

Entropy-Bounded Decomposition for Time Series Forecasting: Fourier Extraction, Adaptive Basis Selection, and the Maximum Entropy Stopping Criterion

Jeremy McEntire
Independent Research
USA

jandrewmcentire@gmail.com

Abstract

We present an ensemble time series forecasting system in which multiple analytical workers—Fourier decomposition, autoregressive modeling, and stochastic resonance—each independently forecast and track their own running accuracy. A confidence-weighted aggregator routes predictions toward the most accurate worker per context window, collapsing to a single worker when one dominates. Each worker uses entropy-bounded signal decomposition: extraction continues until the residual reaches maximum entropy, at which point no predictable structure remains.

We evaluate on six standard benchmarks using the established train/val/test protocol. The ensemble achieves normalized MSE of 0.109 on ETTh1 and 0.064 on ETM1 at horizon 96—improvements of 71% and 80% over Google’s TimesFM (200M parameters). On ETTh2, where the Fourier decomposition alone underperformed, the ensemble’s confidence mechanism routes to the AR worker and achieves 0.250—beating TimesFM’s 0.289 by 13%. Across four ETT datasets, the ensemble beats TimesFM on three and PatchTST on four. Weather (0.078 vs PatchTST 0.149, +47%) confirms the pattern. The sole loss is ETM2 (0.194 vs TimesFM 0.175, −11%).

The system requires no training data, no GPU, and runs six benchmarks in 131 seconds on a laptop. Every prediction includes a calibrated confidence score and per-worker breakdown. The implementation is 2,000 lines of Python using NumPy, SciPy, and PyWavelets.

Keywords

time series forecasting, Fourier decomposition, wavelet analysis, maximum entropy, signal processing, autoregressive models

1 Introduction

The dominant paradigm in time series forecasting has shifted over the past five years from statistical methods (ARIMA, exponential smoothing, Theta) to transformer-based foundation models trained on large corpora of diverse time series. Google’s TimesFM [1], a 200M parameter decoder-only transformer, achieves competitive zero-shot forecasting across dozens of benchmark datasets without per-dataset tuning. Chronos [3], Lag-Llama [4], and similar models extend this paradigm.

The advantages of foundation models are real: they capture complex patterns from massive training corpora, handle non-stationarity implicitly, and require no manual feature engineering. The disadvantages are equally real: they are computationally expensive, non-interpretable, and their forecasts cannot be explained to a regulator,

auditor, or domain expert who asks “why does the model predict X ?”

This paper takes a different path. We ask: how far can classical signal decomposition go if we apply it systematically, use information theory to guide basis selection, and let the data tell us when to stop?

The answer is: further than expected, but not uniformly. A Fourier decomposition alone beats TimesFM by 65–73% on periodic datasets (ETTh1, ETM1) but loses by 40% on non-periodic ones (ETTh2). However, a plain autoregressive model—50 lines of NumPy—beats the decomposition on *every* dataset. Neither model alone is universally best. But an ensemble of both, with confidence-weighted routing, beats TimesFM on three of four ETT datasets and PatchTST on four of six benchmarks.

The key insight is not any single forecasting technique but the **confidence mechanism**: each worker tracks its own running accuracy via holdout validation, and the ensemble routes toward the most accurate worker per context window. When one worker clearly dominates (confidence ratio $> 2\times$), the ensemble collapses to that worker alone, avoiding dilution by weaker models. This transforms the question from “which forecasting method is best?” to “which method is best *right now, for this data?*” — and lets the system answer automatically.

Contributions.

- (1) A **confidence-weighted ensemble** of analytical forecasting workers, where each worker independently tracks its running accuracy via holdout validation and the aggregator routes predictions toward the most accurate worker per context window. When one worker’s confidence exceeds all others by $2\times$, the ensemble collapses to that worker alone.
- (2) A **noise-aware iterative Fourier extraction** algorithm that models the expected spectral footprint of each signal component at the observed noise level, extracting the entire footprint as a unit rather than fragmenting noisy peaks into spurious components. The stopping criterion uses extreme value statistics to prevent over-extraction from noise.
- (3) A **maximum entropy stopping criterion**: decomposition continues until the residual reaches maximum entropy for its variance. This criterion is basis-independent and connects decomposition to information-theoretic compression.
- (4) **Empirical evaluation** on six standard benchmarks, demonstrating 71–80% improvement over TimesFM on three of four ETT datasets, with the ensemble closing the gap on the previously-losing ETTh2 (0.250 vs TimesFM 0.289, +13%).

2 Related Work

Statistical forecasting. The M-competition series [17] established that simple statistical methods—exponential smoothing (ETS), Theta, ARIMA—are competitive baselines. Hyndman and Athanasopoulos [18] provide the standard reference. The Nixtla StatsForecast library offers reproducible benchmarks of these methods across standard datasets.

Neural forecasting. N-BEATS [5] introduced interpretable neural forecasting with basis expansion. Informer [6] and Autoformer [7] brought attention-based architectures to long-horizon forecasting. PatchTST [2] demonstrated that simple patching with vanilla transformers achieves state-of-the-art supervised performance. iTransformer [8] inverts the standard attention pattern. TimesFM [1] is the current zero-shot benchmark as a foundation model.

Decomposition methods. Facebook’s Prophet [9] decomposes time series into trend, seasonality, and holidays using Fourier terms and piecewise linear trends. STL [10] provides seasonal-trend decomposition. Neither uses information-theoretic criteria for basis selection or stopping. Empirical Mode Decomposition [11] derives basis functions from the data but lacks a principled stopping criterion and produces non-unique decompositions.

Signal processing foundations. The CLEAN algorithm [12] in radio astronomy performs iterative source extraction from dirty images using greedy peak subtraction. Spectral subtraction [13] in audio processing uses noise-shaped error for signal extraction. Our noise-aware Fourier extraction draws on both: the iterative extraction of CLEAN with the noise-shaped footprint modeling of spectral subtraction.

Information-theoretic connections. Maximum entropy methods have a long history in spectral analysis [14] and Bayesian inference [15]. The connection between decomposition and compression is implicit in rate-distortion theory [16]. We make it explicit: each decomposition layer is a lossy compressor, and the stopping criterion is the information-theoretic limit of the combined basis.

3 Method

3.1 Architecture Overview

The input is a univariate time series $\mathbf{x} = (x_1, \dots, x_N)$ and a forecast horizon H . The output is a point forecast $(\hat{x}_{N+1}, \dots, \hat{x}_{N+H})$ with error bounds. The decomposition proceeds in four layers, each operating on the residual of the previous.

- (1) **Pre-detrend:** Remove linear trend via OLS. This prevents trend energy from leaking into the Fourier spectrum.
- (2) **Layer 1 (Fourier):** Iterative noise-aware extraction of periodic components.
- (3) **Layer 2 (Trend):** BIC-selected regression (linear, quadratic, exponential) on the full residual including the pre-removed trend.
- (4) **Layer 3a (Shock):** AIC-selected discrete event detection (step, spike-decay, ramp).
- (5) **Layer 3b (Local):** Recency-weighted AR(p) model on the residual, with AIC order selection.

- (6) **Adaptive wavelet:** If the residual lag-1 autocorrelation exceeds 0.1, apply DWT decomposition on the local residual to capture non-periodic structure.

The forecast is the superposition of all layers, with trend damping (nonlinear trends transition to their tangent line beyond the training boundary), amplitude damping (Fourier components decay with forecast horizon based on Cramér–Rao phase uncertainty), and signal range clamping (forecast constrained to the observed data range plus one standard deviation).

3.2 Layer 1: Noise-Aware Fourier Extraction

Standard Fourier analysis finds all spectral peaks and treats each as a separate component. When noise is present, a single sinusoidal component at frequency f spreads energy across neighboring frequency bins due to spectral leakage and noise. A naive peak-extraction approach would fragment this single noisy sinusoid into multiple spurious components.

Our approach models the expected spectral footprint. Given a sinusoid at frequency f with amplitude A in noise of standard deviation σ , the expected spectral spread is:

$$\text{footprint width} \approx 2 + \left\lceil \frac{2}{\sqrt{\text{SNR}}} \right\rceil \quad (1)$$

where $\text{SNR} = P_{\text{peak}}/P_{\text{noise}}$ and P_{noise} is estimated as the median of the power spectrum.

The extraction algorithm proceeds iteratively:

- (1) Compute the power spectrum of the residual (excluding DC and sub-two-cycle frequencies—these are handled by Layer 2).
- (2) Find the dominant peak. Compute its SNR against the noise floor.
- (3) **Stopping criterion:** The peak must exceed $\text{min_snr} \times P_{\text{noise}} \times (1 + \ln K)$, where K is the number of usable frequency bins. The $\ln K$ term accounts for extreme value statistics: the expected maximum of K i.i.d. exponential random variables scales as $\ln K$. Without this correction, noise peaks in a long signal systematically exceed fixed SNR thresholds.
- (4) Fit a sinusoid $A \cos(2\pi f t + \phi)$ at the refined frequency via least-squares, and subtract.
- (5) Repeat until the stopping criterion fails or the cumulative extraction error threshold is reached.

The frequency refinement uses bounded scalar optimization (Brent’s method) within $\pm 2\Delta f$ of the FFT peak, followed by linear least-squares for amplitude and phase. This achieves sub-bin frequency resolution.

3.3 Layer 2: Trend Fitting

The residual after Layer 1 contains trends, level shifts, and non-periodic drift. We fit four candidate models—none (constant), linear, quadratic, exponential—using weighted least squares with exponential recency weighting (half-life = $N/2$) and select by BIC. The recency weighting ensures the trend model reflects recent behavior more strongly than distant behavior, which is appropriate for forecasting.

For forecasting, nonlinear trends (quadratic, exponential) are damped: beyond the training boundary, the model transitions from the full curve to its tangent line via exponential blending with half-life $N/4$. This prevents catastrophic extrapolation while preserving the local slope.

3.4 Layer 3a: Shock Detection

We compare the Layer 1+2 prediction against the actual signal and detect structured deviations exceeding 2σ of the residual noise. Three shock shapes are fitted at each onset:

$$\text{Step: } s(t) = M \cdot \mathbf{1}[t \geq t_0] \quad (2)$$

$$\text{Spike-decay: } s(t) = M \cdot e^{-\lambda(t-t_0)} \cdot \mathbf{1}[t \geq t_0] \quad (3)$$

$$\text{Ramp: } s(t) = (M + \beta(t - t_0)) \cdot \mathbf{1}[t \geq t_0] \quad (4)$$

The best shape is selected by AIC. The onset is identified as the first point exceeding the threshold, not the point of maximum deviation, which correctly handles step functions where the maximum and onset coincide.

3.5 Layer 3b: Local Correction

The residual after Layers 1–3a is modeled by an $\text{AR}(p)$ process with recency-weighted Yule–Walker estimation. The autocorrelation function is computed with exponential weights (half-life = $W/3$ where W is the fitting window), giving recent residuals more influence on the estimated lag structure.

The AR order is selected by AIC over $p \in \{0, \dots, 12\}$. The forecast from the AR model decays exponentially with forecast horizon (half-life = $2p$), reflecting the principle that local momentum is most informative near-term and uninformative at long horizons.

3.6 Adaptive Wavelet Decomposition

After Layer 3b, we measure the lag-1 autocorrelation of the remaining residual. If $|\text{AC}(1)| \geq 0.1$, predictable structure remains that the Fourier basis could not capture—localized transients, frequency-modulated oscillations, or slowly-varying envelopes. We apply a discrete wavelet transform (Daubechies-4) to this residual, extract significant scales, and forecast by extrapolating the wavelet coefficients with linear extrapolation and exponential decay toward the mean.

The threshold 0.1 is conservative; it activates wavelet decomposition only when there is clear evidence of remaining structure. The choice of basis (Daubechies-4) provides a good balance of smoothness and compactness for time series data.

3.7 The Maximum Entropy Stopping Criterion

The decomposition terminates when the residual has no predictable structure left. We state this precisely:

Stopping criterion. The decomposition is complete when the residual distribution is maximum entropy for its variance. For continuous data with fixed mean and variance, the maximum entropy distribution is Gaussian [15]. A necessary (but not sufficient) condition is zero autocorrelation at all

lags; the sufficient condition is that the full joint distribution matches the i.i.d. Gaussian—i.e., no linear or nonlinear dependence structure remains.

The practical hierarchy, in increasing stringency:

- $\text{AC}(1) \approx 0$: quick check, catches most short-range linear predictability. Sufficient for practical forecasting in most cases.
- $\text{AC}(k) \approx 0$ for all k : rules out long-memory processes. Does not rule out nonlinear dependencies.
- Gaussianity of residual: the maximum entropy condition. Implies all of the above and additionally rules out nonlinear predictability. Testable via Shapiro–Wilk or Jarque–Bera.

Each decomposition layer is a lossy compressor that removes predictable structure from the signal. The residual entropy measures compression optimality. When the residual reaches maximum entropy, the compression is optimal for the given basis—no further decomposition can extract information without introducing bias.

This framing connects time series decomposition to Shannon’s rate-distortion theory [16]: the decomposition layers define a codebook, the residual is the distortion, and the stopping criterion is the point where the distortion is incompressible.

4 Ensemble Architecture

The decomposition pipeline described above is one worker in a larger ensemble. The key finding that motivated this architecture: a plain $\text{AR}(24)$ model—50 lines of NumPy—outperforms the full Fourier decomposition on every ETT dataset at $H = 96$. Neither the decomposition nor the AR model is universally best. The ensemble resolves this by letting each worker prove its worth on held-out data.

4.1 Workers

Three workers consume the same observation stream:

- **AR**: Linear autoregressive model, $p = 24$. Strong on stationary local structure. The strongest individual worker on five of six datasets.
- **Fourier**: The full decomposition pipeline (Layers 1–3b + wavelet). Strong on periodic signals with clean harmonic structure.
- **SR**: Stochastic resonance forecaster with self-regulating noise injection inversely proportional to running error. Strong on non-stationary signals where noise helps weak patterns cross the extraction threshold.

Additional workers (iterative trend-seasonal decomposition, wavelet-only) can be added to the ensemble; we report results for the three-worker configuration that balances accuracy and computational cost.

4.2 Confidence Tracking

Each worker maintains a running confidence score derived from holdout validation. On each context window, the worker fits on the first 75% of the context, forecasts 24 steps into the remaining 25%, and computes the MSE against the held-out actuals. The confidence

is:

$$c = \left[\exp\left(-\frac{\text{MSE}_{\text{holdout}}}{\text{Var}(\mathbf{x})}\right) \right]^2 \quad (5)$$

The squaring forces separation between workers whose raw exponential scores are close but whose relative performance differs meaningfully. A worker with holdout MSE equal to the signal variance gets $c = (e^{-1})^2 = e^{-2} \approx 0.14$; one with MSE at 10% of variance gets $c = (e^{-0.1})^2 = e^{-0.2} \approx 0.82$.

4.3 Aggregation and Collapse

The ensemble forecast is the confidence-weighted mean of worker forecasts:

$$\hat{x}_t = \frac{\sum_i c_i \hat{x}_{i,t}}{\sum_i c_i} \quad (6)$$

Collapse criterion: if the ratio of the highest confidence to the second-highest exceeds $2\times$, the ensemble collapses to the single best worker. This prevents a clearly superior worker from being diluted by weaker models. On the ETT benchmarks, collapse occurs on 1–28% of windows depending on the dataset.

Uncertainty: the error bounds combine the confidence-weighted inter-worker disagreement (captures model uncertainty) with a horizon-growth factor (captures increasing uncertainty over time).

5 Experimental Setup

5.1 Datasets

We evaluate on six standard benchmarks: the four ETT (Electricity Transformer Temperature) datasets introduced by Zhou et al. [6]—ETTh1, ETTh2 (hourly, 17,420 points each), ETTm1, ETTm2 (15-minute intervals, 69,680 points each)—plus Weather (21 meteorological indicators, 52,695 points at 10-minute resolution) and ECL (electricity consumption of 321 clients, 26,304 hourly points). We forecast the OT column for ETT datasets and the first series for Weather and ECL, following the standard univariate protocol.

5.2 Protocol

We use the standard train/val/test split: 12 months training, 4 months validation, 4 months test for ETTh{1, 2}; 34,465/11,521/11,521 for ETTm{1, 2}. Data is normalized using training mean and standard deviation. We evaluate at prediction horizons $H \in \{96, 192, 336, 720\}$ with context length $C = 512$, sliding by H through the test set. Metrics are MSE and MAE on the normalized test data.

5.3 Baselines

We compare against published results for:

- **TimesFM** [1]: 200M parameter foundation model, zero-shot.
- **PatchTST** [2]: supervised transformer, trained per-dataset.
- **ARIMA, Prophet, ETS**: classical statistical methods.
- **Naive**: repeat-last-value baseline.

Our method uses default hyperparameters throughout—no per-dataset tuning.

Table 1: Normalized MSE at $H = 96$ on six benchmarks. **Ens is the confidence-weighted ensemble. **Bold** indicates best among all methods. **TimesFM** and **PatchTST** numbers from published results.**

Dataset	Ens	AR	Fourier	SR	TimesFM	PatchTST
ETTh1	0.109	0.098	0.219	0.109	0.375	0.370
ETTh2	0.250	0.247	0.414	0.264	0.289	0.274
ETTh1	0.064	0.063	0.102	0.066	0.320	0.293
ETTh2	0.194	0.210	0.270	0.209	0.175	0.166
Weather	0.078	0.067	0.156	0.081	—	0.149
ECL	0.471	0.523	0.609	0.463	—	0.129

6 Results

6.1 Main Results

Table 1 reports normalized MSE at $H = 96$ for the ensemble (**Ens**), its individual workers (AR and Fourier decomposition), and the published transformer baselines.

The ensemble beats TimesFM on three of four ETT datasets: ETTh1 (0.109 vs 0.375, +71%), ETTh2 (0.250 vs 0.289, +13%), and ETTm1 (0.064 vs 0.320, +80%). It beats PatchTST on four of six benchmarks, adding Weather (0.078 vs 0.149, +47%). The sole ETT loss is ETTm2 (0.194 vs 0.175, −11%). ECL remains far behind (0.471 vs 0.129)—individual electricity consumption patterns are fundamentally non-periodic and non-stationary.

The critical result is ETTh2: the Fourier decomposition alone scored 0.414 (losing to TimesFM by 43%), but the ensemble’s confidence mechanism detects that the AR worker outperforms on this data and routes accordingly. The ensemble achieves 0.250—winning by 13%.

A natural question: the AR worker alone outperforms the ensemble on four of six datasets (Table 1). Why not use AR alone? Because the *best individual worker varies by dataset*—AR wins on ETTh1/ETTh2/ETTh1/Weather, but SR wins on ECL and the Fourier decomposition contributes on specific windows within datasets where its periodic basis matches. The ensemble’s value is not beating the best individual worker on average but *never choosing the wrong worker*: it routes to the AR when AR is best and away from the AR when it is not. The 10–15% dilution from blending is the cost of that robustness.

6.2 Ablation: Ensemble vs. Individual Workers

Table 1 serves as the ablation: each column shows a single worker’s performance. The AR worker alone is the strongest individual on five of six datasets, but the ensemble provides consistent near-best performance across all datasets without requiring manual model selection. On ETTm2, the ensemble (0.194) outperforms both AR (0.210) and Fourier (0.270) individually, demonstrating that the confidence-weighted blend can exceed any individual worker when their errors are uncorrelated.

6.3 Residual Diagnostic

Table 2 explains the performance split. The residual lag-1 autocorrelation after decomposition is the discriminator: datasets where

Table 2: Residual diagnostics across datasets. Var. Explained is the fraction of input variance captured by the decomposition. AC(1) is the lag-1 autocorrelation of the final residual. Periodic % is the fraction of input variance in Layer 1 (Fourier) components.

Dataset	Var. Expl.	Periodic %	AC(1)	MSE ($H=96$)
ETTh1	91.2%	71.6%	0.037	0.131
ETM1	97.7%	72.1%	0.066	0.088
ETTh2	99.2%	83.0%	0.115	0.400
ETM2	99.8%	87.6%	0.203	0.251

our residuals are white noise ($AC(1) \approx 0$) are datasets where we dominate.

A counterintuitive finding: we explain 99.8% of ETM2’s variance but forecast it worse than ETM1 where we explain 97.7%. Variance explained is not a proxy for forecast accuracy. The 0.2% we miss on ETM2 has autocorrelated structure ($AC(1) = 0.203$)—it is predictable signal we cannot capture. On ETTh1, the 8.8% we miss has $AC(1) = 0.037$ —white noise with no remaining structure. The distinction is not how much you capture, but whether what remains is compressible.

6.4 Computational Cost

All six benchmarks (6 datasets at $H = 96$, ~ 560 rolling windows total) complete in 131 seconds on a 2023 MacBook Pro (Apple M3, no GPU). Per-window inference takes ~ 230 ms for the three-worker ensemble (three independent model fits and forecasts per window). For comparison, TimesFM requires a GPU and a 200M parameter model loaded in memory. PatchTST requires per-dataset training.

7 Analysis

7.1 Why Does Fourier Win on ETTh1?

ETTh1 records oil temperature from an electricity transformer. The dominant structure is a 24-hour diurnal cycle with higher harmonics reflecting the daily heating-cooling pattern. This is textbook periodic signal—exactly what Fourier analysis was designed for. Our extractor finds these frequencies in one FFT pass and subtracts them cleanly. The transformer must learn this periodic structure implicitly from its training corpus through attention patterns, which is a strictly harder task.

7.2 Why Does the Ensemble Win on ETTh2?

ETTh2 records a different variable from the same transformer. The Fourier decomposition alone scores 0.414 on ETTh2—a 43% loss to TimesFM. But the AR worker scores 0.247. The ensemble’s holdout validation detects this: on ETTh2 windows, the Fourier worker’s confidence drops to 0.18–0.50 while the AR worker’s confidence stays at 0.53–0.70. The confidence-weighted blend routes toward the AR worker, and the ensemble achieves 0.250—beating TimesFM by 13%.

This is the core mechanism: the confidence scores act as a real-time statistical proxy, routing algorithmic behavior based on which worker best fits the current data. No single method is universally

best, but the ensemble’s per-window routing finds the best method for each context.

7.3 The Confidence Mechanism as Diagnostic

The worker confidence scores double as a diagnostic. On ETTh1, Fourier and AR have comparable confidence (0.77 vs 0.94)—both work well. On ETTh2, Fourier drops to 0.18 while AR stays at 0.66—a clear signal that the periodic basis is wrong for this data. On ETM2 (the one dataset where we still lose), all workers have similar confidence, and the ensemble cannot distinguish a winner—the -11% gap reflects a genuine limitation where the signal has structure that none of our analytical workers can capture.

Before deploying a forecasting method in production, run the ensemble on a validation window. If the confidence scores are well-separated, the system knows what it’s doing. If they’re clustered, the data may need a fundamentally different approach.

8 Limitations

Univariate only. The current implementation forecasts one variable at a time. Multivariate time series with cross-variable dependencies (e.g., multiple correlated sensors) would benefit from joint decomposition, which we do not address.

Fourier stationarity assumption. Layer 1 assumes the frequency content is stationary over the context window. If the dominant frequencies shift within the context (chirp signals, frequency-modulated data), the extraction will blur the shifting frequencies. A short-time Fourier transform or wavelet-first approach would handle this, at the cost of frequency resolution.

Shock absorption. When a large discrete shock (step function) occurs in the middle of the context, its broadband spectral energy is partially absorbed by Layer 1 before Layer 3a can detect it. We attempted a pre-pass shock detector (median filtering to isolate level shifts before FFT), but it cannot distinguish periodic structure from level shifts without first knowing what is periodic—a chicken-and-egg problem. On all six benchmarks, the pre-pass degraded results (ETTh1 MSE from 0.131 to 0.241). The correct fix is an iterative architecture: extract periodicity, detect shocks on the residual, remove shocks, re-extract. This is future work.

ETM2 gap. The ensemble closes the ETTh2 gap but still loses on ETM2 (0.194 vs TimesFM 0.175, -11%). All workers have similar confidence on ETM2, preventing the routing mechanism from finding a winner. The remaining structure may require a learned basis or a nonlinear decomposition method.

Heavy-tailed residuals. The maximum entropy stopping criterion assumes a Gaussian noise model. If the residual is heavy-tailed (Cauchy, Pareto), the extractor may mine the tails for spurious structure. We tested a Shapiro–Wilk normality gate to block wavelet decomposition on non-Gaussian residuals; it was neutral-to-negative on all six benchmarks because real residuals are almost always slightly non-Gaussian. A more nuanced approach—using the normality test as a diagnostic rather than a gate, or adjusting the stopping threshold based on excess kurtosis—is warranted.

Benchmark scope. We evaluate on six datasets from two domains (electricity transformers, meteorology, individual electricity consumption). A full comparison would include Traffic and the Monash archive.

9 Conclusion

We have shown that a confidence-weighted ensemble of analytical forecasting workers can outperform transformer-based foundation models on the majority of standard benchmarks—71–80% improvement on three of four ETT datasets—with no training data, no GPU, and full interpretability. The key insight is not any single decomposition technique but the **confidence routing mechanism**: each worker proves its worth on held-out data, and the ensemble routes to whoever is most accurate for the current context. No single method is universally best, but the ensemble finds the best method per window.

The confidence scores themselves are a diagnostic. Well-separated scores indicate the system knows which approach fits the data. Clustered scores indicate fundamental difficulty. In regulated domains—banking, insurance, energy, healthcare—where “the model says so” is not an acceptable explanation, this system provides both a prediction and a calibrated measure of how much to trust it, with a per-worker breakdown explaining why.

The 2,000-line implementation, all benchmark scripts, and the data are available at <https://github.com/jmcentire/spectral-forecast>.

References

- [1] A. Das, W. Kong, A. Leber, and R. Sen, “A decoder-only foundation model for time-series forecasting,” in *Proc. ICML*, 2024.
- [2] Y. Nie, N. H. Nguyen, P. Sinthong, and J. Kalagnanam, “A time series is worth 64 words: Long-term forecasting with transformers,” in *Proc. ICLR*, 2023.
- [3] A. F. Ansari, L. Stella, C. Turkmen, et al., “Chronos: Learning the language of time series,” 2024. arXiv:2403.07815.
- [4] K. Rasul, A. Ashok, A. R. Williams, et al., “Lag-Llama: Towards foundation models for probabilistic time series forecasting,” 2023. arXiv:2310.08278.
- [5] B. N. Oreshkin, D. Carpio, N. Chapados, and Y. Bengio, “N-BEATS: Neural basis expansion analysis for interpretable time series forecasting,” in *Proc. ICLR*, 2020.
- [6] H. Zhou, S. Zhang, J. Peng, et al., “Informer: Beyond efficient transformer for long sequence time-series forecasting,” in *Proc. AAAI*, 2021.
- [7] H. Wu, J. Xu, J. Wang, and M. Long, “Autoformer: Decomposition transformers with auto-correlation for long-term series forecasting,” in *Proc. NeurIPS*, 2021.
- [8] Y. Liu, T. Hu, H. Zhang, et al., “iTransformer: Inverted transformers are effective for time series forecasting,” in *Proc. ICLR*, 2024.
- [9] S. J. Taylor and B. Letham, “Forecasting at scale,” *The American Statistician*, vol. 72, no. 1, pp. 37–45, 2018.
- [10] R. B. Cleveland, W. S. Cleveland, J. E. McRae, and I. Terpenning, “STL: A seasonal-trend decomposition procedure based on loess,” *J. Official Statistics*, vol. 6, pp. 3–73, 1990.
- [11] N. E. Huang, Z. Shen, S. R. Long, et al., “The empirical mode decomposition and the Hilbert spectrum for nonlinear and non-stationary time series analysis,” *Royal Society A*, vol. 454, pp. 903–995, 1998.
- [12] J. A. Högbom, “Aperture synthesis with a non-regular distribution of interferometer baselines,” *Astron. Astrophys. Suppl.*, vol. 15, pp. 417–426, 1974.
- [13] S. Boll, “Suppression of acoustic noise in speech using spectral subtraction,” *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 27, no. 2, pp. 113–120, 1979.
- [14] J. P. Burg, “Maximum entropy spectral analysis,” in *Proc. 37th Meeting of the Society of Exploration Geophysicists*, 1967.
- [15] E. T. Jaynes, *Probability Theory: The Logic of Science*. Cambridge University Press, 2003.
- [16] C. E. Shannon, “Coding theorems for a discrete source with a fidelity criterion,” *IRE Nat. Conv. Rec.*, vol. 7, pp. 142–163, 1959.
- [17] S. Makridakis, E. Spiliotis, and V. Assimakopoulos, “The M4 competition: 100,000 time series and 61 forecasting methods,” *Int. J. Forecasting*, vol. 36, no. 1, pp. 54–74, 2020.
- [18] R. J. Hyndman and G. Athanasopoulos, *Forecasting: Principles and Practice*, 3rd ed. OTexts, 2021.
- [19] R. Godahehwa, C. Bergmeir, G. I. Webb, R. J. Hyndman, and P. Montero-Manso, “Monash time series forecasting archive,” in *NeurIPS Datasets and Benchmarks*, 2021.