

# Machine Learning for CO<sub>2</sub> Emission Prediction using Parallel Computing

Hafizah Farhah Saipan Saipol, Syarifah Zyurina Nordin  
Malaysia-Japan International Institute of Technology, Universiti Teknologi Malaysia,  
54100 Kuala Lumpur  
Email: hafizah.farhah@utm.my, szyurina@utm.my



## INTRODUCTION

Carbon dioxide (CO<sub>2</sub>) emissions are the primary driver of global climate change. Human activities such as burning fossil fuels and land-use changes have raised atmospheric CO<sub>2</sub> from 278 ppm in the 18<sup>th</sup> century to 420 ppm in 2024 which leads to a 1.1°C rise in global temperature since 1880. The ongoing increase aggravates global warming, heightening climate-related hazards such as high temperatures, rising sea levels, and loss of biodiversity. Climate change mitigation requires accurate and timely CO<sub>2</sub> emission forecasting for policy decisions. Machine learning models have emerged as powerful tools for emissions prediction, which offer superior accuracy compared to traditional statistical methods [1]. However, computational efficiency remains a critical bottleneck in operational deployment, particularly when evaluating multiple algorithms or conducting hyperparameter optimization [2].

## PROBLEM STATEMENT

Sequential training of multiple ML models for emissions prediction is computationally expensive and time intensive. With increasing availability of multi-core processors, there exists significant potential for performance optimization through parallel computing architectures. Limited works specifically addresses parallel optimization for CO<sub>2</sub> emissions prediction models.

## RESEARCH OBJECTIVES

- The objectives of this research is to:
- a) Implement parallel processing framework for training ML models
  - b) Evaluate parallel performance compared to sequential performance
  - c) Validate parallelization maintains prediction accuracy.

## METHODOLOGY

1. Dataset	<div>1. Source: OWID CO<sub>2</sub> emission dataset</div> <div>2. Temporal coverage: 1990 – 2023 (34 years)</div> <div>3. Geographic focus: 8 Southeast Asian Countries</div> <div>4. Features used: 72 variables including emissions metrics, energy consumption, economic indicators</div> <div>5. Target variable: Total CO<sub>2</sub> emissions (Mt)</div>
2. Data Preprocessing	<div>1. Missing Value treatment: Dropped columns with &gt;20% missing data; applied median imputation for remaining gaps</div> <div>2. Feature Selection: Random Forest filtering (threshold &gt;0.01), reducing to 23 features.</div> <div>3. Normalization: StandardScaler for distance-based algorithms (SVR)</div>
3. Machine Learning Models	<div>Five supervised models were developed</div> <div>i) Random Forest: Ensemble robustness</div> <div>ii) Gradient Boosting: Sequential error correction</div> <div>iii) Neural Network: Baseline linear model with regularization</div> <div>iv) Extra Trees: Randomized ensemble approach and feature importance capability</div> <div>v) Support Vector Regression: non-ensemble baseline, allowing comparison with tree-based and neural models</div>
4. Parallel Implementation	<div>Framework: Python 3.10 with joblib 1.3.2</div> <div>Backend: Loky (process-based parallelism via joblib)</div> <div>Hardware: 4-core CPU, 16GB RAM (Kaggle cloud-based platform)</div>
5. Performance Metrics	<div>Computational efficiency:</div> <div>i) Execution time (T): Total training time (sec)</div> <div>ii) Speedup (S): <math>S = \frac{T_{seq}}{T_{parallel}}</math></div> <div>iii) Parallel efficiency (Eff): <math>Eff = \frac{S}{p} * 100\%</math></div> <div>Prediction accuracy are based on:</div> <div>• Mean Absolute Error (MAE)</div> <div>• Root Mean Squared Error (RMSE)</div> <div>• Coefficient of Determinant (R<sup>2</sup>)</div> <div>• Mean Absolute Percentage Error (MAPE)</div>

## RESULTS

Table 1. Model performance metrics

Model	Parameters (Tuned)	MAE	RMSE	R <sup>2</sup>	MAPE
Neural Network (Large)	Hidden_layers = (200, 100, 50) Max_iteration=1000 Learning rate=0.001	3.24	5.38	0.9989	1.77%
Extra Trees (High)	n_estimators=1000, max_depth=20	5.18	9.68	0.9964	2.47%
Random Forest (High)	n_estimators=1000, max_depth=20	7.25	12.33	0.9941	4.85%
Random Forest (Med)	n_estimators=500, max_depth=20	8.05	14.06	0.9923	5.58%
Gradient Boosting (High)	n_estimators=1000, max_depth=7	8.48	14.10	0.9923	7.23%

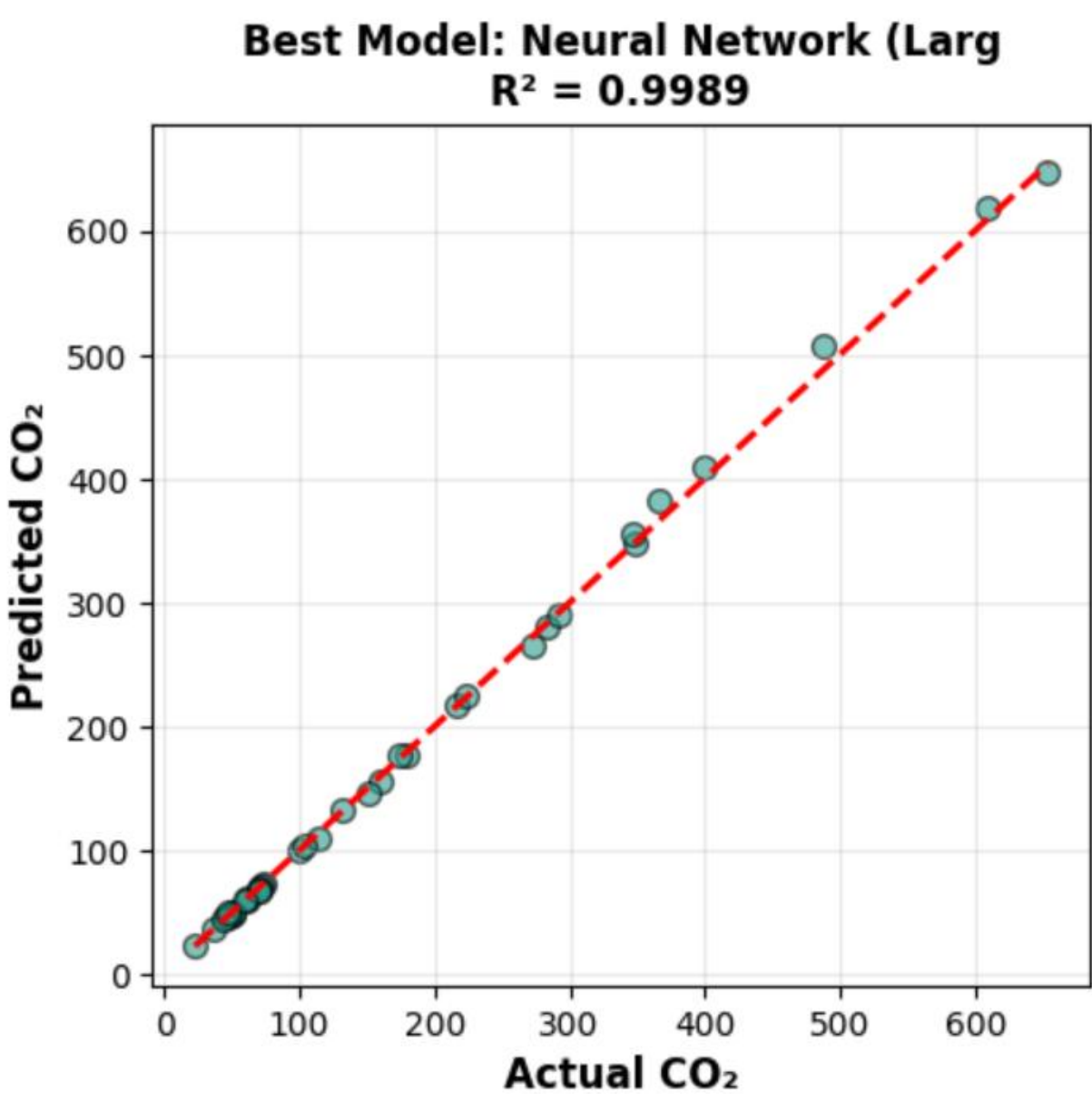


Fig. 1. Actual vs predicted CO<sub>2</sub>



Fig. 2. Sequential and parallel performance

Table 2. Performance comparison

Details	Performance
Sequential time	10.25 sec
Parallel time	5.45 sec
Time saved	4.80 sec (46.8%)
Speedup	1.88x
Parallel efficiency	47.0%
CPU cores used	4
Dataset size	148 training samples
Models trained	9

## CONCLUSION

This study demonstrates that parallel computing significantly enhances the efficiency of machine learning models for CO<sub>2</sub> emissions prediction without compromising model accuracy. By utilizing Kaggle-based multi-core environment with four CPUs, the parallel execution achieved a 1.88x speedup and 47% parallel efficiency which reduce total computation time by 46.8% compared to sequential execution. Among the nine trained models, the Neural Network (Large) achieved the best predictive accuracy with R<sup>2</sup> = 0.9889 and MAPE = 1.77%, followed closely by ensemble-based models such as Extra Trees, and Random Forest. To conclude, as climate modeling demands increase and datasets grow larger, parallel optimization will become essential rather than optional. This research could provide a solution for accelerating environmental ML applications while maintaining scientific rigor.

## REFERENCES

[1] Jin, Y. and Sharifi, A., 2025. Machine learning for predicting urban greenhouse gas emissions: A systematic literature review. Renewable and Sustainable Energy Reviews, 215, p.115625.  
[2] Hassanpouri Baesmat, K., Farrokhi, Z., Chmaj, G. and Regentova, E.E., 2025. Parallel Multi-Model Energy Demand Forecasting with Cloud Redundancy: Leveraging Trend Correction, Feature Selection, and Machine Learning. Forecasting, 7(2), p.25.

## ACKNOWLEDGEMENTS

We would like to thank MJLIT for the continuous support.