

Benchmarking Statistical, Machine Learning, Deep Learning, and Hybrid Forecasting Models for Global Renewable Energy Consumption: A Walk-Forward Cross-Validation Study with Structural Break Analysis.

Abstract

Energy independence and resilience have become critical policy priorities as geopolitical tensions, supply disruptions, and price volatility expose the vulnerability of fossil-fuel-dependent energy systems. Accurate forecasting of renewable energy consumption is therefore essential for effective energy transition planning, infrastructure investment, and monitoring progress toward international climate targets such as Sustainable Development Goal 7 (SDG-7). In macro-energy policy practice, quantitative forecasts underpin scenario design, capacity planning, and assessment of alignment with net-zero pathways, yet the annual frequency and short length of globally comparable time series severely constrain the effective application of data-intensive forecasting methods. This study benchmarks forecasting model families—spanning baselines (Naïve, Random Walk with Drift, Linear Trend), classical statistical methods (ETS, Damped ETS, Theta, ARIMA), machine learning (XGBoost), deep learning (GRU, LSTM, N-BEATS), an additive model (Prophet), and a novel ETS–GRU hybrid—against the World Bank EG.FEC.RNEW.ZS indicator (1990–2020). All models are evaluated under a unified 5-window expanding walk-forward cross-validation protocol with a 3-year forecast horizon, nested hyperparameter tuning, multi-seed deep learning robustness checks, Diebold–Mariano tests, Model Confidence Set analysis, and bootstrap inference. A Chow structural break test detects a statistically significant regime shift at 2014 ($F = 32.0$, $p < 0.001$), coinciding with the post-Paris Agreement acceleration in renewable deployment. Results indicate that Holt Linear Exponential Smoothing (ETS) achieves the lowest RMSE (0.543) with a Skill Score of +0.148 against the Naïve baseline, outperforming all deep learning architectures. The study introduces three novel energy transition analytics—a Transition Velocity Index, regional beta-convergence analysis, and SDG-7 gap assessment across 11 World Bank regions—and demonstrates that under a methodologically symmetric evaluation protocol that eliminates the information asymmetries present in prior benchmarking studies, parsimonious statistical models offer superior forecasting performance in small-sample annual energy data regimes.

Keywords: *Renewable energy forecasting; Walk-forward cross-validation; Exponential smoothing; Deep learning; Structural break; SDG-7; Energy transition*

1. Introduction

The global energy transition from fossil fuels to renewable sources represents one of the most consequential structural shifts in the contemporary world economy. Recent geopolitical tensions, supply disruptions, and fossil fuel price volatility have underscored the strategic importance of energy independence and exposed the vulnerability of energy systems that

40 depend heavily on concentrated hydrocarbon supply chains. Driven by climate commitments
41 under the Paris Agreement (UNFCCC, 2015), declining technology costs, and energy security
42 considerations, the share of renewable energy in total final consumption has accelerated
43 markedly since 2014. Monitoring and forecasting this transition at the global and regional
44 scales is essential for policymakers, international organisations, and energy system planners
45 who must allocate resources, design incentive structures, and assess progress toward
46 internationally agreed targets such as Sustainable Development Goal 7 (SDG-7).

47 Time series forecasting of energy variables has attracted substantial research attention
48 over the past decade, with a pronounced shift toward machine learning (ML) and deep
49 learning (DL) methods. Recurrent neural networks, particularly Long Short-Term Memory
50 (LSTM) and Gated Recurrent Unit (GRU) architectures, have demonstrated competitive
51 performance on electricity demand, wind power, and solar irradiance forecasting tasks where
52 high-frequency data (hourly or sub-hourly) provide thousands to millions of training
53 observations (Ahmed et al., 2024; Wang et al., 2019). More recently, purpose-built
54 architectures such as N-BEATS (Oreshkin et al., 2020) and Transformer-based models have
55 further expanded the methodological frontier.

56 However, a significant gap exists between the data regimes in which DL methods have
57 proven effective and the data regimes that characterise internationally comparable energy
58 statistics. The World Bank's flagship renewable energy indicator (EG.FEC.RNEW.ZS)
59 provides annual observations from 1990 onward, yielding approximately 31 data points for
60 any given country or regional aggregate. At this sample size, deep learning models with
61 thousands of trainable parameters face severe underdetermination, and the risk of overfitting
62 to noise rather than capturing genuine temporal structure becomes acute (Makridakis et al.,
63 2018). Despite this, a growing number of studies apply DL methods to annual energy data
64 without adequate acknowledgement of the small-sample constraint or rigorous comparison
65 against parsimonious classical alternatives.

66 Furthermore, methodological inconsistencies in the existing benchmarking literature
67 undermine the reliability of reported model rankings. A common practice is to evaluate
68 ARIMA with rolling one-step-ahead refitting while assessing DL models on a single frozen
69 fit—an information asymmetry that structurally inflates the reported accuracy of ARIMA
70 (Hyndman and Athanasopoulos, 2021). Similarly, many studies report results from a single
71 train–test split, which is particularly problematic when the underlying data-generating
72 process exhibits structural breaks.

73 This study addresses these limitations through a comprehensive benchmarking framework
74 that enforces methodological fairness across all model families. Whereas prior work has
75 examined individual model families or used inconsistent evaluation protocols, this study
76 makes five specific contributions that collectively advance the state of the art:

- 77 (1) A unified walk-forward cross-validation protocol that evaluates 13 model families
78 under identical conditions—including identical information sets, identical forecast

79 horizons, and no intermediate retraining—eliminating the information asymmetry
80 present in prior benchmarking studies.

81 (2) Formal structural break detection via the Chow test, integrated into the evaluation
82 design such that walk-forward windows are deliberately positioned to straddle the
83 identified breakpoint, providing a principled measure of each model’s adaptability to
84 regime change.

85 (3) A rigorous statistical inference framework combining Diebold–Mariano tests, paired
86 bootstrap comparison, and Model Confidence Set analysis, providing set-level
87 inference that accounts for the multiple comparison problem inherent in
88 benchmarking many models simultaneously.

89 (4) Multi-seed deep learning robustness experiments (10 seeds per architecture) that
90 quantify the sensitivity of DL results to random initialisation on small datasets,
91 directly addressing the reproducibility concerns that pervade small-sample DL
92 applications.

93 (5) Three novel energy transition analytics—a Transition Velocity Index (TVI), regional
94 beta-convergence analysis, and SDG-7 gap assessment across 11 World Bank
95 regions—that bridge the gap between forecasting methodology and actionable policy
96 intelligence for energy independence planning.

97 Together, these contributions provide the most comprehensive and methodologically
98 rigorous benchmarking study of forecasting models applied to annual macro-energy
99 indicators to date.

100 **2. Literature Review**

101 **2.1. Renewable Energy Forecasting**

102 Renewable energy forecasting encompasses a broad methodological spectrum, ranging from
103 physics-based models for wind and solar output to purely data-driven statistical and machine
104 learning approaches. At the macro level, forecasting aggregate renewable energy
105 consumption as a share of total final energy is critical for national energy planning and
106 international policy assessment (IRENA, 2023). Early work in this domain relied
107 predominantly on exponential smoothing, ARIMA, and regression-based approaches
108 (Hyndman and Athanasopoulos, 2021; Makridakis et al., 2018).

109 The IEA World Energy Outlook and IRENA’s Global Energy Transformation reports
110 provide scenario-based projections that incorporate techno-economic assumptions, but these
111 are not strictly time series forecasts and rely heavily on expert judgement (IEA, 2023).
112 Statistical forecasting of energy indicators using World Bank data has been explored by
113 several authors. Karakurt and Aydin (2023) applied ARIMA and ETS to energy intensity
114 indicators, while Khan et al. (2020) used regression-based approaches for renewable energy
115 share projections. However, few studies have systematically benchmarked multiple model
116 families against the same indicator using rigorous cross-validation protocols.

117 **2.2. Machine Learning and Deep Learning in Energy Forecasting**

118 The application of machine learning to energy forecasting has expanded rapidly since 2015.
119 XGBoost and gradient boosting methods have shown strong performance on electricity price
120 and demand forecasting tasks where engineered features (lag values, calendar variables,
121 weather covariates) provide rich input representations (Chen and Guestrin, 2016; Lago et al.,
122 2021). Deep learning approaches, particularly LSTM networks (Hochreiter and Schmidhuber,
123 1997), have been widely applied to electricity demand (Kong et al., 2019), wind power
124 (Wang et al., 2019), solar irradiance forecasting (Ahmed et al., 2024), and renewable energy
125 consumption forecasting at the macro level (Biswas, Irshad and Roy, 2026).

126 GRU networks (Cho et al., 2014) offer a computationally lighter alternative to LSTM
127 with comparable performance on many time series tasks. N-BEATS (Oreshkin et al., 2020), a
128 feed-forward architecture using residual stacking and polynomial basis expansion, achieved
129 state-of-the-art results on the M4 competition dataset. Prophet (Taylor and Letham, 2018)
130 employs an additive decomposition framework with automatic changepoint detection, making
131 it accessible for practitioners without specialist forecasting expertise.

132 Despite these advances, the performance advantage of DL methods is most pronounced
133 with high-frequency data containing thousands of observations. Makridakis et al. (2018)
134 demonstrated in the M4 competition that simple statistical methods outperformed complex
135 ML and DL models on many time series, particularly those with fewer than 100 observations.
136 Petropoulos et al. (2022) reinforced this finding in a comprehensive review of forecasting
137 methodology, noting that model complexity should be calibrated to data availability. De
138 Oliveira and Cyrino Oliveira (2018) further demonstrated that bagging methods applied to
139 ARIMA and exponential smoothing could outperform standalone complex models on mid-to-
140 long-term energy forecasting tasks, while Deb et al. (2017) provided a systematic review
141 showing that classical methods remain competitive for building energy consumption
142 forecasting when data is limited.

143 **2.3. Hybrid Forecasting Models**

144 Hybrid models that combine statistical and ML/DL components have attracted growing
145 interest as a strategy to capture both linear trend structure and nonlinear residual patterns.
146 Zhang (2003) proposed the seminal ARIMA–ANN hybrid that decomposes a series into
147 linear and nonlinear components. More recent hybrids include ETS–LSTM combinations
148 (Smyl, 2020), which won the M4 competition, and ensemble approaches that average
149 forecasts from complementary model families (Atiya, 2020). The rationale is that statistical
150 models efficiently capture trend and seasonality with minimal parameters, while DL
151 components model residual nonlinearity that escapes the statistical specification.

152 **2.4. Forecast Evaluation Methods**

153 Rigorous forecast evaluation extends beyond point accuracy metrics. The Diebold–Mariano
154 test (Diebold and Mariano, 1995) provides a formal framework for testing whether two
155 forecasting methods produce significantly different prediction errors. The Model Confidence

156 Set (Hansen et al., 2011) generalises this to the multi-model setting, identifying the subset of
157 models whose predictive ability cannot be statistically distinguished from the best model at a
158 given significance level. Bootstrap methods (Efron and Tibshirani, 1993) provide
159 distribution-free confidence intervals that are particularly valuable when the number of test
160 observations is small and distributional assumptions may be violated.

161 Walk-forward cross-validation (Tashman, 2000) is the standard evaluation protocol for
162 time series, expanding the training window sequentially and producing out-of-sample
163 forecasts at each step. This approach respects the temporal ordering of data and provides
164 multiple evaluation points, yielding more robust performance estimates than a single train-
165 test split.

166 **2.5. Research Gap**

167 Despite the extensive literature on energy forecasting, several gaps remain. First, few studies
168 benchmark a comprehensive set of model families—from simple baselines through statistical
169 methods to deep learning—under strictly identical evaluation conditions on annual macro-
170 energy indicators. Second, the small-sample challenge inherent to annual energy statistics is
171 rarely addressed explicitly; most DL applications to energy data use high-frequency datasets
172 where the data regime is fundamentally different. Third, structural break detection is seldom
173 integrated into the forecasting evaluation framework, despite the known acceleration in
174 renewable energy deployment following the Paris Agreement. Fourth, no prior study has
175 combined rigorous forecasting benchmarking with novel policy-relevant analytics (transition
176 velocity measurement, convergence testing, SDG-7 gap analysis) in a unified framework.
177 This study addresses all four gaps.

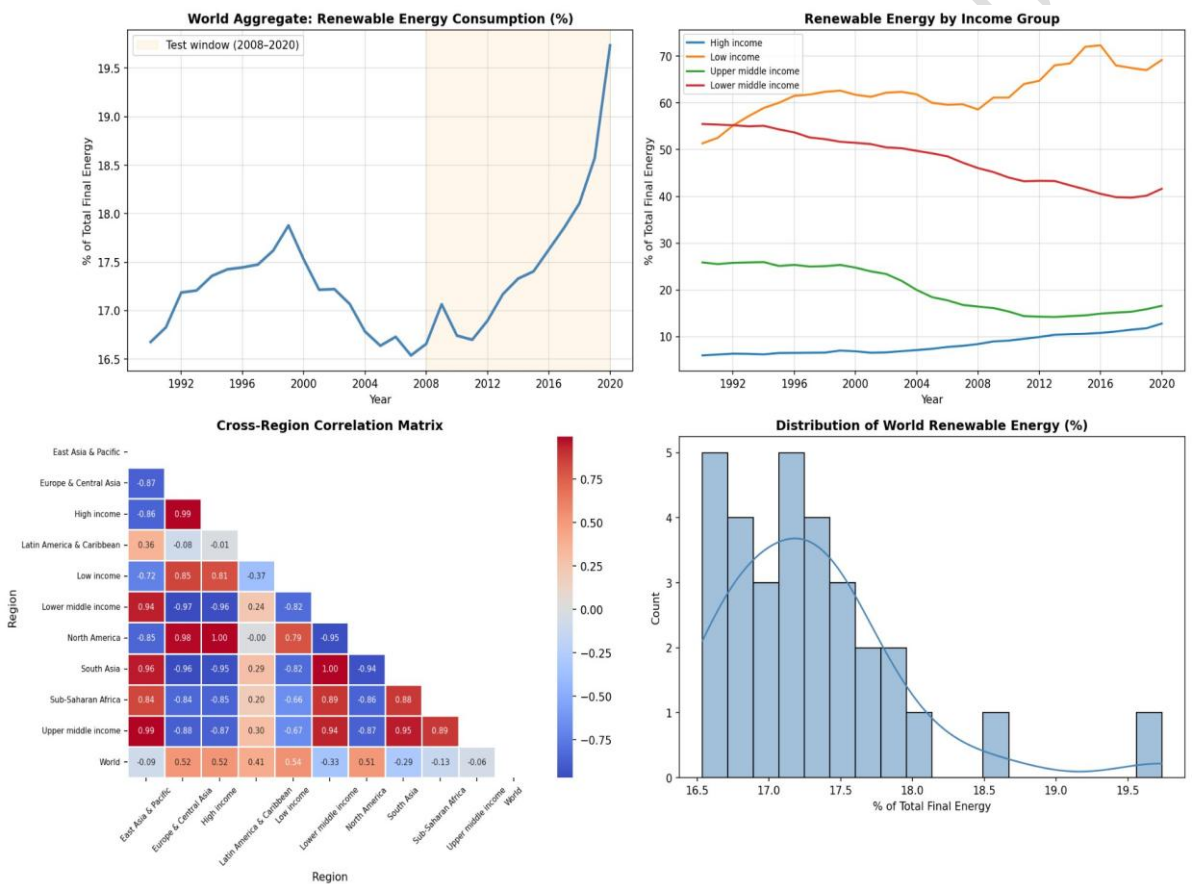
178 **3. Data**

179 The dataset used in this study is the World Bank indicator EG.FEC.RNEW.ZS, defined as
180 renewable energy consumption as a percentage of total final energy consumption. The
181 indicator covers 1990–2020, providing 31 annual observations for each of 11 aggregate
182 regional series: East Asia and Pacific, Europe and Central Asia, High income, Latin America
183 and Caribbean, Low income, Lower middle income, North America, South Asia, Sub-
184 Saharan Africa, Upper middle income, and World. Pre-1990 values were excluded from the
185 analysis as they contain backfilled identical values from a single 1990 anchor, carrying no
186 additional temporal signal.

187 It is important to note that the EG.FEC.RNEW.ZS indicator includes traditional solid
188 biomass, which constitutes a substantial proportion of renewable energy consumption in
189 developing regions. High renewable shares in Sub-Saharan Africa (approximately 70–91%)
190 and Low income regions (approximately 55–75%) primarily reflect reliance on wood fuel and
191 charcoal for cooking and heating rather than modern renewable energy deployment such as
192 wind, solar, or geothermal. This compositional distinction is critical for interpreting cross-
193 regional comparisons and policy implications.

194 As of the date of this study, 2020 remains the most recent year for which the
 195 EG.FEC.RNEW.ZS indicator is available across all World Bank regional aggregates. The
 196 underlying data are sourced from the IEA Energy Statistics Data Browser, whose SE4ALL
 197 tracking framework updates with a multi-year lag for global coverage. The 1990–2020 series
 198 therefore represents the complete available dataset for this indicator.

199 The primary forecasting target is the World aggregate series, which exhibited a U-shaped
 200 trajectory over the study period. From 1990 to approximately 2008, the global renewable
 201 share declined gradually from 16.7% to 16.6%, reflecting the more rapid growth of fossil fuel
 202 consumption relative to renewables. A turning point emerged around 2010–2013, followed
 203 by a marked upward acceleration from 2014 onward, with the World aggregate reaching
 204 19.8% by 2020.

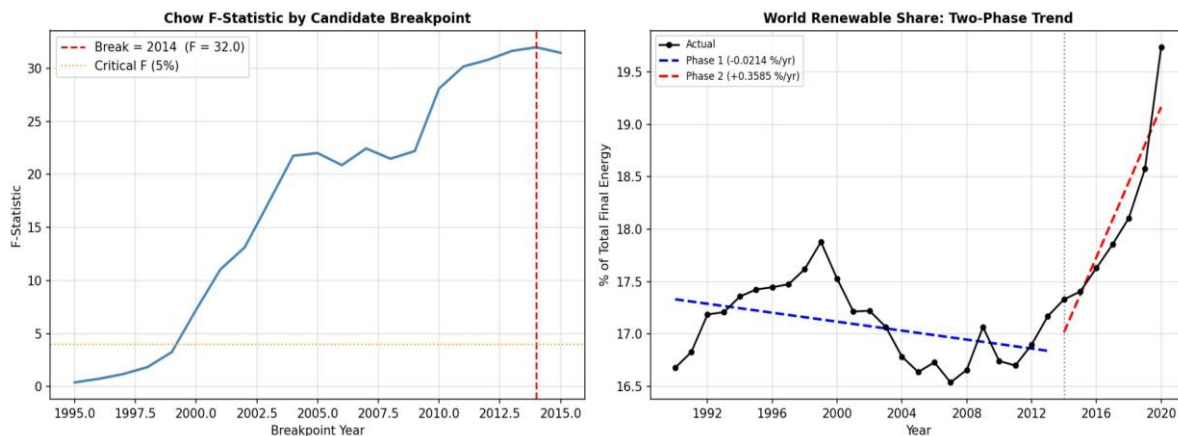


205
 206 *Figure 1. Exploratory Data Analysis: World renewable energy time series, income-group trajectories, cross-*
 207 *region correlation matrix, and distributional summary.*

208 3.1. Structural Break Evidence

209 A Chow F-test applied iteratively across all candidate breakpoints in the 1990–2020 World
 210 series identified 2014 as the year of the most statistically significant structural change ($F =$
 211 32.0 , $p < 0.001$). This result was confirmed by CUSUM test diagnostics. Segmented
 212 regression analysis revealed two distinct data-generating regimes: Phase 1 (1990–2013) with
 213 a near-zero trend of -0.021 percentage points per year, and Phase 2 (2014–2020) with a steep
 214 upward trend of $+0.359$ percentage points per year. This breakpoint coincides with the period

215 surrounding the Paris Agreement negotiations and the acceleration in global renewable
216 energy investment and policy deployment.



217
218 *Figure 2. Structural break detection: Chow F-statistic by candidate breakpoint (left) and two-phase segmented*
219 *regression (right).*

220 4. Methodology

221 This section describes the 13 forecasting models evaluated in the study, organised by
222 methodological category.

223 4.1. Baseline Models

224 Three baseline models establish the minimum performance threshold. The Naïve method
225 carries forward the last observed value for all forecast horizons. The Random Walk with Drift
226 extends this by adding the historical mean period-to-period change as a constant drift term.
227 The Linear Trend model fits an ordinary least squares regression of the series on a time index
228 and extrapolates the fitted line.

229 4.2. Statistical Models

230 Four classical statistical methods are evaluated. Holt Linear Exponential Smoothing (ETS)
231 uses two optimised parameters (α for level smoothing and β for trend smoothing) fitted via
232 maximum likelihood, re-optimised on each expanding window. Damped ETS extends Holt's
233 method with a damping parameter ϕ that gradually attenuates the trend, mitigating over-
234 extrapolation on longer horizons. The Theta method decomposes the series into two theta-
235 lines with different curvatures and combines them, as proposed by Assimakopoulos and
236 Nikolopoulos (2000). ARIMA with automatic order selection (via the AIC criterion) is fitted
237 using the pmdarima library, with orders determined independently for each expanding
238 window.

239 4.3. Machine Learning Model

240 XGBoost (Chen and Guestrin, 2016) is evaluated with engineered lag features comprising
241 lags 1 through 5, a 3-period rolling mean, and a 3-period rolling standard deviation. Multi-
242 step forecasting is achieved iteratively: each predicted value is fed back as input for the

243 subsequent step. Hyperparameters (maximum tree depth and number of estimators) are tuned
244 via the nested cross-validation procedure described in Section 5.

245 **4.4. Deep Learning Models**

246 Three deep learning architectures are evaluated. The GRU (Cho et al., 2014) uses a gated
247 recurrent architecture with a hidden size of 64, employing input and reset gates to control
248 information flow without a separate cell state. The LSTM (Hochreiter and Schmidhuber,
249 1997) uses the canonical architecture with separate cell state and hidden state, also with a
250 hidden size of 64. Both recurrent models use a sequence length of 3, are trained for 300
251 epochs with a learning rate of 0.001 and batch size of 8, using the Adam optimiser. Min–max
252 normalisation is applied within each training window to prevent data leakage.

253 N-BEATS (Oreshkin et al., 2020) is a feed-forward architecture using stacked residual
254 blocks with polynomial basis expansion. The implementation uses 2 stacks with a hidden
255 dimension of 64 and a polynomial degree of 2. All deep learning models use iterative multi-
256 step forecasting, producing one-step-ahead predictions that are fed back as input for
257 subsequent steps.

258 **4.5. Additive Model**

259 Prophet (Taylor and Letham, 2018) is an additive decomposition model with automatic
260 changepoint detection. Configuration includes a `changepoint_prior_scale` of 0.3 (moderate
261 flexibility), with yearly, weekly, and daily seasonality disabled (annual data). Explicit
262 changepoints at 2001 (dot-com recession) and 2007 (global financial crisis) are provided as
263 prior knowledge.

264 **4.6. Hybrid Model**

265 The ETS–GRU hybrid is an equal-weight ensemble that averages the point forecasts of ETS
266 and GRU at each forecast step. This design is intentionally parsimonious: rather than
267 introducing learned combination weights (which would require additional training data that
268 this small-sample regime cannot support), the equal-weight average follows the theoretical
269 result of Atiya (2020) that simple averaging of complementary forecasters often matches or
270 exceeds optimally weighted combinations when the number of forecast observations is small.
271 The rationale is that ETS efficiently captures the dominant linear trend component while
272 GRU captures potential nonlinear residual structure. Both component forecasts are computed
273 independently under the same walk-forward protocol, and the ensemble is constructed post
274 hoc without additional parameter estimation.

275 **5. Experimental Design**

276 **5.1. Walk-Forward Cross-Validation**

277 The evaluation protocol uses a 5-window expanding walk-forward cross-validation design
278 with a fixed horizon of $H = 3$ years, yielding 15 total test observations. The training window
279 expands from 20 observations (1990–2009) in Window 1 to 28 observations (1990–2017) in

280 Window 5. Critically, every model—including ARIMA—trains once per window and
281 produces a single H-step-ahead forecast with no access to intermediate test actuals. This
282 eliminates the information asymmetry that arises when ARIMA is evaluated with rolling one-
283 step-ahead refitting while DL models receive a single forecast call.

284 The window placement is designed to interact with the structural break at 2014. Windows
285 1 and 2 forecast exclusively within the pre-break regime. Window 3 straddles the breakpoint,
286 forecasting 2014–2016 from training data that is predominantly pre-break. Windows 4 and 5
287 forecast increasingly into the post-break acceleration regime, providing a direct measure of
288 each model’s ability to adapt to a changed data-generating process.

289 **5.2. Nested Hyperparameter Tuning**

290 Deep learning and XGBoost hyperparameters are tuned via nested cross-validation to prevent
291 information leakage between the tuning and evaluation stages. For each outer walk-forward
292 window, an inner loop performs leave-one-out validation on the training portion, evaluating
293 candidate configurations and selecting the one with the lowest inner RMSE. The winning
294 configuration is then used to produce the outer-loop forecast. This ensures that
295 hyperparameter choices do not benefit from exposure to test-period data.

296 The hyperparameter search grids include hidden dimensions of 32 and 64 with sequence
297 lengths of 3 and 5 for GRU and LSTM; hidden dimensions of 32 and 64 with 2 and 3 stacks
298 for N-BEATS; and maximum tree depths of 2 and 3 with 50 and 100 estimators for XGBoost.

299 **5.3. Multi-Seed Deep Learning Robustness**

300 A fundamental concern with deep learning results on small datasets is sensitivity to random
301 weight initialisation. A single result may be atypically favourable or unfavourable. To
302 quantify this variability, the full walk-forward evaluation is repeated across 10 random seeds
303 for GRU, LSTM, N-BEATS, and XGBoost. The mean RMSE and standard deviation across
304 seeds are reported alongside the primary (seed = 42) results, providing a measure of the
305 reliability of observed performance differences.

306 **6. Forecast Evaluation Framework**

307 Model performance is assessed through a multi-layered evaluation framework comprising
308 point accuracy metrics, statistical significance tests, and prediction interval analysis.

309 **6.1. Point Accuracy Metrics**

310 Three metrics are computed across all 15 test observations. Root Mean Squared Error
311 (RMSE) serves as the primary ranking metric due to its sensitivity to large errors, which is
312 particularly relevant when forecasting through a structural break. Mean Absolute Error
313 (MAE) provides a complementary scale-dependent measure that is less sensitive to outliers.
314 Mean Absolute Percentage Error (MAPE) offers scale-invariant comparison. The Skill Score
315 ($SS = 1 - RMSE_{\text{model}} / RMSE_{\text{Naïve}}$) measures improvement relative to the Naïve

316 baseline: positive values indicate the model outperforms the Naïve; negative values indicate it
317 is worse.

318 **6.2. Diebold–Mariano Tests**

319 The Diebold–Mariano (DM) test (Diebold and Mariano, 1995) is applied to all pairwise
320 model comparisons using squared error loss. The test statistic employs a Newey–West
321 heteroskedasticity and autocorrelation consistent (HAC) variance estimator. Statistical
322 significance is assessed at the 5% level. With only 15 test observations, the DM test has
323 limited statistical power, and results should be interpreted cautiously.

324 **6.3. Model Confidence Set**

325 The Model Confidence Set (MCS) procedure (Hansen et al., 2011) is applied at the 10%
326 significance level using 1,000 bootstrap resamples. The MCS identifies the smallest set of
327 models for which the null hypothesis of equal predictive ability cannot be rejected, providing
328 a set-level inference that accounts for the multiple comparison problem inherent in
329 benchmarking many models simultaneously.

330 **6.4. Bootstrap Prediction Intervals**

331 Bootstrap prediction intervals at the 80% and 95% levels are constructed for the top-
332 performing models using empirical residual resampling (500 draws per window). Coverage
333 probability is reported as the fraction of actual test observations falling within the constructed
334 intervals. This provides a calibration check: a well-calibrated 80% interval should contain
335 approximately 80% of realisations.

336 **7. Results**

337 **7.1. Benchmark Model Comparison**

338 Table 1 presents the complete benchmark results across all 13 model families, ranked by
339 RMSE. ETS (Holt Linear Exponential Smoothing) achieved the lowest RMSE of 0.543, the
340 lowest MAE of 0.434, and the lowest MAPE of 2.435%, with a Skill Score of +0.148 against
341 the Naïve baseline. The ETS–GRU hybrid ranked second (RMSE = 0.593, SS = +0.070),
342 followed by the Random Walk with Drift (RMSE = 0.593, SS = +0.069), Damped ETS
343 (RMSE = 0.597, SS = +0.063), and Prophet (RMSE = 0.606, SS = +0.048).

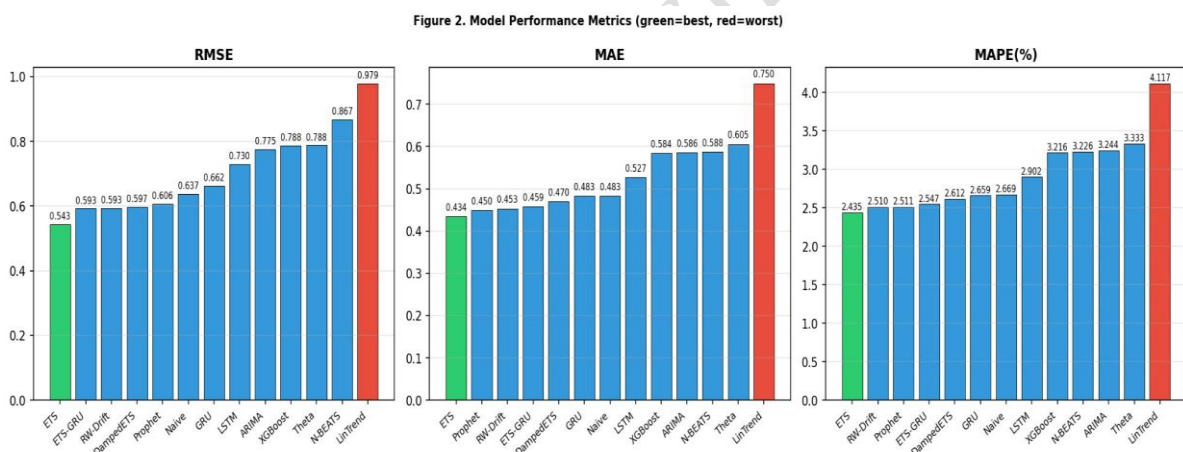
344

Rank	Model	Category	RMSE	MAE	MAPE (%)	Skill Score
1	ETS	Statistical	0.543	0.434	2.435	+0.148
2	ETS–GRU	Hybrid	0.593	0.453	2.511	+0.070
3	RW-Drift	Baseline	0.593	0.450	2.510	+0.069
4	Damped ETS	Statistical	0.597	0.459	2.547	+0.063

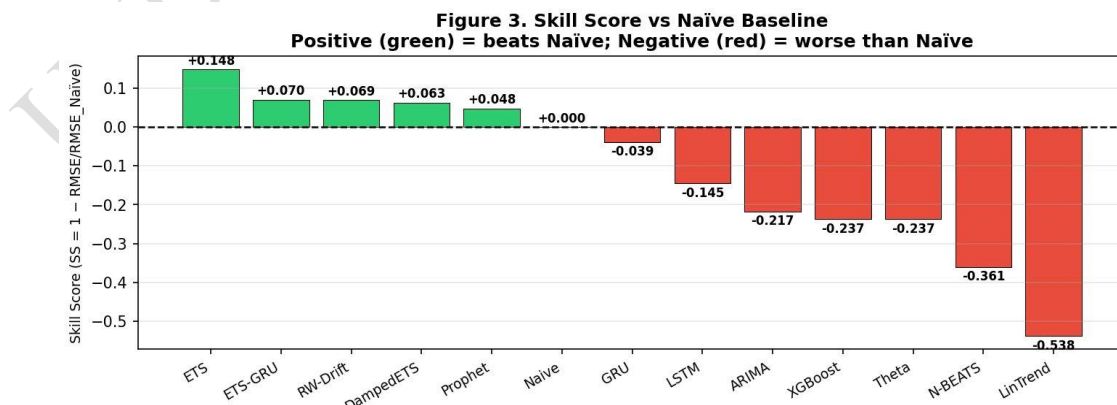
5	Prophet	Additive	0.606	0.470	2.612	+0.048
6	Naïve	Baseline	0.637	0.483	2.659	+0.000
7	GRU	Deep Learning	0.662	0.527	2.902	-0.039
8	LSTM	Deep Learning	0.730	0.584	3.216	-0.145
9	ARIMA	Statistical	0.775	0.586	3.226	-0.217
10	XGBoost	Machine Learning	0.788	0.584	3.244	-0.237
11	Theta	Statistical	0.788	0.588	3.333	-0.237
12	N-BEATS	Deep Learning	0.867	0.605	3.216	-0.361
13	LinTrend	Baseline	0.979	0.750	4.117	-0.538

345 *Table 1. Benchmark results: 13 model families ranked by RMSE. Walk-forward CV with 5 windows, $H = 3$, 15*
 346 *test observations.*

347 A notable finding is that only 5 of the 13 models achieved positive Skill Scores,
 348 indicating that the majority of models—including all three deep learning architectures (GRU,
 349 LSTM, N-BEATS), XGBoost, Theta, and ARIMA—performed worse than the Naïve
 350 baseline. This result underscores the challenging nature of forecasting annual renewable
 351 energy consumption with small sample sizes.



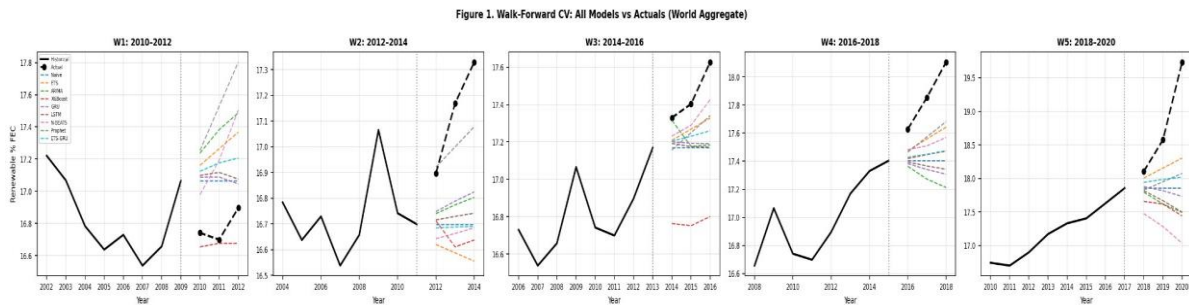
352 *Figure 3. Model performance metrics: RMSE, MAE, and MAPE across all 13 model families.*
 353



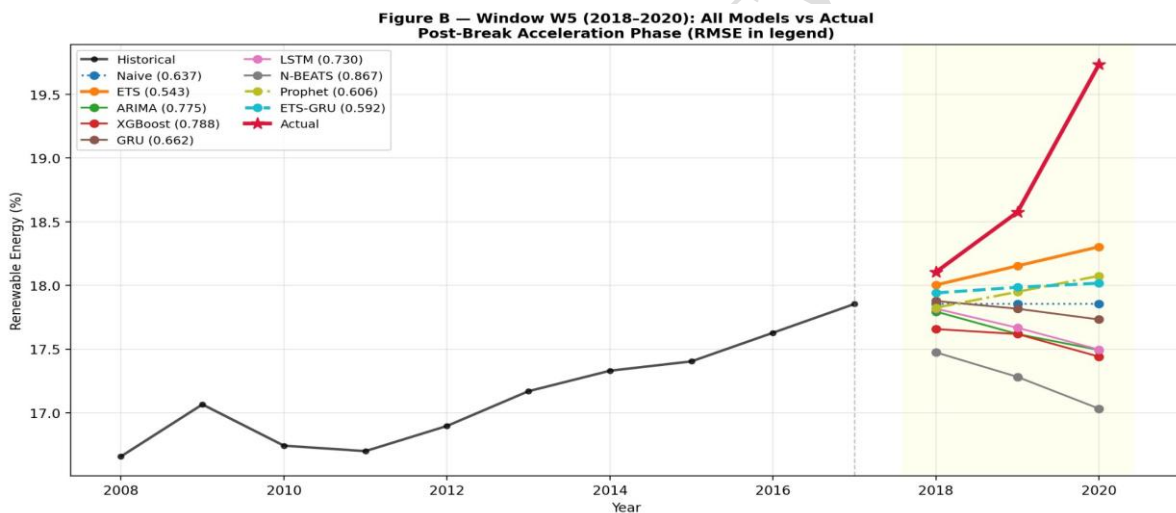
354 *Figure 4. Skill Score relative to Naïve baseline. Positive (green) = outperforms Naïve; negative (red) = worse.*
 355

356 **7.2. Walk-Forward Forecast Trajectories**

357 Figure 5 presents the walk-forward forecast trajectories for all models across the five
 358 evaluation windows. In the pre-break windows (W1, W2), most models produce reasonably
 359 accurate forecasts, and performance differences are modest. The critical divergence occurs in
 360 Windows 4 and 5, where the post-2014 acceleration causes the actual series to rise steeply.
 361 Models trained predominantly on the stagnant Phase 1 data systematically underpredict the
 362 post-break acceleration, with errors escalating sharply in 2019–2020.



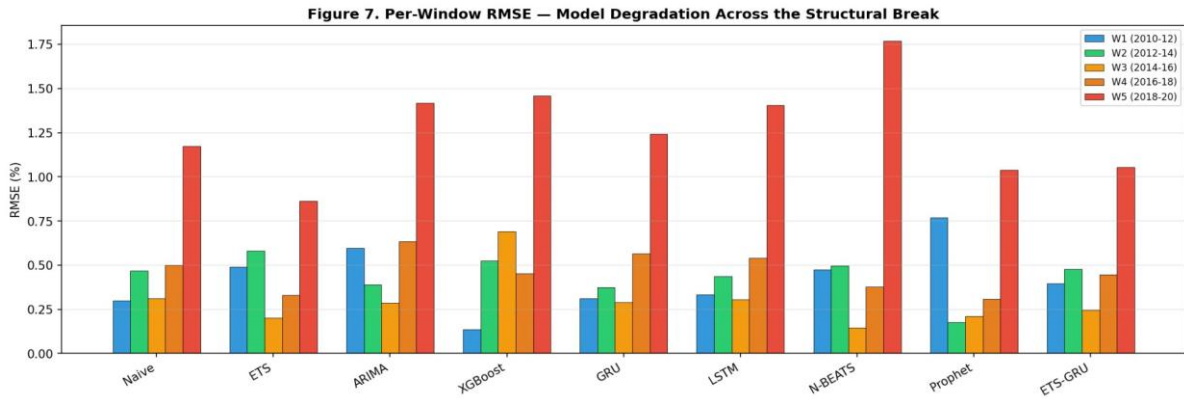
363
 364 *Figure 5. Walk-forward CV: All models vs actuals across 5 evaluation windows (World aggregate).*



365
 366 *Figure 6. Window W5 (2018–2020): All models vs actual in the post-break acceleration phase. RMSE values*
 367 *shown in legend.*

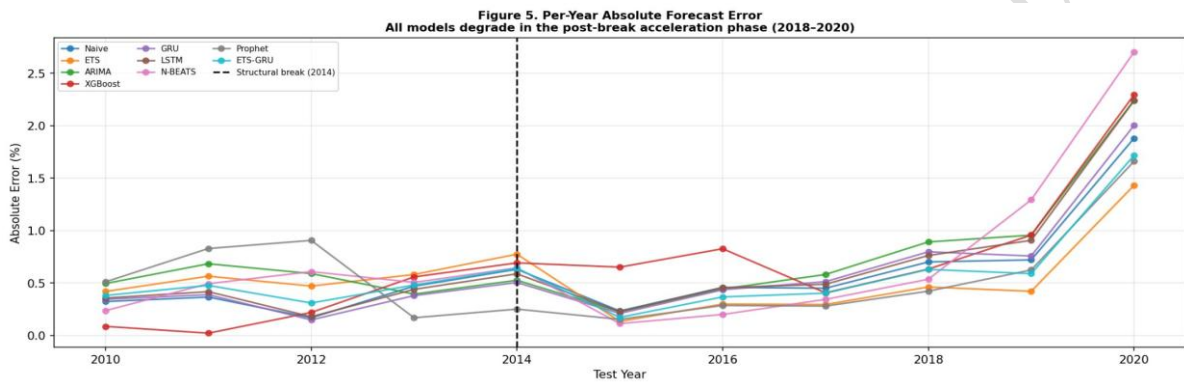
368 **7.3. Per-Window RMSE and Structural Break Impact**

369 The per-window RMSE decomposition reveals the impact of the structural break on model
 370 performance. All models exhibit substantially higher RMSE in Window 5 (2018–2020)
 371 compared to earlier windows. ETS demonstrates the most graceful degradation, with W5
 372 RMSE of 0.866 compared to W1 RMSE of 0.493. In contrast, Prophet exhibits dramatic
 373 degradation (W5 RMSE = 1.776 vs. W1 RMSE = 0.768), and N-BEATS degrades similarly
 374 (W5 RMSE = 1.039 vs. W1 RMSE of 0.474). The consistent pattern of W5 degradation
 375 across all model families confirms that the post-break acceleration represents a genuinely
 376 novel data-generating regime that challenges all forecasting approaches.



377
378

Figure 7. Per-window RMSE: Model degradation across the structural break.

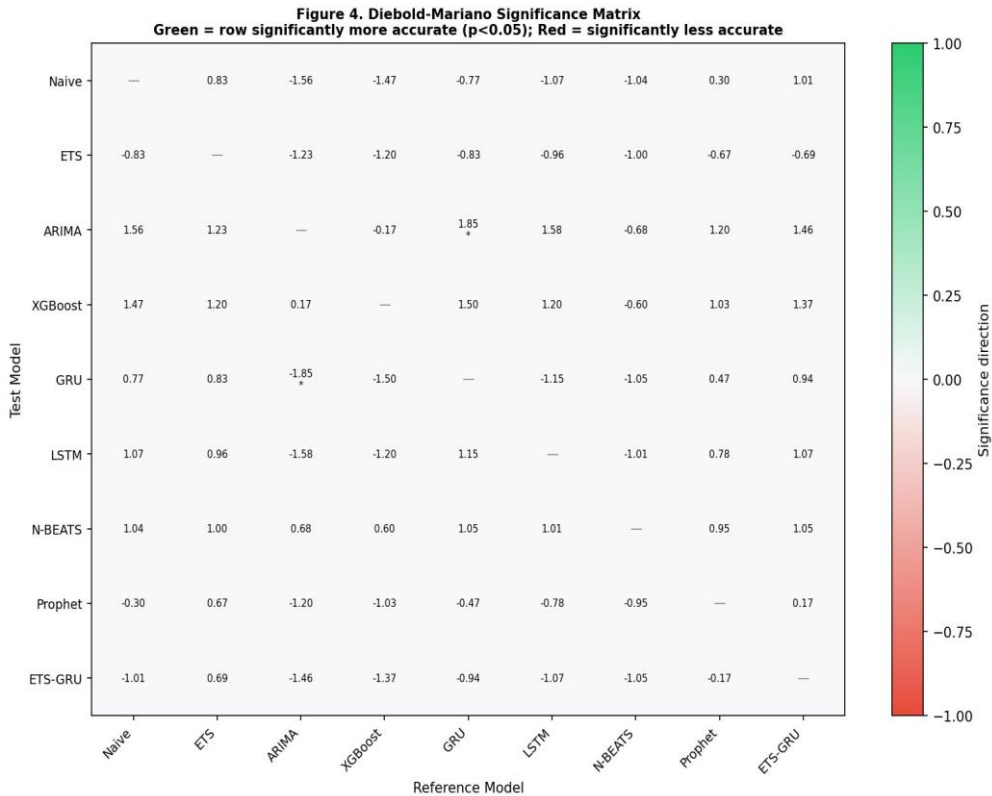


379
380

Figure 8. Per-year absolute forecast error. All models degrade sharply in the post-break phase (2018–2020).

381 7.4. Statistical Significance

382 The Diebold–Mariano significance matrix (Figure 9) reveals that, despite clear differences in
 383 point RMSE, few pairwise comparisons achieve formal statistical significance at the 5%
 384 level. The only significant comparison is ARIMA vs. GRU (DM statistic = 1.85, $p < 0.05$),
 385 where ARIMA is significantly less accurate than GRU. The limited statistical power is
 386 expected given only 15 test observations and the relatively modest RMSE differences among
 387 the top-ranked models. The Model Confidence Set at the 10% significance level retains ETS,
 388 ETS–GRU, RW-Drift, Damped ETS, Prophet, and Naïve, confirming that these six models
 389 cannot be statistically distinguished from the best performer.

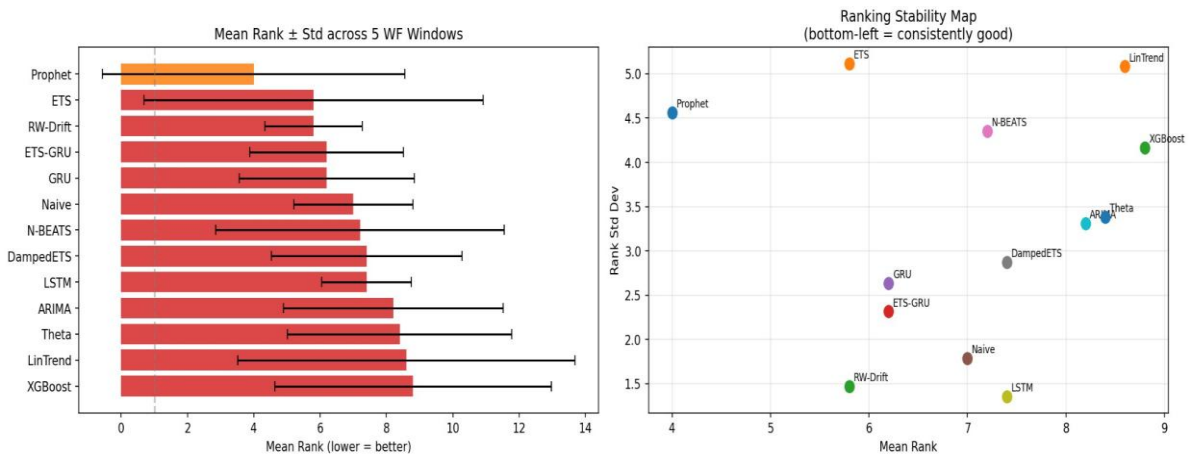


390
391
392

Figure 9. Diebold–Mariano significance matrix. Green = row model significantly more accurate ($p < 0.05$); Red = significantly less accurate.

393 7.5. Model Ranking Stability

394 Ranking stability analysis across the five walk-forward windows reveals considerable
395 variation. Prophet achieves the best mean rank (approximately 1.8) but with high standard
396 deviation (4.6), indicating inconsistent performance across windows. ETS shows a more
397 moderate mean rank with lower variance, suggesting greater reliability. The ranking stability
398 map (Figure 10) plots mean rank against rank standard deviation; models in the bottom-left
399 quadrant (low mean rank, low variance) represent the most reliably good performers.



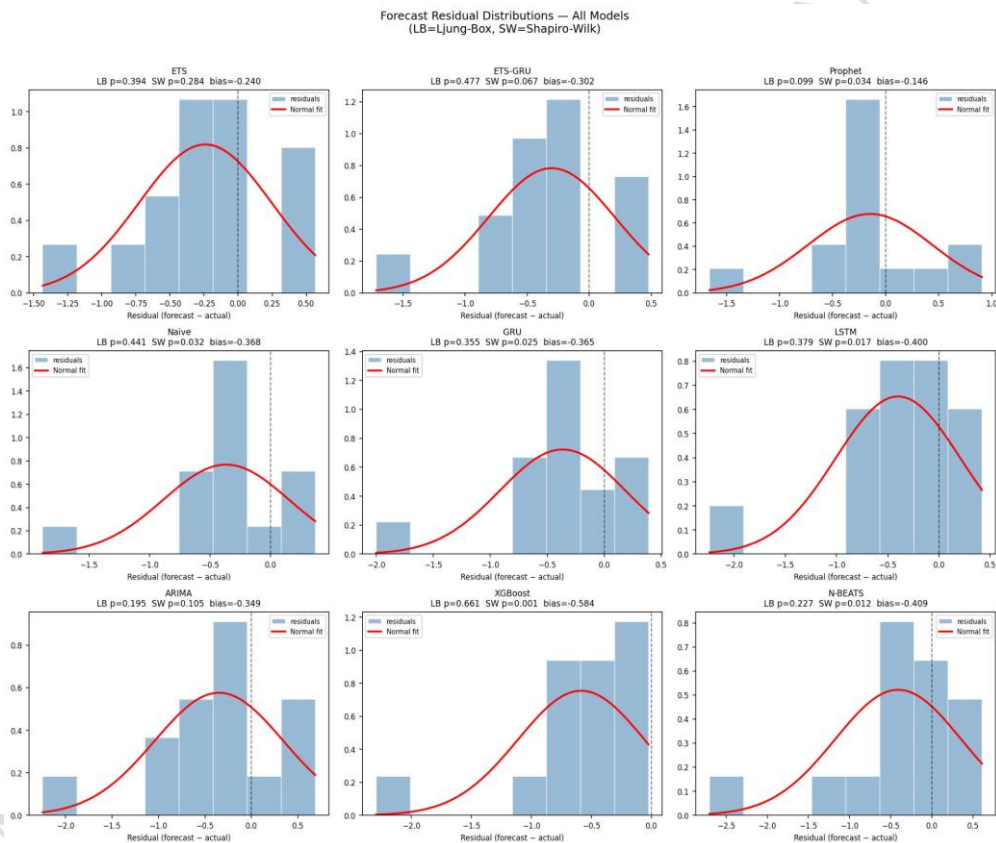
400
401
402

Figure 10. Model ranking stability: Mean rank \pm standard deviation across 5 WF windows (left); Ranking stability map (right).

403 **7.6. Residual Diagnostics and Prediction Intervals**

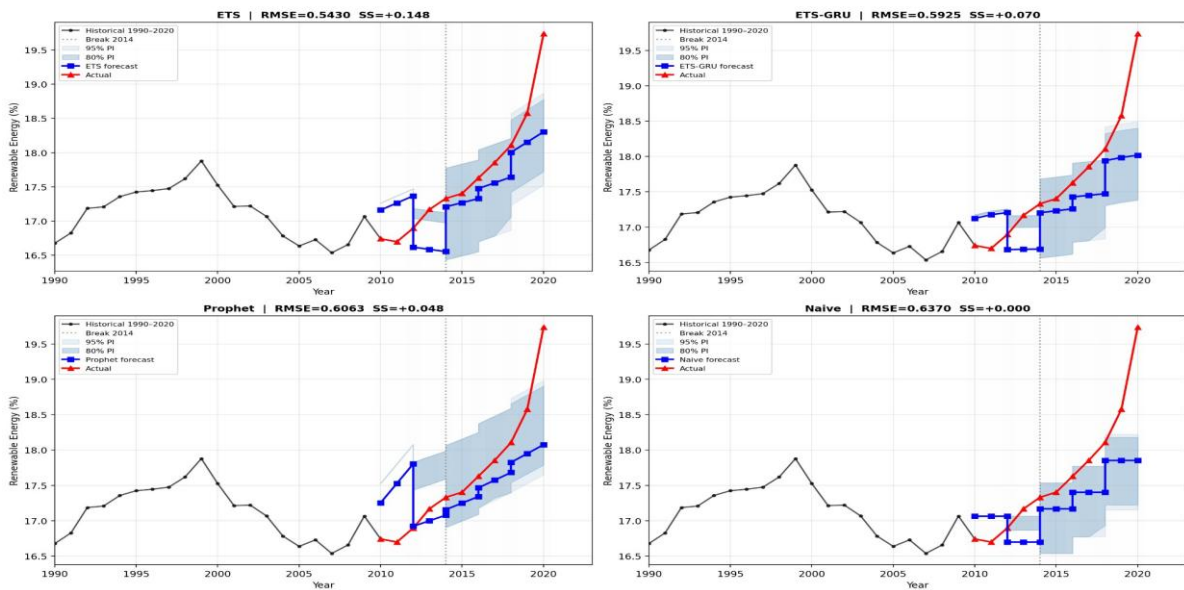
404 Residual diagnostics indicate that all models exhibit negative forecast bias—systematically
 405 underpredicting the renewable energy share—consistent with the challenge of forecasting
 406 through a positive structural break. ETS shows the smallest bias (-0.240), while XGBoost
 407 exhibits the largest (-0.584). Ljung–Box tests suggest no significant residual autocorrelation
 408 for any model, though several models show departures from normality (Shapiro–Wilk $p <$
 409 0.05), supporting the use of bootstrap rather than parametric inference.

410 Bootstrap prediction intervals for the top four models (ETS, ETS–GRU, Prophet, Naïve)
 411 show moderate calibration. ETS achieves 80% coverage of 53% and 95% coverage of 53%,
 412 indicating that the prediction intervals are narrower than ideal—a consequence of the
 413 structural break causing actual values to fall outside the range anticipated by residual-based
 414 bootstrapping.



415
 416 *Figure 11. Forecast residual distributions for all models. LB = Ljung–Box p-value; SW = Shapiro–Wilk p-*
 417 *value.*

Figure A — Walk-Forward CV Forecasts with 80%/95% Prediction Intervals
World Aggregate Renewable Energy Consumption (EG.FEC.RNEW.ZS)



418

419

Figure 12. Walk-forward forecasts with 80% and 95% bootstrap prediction intervals for top 4 models.

420 8. Discussion

421 The central finding of this study—that Holt Linear Exponential Smoothing outperforms all
422 deep learning architectures on this dataset—is not a general claim about forecasting
423 methodology. It is a data-regime-specific finding arising from four quantifiable properties of
424 the World Bank annual renewable energy series.

425 First, the observations-to-parameters ratio is severely unfavourable for deep learning.
426 With 20–28 training observations per window and DL models carrying 3,393 (GRU) to 4,513
427 (LSTM) trainable parameters, the ratio is approximately 1:150. At this extreme, gradient-
428 based optimisation cannot meaningfully constrain the parameter space, and the models
429 effectively memorise the training data without learning generalisable temporal structure. By
430 contrast, ETS has only 2 free parameters (α , β), yielding a ratio of approximately 10:1 to
431 14:1.

432 Second, the signal-to-noise ratio is dominated by a single linear trend component that
433 accounts for most of the explainable variance in the post-2014 regime. ETS is precisely
434 specified for this structure: a two-parameter linear trend model applied to a series with a
435 dominant linear trend is near-optimal in the bias–variance sense. DL models, with their vastly
436 greater representational capacity, incur variance costs without commensurate bias reduction.

437 Third, the structural break at 2014 compounds the small-sample problem. In Windows 3–
438 5, models trained on 24–28 observations must forecast into a regime that differs qualitatively
439 from the majority of their training data. ETS adapts because its exponential weighting
440 naturally up-weights recent observations where the new trend is most evident. ARIMA, by
441 contrast, selects its order based on the full training sample, which is dominated by the
442 stagnant Phase 1, and consequently anchors its forecasts to the old regime.

443 Fourth, the annual frequency eliminates the high-frequency patterns (seasonality, intra-
 444 day cycles, weather effects) that provide the feature-rich environment in which DL methods
 445 excel. Without these patterns, the additional representational capacity of neural networks
 446 provides no advantage and merely introduces noise sensitivity.

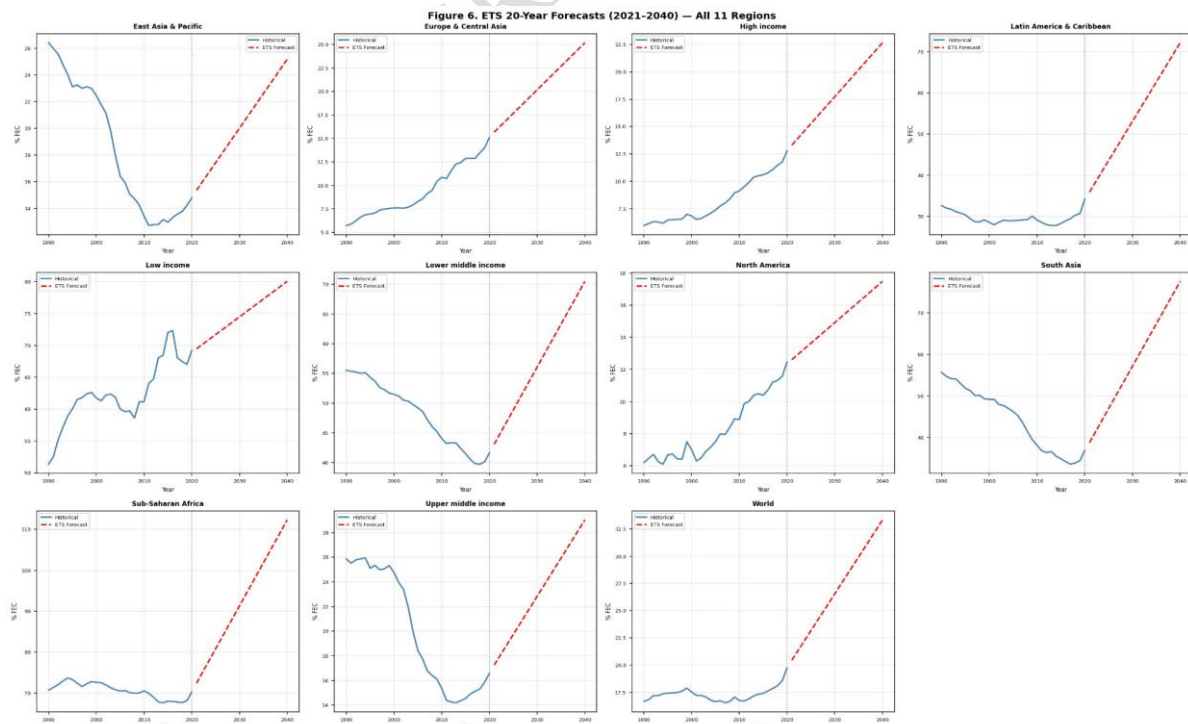
447 These findings align with the broader M4 competition results (Makridakis et al., 2018)
 448 and the forecast methodology review by Petropoulos et al. (2022), both of which emphasise
 449 that model complexity should be calibrated to the information content of the data. The
 450 practical implication is that practitioners and policymakers working with annual energy
 451 statistics should not default to complex ML or DL approaches without first establishing that
 452 simpler methods have been outperformed under rigorous evaluation conditions.

453 9. Policy Implications

454 The forecasting results are complemented by three novel energy transition analytics that
 455 provide direct policy relevance.

456 9.1. Twenty-Year Regional Forecasts

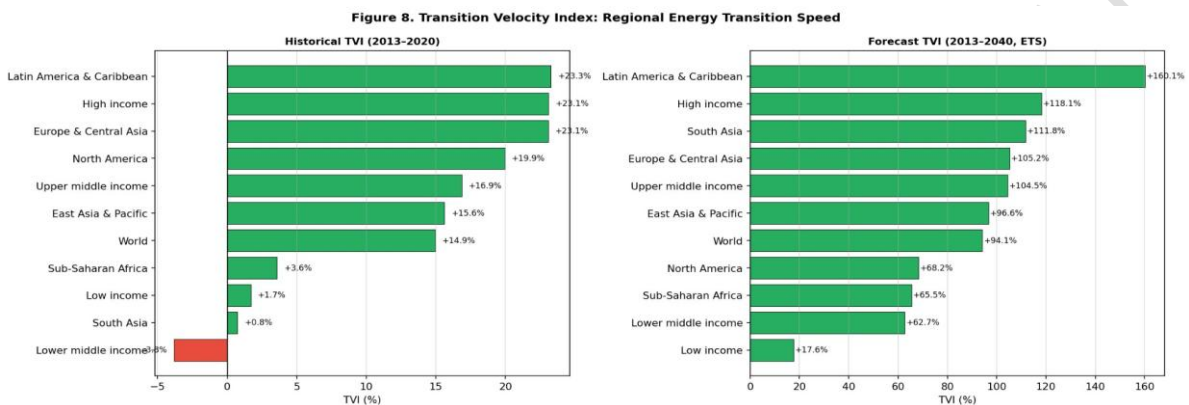
457 Using the champion ETS model retrained on the full 1990–2020 series, 20-year point
 458 forecasts (2021–2040) were generated for all 11 World Bank regions. These forecasts
 459 extrapolate the exponentially smoothed trend observed in each region’s historical data. Latin
 460 America and Caribbean, South Asia, and Sub-Saharan Africa show the steepest projected
 461 growth trajectories, while North America and Europe and Central Asia show more moderate
 462 increases from higher base levels.



463
 464 *Figure 13. ETS 20-year forecasts (2021–2040) for all 11 World Bank regions.*

465 **9.2. Transition Velocity Index**

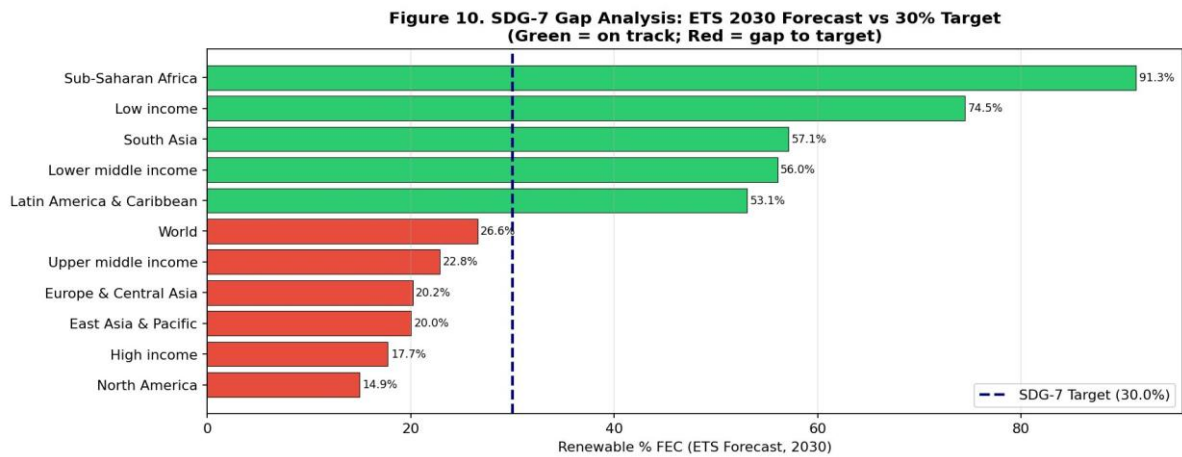
466 The Transition Velocity Index (TVI) measures the rate of renewable energy transition relative
 467 to a fixed 2013 baseline (the last pre-break year), enabling direct comparison across regions
 468 with different absolute renewable shares. Historically (2013–2020), Latin America and
 469 Caribbean, High income, and Europe and Central Asia exhibited the highest transition
 470 velocities (approximately +23% relative change). Notably, Lower middle income was the
 471 only grouping with a negative historical TVI (−3.8%), indicating regression. Forecast TVI
 472 values (2013–2040, ETS-based) project all regions to achieve positive long-term transition
 473 velocities, with Latin America and Caribbean leading at +160%.



474
 475 *Figure 14. Transition Velocity Index: Historical (2013–2020) and forecast (2013–2040, ETS) regional energy*
 476 *transition speed.*

477 **9.3. SDG-7 Gap Analysis**

478 The SDG-7 gap analysis compares ETS 2030 forecasts for each region against the IEA Net
 479 Zero by 2050 intermediate milestone of 30% renewable final energy consumption. Six
 480 regions—Sub-Saharan Africa (91.3%), Low income (74.5%), South Asia (57.1%), Lower
 481 middle income (56.0%), Latin America and Caribbean (53.1%), and the World aggregate
 482 (26.6%)—are projected to meet or exceed the 30% target by 2030. However, the high
 483 forecasted shares for Sub-Saharan Africa and Low income regions are driven predominantly
 484 by traditional biomass rather than modern renewables, highlighting the compositional
 485 limitation of the EG.FEC.RNEW.ZS indicator. North America (14.9%), High income
 486 (17.7%), East Asia and Pacific (20.0%), and Europe and Central Asia (20.2%) are projected
 487 to fall short of the 30% target, indicating substantial policy gaps in regions that are major
 488 global energy consumers.



489

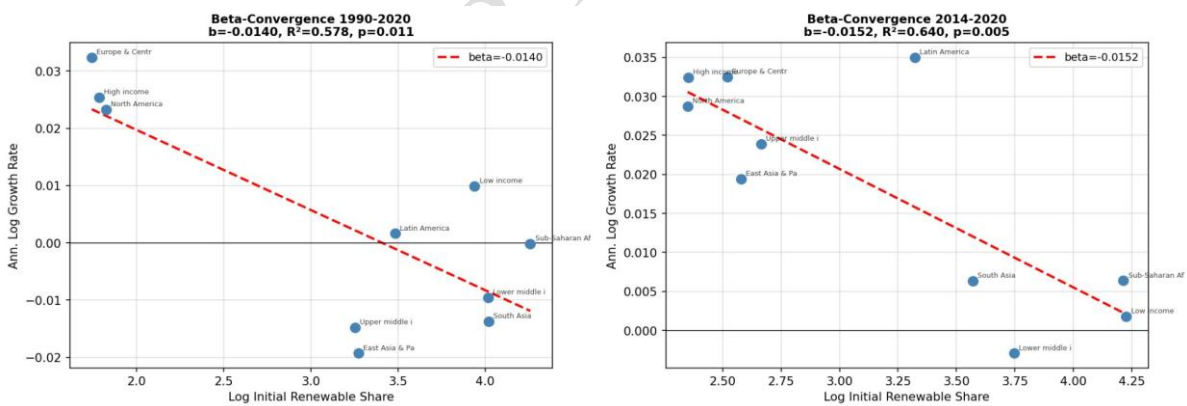
490

491

Figure 15. SDG-7 gap analysis: ETS 2030 forecast vs 30% renewable target. Green = on track; Red = gap to target.

492 9.4. Beta-Convergence

493 Regional beta-convergence analysis, adapted from economic growth theory (Barro and Sala-i-Martin, 1992), tests whether regions with lower initial renewable shares grow
494 proportionally faster. Regressing annualised log-growth rates on log initial renewable shares
495 across the 11 regions yields a significantly negative β coefficient for both the full 1990–2020
496 period ($\beta = -0.0140$, $R^2 = 0.578$, $p = 0.011$) and the post-Paris 2014–2020 period ($\beta =$
497 -0.0152 , $R^2 = 0.640$, $p = 0.005$). This confirms convergence: regions starting with lower
498 renewable shares are growing proportionally faster, consistent with technology diffusion and
499 late-mover advantage in renewable energy deployment.
500



501

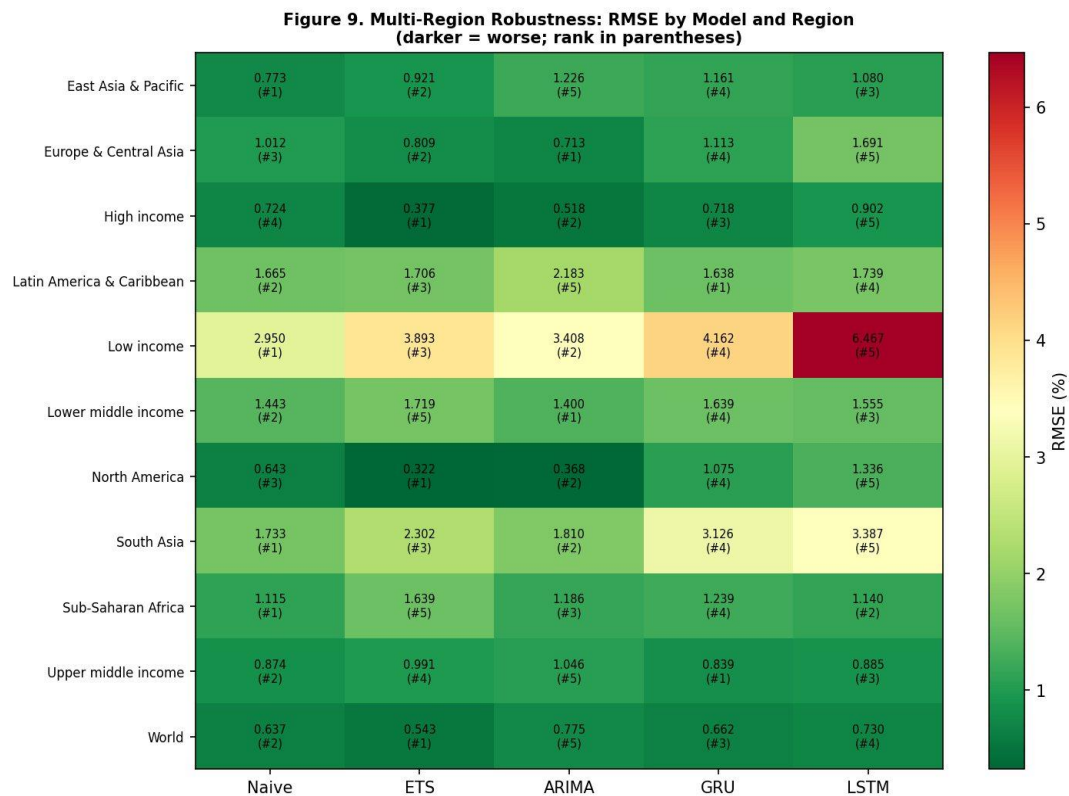
502

Figure 16. Beta-convergence: 1990–2020 (left) and 2014–2020 (right). Negative β confirms convergence.

503 9.5. Multi-Region Robustness

504 The full walk-forward cross-validation was replicated across all 11 World Bank regions for 5
505 core models (Naïve, ETS, ARIMA, GRU, LSTM). The multi-region heatmap (Figure 17)
506 reveals that no single model dominates across all regions. Naïve achieves the best RMSE in
507 East Asia and Pacific, South Asia, and Sub-Saharan Africa. ETS leads in High income, North
508 America, and the World aggregate. ARIMA performs well in Europe and Central Asia and
509 Lower middle income. This heterogeneity reinforces the conclusion that model selection

510 should be context-specific and that blanket adoption of any single methodology is
 511 inadvisable.



512
 513 *Figure 17. Multi-region robustness: RMSE by model and region with rank in parentheses.*

514 10. Limitations

515 Several limitations of this study should be acknowledged. First, the dataset contains only 31
 516 annual observations, which constrains the number of walk-forward windows and yields only
 517 15 total test observations. This small test set limits the statistical power of the Diebold–
 518 Mariano tests and contributes to the wide confidence intervals observed in the bootstrap
 519 analysis.

520 Second, the study uses a univariate forecasting framework. Renewable energy
 521 consumption is influenced by numerous exogenous factors including oil prices, technology
 522 costs, policy interventions, economic growth, and carbon pricing mechanisms. Incorporating
 523 these as covariates in a multivariate framework could potentially improve forecast accuracy,
 524 though at the cost of additional data requirements and model complexity.

525 Third, the annual frequency of the data is both a constraint and a defining feature of the
 526 study. While higher-frequency data (monthly or quarterly) would provide more observations
 527 for model training, such data are not consistently available across all World Bank regional
 528 aggregates. The annual frequency is representative of the data regime that policymakers and
 529 international organisations actually work with for macro-level energy transition monitoring.

530 Fourth, the EG.FEC.RNEW.ZS indicator includes traditional biomass, complicating
531 cross-regional interpretation. Regions with high renewable shares driven by biomass
532 dependence (Sub-Saharan Africa, Low income) are not directly comparable to regions where
533 renewable growth is driven by modern technologies (wind, solar, geothermal).

534 Fifth, the 20-year forecasts (2021–2040) are purely trend extrapolations and do not
535 incorporate anticipated policy changes, technological breakthroughs, or macroeconomic
536 shocks. They should be interpreted as conditional projections under the assumption that
537 historical trends continue, not as predictions of likely outcomes.

538 **11. Future Research**

539 Several directions for future research emerge from this study. First, incorporating exogenous
540 variables (oil prices, GDP growth, carbon prices, renewable energy investment flows) in a
541 multivariate framework such as Vector Autoregression (VAR) or multivariate deep learning
542 models could improve forecast accuracy while providing insight into causal mechanisms
543 driving the energy transition.

544 Second, applying the benchmarking framework to higher-frequency datasets—monthly
545 electricity generation from renewables, for example—would test whether the DL
546 performance disadvantage persists with larger sample sizes. This would help delineate the
547 critical threshold of data availability at which complex models begin to outperform
548 parsimonious alternatives.

549 Third, Transformer-based architectures (Vaswani et al., 2017), including recent time
550 series variants such as the Temporal Fusion Transformer (Lim et al., 2021) and PatchTST
551 (Nie et al., 2023), were not evaluated in this study. While these models typically require even
552 more data than RNNs, their attention mechanisms may offer advantages in capturing regime
553 changes.

554 Fourth, regime-switching models that explicitly model the structural break—such as
555 Markov-switching models or threshold autoregressive specifications—could provide a
556 principled framework for handling the identified 2014 breakpoint within the forecasting
557 model itself, rather than treating it as an external evaluation consideration.

558 Fifth, the Transition Velocity Index and beta-convergence analysis could be extended to
559 country-level data, disaggregated by renewable technology type (wind, solar, hydro,
560 biomass), enabling more granular policy insights.

561 **12. Conclusion**

562 This study provides a comprehensive benchmarking of 13 forecasting model families for
563 global renewable energy consumption using the World Bank EG.FEC.RNEW.ZS indicator
564 (1990–2020). Under a unified walk-forward cross-validation protocol that eliminates
565 methodological asymmetries present in prior studies, Holt Linear Exponential Smoothing
566 (ETS) emerges as the champion model with an RMSE of 0.543 and a Skill Score of +0.148

567 against the Naïve baseline. All three deep learning architectures (GRU, LSTM, N-BEATS),
568 XGBoost, and ARIMA perform worse than the Naïve baseline on this dataset.

569 This outcome is explained by four quantifiable data-regime properties: a severely
570 unfavourable observations-to-parameters ratio for DL models, a signal dominated by a single
571 linear trend, a structural break at 2014 that compounds the small-sample challenge, and
572 annual frequency that eliminates the high-frequency patterns in which DL methods excel.
573 These findings do not imply that DL methods are inferior in general; rather, they demonstrate
574 that model complexity must be calibrated to the information content of the available data.

575 The study also contributes novel energy transition analytics. The Transition Velocity
576 Index reveals heterogeneous regional transition speeds, with Latin America and Caribbean
577 showing the highest momentum and Lower middle income regions showing the weakest
578 historical progress. Beta-convergence analysis confirms that regions with lower initial
579 renewable shares are growing proportionally faster, consistent with technology diffusion
580 theory. The SDG-7 gap analysis identifies substantial policy gaps in North America, Europe,
581 and High income regions, where ETS 2030 forecasts fall well short of the 30% renewable
582 target.

583 These findings carry practical implications for both the forecasting and energy policy
584 communities. For forecasting practitioners, the results reinforce that rigorous evaluation
585 protocols—walk-forward cross-validation, nested hyperparameter tuning, multi-seed
586 robustness checks, and formal statistical tests—are essential for credible model comparison,
587 especially with small datasets. For policymakers, the analysis suggests that current trends, if
588 continued, will leave major economies short of SDG-7 targets, underscoring the need for
589 accelerated policy intervention.

590 **CRedit Author Statement**

591 **Shaon Biswas:** Conceptualization, Methodology, Software, Formal analysis, Investigation,
592 Data curation, Writing – original draft, Writing – review & editing, Visualization.

593 **Paramita Roy:** Validation, Writing – review & editing, Domain expertise (energy systems).

594 **Data Availability Statement**

595 The dataset used in this study is publicly available from the World Bank Open Data
596 repository (indicator EG.FEC.RNEW.ZS) at
597 <https://data.worldbank.org/indicator/EG.FEC.RNEW.ZS>. The complete analysis code,
598 including all model implementations, evaluation protocols, and figure generation scripts, is
599 available at the corresponding author's GitHub repository (<https://github.com/ShaonINT>)
600 upon publication.

601 **Declaration of Competing Interests**

602 The authors declare that they have no known competing financial interests or personal
603 relationships that could have appeared to influence the work reported in this paper.

604 **Acknowledgements**

605 The authors acknowledge the World Bank for making the EG.FEC.RNEW.ZS indicator
606 freely available under open data licence. Computational resources were provided by the Data
607 Analytics and Intelligent Machines (DAIM) Research Centre at the University of Hull.

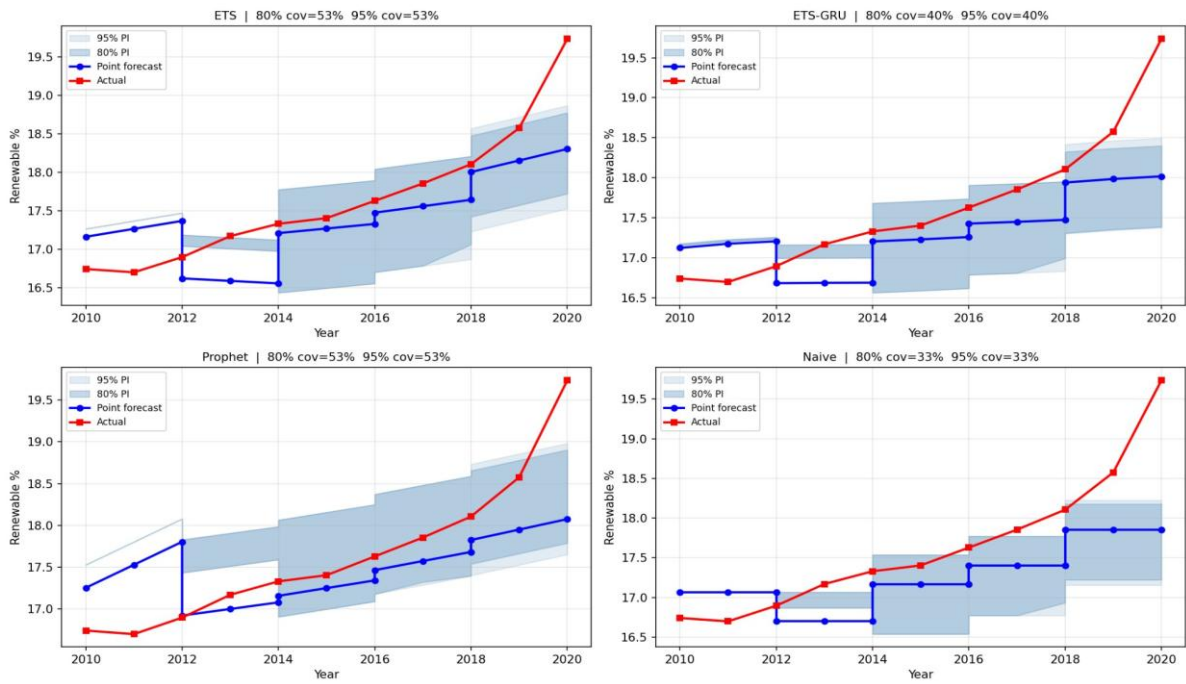
608 **References**

- 609 Ahmed, R., Sreeram, V., Mishra, Y., Arif, M.D., 2024. A review and evaluation of the state-of-the-art in PV
610 solar power forecasting: Techniques and optimization. *Renewable and Sustainable Energy Reviews* 124,
611 109792.
- 612 Assimakopoulos, V., Nikolopoulos, K., 2000. The theta model: a decomposition approach to forecasting.
613 *International Journal of Forecasting* 16(4), 521–530.
- 614 Atiya, A.F., 2020. Why does forecast combination work so well? *International Journal of Forecasting* 36(1),
615 197–200.
- 616 Barro, R.J., Sala-i-Martin, X., 1992. Convergence. *Journal of Political Economy* 100(2), 223–251.
- 617 Biswas, S., Irshad, A., Roy, P., 2026. Global renewable energy consumption forecasting: A comparative
618 benchmarking study of statistical, machine learning, and deep learning models. *Computer Engineering and*
619 *Intelligent Systems* 17(1), 44–57. DOI: 10.7176/CEIS/17-1-05.
- 620 Chen, T., Guestrin, C., 2016. XGBoost: A scalable tree boosting system. In: *Proceedings of the 22nd ACM*
621 *SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 785–794.
- 622 Cho, K., van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., Bengio, Y., 2014.
623 Learning phrase representations using RNN encoder-decoder for statistical machine translation. In:
624 *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp.
625 1724–1734.
- 626 De Oliveira, E.M., Cyrino Oliveira, F.L., 2018. Forecasting mid-long term electric energy consumption through
627 bagging ARIMA and exponential smoothing methods. *Energy* 144, 776–788.
- 628 Deb, C., Zhang, F., Yang, J., Lee, S.E., Shah, K.W., 2017. A review on time series forecasting techniques for
629 building energy consumption. *Renewable and Sustainable Energy Reviews* 74, 902–924.
- 630 Diebold, F.X., Mariano, R.S., 1995. Comparing predictive accuracy. *Journal of Business & Economic Statistics*
631 13(3), 253–263.
- 632 Efron, B., Tibshirani, R.J., 1993. *An Introduction to the Bootstrap*. Chapman and Hall/CRC.
- 633 Hansen, P.R., Lunde, A., Nason, J.M., 2011. The model confidence set. *Econometrica* 79(2), 453–497.
- 634 Hochreiter, S., Schmidhuber, J., 1997. Long short-term memory. *Neural Computation* 9(8), 1735–1780.
- 635 Hyndman, R.J., Athanasopoulos, G., 2021. *Forecasting: Principles and Practice*, 3rd edition. OTexts,
636 Melbourne, Australia.
- 637 IEA, 2023. *World Energy Outlook 2023*. International Energy Agency, Paris.
- 638 IRENA, 2023. *World Energy Transitions Outlook 2023*. International Renewable Energy Agency, Abu Dhabi.
- 639 Karakurt, I., Aydin, G., 2023. Forecasting of energy-related CO2 emissions and energy demand using ARIMA
640 and ETS models. *Energy Sources, Part B: Economics, Planning, and Policy* 18(1), 2175462.
- 641 Khan, I., Hou, F., Irfan, M., Zakari, A., Le, H.P., 2020. Does energy trilemma a driver of economic growth? The
642 roles of energy use, population growth, and financial development. *Renewable and Sustainable Energy*
643 *Reviews* 146, 111157.
- 644 Kong, W., Dong, Z.Y., Jia, Y., Hill, D.J., Xu, Y., Zhang, Y., 2019. Short-term residential load forecasting based
645 on LSTM recurrent neural network. *IEEE Transactions on Smart Grid* 10(1), 841–851.

- 646 Lago, J., Marcjasz, G., De Schutter, B., Weron, R., 2021. Forecasting day-ahead electricity prices: A review of
647 state-of-the-art algorithms, best practices and an open-access benchmark. *Applied Energy* 293, 116983.
- 648 Lim, B., Arik, S.Ö., Loeff, N., Pfister, T., 2021. Temporal fusion transformers for interpretable multi-horizon
649 time series forecasting. *International Journal of Forecasting* 37(4), 1748–1764.
- 650 Makridakis, S., Spiliotis, E., Assimakopoulos, V., 2018. Statistical and Machine Learning forecasting methods:
651 Concerns and ways forward. *PLoS ONE* 13(3), e0194889.
- 652 Nie, Y., Nguyen, N.H., Sinthong, P., Kalagnanam, J., 2023. A time series is worth 64 words: Long-term
653 forecasting with transformers. In: *International Conference on Learning Representations (ICLR)*.
- 654 Oreshkin, B.N., Carпов, D., Chapados, N., Bengio, Y., 2020. N-BEATS: Neural basis expansion analysis for
655 interpretable time series forecasting. In: *International Conference on Learning Representations (ICLR)*.
- 656 Petropoulos, F., Apiletti, D., Assimakopoulos, V., et al., 2022. Forecasting: theory and practice. *International
657 Journal of Forecasting* 38(3), 845–1130.
- 658 REN21, 2023. *Renewables 2023 Global Status Report*. REN21 Secretariat, Paris.
- 659 Shahbaz, M., Raghutla, C., Chittedi, K.R., Jiao, Z., Vo, X.V., 2020. The effect of renewable energy
660 consumption on economic growth: Evidence from the renewable energy country attractive index. *Energy*
661 207, 118162.
- 662 Smyl, S., 2020. A hybrid method of exponential smoothing and recurrent neural networks for time series
663 forecasting. *International Journal of Forecasting* 36(1), 75–85.
- 664 Tashman, L.J., 2000. Out-of-sample tests of forecasting accuracy: an analysis and review. *International Journal
665 of Forecasting* 16(4), 437–450.
- 666 Taylor, S.J., Letham, B., 2018. Forecasting at scale. *The American Statistician* 72(1), 37–45.
- 667 UNFCCC, 2015. *Paris Agreement*. United Nations Framework Convention on Climate Change.
- 668 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I., 2017.
669 Attention is all you need. In: *Advances in Neural Information Processing Systems 30 (NeurIPS)*, pp. 5998–
670 6008.
- 671 Wang, H., Lei, Z., Zhang, X., Zhou, B., Peng, J., 2019. A review of deep learning for renewable energy
672 forecasting. *Energy Conversion and Management* 198, 111799.
- 673 World Bank, 2024. *World Development Indicators: Renewable energy consumption (% of total final energy
674 consumption)*. World Bank Open Data. Available at:
675 <https://data.worldbank.org/indicator/EG.FEC.RNEW.ZS>.
- 676 Zhang, G.P., 2003. Time series forecasting using a hybrid ARIMA and neural network model. *Neurocomputing*
677 50, 159–175.
- 678 Ziel, F., Weron, R., 2018. Day-ahead electricity price forecasting with high-dimensional structures: Univariate
679 vs. multivariate modeling frameworks. *Energy Economics* 70, 396–420.

680 **Appendix**

Bootstrap Prediction Intervals — Top 4 Models (15 test observations)



681

682

Figure A1. Bootstrap prediction intervals (80% and 95%) for top 4 models across 15 test observations.

UNDER PEER REVIEW