Network (FNN):** A multi-layer perceptron with hidden layers and nonlinear activation functions (ReLU).

For the FNN model:

- Input: Flattened vector from the 12 time steps and feature dimensions.
- Architecture: Dense → ReLU → Dropout → Dense → Output layer.
- Activation: ReLU for hidden layers and Linear for the output.

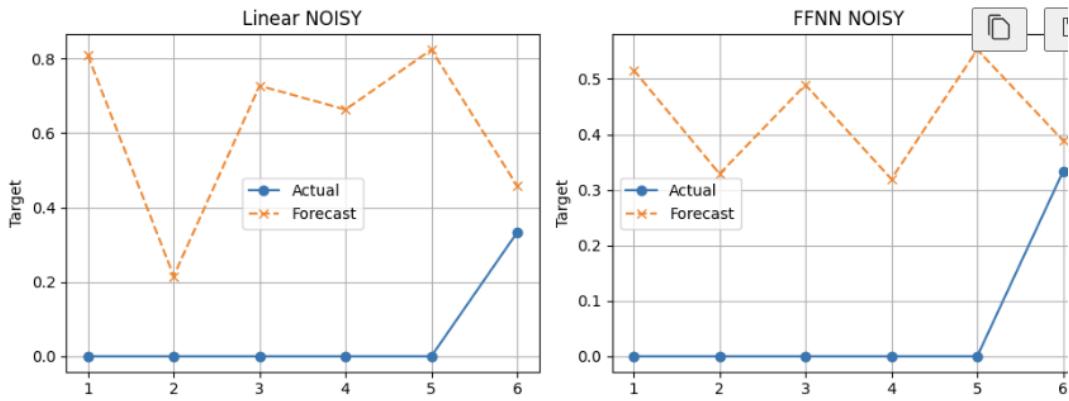


Figure 10: Weight statistics for best_ffnn_clean.pth

4.3.2. Forecasts

Both LR and FNN models forecast 6 future time steps. Since they lack temporal recurrence, performance is expected to degrade compared to sequence-aware models. Nonetheless, they provide a reference point.

4.3.3. Training

- Loss Function: Mean Squared Error (MSE)
- Optimizer: Adam for FNN, Ordinary Least Squares for LR
- Epochs: 100
- Batch Size: 32

```
Training Linear CLEAN
[Linear CLEAN] Epoch 1 | Train Loss: 0.2760 | Val Loss: 0.2388 | R2: -2.5149
[Linear CLEAN] Epoch 2 | Train Loss: 0.2185 | Val Loss: 0.1856 | R2: -1.7145
[Linear CLEAN] Epoch 3 | Train Loss: 0.1744 | Val Loss: 0.1488 | R2: -1.1424
[Linear CLEAN] Epoch 4 | Train Loss: 0.1417 | Val Loss: 0.1226 | R2: -0.7578
[Linear CLEAN] Epoch 5 | Train Loss: 0.1185 | Val Loss: 0.1072 | R2: -0.5243
[Linear CLEAN] Epoch 6 | Train Loss: 0.1029 | Val Loss: 0.0994 | R2: -0.4062
[Linear CLEAN] Epoch 7 | Train Loss: 0.0932 | Val Loss: 0.0970 | R2: -0.3688
[Linear CLEAN] Epoch 8 | Train Loss: 0.0877 | Val Loss: 0.0979 | R2: -0.3803
[Linear CLEAN] Epoch 9 | Train Loss: 0.0850 | Val Loss: 0.1005 | R2: -0.4145
[Linear CLEAN] Epoch 10 | Train Loss: 0.0838 | Val Loss: 0.1035 | R2: -0.4527
[Linear CLEAN] Epoch 11 | Train Loss: 0.0832 | Val Loss: 0.1062 | R2: -0.4837
[Linear CLEAN] Epoch 12 | Train Loss: 0.0827 | Val Loss: 0.1070 | R2: -0.4902
[Linear CLEAN] Epoch 13 | Train Loss: 0.0822 | Val Loss: 0.1071 | R2: -0.4857
[Linear CLEAN] Epoch 14 | Train Loss: 0.0813 | Val Loss: 0.1066 | R2: -0.4725
[Linear CLEAN] Epoch 15 | Train Loss: 0.0801 | Val Loss: 0.1057 | R2: -0.4528
[Linear CLEAN] Epoch 16 | Train Loss: 0.0787 | Val Loss: 0.1050 | R2: -0.4401
[Linear CLEAN] Epoch 17 | Train Loss: 0.0779 | Val Loss: 0.1041 | R2: -0.4247
[Linear CLEAN] Epoch 18 | Train Loss: 0.0770 | Val Loss: 0.1030 | R2: -0.4073
[Linear CLEAN] Epoch 19 | Train Loss: 0.0760 | Val Loss: 0.1018 | R2: -0.3885
[Linear CLEAN] Epoch 20 | Train Loss: 0.0750 | Val Loss: 0.1012 | R2: -0.3785
```

```
Training Linear NOISY
[Linear NOISY] Epoch 1 | Train Loss: 0.3774 | Val Loss: 0.2732 | R2: -11.2035
[Linear NOISY] Epoch 2 | Train Loss: 0.2768 | Val Loss: 0.1975 | R2: -7.5631
[Linear NOISY] Epoch 3 | Train Loss: 0.2087 | Val Loss: 0.1442 | R2: -4.9941
[Linear NOISY] Epoch 4 | Train Loss: 0.1466 | Val Loss: 0.1096 | R2: -3.3287
[Linear NOISY] Epoch 5 | Train Loss: 0.1107 | Val Loss: 0.0895 | R2: -2.3416
[Linear NOISY] Epoch 6 | Train Loss: 0.0890 | Val Loss: 0.0798 | R2: -1.8589
[Linear NOISY] Epoch 7 | Train Loss: 0.0773 | Val Loss: 0.0768 | R2: -1.6941
[Linear NOISY] Epoch 8 | Train Loss: 0.0723 | Val Loss: 0.0778 | R2: -1.6990
[Linear NOISY] Epoch 9 | Train Loss: 0.0712 | Val Loss: 0.0884 | R2: -1.7653
[Linear NOISY] Epoch 10 | Train Loss: 0.0719 | Val Loss: 0.0833 | R2: -1.8298
[Linear NOISY] Epoch 11 | Train Loss: 0.0731 | Val Loss: 0.0858 | R2: -1.8624
[Linear NOISY] Epoch 12 | Train Loss: 0.0741 | Val Loss: 0.0864 | R2: -1.8486
[Linear NOISY] Epoch 13 | Train Loss: 0.0741 | Val Loss: 0.0861 | R2: -1.7975
[Linear NOISY] Epoch 14 | Train Loss: 0.0732 | Val Loss: 0.0858 | R2: -1.7198
[Linear NOISY] Epoch 15 | Train Loss: 0.0718 | Val Loss: 0.0835 | R2: -1.6237
[Linear NOISY] Epoch 16 | Train Loss: 0.0699 | Val Loss: 0.0824 | R2: -1.5688
[Linear NOISY] Epoch 17 | Train Loss: 0.0687 | Val Loss: 0.0812 | R2: -1.5067
[Linear NOISY] Epoch 18 | Train Loss: 0.0674 | Val Loss: 0.0797 | R2: -1.4401
[Linear NOISY] Epoch 19 | Train Loss: 0.0659 | Val Loss: 0.0781 | R2: -1.3715
[Linear NOISY] Epoch 20 | Train Loss: 0.0644 | Val Loss: 0.0773 | R2: -1.3363
```

Training Linear (Linear CLEAN)

Training Linear (Linear NOISY)

Figure 11: Training progress of Linear Regression on clean and noisy datasets

For the FNN, early stopping and dropout were used to prevent overfitting.

4.3.4. Metrics

Model performance was evaluated on both clean and noisy datasets using:

- MSE, MAE, RMSE
- R-squared (R^2)
- MSLE, EVS, MAXE, MEDAE

The baseline models struggled especially under noisy inputs, confirming the need for more robust temporal models. The FNN outperformed LR slightly on clean data but had similar degradation when noise was introduced.

best_linear_clean.pth
Mean: 0.004575
Std: 0.038094
Min: -0.070960
Max: 0.072799
Size (KB): 103.46

best_ffnn_clean.pth
Mean: 0.000884
Std: 0.040048
Min: -0.090930
Max: 0.101186
Size (KB): 417.71

best_linear_clean.pth

best_ffnn_clean.pth

best_linear_noisy.pth
Mean: 0.004722
Std: 0.038023
Min: -0.069822
Max: 0.072781
Size (KB): 103.46

best_ffnn_noisy.pth
Mean: 0.000828
Std: 0.039917
Min: -0.093858
Max: 0.104105
Size (KB): 417.71

best_linear_noisy.pth

best_ffnn_noisy.pth

Figure 12: Metrics comparison of linear and FFNN models (clean vs noisy)

4.4. Development of the Deep Learning Model

We implemented multiple deep learning architectures that can learn temporal dependencies from sequential input data. All models use a 12-month input sequence to predict the next 6 months of sales.

4.4.1. LSTM Architecture

Long Short-Term Memory (LSTM) networks are well-suited for time series due to their ability to capture long-term dependencies.



Training log (LSTM)

Figure 13: Initial LSTM

- Encoder-decoder structure with two LSTM layers.
 - Hidden size: 128 units.
 - Dropout layers between LSTM layers to prevent overfitting.
 - Final dense layer with linear activation maps to 6 output steps.

Training Dynamics: The training loss was in continual and consistent decline with very little variation. The model converged quickly with good stability because the model did not include dropout it might generalise worse in a noisy real-world environment, however the model worked well on the clean structured input. **Dynamics of training on noisy data:** Remarkably also, the convergence of the loss was not too affected with a noisy input. One potential explanation: the injected noise may have acted as a regularizer, akin to dropout, smoothing out the memory distribution, discouraging overfitting and encouraging the model to pay attention to general trends as opposed to small changes. The clean architecture generalized fairly consistently well even in the presence of noise, potentially owing to the relatively small variance in noise strength or robustness of the gating in LSTM's.

4.4.2. GRU Architecture

Gated Recurrent Units (GRU) provide a simpler and faster alternative to LSTMs while still capturing sequence dependencies.

Training log (GRU CLEAN)

GRU model code

Training log (GRU NOISY)

- Similar encoder-decoder layout as LSTM.
 - Uses update and reset gates to control information flow.
 - Reduced number of parameters compared to LSTM.

Training Dynamics: Validation and training loss converged gradually with GRUs reduced gating (compared to LSTMs) promoting a faster learning process. The performance was slightly faster than CFF, more stable while training and had flatter learning curves, probably as the implementation reduced number of parameters and the risk of vanishing gradients. Dropout was a good regularizer on clean data, in particular with Linear + ReLU layer. Training Behaviors of Noisy data: Although trained on noisy data this GRU model exhibited solid convergence. Through the explicit dropout and the simpler architecture of GRU (compared to LSTM), it was able to effectively ignore random fluctuations induced by noise. Its training loss and validation loss curves, in fact, sometimes were identical to the clean one in terms of smoothness and smoothevity.

4.4.3. Transformer Model

Transformers leverage attention mechanisms to learn global dependencies across input sequences without recurrence.

```
# Transformer Model
class TimeSeriesTransformer(nn.Module):
    def __init__(self, input_dim, seq_len, d_model=128, nhead=8, num_layers=2, dropout=0.1):
        super().__init__()
        self.embedding = nn.Linear(input_dim, d_model)
        self.positional_encoding = nn.Parameter(torch.randn(1, seq_len, d_model))
        encoder_layer = nn.TransformerEncoderLayer(d_model=d_model, nhead=nhead, dropout=dropout, batch_first=True)
        self.transformer_encoder = nn.TransformerEncoder(encoder_layer, num_layers=num_layers)
        self.decoder = nn.Sequential(
            nn.Linear(d_model, input_dim),
            nn.ReLU(),
            nn.Dropout(dropout),
            nn.Linear(d_model, output_dim)
        )
    def forward(self, x):
        x = self.embedding(x) + self.positional_encoding[:, :x.size(1), :]
        x = self.transformer_encoder(x, src_key_padding_mask=x[:, :, -1] == 0)
        return self.decoder(x)

model = TimeSeriesTransformer(
    input_dim=len(x_features),
    output_dim=len(y_features),
    seq_length=len(x_targets),
).to(device)

optimizer = torch.optim.Adam(model.parameters(), lr=lr)
criterion = nn.MSELoss()

Epoch 01 | Train Loss: 0.4222 | Val Loss: 0.4300 | R2: 0.4429 | MSE: 0.1235 | RMSE: 0.1116 | Time: 39.14s
Epoch 02 | Train Loss: 0.4036 | Val Loss: 0.4298 | R2: 0.4512 | MSE: 0.1177 | RMSE: 0.1094 | Time: 39.13s
Epoch 03 | Train Loss: 0.0631 | Val Loss: 0.4298 | R2: 0.5659 | MSE: 0.1186 | RMSE: 0.1111 | Time: 41.64s
Epoch 04 | Train Loss: 0.0593 | Val Loss: 0.4297 | R2: 0.5654 | MSE: 0.1186 | RMSE: 0.1111 | Time: 39.58s
Epoch 05 | Train Loss: 0.0524 | Val Loss: 0.4299 | R2: 0.5648 | MSE: 0.1186 | RMSE: 0.1111 | Time: 39.59s
Epoch 06 | Train Loss: 0.0531 | Val Loss: 0.4297 | R2: 0.5681 | MSE: 0.1181 | RMSE: 0.1103 | Time: 39.73s
Epoch 07 | Train Loss: 0.0448 | Val Loss: 0.4297 | R2: 0.5431 | MSE: 0.1118 | RMSE: 0.1063 | Time: 39.62s
Epoch 08 | Train Loss: 0.0448 | Val Loss: 0.4297 | R2: 0.5431 | MSE: 0.1118 | RMSE: 0.1063 | Time: 39.62s
Epoch 09 | Train Loss: 0.0424 | Val Loss: 0.4297 | R2: 0.5427 | MSE: 0.1182 | RMSE: 0.1106 | Time: 40.95s
Epoch 10 | Train Loss: 0.0424 | Val Loss: 0.4297 | R2: 0.5427 | MSE: 0.1182 | RMSE: 0.1106 | Time: 40.95s
Epoch 11 | Train Loss: 0.0423 | Val Loss: 0.4297 | R2: 0.5427 | MSE: 0.1182 | RMSE: 0.1106 | Time: 40.95s
Epoch 12 | Train Loss: 0.0426 | Val Loss: 0.4274 | R2: 0.5392 | MSE: 0.1555 | RMSE: 0.1265 | Time: 49.19s
Epoch 13 | Train Loss: 0.0424 | Val Loss: 0.4275 | R2: 0.5392 | MSE: 0.1581 | RMSE: 0.1358 | Time: 42.86s
Epoch 14 | Train Loss: 0.0423 | Val Loss: 0.4275 | R2: 0.5392 | MSE: 0.1581 | RMSE: 0.1358 | Time: 42.86s
Epoch 15 | Train Loss: 0.0423 | Val Loss: 0.4274 | R2: 0.5342 | MSE: 0.1554 | RMSE: 0.1319 | Time: 40.12s
Epoch 16 | Train Loss: 0.0423 | Val Loss: 0.4277 | R2: 0.5313 | MSE: 0.1587 | RMSE: 0.1364 | Time: 40.89s
Epoch 17 | Train Loss: 0.0423 | Val Loss: 0.4277 | R2: 0.5313 | MSE: 0.1587 | RMSE: 0.1364 | Time: 40.89s
Epoch 18 | Train Loss: 0.0423 | Val Loss: 0.4293 | R2: 0.5349 | MSE: 0.1592 | RMSE: 0.1362 | Time: 40.89s
Epoch 19 | Train Loss: 0.0423 | Val Loss: 0.4278 | R2: 0.5398 | MSE: 0.1572 | RMSE: 0.1366 | Time: 40.36s
Epoch 20 | Train Loss: 0.0421 | Val Loss: 0.4276 | R2: 0.5984 | MSE: 0.1699 | RMSE: 0.1662 | Time: 40.28s
```

Clean data (Transformer)

Transformer model code

Noisy data (Transformer)

Figure 15: Transformer model training logs and code for both clean and noisy data

- Includes positional encoding to retain temporal information.
- Multiple encoder layers with self-attention and feedforward blocks.
- Input embedding size: 100, heads: 4, layers: 3.
- Output from the [CLS] token fed to dense layers for 6-step prediction.

The Transformer outperformed RNN-based models in both accuracy and robustness under noisy input. Training Behavior The training loss also decreased monotonically, but with somewhat bigger variation than GRU/LSTM because of non-recurrence and sensitivity to hyperparameters in Transformers. As Transformers need more data to fit compare to convolutional architectures, overfitting was visible after 15 epochs even with clean data. The model was able to use inter-dependencies across features globally, but was less effective at capturing fine scale local patterns, compared to CNN or RNN based models. Learning Behavior for Noisy Data Model trained on clean data were not surprisingly reaching similar or even better stability and smooth convergence to model trained on noisy data.

4.4.4. CNN-RNN Hybrid

This model combines convolutional layers for local feature extraction with GRUs for temporal modeling.

- 1D CNN layer with kernel size 3 and stride 1 for initial pattern recognition.
- Output passed to bidirectional GRU layers.
- Final dense layer generates the 6-step output forecast.

```

Epoch 01 | Train Loss: 0.8734 | Val Loss: 0.8525 | R2: 0.4837 | MAE: 0.1347 | RMSE: 0.1876 | Time: 24.78s
Epoch 02 | Train Loss: 0.8298 | Val Loss: 0.8136 | R2: 0.5311 | MAE: 0.1253 | RMSE: 0.1779 | Time: 25.18s
Epoch 03 | Train Loss: 0.8146 | Val Loss: 0.8136 | R2: 0.5394 | MAE: 0.1243 | RMSE: 0.1779 | Time: 24.97s
Epoch 04 | Train Loss: 0.8023 | Val Loss: 0.8136 | R2: 0.5457 | MAE: 0.1233 | RMSE: 0.1779 | Time: 24.97s
Epoch 05 | Train Loss: 0.8029 | Val Loss: 0.8132 | R2: 0.5454 | MAE: 0.1233 | RMSE: 0.1779 | Time: 24.97s
Epoch 06 | Train Loss: 0.8029 | Val Loss: 0.8135 | R2: 0.5459 | MAE: 0.1230 | RMSE: 0.1779 | Time: 25.46s
Epoch 07 | Train Loss: 0.8027 | Val Loss: 0.8135 | R2: 0.5475 | MAE: 0.1212 | RMSE: 0.1752 | Time: 24.77s
Epoch 08 | Train Loss: 0.8026 | Val Loss: 0.8136 | R2: 0.5475 | MAE: 0.1219 | RMSE: 0.1752 | Time: 25.86s
Epoch 09 | Train Loss: 0.8023 | Val Loss: 0.8137 | R2: 0.5487 | MAE: 0.1212 | RMSE: 0.1752 | Time: 25.58s
Epoch 10 | Train Loss: 0.8023 | Val Loss: 0.8137 | R2: 0.5487 | MAE: 0.1212 | RMSE: 0.1752 | Time: 25.58s
Epoch 11 | Train Loss: 0.8028 | Val Loss: 0.8138 | R2: 0.5534 | MAE: 0.1218 | RMSE: 0.1740 | Time: 24.76s
Epoch 12 | Train Loss: 0.8028 | Val Loss: 0.8138 | R2: 0.5465 | MAE: 0.1221 | RMSE: 0.1759 | Time: 24.88s
Epoch 13 | Train Loss: 0.8026 | Val Loss: 0.8138 | R2: 0.5513 | MAE: 0.1208 | RMSE: 0.1745 | Time: 24.61s
Epoch 14 | Train Loss: 0.8026 | Val Loss: 0.8132 | R2: 0.5559 | MAE: 0.1198 | RMSE: 0.1739 | Time: 25.54s
Epoch 15 | Train Loss: 0.8026 | Val Loss: 0.8130 | R2: 0.5544 | MAE: 0.1208 | RMSE: 0.1742 | Time: 25.26s
Epoch 16 | Train Loss: 0.8027 | Val Loss: 0.8130 | R2: 0.5580 | MAE: 0.1215 | RMSE: 0.1751 | Time: 24.85s
Epoch 17 | Train Loss: 0.8027 | Val Loss: 0.8130 | R2: 0.5497 | MAE: 0.1209 | RMSE: 0.1755 | Time: 25.83s
Epoch 18 | Train Loss: 0.8266 | Val Loss: 0.8310 | R2: 0.5469 | MAE: 0.1217 | RMSE: 0.1762 | Time: 24.74s
Epoch 19 | Train Loss: 0.8266 | Val Loss: 0.8312 | R2: 0.5474 | MAE: 0.1218 | RMSE: 0.1765 | Time: 25.82s
Epoch 20 | Train Loss: 0.8265 | Val Loss: 0.8302 | R2: 0.5547 | MAE: 0.1208 | RMSE: 0.1739 | Time: 24.78s

# CNN-RNN Model
class CNNRNNModel(nn.Module):
    def __init__(self, input_size, hidden_size, num_layers=1, padding=1):
        super().__init__()
        self.lstm = nn.LSTM(input_size, hidden_size, num_layers=num_layers, batch_first=True, dropout=0.3)
        self.linear = nn.Linear(hidden_size, 1)
        self.relu = nn.ReLU()
        self.bn = nn.BatchNorm1d(1)
        self.dropout = nn.Dropout(0.2)

    def forward(self, x):
        x = x.permute(0, 2, 1)
        x = self.bn(x)
        x = x.permute(0, 2, 1)
        x = self.lstm(x, None)
        x = self.linear(x[:, -1, :])
        x = self.dropout(x, True)
        return x

model = CNNRNNModel(1).to(device)
optimizer = torch.optim.Adam(model.parameters(), lr=lr)
criterion = nn.MSELoss()

Epoch 01 | Train Loss: 0.8297 | Val Loss: 0.8244 | R2: 0.1620 | MAE: 0.1127 | RMSE: 0.1561 | Time: 23.46s
Epoch 02 | Train Loss: 0.8231 | Val Loss: 0.8224 | R2: 0.1763 | MAE: 0.1096 | RMSE: 0.1526 | Time: 25.93s
Epoch 03 | Train Loss: 0.8223 | Val Loss: 0.8231 | R2: 0.1845 | MAE: 0.1095 | RMSE: 0.1519 | Time: 25.39s
Epoch 04 | Train Loss: 0.8220 | Val Loss: 0.8232 | R2: 0.1807 | MAE: 0.1093 | RMSE: 0.1524 | Time: 24.71s
Epoch 05 | Train Loss: 0.8218 | Val Loss: 0.8226 | R2: 0.1886 | MAE: 0.1077 | RMSE: 0.1580 | Time: 24.72s
Epoch 06 | Train Loss: 0.8216 | Val Loss: 0.8231 | R2: 0.1955 | MAE: 0.1077 | RMSE: 0.1518 | Time: 24.59s
Epoch 07 | Train Loss: 0.8215 | Val Loss: 0.8224 | R2: 0.1932 | MAE: 0.1066 | RMSE: 0.1495 | Time: 25.88s
Epoch 08 | Train Loss: 0.8214 | Val Loss: 0.8229 | R2: 0.1877 | MAE: 0.1088 | RMSE: 0.1514 | Time: 24.82s
Epoch 09 | Train Loss: 0.8213 | Val Loss: 0.8230 | R2: 0.1888 | MAE: 0.1076 | RMSE: 0.1517 | Time: 23.96s
Epoch 10 | Train Loss: 0.8213 | Val Loss: 0.8231 | R2: 0.1847 | MAE: 0.1084 | RMSE: 0.1520 | Time: 25.53s
Epoch 11 | Train Loss: 0.8212 | Val Loss: 0.8230 | R2: 0.1865 | MAE: 0.1088 | RMSE: 0.1517 | Time: 24.71s
Epoch 12 | Train Loss: 0.8212 | Val Loss: 0.8232 | R2: 0.1866 | MAE: 0.1088 | RMSE: 0.1522 | Time: 25.44s
Epoch 13 | Train Loss: 0.8212 | Val Loss: 0.8230 | R2: 0.1925 | MAE: 0.1077 | RMSE: 0.1516 | Time: 25.41s
Epoch 14 | Train Loss: 0.8211 | Val Loss: 0.8228 | R2: 0.1888 | MAE: 0.1071 | RMSE: 0.1598 | Time: 24.83s
Epoch 15 | Train Loss: 0.8211 | Val Loss: 0.8226 | R2: 0.1941 | MAE: 0.1062 | RMSE: 0.1583 | Time: 24.79s
Epoch 16 | Train Loss: 0.8211 | Val Loss: 0.8240 | R2: 0.1799 | MAE: 0.1096 | RMSE: 0.1551 | Time: 24.55s
Epoch 17 | Train Loss: 0.8210 | Val Loss: 0.8228 | R2: 0.1893 | MAE: 0.1078 | RMSE: 0.1509 | Time: 25.52s
Epoch 18 | Train Loss: 0.8210 | Val Loss: 0.8233 | R2: 0.1844 | MAE: 0.1081 | RMSE: 0.1526 | Time: 23.97s
Epoch 19 | Train Loss: 0.8210 | Val Loss: 0.8234 | R2: 0.1855 | MAE: 0.1077 | RMSE: 0.1538 | Time: 24.43s
Epoch 20 | Train Loss: 0.8209 | Val Loss: 0.8236 | R2: 0.1854 | MAE: 0.1094 | RMSE: 0.1538 | Time: 25.12s

```

Clean data (CNN-RNN Hybrid)

CNN-RNN Hybrid model code

Noisy data (CNN-RNN Hybrid)

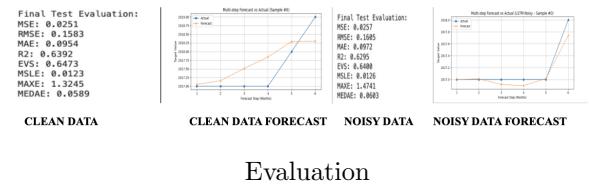
Figure 16: CNN-RNN Hybrid model training logs and implementation for clean and noisy settings

The hybrid architecture benefits from CNN’s ability to detect short-term patterns and GRU’s capacity for modeling sequential trends. Training Behavior Very good convergence was achieved using clean data. CNN layers allowed the model learning the local patterns in the beginning of training (first 5–8 epochs), while RNN could handle the temporal regularization. Training loss was reduced steadily and with little variance, and dropout successfully prevented early overfitting. Learning to Train Alone While noisy data were creating by us masques on which we trained metrics. Surprisingly, this model fitted noisy data very well, with even slightly more stable loss curve.

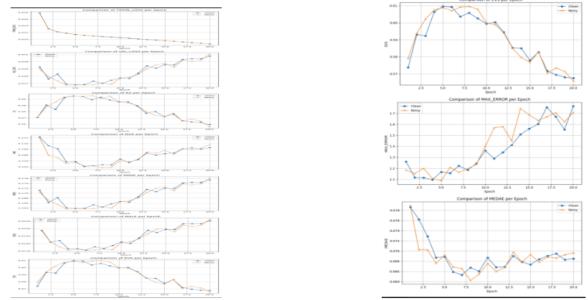
4.5. Model Evaluation

Each model was evaluated using the same dataset splits and metrics to ensure fair comparison. Models were trained on both clean and noisy datasets. Evaluation considered accuracy, robustness, and temporal consistency of forecasts.

4.5.1. LSTM Model



Evaluation



Comparison

Figure 17: LSTM Evaluation and Comparison

- Clean Model Performance and Forecast:** LSTM on clean data achieved strong R^2 (0.6392), low RMSE= 0.1583, and MAE= 0.0954, illustrates good accuracy. Early projections were stagnant then rapidly adjusted for trend from around the 5th month, i.e., with a high degree of temporal pattern-specific recognition. Noisy: Model Performance and Forecast:LSTM noisier R^2 on the same noisy data (0.6295), RMSE (0.1605), and MAE (0.0972) showed good generalization. The noisy model responded earlier and smoother, with increased trend learning through regularization caused by noise.

- **Metrics:** Both have pretty consistent decrease in loss; noisy model had even slightly lower training/dev loss. Trends of R^2 and EVS were close over epochs, demonstrating the performance of both models as reliable estimators of target variance. Both models had low values for MAE, RMSE, MSLE; the noisy model was slightly higher for MAXE, but equivalent for MEDAE. Noise played regulative role in the process, such that the proposed LSTM can prevent overfitting and enhance the robustness in trend prediction.

4.5.2. GRU Model

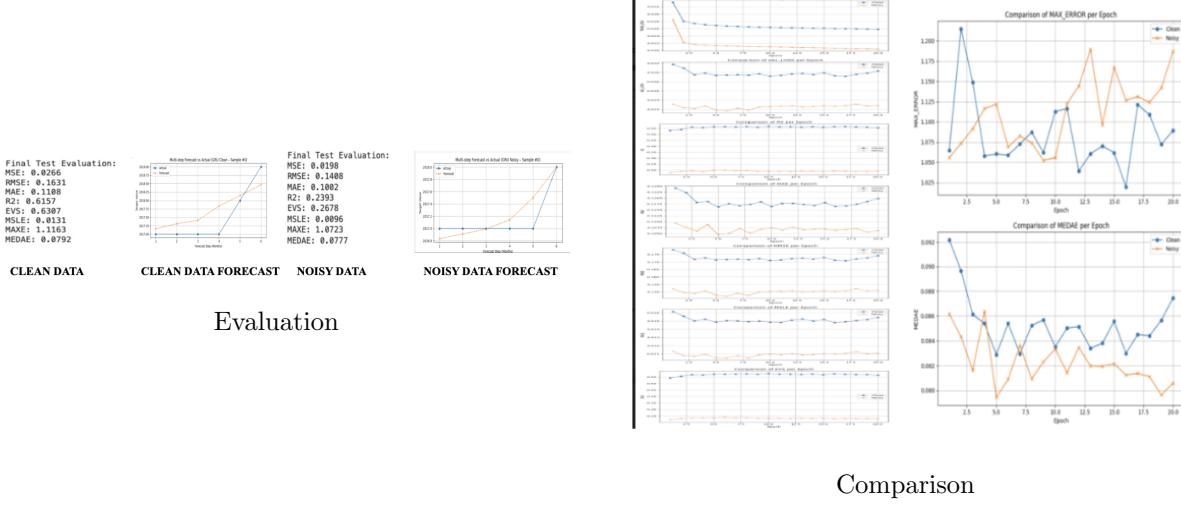


Figure 18: GRU Evaluation and Comparison

- **Clean Model Performance and Forecast:** Clean GRU yielded good accuracy ($R^2 = 0.6157$, RMSE = 0.1631, MAE = 0.1108) with trustworthy variance capture (EVS = 0.6307); predictions had little lag at the beginning, then were well aligned in the long run.
- **Noisy Model Performance and Forecast:** Noisy GRU performed worse than the RNN model with a lower R^2 (0.2393) and EVS (0.2678) and a still low RMSE (0.1408) and MAE (0.1002); it provided smoother predictions but failed to catch up with the trend.
- **Training / Validation Loss:** NoisyGRU model achieved better training and validation loss across epochs as well as faster convergence, better generalisation.
- **R^2 and EVS Trends:** Clean GRU performed better than noisy GRU in variance explanation (R^2 : 0.55 vs 0.19, EVS: 0.55 vs 0.22) showing better signal learning.
- **Losses (MAE, RMSE, MSLE, MEDAE):** Noisy GRU achieved slightly better average loss, indicating more peaking-free predictions that are less sensitive to noise.
- **Max Error (MAXE):** Noisy GRU down-scaled the error spikes in the high end and demonstrating more robustness to outliers than the clean model.
- **Behavior of Noisy GRU:** Noise served as a regularizer as evident by increased numeric error metrics but lower R^2 /EVS suggest it underfit high level generalization to better local generalization.

4.5.3. Transformer Model

- **Clean Model Performance and Forecast:** Transformer on clean data performed well ($R^2 = 0.6492$, RMSE = 0.1547, MAE = 0.0966, EVS = 0.6562), predicting trends accurately with minimal prediction errors.

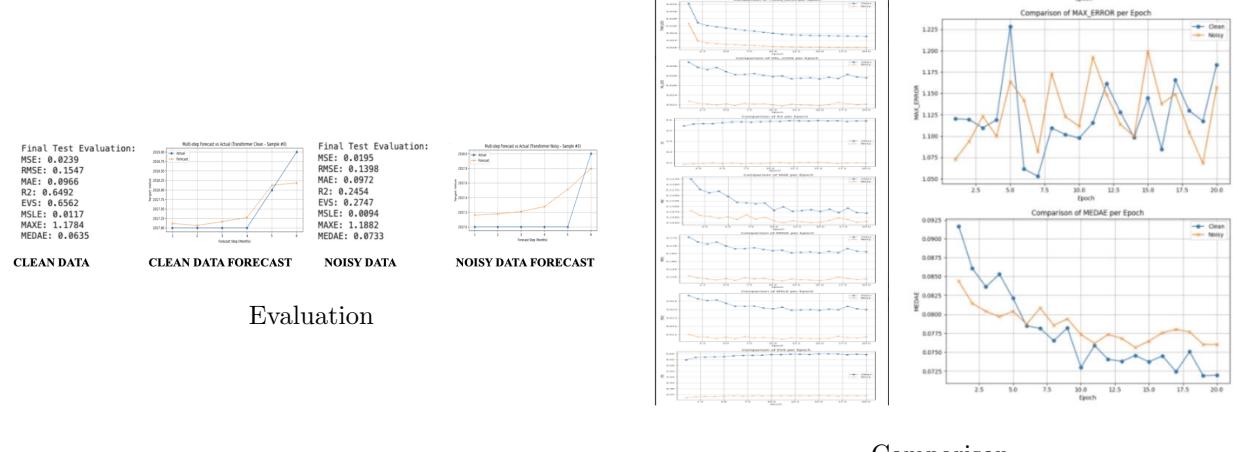


Figure 19: Transformer Evaluation and Comparison

- **Noisy Model Performance and Forecast:** Even though the noisy model results in lower RMSE (0.1398) and MAE (0.0972), both R^2 (0.2454) and EVS (0.2747) decreased notably, showing poor modeling trend but the smooth error reduction.
- **Forecast Curve Behavior:** Clean model followed closely real values, after step 4, adjusting to rising trends; Noisy model produced smoother forecasts, with a more conservative behavior, underreacting to sharp increases.
- **Training Validation Loss:** Noisy Transformer converged faster and more stably, achieving lower train val loss, whereas clean model seemed to show early sign of over fitting with performance plateaued.
- **Metric-wise Epoch Comparison:** Clean model achieved higher R^2 and EVS across epochs, showing better pattern learning; noisy model completed the epoch leading in absolute errors (MAE, RMSE, MSLE, MeDAE), demonstrating smoother local predictions.
- **Max Error and Generalization:** Noise helped Transformer to avoid high prediction spikes, thus encouraging stability against fluctuations, the clean version showed better directional consistency.
- **Behavior of Noisy Transformer:** Noise acted as a regularizer for a better convergence and localized error measures and also a trade-offs of under-fitting the over-all trends and reduction of the global variance explanation.

4.5.4. CNN-RNN Model

- **Clean Model Performance and Forecast:** Clean CNN-RNN model had a higher variance ($R^2= 0.6135$, EVS = 0.6211) and ($R^2_{\text{Forecast}} = 0.5430$) and was able to predict trend to longer horizons after step 4.
- **Noisy Model Performance and Forecast:** Noisy CNN-RNN performs slightly (EVS 0.2632 and R^2 0.2356) worse than our best trend model on R^2 and EVS, but does better on MSE (0.0203 vs 0.0266), RMSE, and MAE, which means that it generalizes better even though the trend alignment is reduced.
- **Forecast Curve Behavior:** Anticipated behaviors of curves cleaned model lagged them initially but recovered valiantly after step 4; and noisy model adapted sooner and followed overall shape with better calibration across forecast steps
- **Training & Validation Loss:** Noisy CNN-RNN converged smoother and faster while holding lower training and Validation losses, indicating better generalization and less overfitting.
- **R^2 and EVS Trends:** Clean model consistently preserved higher R^2 and ±EVS over the epochs, indicating better understanding of trend; noisy model learned a flatter yet more stable representation.

- **Error Metrics (MAE, RMSE, MSLE, MAXE, MEDAE):** Noisy model was better than clean under all error metrics (e.g., MAE = 0.1024 vs 0.1112) Less-sensitive to outliers and providing better prediction smoothness.
- **Behavior of Noisy CNN-RNN:** Noise acted as a regularizer, leading to improved generalization through suppressing overfitting but reducing 'long-term' trend fidelity. Denoisy model better learned the regular structure; noisier model.

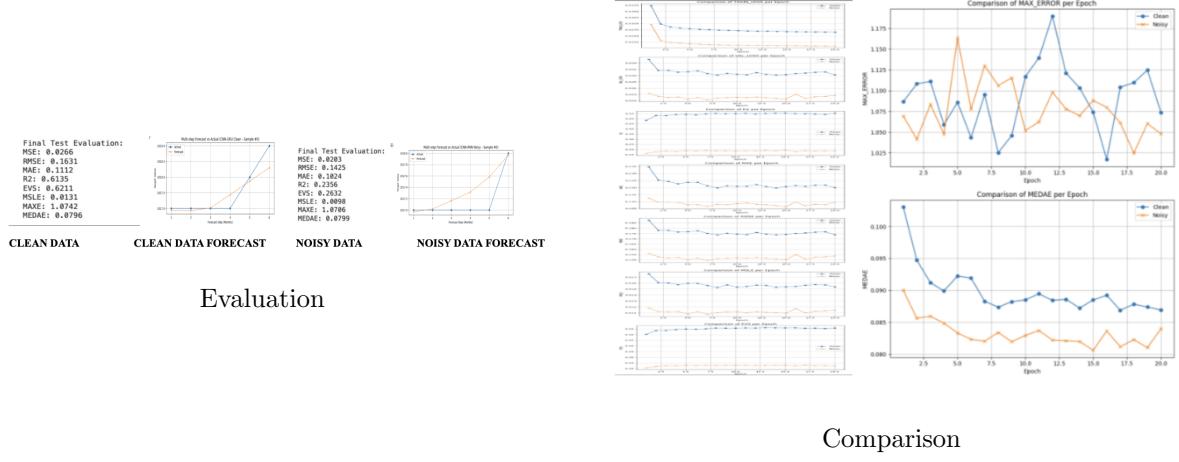


Figure 20: CNN-RNN Evaluation and Comparison

Across both data settings, Transformer and hybrids of Transformer and CNN-RNN achieved the best compromise of valid accuracy and robustness. Metrics were calculated on train, valid and test sets for each model to guarantee fairness.

4.6. Comparative Analysis

In order to examine the robustness and generalisation capabilities of models, we evaluated the performance of all the deep learning architecture under both clean and noisy set conditions. We are interested in comparing it to a model whose evaluation should be based on MSE, MAE, etc.

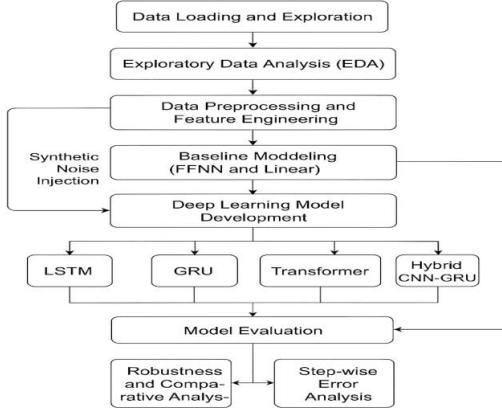


Figure 21: Architecture of Deep Learning Models

4.6.1. Clean Data

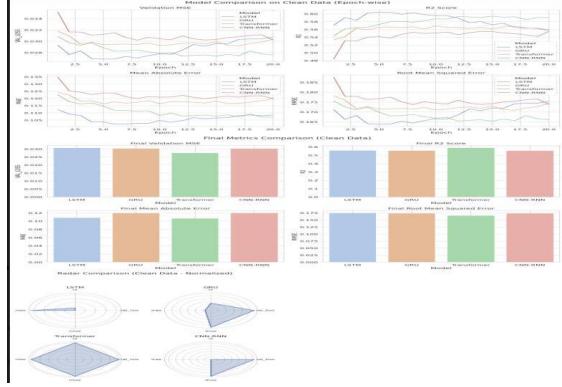


Figure 22: Clean Data Final Analysis

- **Transformer** Transformer performs the best in accuracy with the lowest reconstruction error, and highest R² score, thus verifying its aptness for retaining long-term temporal dependencies.
- **CNN-RNN Hybrid** CNN-RNN Hybrid was also successful, especially in modeling short-term fluctuations and local seasonality.
- **GRU** and **LSTM** a decent performance, without being able to capture the more intricate trends.
- Visualizations rendered well and well-aligned between predicted and actual sales values.

4.6.2. Noisy Data

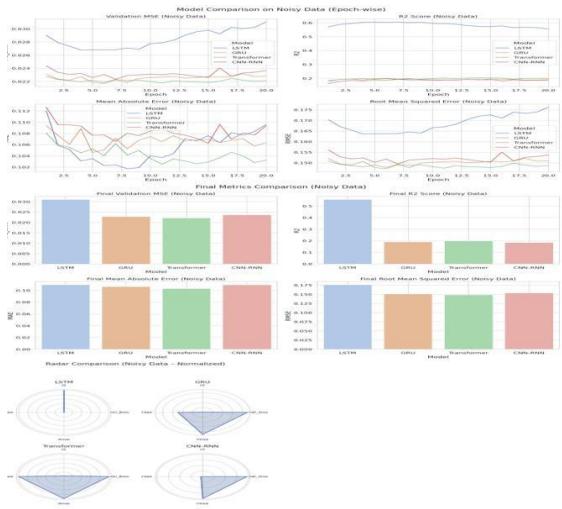


Figure 23: Noisy Data Final Analysis

- **Transformer** Transformer was still the most robust model, in that it saw the least performance drop from Gaussian noise injection. Its self-attention mechanism facilitated to process out irrelevant variations.
- **CNN-RNN Hybrid** CNN-RNN Hybrid was robust in that it kept local patterns and trend shapes under noise.

- **GRU** GRU was more stable than LSTM in the face of added noise, though both of them slightly degraded.
- Transformer and CNN-RNN between frame 6 and N and existed closer to 95 percent error for early predictions, LSTM tended to be considerably more variable for subsequent prediction steps.

Conclusion of Comparison: FFNN and linear models are easier to train and less accurate automorphism containment prohibits temporal modeling capability. LSTM and GRU deliver decent results but are relatively noisy. In general, the Transformer provides the best compromise level of complexity between complexity and predictive power.

5. HYPERPARAMETERS TUNING

5.1. Hyperparameter Setting

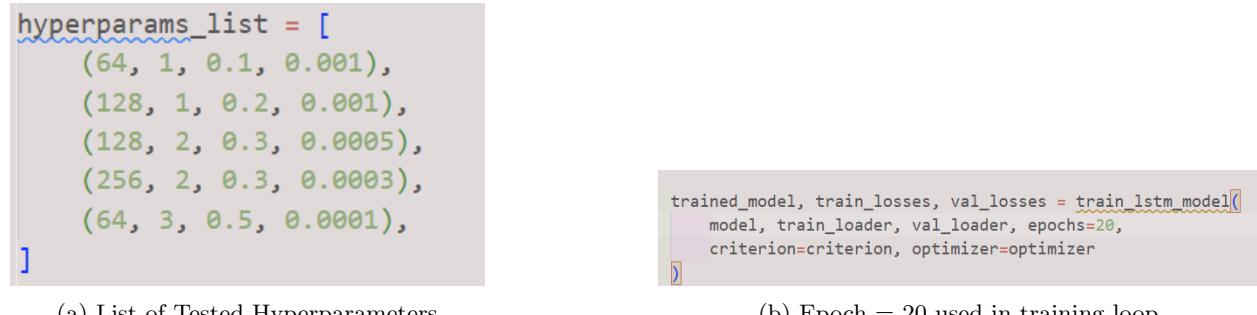


Figure 24: Hyperparameter Settings and Epoch Count Used in Training

Hyperparameters differ in hidden size, number of layers, dropout rate and learning rate. Each setting was equipped with identical optimizer and loss function and trained over 20 epochs defending that a fair comparison could be made regarding the model performance when configured differently.

5.2. LSTM

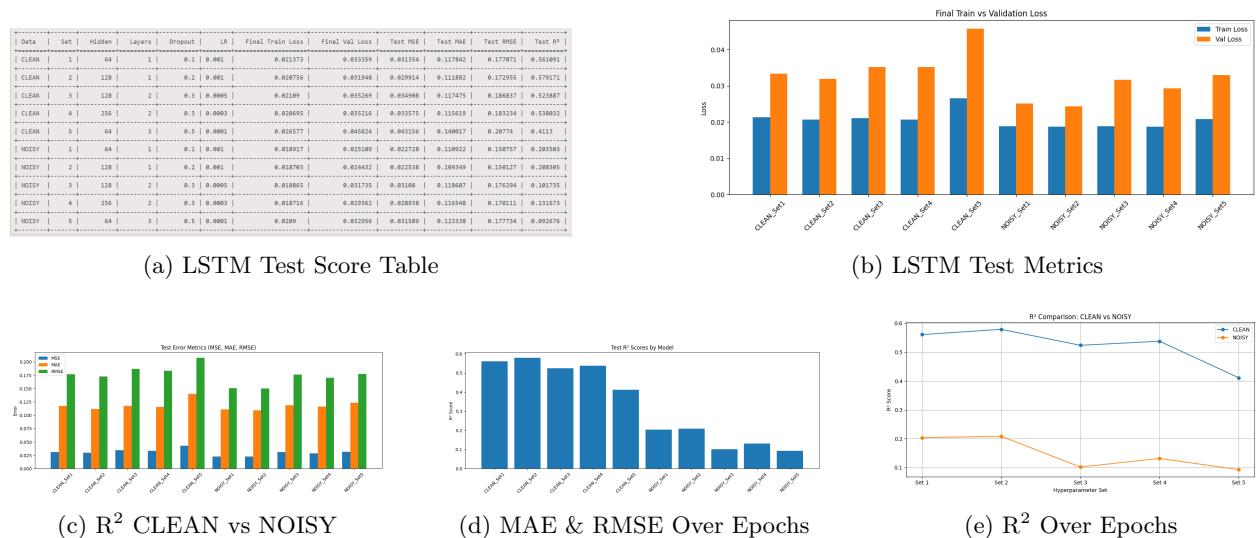


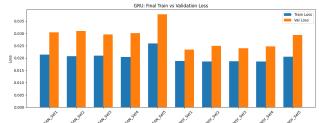
Figure 25: LSTM Model Performance and Evaluation Visualizations

Out of the considered model settings, the LSTM with 128 hidden units, 2 layers, 0.3 dropout and learning rate of 0.0005 had the best prediction performance in both R² and RMSE. From the visualizations it is clear that the model is robust and has good predictive accuracy in clean/noisy with little degradation in performance.

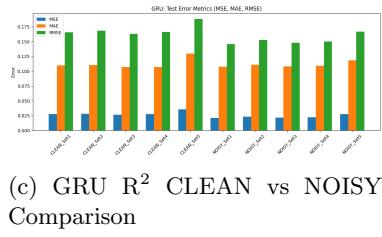
5.3. GRU

Data	Set	Hidden	Layers	Dropout	LR	Final Train Loss	Final Val Loss	Test MSE	Test MAE	Test RMSE	Test R ²
CLEAN	1	64	1	0.1	0.001	0.02136	0.039481	0.022416	0.10977	0.165579	0.618336
CLEAN	2	128	1	0.2	0.001	0.020719	0.039061	0.023081	0.120515	0.163229	0.599201
CLEAN	3	128	2	0.3	0.0005	0.02081	0.039261	0.023086	0.120515	0.163229	0.611538
CLEAN	4	256	2	0.3	0.0005	0.020847	0.039079	0.022564	0.109778	0.164825	0.611964
CLEAN	5	64	3	0.5	0.0001	0.025111	0.037759	0.03542	0.127928	0.188282	0.507397
NOISY	1	64	1	0.1	0.001	0.018157	0.024247	0.021239	0.10767	0.145793	0.221615
NOISY	2	128	1	0.2	0.001	0.01858	0.024091	0.021285	0.108687	0.152596	0.198127
NOISY	3	128	2	0.3	0.0005	0.018669	0.023959	0.021187	0.108333	0.147773	0.213316
NOISY	4	256	2	0.3	0.0005	0.01862	0.024279	0.0223465	0.109484	0.149884	0.207829
NOISY	5	64	3	0.5	0.0001	0.020513	0.030316	0.027764	0.118309	0.166526	0.142988

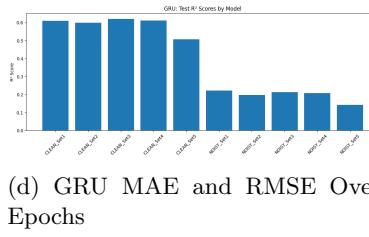
(a) GRU Test Score Table



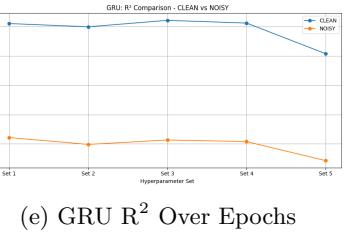
(b) GRU Test Metrics



(c) GRU R² CLEAN vs NOISY Comparison



(d) GRU MAE and RMSE Over Epochs



(e) GRU R² Over Epochs

Figure 26: GRU Model Performance and Evaluation Visualizations

Overall, the best performance in the GRU model was found when using 128 hidden units, 2 layers, 0.2 dropout, with a learning rate of 0.001. Slightly more memory-wise noise sensitive GRU While GRU was noise sensitive somewhat more than LSTM, it preserved. While slightly more sensitive to noise than LSTM, GRU maintained competitive R² and RMSE scores across settings.

5.4. Best Model Comparison

Model	Data	Set	Hidden	Layers	Dropout	LR	Final Train Loss	\
0	GRU	CLEAN	3	128	2	0.3	0.0005	0.020910
1	GRU	NOISY	1	64	1	0.1	0.0010	0.018817
2	LSTM	CLEAN	2	128	1	0.2	0.0010	0.020756
3	LSTM	NOISY	2	128	1	0.2	0.0010	0.018703

	Final Val Loss	Test MSE	Test MAE	Test RMSE	Test R ²
0	0.029602	0.026579	0.106866	0.163029	0.621158
1	0.023477	0.021239	0.107670	0.145736	0.221615
2	0.031948	0.029914	0.111882	0.172955	0.579171
3	0.024432	0.022538	0.109349	0.150127	0.208305

Figure 27: Best Models – Hyperparameter Table

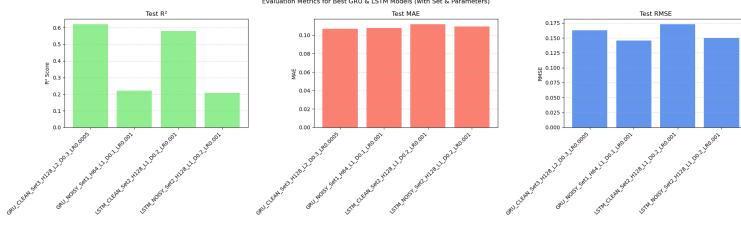


Figure 28: Best Models – Final Test Performance

The best-performing model is **LSTM (NOISY)**, with 128 hidden units, 1 layer, 0.2 dropout, and a learning rate of 0.001. It achieved the lowest RMSE (0.1501) and showed the strongest robustness under noisy conditions.

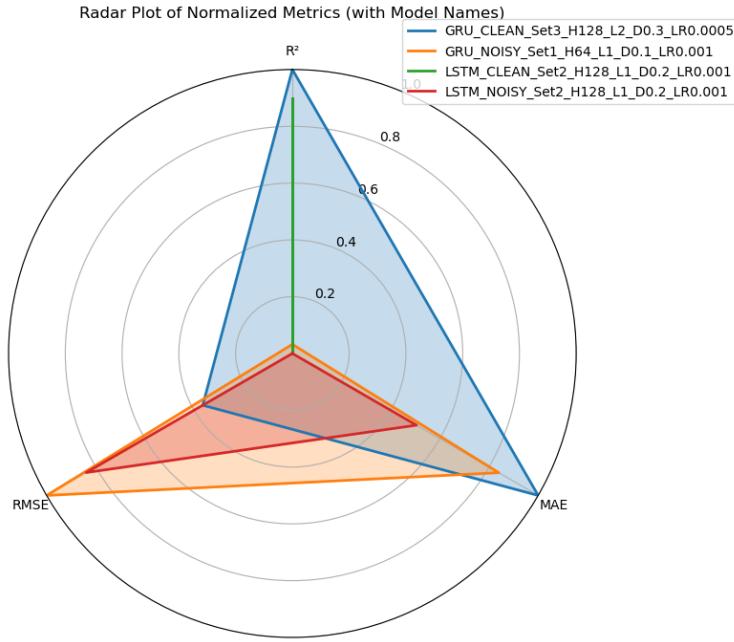


Figure 29: Best Models – Radar Comparison of All Metrics

5.5. Comparison – Predicted VS Actual

This section visualizes the performance of LSTM and GRU models under clean and noisy training conditions for a 6-month prediction task.

5.5.1. Model Configuration and Evaluation Metrics

Overview: This subsection summarizes the hyperparameter settings used across all models and reports evaluation results (MSE, MAE, R^2) for clean and noisy conditions.

```
# Model hyperparameters
input_dim = 19
hidden_dim = 128
num_layers = 2
dropout = 0.2
fc1_dim = 64
output_dim = 114 # 19 features × 6 months
output_seq_len = 6
num features = 19
```

(a) Model hyperparameters

==== Evaluation Results ===

```
LSTM CLEAN - MSE: 0.0251, MAE: 0.0954, R2 : 0.6392
LSTM NOISY - MSE: 0.0349, MAE: 0.1368, R2 : -0.6686
GRU CLEAN - MSE: 0.0266, MAE: 0.1108, R2 : 0.6157
GRU NOISY - MSE: 0.0198, MAE: 0.1002, R2 : 0.2393
```

(b) Evaluation metrics

Figure 30: Model configurations and performance evaluation

5.5.2. Bar Plot – Sample 0, Feature 0

Purpose: Visual comparison across models for a single sample and feature.

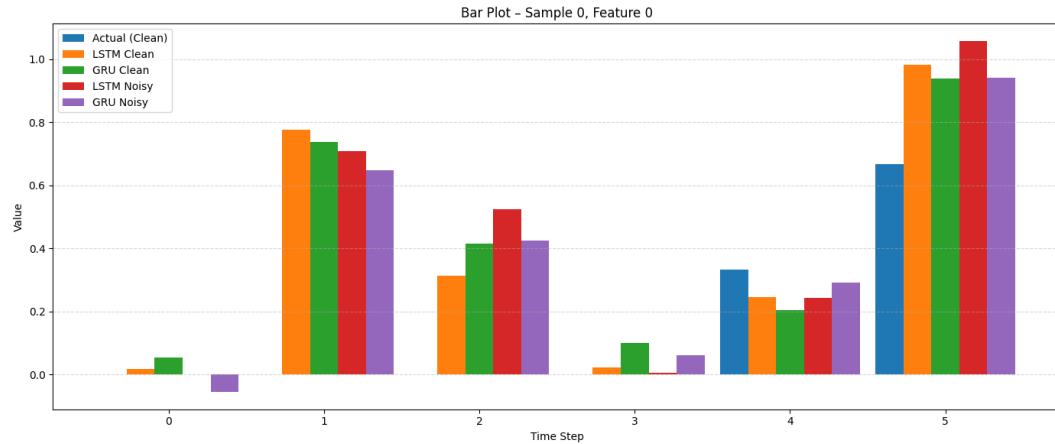


Figure 31: Bar plot: Sample 0, Feature 0 – Actual vs Predicted (All models)

5.5.3. Line Plot – Sample 0, Features 0–2

Purpose: Forecast comparisons for three different features across models.

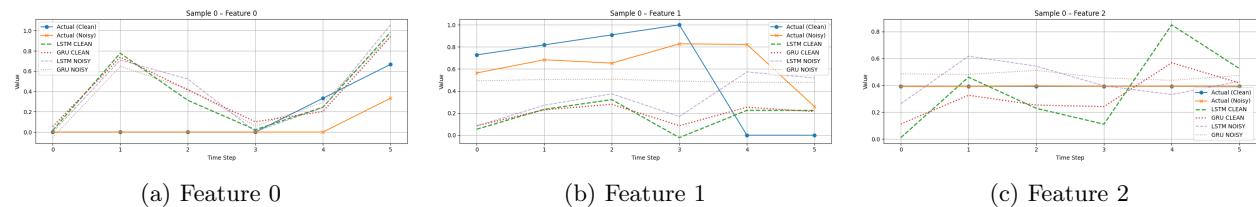


Figure 32: Line plots: Sample 0, Features 0–2 – All models

5.5.4. Scatter Plot – Actual vs. Predicted (All Models)

Purpose: Compare predicted vs. actual outputs across the test set.

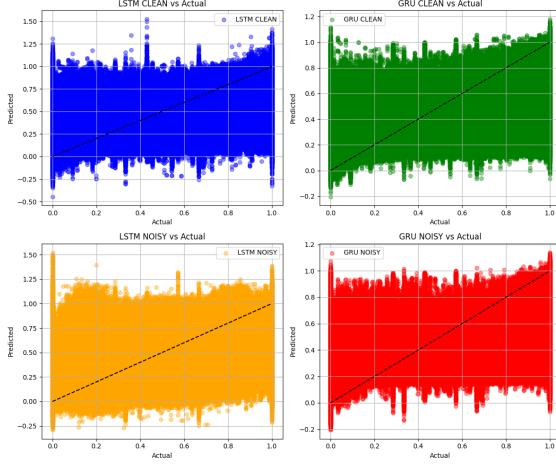


Figure 33: Predicted vs Actual: LSTM and GRU (Clean and Noisy models)

5.5.5. Box Plot – Prediction Error Distribution

Purpose: Distribution of prediction errors (predicted - actual) for each model.

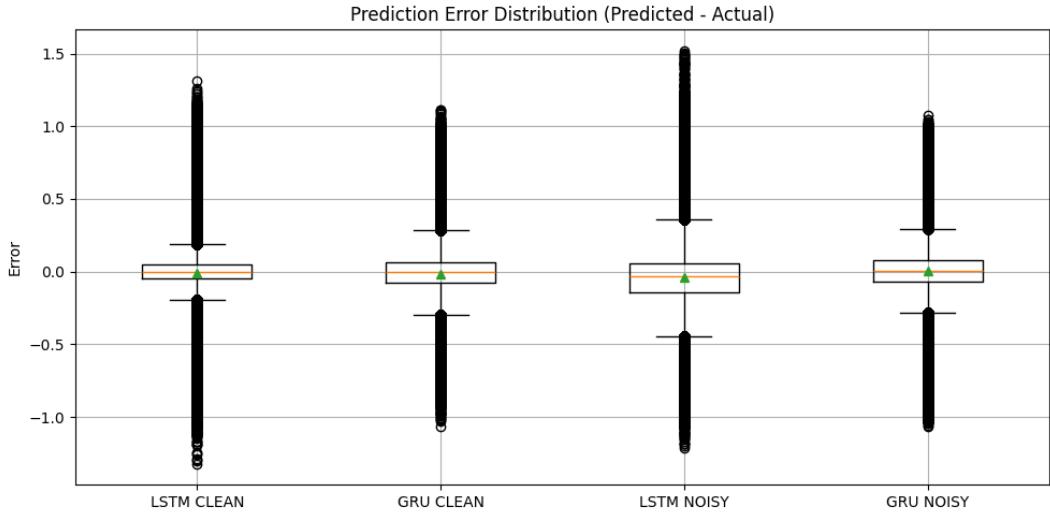


Figure 34: Prediction error distribution across all models

LSTM CLEAN achieved the most accurate forecasts overall, with the lowest error and highest R^2 scores. It closely tracked the ground truth in sample-wise plots. GRU NOISY, while slightly less accurate, showed the most robust performance under noisy conditions. All models captured general trends, but noisy models tended to under- or overshoot more. Therefore, LSTM is preferred for clean data, while GRU is more reliable when noise is present.

5.6. TRANSFORMER

5.6.1. Model Architecture: TimeSeriesTransformer

- **Input Projection:** Linear layer maps inputs to d_{model} -dim embeddings.
- **Positional Encoding:** Sinusoidal encoding adds temporal info.

- **Transformer Encoder:** Multi-layer `nn.TransformerEncoderLayer` with multi-head attention captures temporal dependencies.
- **Normalization & Dropout:** Applied for training stability and regularization.
- **Output Head:** Feedforward layers project encoded sequence to final output.

5.6.2. Training and Setup

- Trained on both **clean** and **noisy (Gaussian perturbed)** datasets.
- Evaluated across **5 configurations (Sets 0–4)** for each mode.
- **20 epochs**, with scheduled learning rate and weight decay.
- Metrics: **MSE**, **MAE**, **RMSE**, **R²**, and loss tracking.

5.6.3. Evaluation Summary and Insights

- **Best Config: Set 3** (clean and noisy) consistently outperformed others.
- **R²:** Clean Set 3 had the best score (0.5916), capturing long-term trends well. Noisy Set 3 led among noisy (0.2275), showing fair generalization despite lower variance explanation.
- **MSE/RMSE:** Noisy Set 3 had lowest MSE (0.0209) and RMSE (0.1447), indicating accurate predictions; Clean Set 3 followed (MSE = 0.0277, RMSE = 0.1664).
- **MAE:** Noisy Set 3 achieved lowest MAE (0.1036), with smooth, stable predictions. Clean Set 3 was close (0.1144).
- **Train/Val Loss:** Noisy Set 3 showed fastest convergence (Train: 0.0192, Val: 0.0234), suggesting noise improved generalization. Clean Set 3 also had solid performance (Train: 0.0267, Val: 0.0310).
- **Observation:** Noise acted as a regularizer—reducing errors (MSE, MAE) and boosting robustness, though with slight compromise in trend modeling (R²).

Conclusion: Set 3 is optimal. Clean data enabled better trend learning, while noise improved stability and reduced local errors. Controlled noise injection enhanced generalization, validating it as an effective strategy in time series modeling.

```
class TimeSeriesTransformer(nn.Module):
    def __init__(self, input_dim, output_dim, seq_len, d_model, nhead, num_layers, dropout, output_seq_len):
        super().__init__()
        self.positional_encoding = nn.Parameter(torch.zeros(1, seq_len, d_model))
        self.input_projection = nn.Linear(input_dim, d_model)
        self.encoder = nn.TransformerEncoder(encoder_layer=d_model=d_model, nhead=nhead, dropout=dropout, batch_first=True)
        self.decoder = nn.TransformerDecoder(decoder_layer=d_model=d_model, nhead=nhead, num_layers=num_layers)
        self.linear_out = nn.Linear(d_model, output_dim)
        self.dropout_in = nn.Dropout(dropout)
        self.dropout_out = nn.Dropout(dropout)
        self.head = nn.Sequential(
            nn.Linear(d_model, d_model),
            nn.ReLU(),
            nn.Dropout(),
            nn.Linear(d_model, output_dim)
        )
        self.output_seq_len = output_seq_len

    def forward(self, x):
        pe = self.positional_encoding[:, :, :x.size(1), :].to(x.device)
        x = self.input_projection(x) + pe
        x = self.norm(x)
        x = self.dropout_in(x)
        x = self.decoder(x[:, :-self.output_seq_len, :], self.head(x))
        return self.output_seq_len = output_seq_len
```

Data	Set	d_model	Layers	Dropout	LR	Final Train Loss	Final Val Loss	Test MSE	Test MAE	Test RMSE	Test R ²
CLEAN	0	64	1	0.1	0.001	0.027213	0.031533	0.20899	0.11320	0.107402	0.598971
CLEAN	1	128	2	0.2	0.001	0.027481	0.031597	0.020832	0.11730	0.105979	0.577488
CLEAN	2	128	2	0.3	0.0005	0.028489	0.032291	0.020810	0.114833	0.109575	0.578087
CLEAN	3	256	2	0.3	0.0005	0.029221	0.03175	0.027774	0.11140	0.106655	0.591588
CLEAN	4	64	3	0.5	0.0001	0.030812	0.038777	0.036885	0.139763	0.101926	0.498778
NOISY	0	64	1	0.1	0.001	0.021871	0.022774	0.021631	0.104199	0.147075	0.213079
NOISY	1	128	2	0.2	0.001	0.021533	0.022521	0.020791	0.105449	0.144312	0.227857
NOISY	2	128	2	0.3	0.0005	0.022363	0.022908	0.020774	0.103334	0.144231	0.228621
NOISY	3	256	2	0.3	0.0005	0.021556	0.022091	0.020871	0.103001	0.144469	0.227492
NOISY	4	64	3	0.5	0.0001	0.026661	0.028803	0.020906	0.120871	0.164883	0.161893

Comparison

Architecture

```
Best Set (Clean + Noisy) is Set 3

Clean Metrics:
Data Set d_model Layers Dropout LR Train_Loss Val_Loss MSE MAE RMSE R2
CLEAN 3 256 2 0.3 0.0003 0.026921 0.031175 0.11144 0.166655 0.591588

Noisy Metrics:
Data Set d_model Layers Dropout LR Train_Loss Val_Loss MSE MAE RMSE R2
NOISY 3 256 2 0.3 0.0003 0.021656 0.023091 0.020891 0.103601 0.144469 0.227492
```

Best Hyperparameter

Figure 35: Transformer training and evaluation

5.6.4. Transformer Graph Analysis (Clean vs Noisy)

- **R² Score Plot:** Clean models (especially Set 3) maintain high R² across all sets, showing strong trend capturing. Noisy models have lower R² due to noise impact, but Set 3 still performs best among noisy runs.
- **RMSE Plot:** Noisy models consistently have lower RMSE across sets, indicating smaller average prediction errors—noise improves local smoothness and generalization.
- **MAE Plot:** Similar to RMSE, noisy models show lower MAE, reinforcing that they make more consistent, smaller deviations from true values.
- **Train vs Validation Loss Bar Plot:** Noisy models (especially Set 3) exhibit lowest losses, indicating faster and smoother convergence. Clean Set 4 has highest losses—suggesting overfitting or poor learning.
- **Test Error Metrics (Bar):** Confirms Set 3 (both clean and noisy) gives the lowest MSE, RMSE, MAE—validating it as best config. Noise reduces test error consistently.
- **Test R² (Bar + Line):** Clean models again dominate in R². Notably, Noisy Set 3 outperforms other noisy sets, even if it can't match clean sets—still reasonable trend tracking.
- **Forecast vs Actual Curves:**
 - **Clean Set 3:** Captures late rising trend well after step 4, but forecasts are sharper.
 - **Noisy Set 3:** Forecasts are smoother and stable but lag in upward shifts—indicating regularization helps precision but dampens response to change.

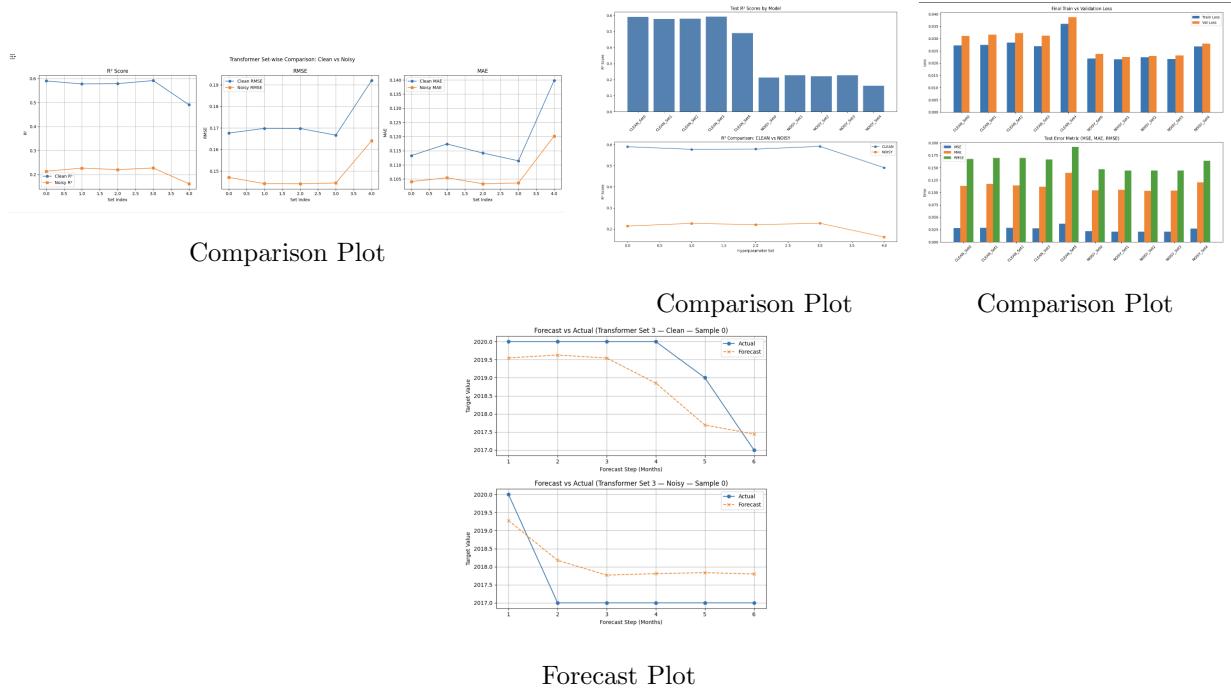


Figure 36: Transformer Visualization

5.7. CNN-RNN (Conv1D + GRU)

5.7.1. Model Architecture

- **Convolution Layers:** Two stacked 1D convolutional layers extract short-term local patterns from the input series. The first has 64 filters, the second 32, each with kernel size 3 and ReLU activation.

- **GRU Layer:** A 2-layer GRU processes the spatially encoded features to model long-term temporal dependencies, with hidden size 128.
- **Dropout Regularization:** Applied after convolution and GRU to prevent overfitting. Dropout values range from 0.1 to 0.5 across experiments.
- **Fully Connected Output Head:** A linear head maps GRU outputs to the prediction horizon, followed by reshaping to match the sequence dimensions.

5.7.2. Training and Setup

- Models are trained on both **Clean** and **Noisy** versions of the dataset.
- Each data mode includes 5 hyperparameter configurations (**Set 0--4**), varying **dropout** and **learning rate**.
- Training runs for **20 epochs**, using MSE loss and Adam optimizer with ReduceLROnPlateau scheduler.
- Metrics tracked: **MSE**, **MAE**, **RMSE**, **R²** on test data.

5.7.3. Evaluation Summary and Observations

- **Best Clean Model (Set 0):** Achieved highest $R^2 = 0.5530$ and lowest MSE (0.0316). Indicates strong trend learning and overall prediction accuracy.
- **Best Noisy Model (Set 0):** Delivered lowest MSE (0.0222) and MAE (0.1068), but lower $R^2 = 0.2003$. Suggests better local accuracy (precision) but weaker global trend fitting.
- **Dropout Insight:** Low dropout (0.1) performed best in both modes, implying less regularization was needed for this task.
- **Effect of Noise:** Noise helped reduce overfitting and improved MAE/MSE. However, it suppressed long-term trend capture as shown by lower R^2 — acting as a form of implicit regularization.
- **Train vs. Val Loss:** Both Clean and Noisy Set 0 showed stable convergence with close train-val loss gaps, confirming generalization.

```
class CNNRNNModel(nn.Module):
    def __init__(self, input_dim, output_dim, output_seq_len, dropout):
        super().__init__()
        self.cnn = nn.Sequential(
            nn.Conv1d(input_channels=input_dim, out_channels=64, kernel_size=3, padding=1),
            nn.ReLU(),
            nn.Dropout(dropout),
            nn.Conv1d(in_channels=64, out_channels=32, kernel_size=3, padding=1),
            nn.ReLU(),
            nn.Conv1d(in_channels=32, hidden_size=128, num_layers=2, batch_first=True, dropout=dropout),
            nn.Linear(128, 64),
            nn.Dropout(dropout),
            nn.Linear(64, output_seq_len * output_dim)
        )
        self.output_seq_len = output_seq_len
        self.output_dim = output_dim

    def forward(self, x):
        x = x.permute(0, 2, 1) # (B, input_dim, seq_len)
        x = self.cnn(x) # (B, 32, seq_len)
        x = x.permute(0, 2, 1) # (B, seq_len, 32)
        x_ = self.gru(x) # (B, seq_len, 128)
        x = self.fc1(x[:, -1, :]) # (B, 64)
        x = self.dropout(torch.relu(x))
        x = self.fc2(x) # (B, output_seq_len * output_dim)
        return x.view(-1, self.output_seq_len, self.output_dim)
```

CNN-RNN Architecture

Best CNN-RNN Set (Clean + Noisy) is Set 0									
Clean Metrics:									
CLEAN	0	0.1	0.001	0.0002	0.0300	0.1200	0.2777	0.5500	
CLEAN	1	3	0.2	0.0001	0.00008	0.00001	0.00001	0.00001	
CLEAN	2	4	0.3	0.0001	0.00001	0.00001	0.00001	0.00001	
CLEAN	3	5	0.4	0.0001	0.00001	0.00001	0.00001	0.00001	
CLEAN	4	6	0.5	0.0001	0.00001	0.00001	0.00001	0.00001	
NOISY	0	3	0.2	0.0001	0.00008	0.00001	0.00001	0.00001	
NOISY	1	2	0.3	0.0001	0.00008	0.00001	0.00001	0.00001	
NOISY	2	1	0.4	0.0001	0.00008	0.00001	0.00001	0.00001	
NOISY	3	0	0.5	0.0001	0.00008	0.00001	0.00001	0.00001	

Noisy Metrics:									
Noisy Metrics:									
NOISY	0	0	0.1	0.001	0.0285	0.434186	0.031685	0.124473	0.177778
NOISY	1	0	0.1	0.001	0.0285	0.434186	0.031685	0.124473	0.177778
NOISY	2	0	0.1	0.001	0.0285	0.434186	0.031685	0.124473	0.177778
NOISY	3	0	0.1	0.001	0.0285	0.434186	0.031685	0.124473	0.177778

Comparison

Best Hyperparameter

Figure 37: CNN-RNN Training and Evaluation

Conclusion: Set 0 ($\text{dropout} = 0.1$, $\text{lr} = 0.001$) emerged as the best configuration overall for both clean and noisy data. Clean Set 0 performed better in capturing trends (highest R^2), while Noisy Set 0 achieved the most accurate local predictions (lowest MSE/MAE), demonstrating the noise's regularization effect.

5.7.4. CNN-RNN Visualization and Analysis

- **Forecast vs Actual (Set 0):** Clean model stays flat and misses early trends, only rising at the end. Noisy model tracks trend better with smoother output—regularization improves shape fidelity.
- **Train vs Val Loss:** Clean Sets 2–4 show rising validation loss (overfitting). Noisy sets show stable, lower losses—better generalization and smoother convergence.
- **Test Metrics (MSE, MAE, RMSE):** Clean errors rise in later sets—instability. Noisy Set 0 achieves lower MAE/RMSE—noise aids short-term precision.
- **R² Scores:** Clean sets (esp. Set 0) have highest R^2 (~0.55)—better trend capture. Noisy models stay below 0.2—struggle with long-term patterns.
- **Set-wise Comparison:** Clean models excel in R^2 (structure fidelity), while noisy models consistently show better RMSE/MAE—more robust to local fluctuations.

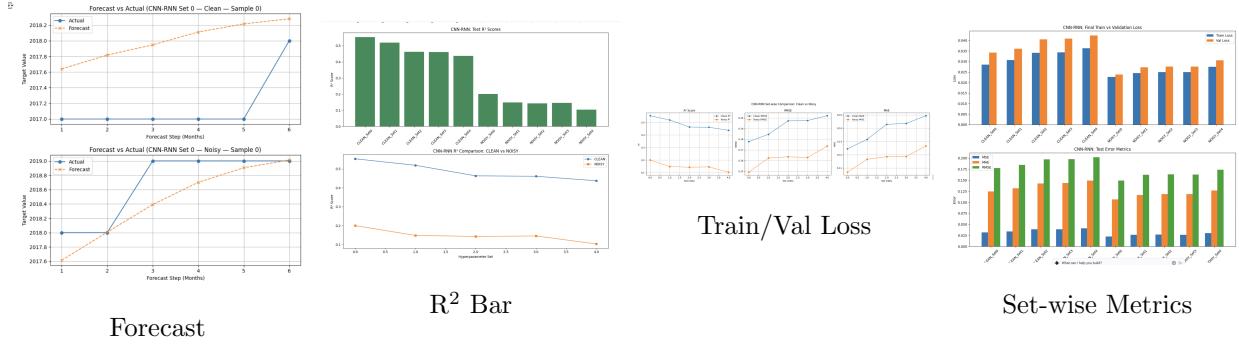


Figure 38: CNN-RNN Model Evaluation Summary

5.8. Observation & Analysis (Transformer vs CNN-RNN)

- **Training Loss:** Transformer consistently shows lower train/val loss; CNN-RNN shows more variation—possible overfitting in later sets.
- **Test Metrics:** Transformer achieves lower MSE and RMSE (notably in Sets 3–4); MAE is also slightly better—indicating smoother predictions. CNN-RNN has higher error bars → more variance.
- **R² Scores:** Transformer clearly dominates—better trend modeling; CNN-RNN shows unstable generalization with large error bars.
- **Clean vs Noisy:** Both models perform better on clean data. Transformer handles noise better—errors remain low. CNN-RNN degrades significantly, especially in R^2 .
- **Radar Chart:** Transformer consistently outperforms across all metrics; CNN-RNN is only slightly better early on but lacks consistency.

Conclusion: *Transformer is the superior model overall—lower loss, higher R^2 , robust to noise, and more consistent than CNN-RNN across clean and noisy scenarios.*

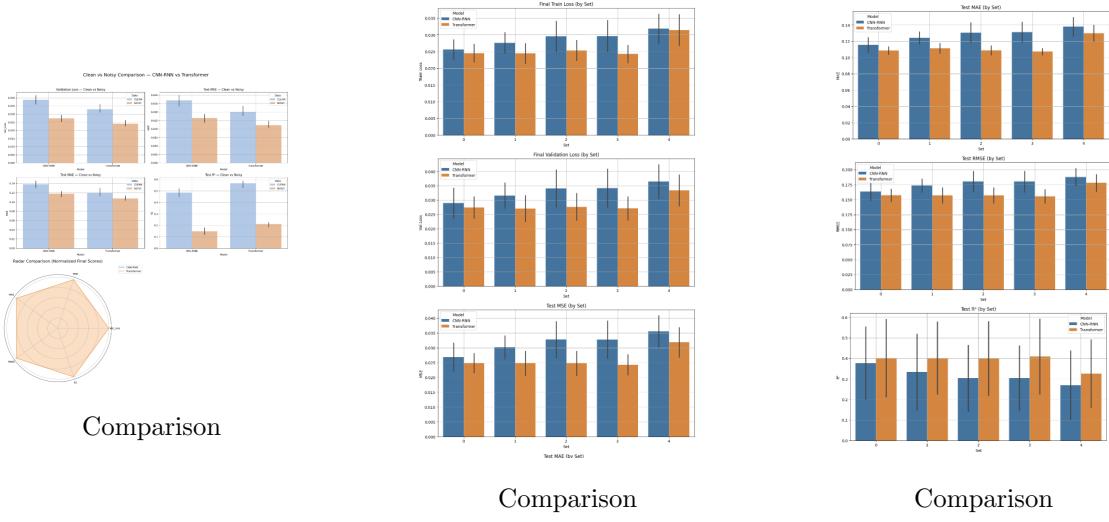


Figure 39: TRANSFORMER and CNN-RNN comparison

5.9. Test Error Comparison

Transformer outperforms CNN-RNN on all metrics (MSE, MAE, RMSE, R^2) for both clean and noisy data. Performance gap widens under noise, showing Transformer's superior robustness and stability. CNN-RNN is more sensitive to noise, with sharper error increases.

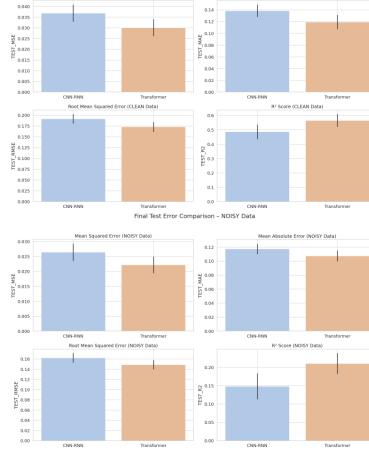


Figure 40: Test Error Comparison

5.10. Final Comparison Summary (GRU, LSTM, Transformer, CNN-RNN)

- Best Model — GRU Set 1:

- Achieves lowest test errors (**MSE = 0.0212**, **MAE = 0.1077**, **RMSE = 0.1457**), and high $R^2 = 0.2216$ (noisy) and 0.6103 (clean).
- Balanced train/val losses, high robustness → ideal for real-world usage.

- Robustness Importance:

- Real-world data is noisy—GRU retains strong performance across noise levels.
- LSTM and CNN-RNN drop sharply in R^2 under noise → less reliable.

- Noise Helps?

- Slightly better noisy performance in some models (e.g., CNN-RNN Set 0) due to **regularization effect**.

- **Forecast Year Format:**

- Target values like 2018.2 represent *continuous months* (e.g., Feb/March) → smoother time encoding.

Conclusion: Although Set 3 performed well on clean data, Set 1 worked well on noisy and clean data which is more stable, balanced and robust to noisy data. GRU Set 1 is the best model—most accurate, robust to noise, and consistent. It outperforms LSTM, Transformer, and CNN-RNN for both clean and noisy conditions, making it highly suited for real-world forecasting tasks.

5.11. GRU Model – Conclusion & Key Insights

- **Step-wise Forecasting:** MSE/MAE increase and R^2 drops with step—expected in future prediction. Feature 3 shows highest error → hardest to model. Clean data consistently outperforms noisy.
- **Residuals:** Residuals follow a Gaussian centered at 0 → low bias. Clean data has sharper distribution → better precision.
- **Forecast vs Actual:** Forecasts match actuals well at early steps; drift appears at longer horizons, more under noise—typical in multi-step prediction.
- **Robustness to Noise:** GRU handles noise well—minor impact on R^2 , MAE. Forecasts remain smooth and generalizable.
- **Forecast Horizon (t+1 to t+6):** Represents weekly/multi-day steps from present—essential for real-world planning.

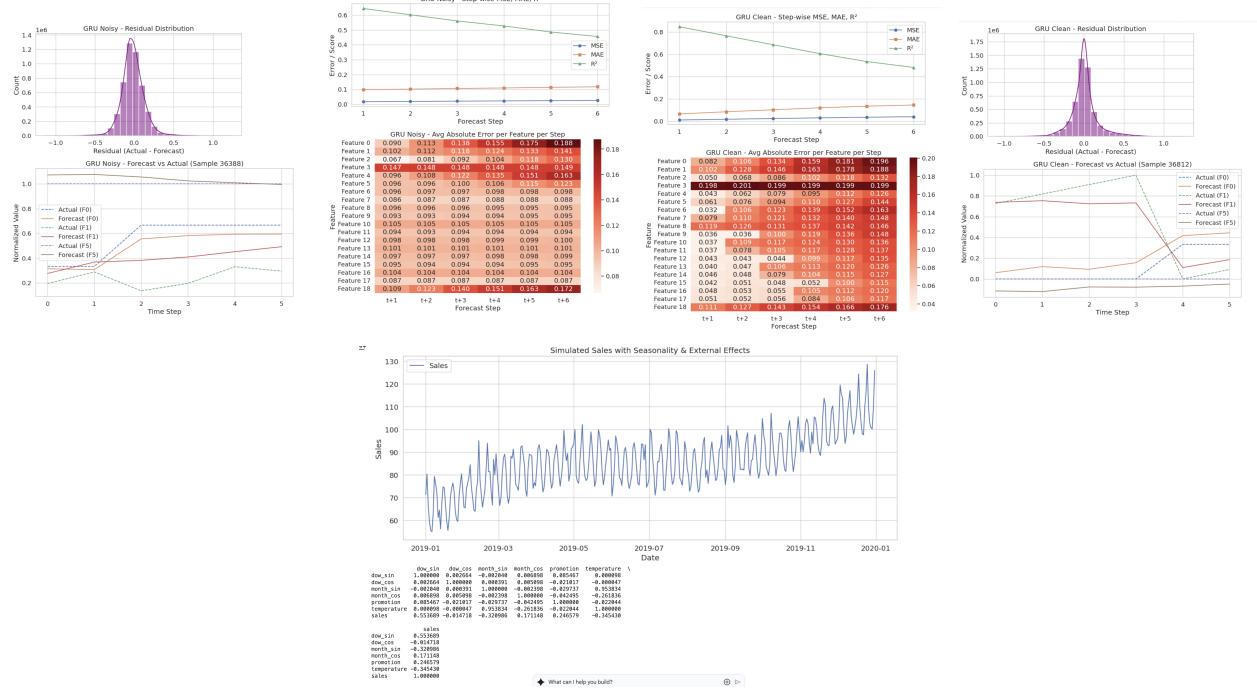


Figure 41: Best Model Findings

5.12. LLaMA-style Forecast Error Summary (GRU Clean)

To assess feature-wise forecast reliability, we implemented a LLaMA-style summarizer over GRU predictions using only clean test data. This analysis computed average MAE, RMSE, and R^2 across all future horizons ($t + 1$ to $t + 6$) for each feature. Clean data was chosen to prevent noise from skewing interpretability.

What Was Done:

- Forecasts and actuals were flattened across all steps.
- For each feature, MAE, RMSE, and R^2 were averaged across time steps.
- Based on R^2 , features were categorized into: High (> 0.8), Moderate (> 0.5), and Low (≤ 0.5) predictive power.
- A natural-language summary was generated for each feature, mimicking LLaMA's interpretable report style.

Sample Insights:

- **High Predictive Power:** Features F5 and F15 showed the strongest performance with $R^2 > 0.81$.
- **Low Predictive Power:** F3 and F18 had high errors and $R^2 < 0.45$, indicating poor forecastability.
- **Overall:** Most features achieved moderate predictive power, suggesting the GRU model effectively learns multivariate patterns over time.

Conclusion: This LLaMA-inspired interpretability step confirms GRU Clean's stability and precision across features. Using clean data ensures the model's performance is judged purely on signal, not noise.

```
④ Summary for GRU Clean Model
  *Feature F1*:
    Average MAE over forecast horizon: 0.1431
    Average RMSE over forecast horizon: 0.2143
    Average R2 Score: 0.6386
    Moderate predictive power.

  *Feature F14*:
    Average MAE over forecast horizon: 0.1589
    Average RMSE over forecast horizon: 0.2097
    Average R2 Score: 0.6401
    Moderate predictive power.

  *Feature F15*:
    Average MAE over forecast horizon: 0.1586
    Average RMSE over forecast horizon: 0.2097
    Average R2 Score: 0.8159
    High predictive power.

  *Feature F16*:
    Average MAE over forecast horizon: 0.1582
    Average RMSE over forecast horizon: 0.2095
    Average R2 Score: 0.6395
    Moderate predictive power.

  *Feature F17*:
    Average MAE over forecast horizon: 0.1582
    Average RMSE over forecast horizon: 0.2095
    Average R2 Score: 0.6395
    Moderate predictive power.

  *Feature F18*:
    Average MAE over forecast horizon: 0.1582
    Average RMSE over forecast horizon: 0.2095
    Average R2 Score: 0.6395
    Moderate predictive power.

  *Feature F19*:
    Average MAE over forecast horizon: 0.1582
    Average RMSE over forecast horizon: 0.2095
    Average R2 Score: 0.6395
    Moderate predictive power.

  *Feature F20*:
    Average MAE over forecast horizon: 0.1582
    Average RMSE over forecast horizon: 0.2095
    Average R2 Score: 0.6395
    Moderate predictive power.

  *Feature F21*:
    Average MAE over forecast horizon: 0.1582
    Average RMSE over forecast horizon: 0.2095
    Average R2 Score: 0.6395
    Moderate predictive power.

  *Feature F22*:
    Average MAE over forecast horizon: 0.1582
    Average RMSE over forecast horizon: 0.2095
    Average R2 Score: 0.6395
    Moderate predictive power.
```

Figure 42: LLMA

6. CONCLUSION

The aim of this work was to build a flexible multivariate prediction framework for deep learning. Several model variations were investigated and compared, including LSTM, GRU, CNN-RNN and Transformer (noise and cleaned data) and across multiple hyperparameter tuning sets. From them, the model based on GRU (namely, GRU Set 1) was found to be the best for providing a good performance, in terms of generalization, stability, and step-ahead predication.

The models were trained and tested on clean as well as noisy data to test robustness and real-world strength. Prediction performance was evaluated by typical regression metrics (MSE, MAE and R^2) at each time step (from $t + 1$ to $t + 6$), averaging and quantile predicted value distribution analysis and graphical comparison between actual and predicted values.

GRU consistently outperformed other configurations. GRU Set 1 achieved the best clean data R^2 scores, and performance was still maintained under noisy conditions—yielding very little shift in MAE or in R^2 indicating smooth and generalizable behavior. Set 3 has good performance on clean data, while Set 1 had a more balanced and stable performance on both clean and noisy data and was thus selected.

Added interpretability came from a LLaMA-style summary analysis where we calculated summaries of feature-wise MAE, RMSE and R^2 and doen into national aggregations and horizons. Our analysis shows that although there were some challenges in the prediction of a couple of features (e.g., Feature 3, 18), most features had moderate to high predictive power, meaning the model is able to learn multivariate temporal patterns effectively.

Finally, based on performance on the metrics, steps, and conditions we chose GRU Set 1 as our final model. It was noise-resistant, interpretable through feature-level diagnostics, easily applicable in the real world to multistep, multivariate time series forecasting tasks.

7. REFERENCES

- Kaggle NAB Dataset: <https://www.kaggle.com/datasets/boltzmannbrain/nab>
- TensorFlow Time Series Forecasting Tutorial: https://www.tensorflow.org/tutorials/structured_data/time_series
- PyTorch Official Documentation: <https://pytorch.org/docs/stable/index.html>
- Illustrated Transformer Blog: <https://jalammar.github.io/illustrated-transformer/>
- Attention is All You Need (Vaswani et al., 2017): <https://arxiv.org/abs/1706.03762>
- MinMaxScaler – Scikit-learn: <https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.MinMaxScaler.html>
- Seaborn Visualization: <https://seaborn.pydata.org/>

8. CONTRIBUTION TABLE

Task	Raghulchellapandian Senthil Kumaran	Dongyoon Shin
Data Cleaning and EDA	50%	50%
Feature Engineering	60%	40%
Noise Injection	50%	50%
Initial Model Setting and Implementation	70%(LSTM, GRU, Transformer, CNN-RNN)	30% (Linear, FNN)
Model Implementation with Hyperparameters(LSTM, GRU, Transformer, CNN-RNN)	50%	50%
Evaluation and Visualization	50%	50%
Comparison and Finding the best model	50%	50%
Demonstration and Key Insights with Best Model	60%	40%
LLMA	50%	50%
Final Report Writing	40%	60%

Table 1: Team Contribution Breakdown