

Lag-Llama: A Foundation Model for Probabilistic Time Series Forecasting

Lag-Llama is a newly released **open-source foundation model** for time series forecasting. It is designed to work on **univariate** series (single-variable) and outputs full probability distributions rather than single-point forecasts ¹ ². In practice, this means Lag-Llama can generate forecasts on entirely new series *without any training on that specific data* (zero-shot forecasting), analogous to how an LLM can generate text on a new topic without fine-tuning ¹ ³. According to Rasul *et al.* (2024), Lag-Llama was pretrained on a **diverse, multi-domain corpus** of time series, giving it strong generalization: it shows “strong zero-shot generalization capabilities” on unseen datasets and, when fine-tuned on a bit of new data, achieves state-of-the-art performance on average ¹ ⁴. Unlike many existing forecasting models, Lag-Llama is explicitly **probabilistic**: its output is a parameterized distribution (e.g. Student’s t) for each future time step, so users automatically get uncertainty intervals on predictions ⁵ ⁶. (In contrast, deterministic forecasters give only point predictions, offering no measure of confidence ⁵.)

Lag-Llama’s **open-source** release (via GitHub/HuggingFace) makes it freely accessible and extensible. It was built using a **decoder-only Transformer** architecture (inspired by Meta’s LLaMA) that attends over past time-series values ⁷ ⁸. Its inputs are tokenized by **concatenating lagged values and time features** into a single vector at each time step. Concretely, for each timestamp t , the model takes the series value at t and a fixed set of past lags (e.g. 1 day ago, 1 week ago, 1 month ago, etc.), plus static date/time covariates (e.g. second-of-minute, hour-of-day, day-of-week, etc.) ⁹ ¹⁰. This token (current value + lags + time features) is linearly projected into the Transformer’s hidden space. The network applies causal self-attention with **pre-normalization (RMSNorm)** and **Rotary Positional Encoding (RoPE)** at each layer ⁷ ¹¹. The diagram below illustrates the high-level architecture of Lag-Llama:

Lag-Llama uses a stack of masked Transformer decoder layers (like LLaMA) to process a sequence of lagged input features and predict the next time step’s distribution. The input token for each time step combines the current value, selected past lags, and time-of-day/week covariates. After projection and layered self-attention, a distribution head maps the final hidden state to parameters of a forecast distribution (here a Student’s t distribution) ⁷ ⁶.

After the Transformer layers, Lag-Llama’s **distribution head** produces parameters of a probabilistic forecast for the next time step ⁶. In the published model, this head outputs three parameters of a Student’s t -distribution (degrees of freedom, mean, scale) with appropriate nonlinearities to keep them valid ⁶. In other words, instead of predicting a single number, Lag-Llama predicts a full distribution (and thus confidence bounds) for the future value. (The authors note that other distribution types could be used instead — for example, one could swap in a Gaussian output or a learnable normalizing-flow or copula head for more flexibility ⁶.) At inference time, given a context window of past observations, the model is applied autoregressively: one obtains the next-step distribution, samples (or takes the median) to get a prediction, feeds that prediction back into the next input, and repeats until the desired horizon. By sampling many such trajectories, one can construct empirical uncertainty intervals ¹² ⁶. In practice, users can

convert Lag-Llama’s outputs into point forecasts (e.g. by taking the median of the predicted distribution) or use the full probabilistic output for risk-aware decision-making ⁵ ¹² .

Lag-Llama was **pretrained** on a massive, heterogeneous corpus of time series. The authors collected thousands of univariate series from six broad domains (energy, transportation, economics, natural environment, air quality, and cloud/operations) with varying frequencies and lengths ¹³ . In total the pretraining data spanned **millions of individual data windows** across all domains ¹³ . During training, each series window is **standardized** in a *robust* way: they subtract the median and divide by the interquartile range (IQR), and also supply the original mean/variance as extra “summary statistic” covariates ¹⁴ . This scaling ensures the model handles series of different magnitudes and outliers robustly. The pretraining objective is simply to minimize negative log-likelihood of the true next values under the predicted distribution. To ensure the model sees diverse patterns, the data sampling is stratified by dataset (so no single dataset dominates) and supplemented with frequency-domain augmentations (Freq-Mix and Freq-Mask) to improve robustness ¹⁵ . In short, Lag-Llama was trained “from scratch” on a large mix of real-world series, learning to encode temporal patterns purely from data.

With this setup, **zero-shot forecasting** becomes possible. The pretrained Lag-Llama can be applied immediately to a new univariate series: one simply tokenizes the recent history and runs the model to predict the future distribution. In experiments, the authors report that Lag-Llama’s *out-of-the-box* forecasts often match or exceed the performance of models trained *from scratch* on that series ⁴ . For example, it “demonstrates strong zero-shot generalization capabilities” across different domains ⁴ . (In engineering terms, one would try a few context lengths — e.g. 32, 64, 128, ... tokens — and likely use RoPE scaling to handle longer contexts than it was originally trained on ¹⁶ .) If even more accuracy is needed, Lag-Llama can be **fine-tuned** on data from the target series. Fine-tuning has been shown to significantly improve its forecasts: with only a modest amount of new data, Lag-Llama outperforms prior deep-learning forecasters and becomes “the best general-purpose model on average” for that task ⁴ . In practice, users fine-tune by training the model (e.g. via the provided GluonTS estimator) on historical data from the target domain. Key hyperparameters to tune include the context (input) length and learning rate ¹⁷ . The Lag-Llama team recommends grid-searching context lengths (e.g. 32–1024) and learning rates (e.g. 1e-2 to 1e-4), using a validation split and early stopping to avoid overfitting ¹⁷ . With sufficient fine-tuning data, Lag-Llama typically beats specialized models – but even without any tuning, its zero-shot performance is often competitive ⁴ ¹⁸ .

Customization and extensions: Being a modular Transformer-based model, Lag-Llama can be adapted in various ways for different use-cases. For instance:

- **Adjusting the input features:** By default, Lag-Llama uses only *internal* features (lags of the series) and fixed time covariates. In custom applications one could add additional external regressors as extra token components (though this would require retraining or modifying the code). More generally, any improvements to tokenization (e.g. including holiday indicators or domain-specific signals) can be incorporated into the input vector as additional covariates.
- **Changing the context length:** The effective memory of the model is controlled by the input “context length”. In zero-shot use or fine-tuning, one should experiment with larger or smaller windows. The developers note that performance often improves as context grows (up to a point) ¹⁶ . If a series has long-range dependencies, one can increase the context length (for example, from 32 up to 1024 tokens) and enable the model’s RoPE scaling to handle it ¹⁶ .

- **Distributional outputs:** By default, Lag-Llama outputs a Student- t distribution. However, its architecture allows the final distribution head to be swapped. For applications that need other uncertainty models, one could implement a Gaussian or even a learned normalizing-flow output. The authors explicitly mention that more expressive heads (copulas, flows) are possible future enhancements ⁶.
- **Fine-tuning vs. zero-shot:** Depending on data availability, one can either rely on the pretrained model or retrain it. Fine-tuning on domain-specific data typically yields much better accuracy (since the model adapts its weights). For very domain-specific series, retraining more layers or the whole model might be justified.
- **Multi-step forecasting:** For longer horizons, Lag-Llama naturally does autoregressive decoding. One can either generate samples repeatedly (sampling from its output distribution at each step) or take the median/mean at each step. The model's built-in uncertainty allows constructing multi-step prediction intervals.
- **Scaling and normalization:** As part of preprocessing, Lag-Llama uses robust scaling (median/IQR). For a new use-case, similar scaling should be applied so the model sees data in a compatible range ¹⁴. If one wanted to deviate (e.g. use min-max scaling), it would require adjusting the input features accordingly.
- **Multivariate extension (future):** Currently Lag-Llama handles one series at a time. Extending it to true multivariate forecasting would require substantial changes. The authors themselves note that "expanding from univariate towards multivariate approaches" is an important future direction ¹⁹.
- **Implementation:** Lag-Llama is implemented in PyTorch/GluonTS. Developers typically load the pretrained checkpoint (via Hugging Face) and use the provided GluonTS `LagLlamaEstimator`. It can integrate with existing time-series pipelines (e.g. GluonTS or other forecasting libraries) to generate predictions ²⁰.

In summary, Lag-Llama is **designed to be flexible**: it excels out-of-the-box on many tasks, but also offers knobs for customization. Typical steps to adapt it include tuning the context window and learning rate, adding any needed covariates to the input token, and (if desired) modifying the output distribution type ¹⁶ ¹⁷. Its open-source nature and use of standard libraries mean practitioners can probe intermediate embeddings or integrate it into larger systems (e.g. feeding embeddings into downstream models, hooking it into a prediction API, etc.).

Potential Applications: Because Lag-Llama is a *general-purpose* forecaster, it has virtually all the usual time-series use-cases – and more – in reach. Any domain that relies on predicting future values from past trends could benefit, especially when uncertainty estimates are important. For example: - **Finance and Economics:** Forecasting stock prices, currency rates, macroeconomic indicators or sales/revenue. Probabilistic outputs help quantify market risk. - **Business and Retail:** Predicting demand, inventory needs, or customer behavior over time. The model's zero-shot ability could handle new products without retraining. - **Supply Chain and Manufacturing:** Anticipating inventory depletion, logistics demand, or production metrics. (IBM specifically mentioned optimizing cement production as a test case for time-series AI ²¹.) - **Energy and Utilities:** Forecasting electricity load, fuel usage, renewable generation. Lag-Llama's ability to learn from diverse historical data could improve accuracy over simpler models. - **Environmental and Climate:** Predicting weather variables, pollution levels, or climate trends. The probabilistic forecasts can guide policy planning or hazard warnings. - **IoT and Operations:** Sensor readings in factories or data centers (e.g. temperature, pressure, resource usage) can be forecasted for anomaly detection or preventive maintenance. - **Science and Research:** Any time series from experiments or simulations (astronomy data, EEG signals, etc.) where patterns are complex. As IBM's blog notes, foundation models might even "explore the mysteries of the universe" by forecasting scientific measurements ²¹. - **General AI Integration:**

Because it's a generative model, Lag-Llama could potentially be wrapped into chatbots or AI assistants for "time series Q&A", e.g. enabling users to query forecasts in natural language (though this would require additional interface layers).

These examples illustrate the breadth of use-cases. In practice, one would use Lag-Llama for any forecasting task that *can be framed as univariate probabilistic time-series prediction*. Its key advantage is that a single pretrained model (with limited fine-tuning) can serve many domains, avoiding the need to build a brand-new model for each dataset. For instance, a data team could use Lag-Llama as a drop-in tool for forecasting weather data one day, sales data the next, and machine metrics thereafter — benefiting from its pretrained knowledge in each case. Moreover, its uncertainty outputs help decision-makers assess confidence: a wide predictive interval could signal caution.

In short, **Lag-Llama works by combining old ideas (lags, Transformers, distributional output) in a new way**: it tokenizes time series into a Transformer-friendly format, leverages massive pretraining for transfer learning, and uses a simple parametric head for probabilistic forecasts ⁹ ⁶. It can be **adapted or "fine-tuned"** for different contexts by adjusting its input length, learning rates, and output distributions ¹⁶ ¹⁷. Early reports already show it outperforms many traditional forecasting methods with far less per-task effort, and ongoing research may extend it to multivariate data and larger scales. Given these capabilities, Lag-Llama has potential to become a go-to model in any field that needs reliable, uncertainty-aware time series forecasts ⁴ ²².

Sources: The above summary is based on the Lag-Llama research paper ¹ ⁹ ⁶, IBM's technical tutorial and blog posts ² ²¹, developer notes and best-practice guides ¹⁶ ¹⁷, and third-party analyses ¹⁸ ²³. All model details (architecture, tokenization, training) come from the original Lag-Llama publication and repository, while application insights are drawn from case examples in IBM and data science community articles.

¹ ⁴ ⁶ ⁷ ⁹ ¹⁰ ¹¹ ¹² ¹³ ¹⁴ ¹⁵ ¹⁹ [2310.08278] Lag-Llama: Towards Foundation Models for Probabilistic Time Series Forecasting
<https://arxiv.org/html/2310.08278>

² ³ ⁵ Time Series Forecasting with Lag-Llama | IBM
<https://www.ibm.com/think/tutorials/lag-llama>

⁸ ²² ²³ Lag-Llama: The Open-Source Time Series Forecasting Champion - EnDevSols
<https://endevsols.com/lag-llama-the-open-source-time-series-forecasting-champion/>

¹⁶ ¹⁷ GitHub - time-series-foundation-models/lag-llama: Lag-Llama: Towards Foundation Models for Probabilistic Time Series Forecasting
<https://github.com/time-series-foundation-models/lag-llama>

¹⁸ Forecasting in the Age of Foundation Models | by Alvaro Corrales Cano | TDS Archive | Medium
<https://medium.com/data-science/forecasting-in-the-age-of-foundation-models-8cd4eea0079d>

²⁰ Lag-Llama: An Open-Source Base Model for Predicting Time Series Data | by Tom Odhiambo | Medium
<https://medium.com/@odhitom09/lag-llama-an-open-source-base-model-for-predicting-time-series-data-2e897fddf005>

²¹ Meet the next-gen in generative AI
<https://community.ibm.com/community/user/blogs/marie-de-groot/2024/02/28/meet-the-next-gen-in-generative-ai>