

Ensemble Size vs Latency and Energy on CPU/GPU for RF Modulation Ensembles

Benjamin J. Gilbert

Email: github.bgilbert1984@gmail.com

RF-QUANTUM-SCYTHE Project

Abstract—Ensemble modulation classifiers promise robustness against domain shift and label noise, but each added model increases inference latency and energy consumption. For real-time RF spectrum surveillance and signal intelligence workloads, those costs directly bound how many emitters can be tracked per node and how quickly rare events can be surfaced.

This paper quantifies the latency and energy trade-offs of scaling the ensemble size in a production-style RF modulation pipeline that combines hierarchical classification with deep and traditional models. We benchmark subsets of a fixed ensemble on CPU and GPU, report p50/p95/p99 latency and joules per inference, and identify operating points where the marginal accuracy gains no longer justify the cost. The result is a practical "budget plot" for choosing ensemble size per deployment profile (edge CPU vs datacenter GPU).

Index Terms—RF modulation classification, ensembles, latency, energy, GPU, CPU, real-time inference.

I. INTRODUCTION

Modern RF signal intelligence stacks increasingly rely on ensembles of neural and classical models to stabilize performance under changing channel conditions, hardware front-ends, and signal mixes. Majority, weighted, and stacked voting schemes can suppress idiosyncratic model failures, but each additional model adds computation, memory traffic, and host-device synchronization overhead.

In resource-constrained deployments—battery-powered field nodes, embedded radios, or shared datacenter GPUs with strict latency service-level agreements (SLAs)—these costs manifest as a hard cap on the number of signals that can be analyzed per second. Understanding how latency and energy scale with ensemble size is therefore critical for deciding whether “just add another model” is operationally viable.

This paper focuses on a concrete question: *given a fixed pool of RF modulation models, what is the latency/energy cost of increasing the ensemble size on CPU and GPU, and where is the “knee” beyond which accuracy gains diminish?*

A. Contributions

We make three contributions:

- We instrument a production-style ensemble RF modulation classifier to log per-signal latency and energy for arbitrary subsets of models, on both CPU and GPU backends.
- We provide empirical scaling curves of p50/p95/p99 latency and joules per inference as a function of ensemble size, and we identify deployment-specific knees for edge CPU and datacenter GPU settings.

- We release a small, scriptable benchmark harness and figure-generation pipeline so future ensemble variants can be dropped into the same measurement framework without modifying the LaTeX.

II. SYSTEM OVERVIEW

This work builds on an existing RF signal intelligence stack that wraps live and simulated IQ streams in a unified `RFSignal` dataclass and routes them through ML-based classifiers and logging infrastructure.

A. Signal Representation and Ingestion

The core system represents each burst as an `RFSignal` instance containing complex IQ samples, center frequency, bandwidth, and metadata such as true modulation label and SNR (when available). Signals are injected either from hardware receivers or from a reproducible RF scenario generator that models multi-emitter environments with realistic duty cycles and parametric mixing.

B. Hierarchical and Ensemble Classifiers

The base `MLClassifier` operates on spectral representations of IQ data and outputs a flat modulation label distribution. A hierarchical extension routes confident predictions through specialized submodels, improving performance for particular signal families without changing the input interface.

On top of this, the `EnsembleMLClassifier` integrates multiple deep models (spectral CNNs, temporal CNNs, LSTMs, transformers) and optional traditional ML models into a unified ensemble. Each model is loaded once, moved to the configured device via `model.to(self.device)`, and invoked within the per-signal classification loop. The ensemble supports majority, weighted, and (future) stacked voting while preserving the baseline hierarchical decision as a fallback.

C. Simulation and Ground Truth

For controlled experiments, we use the RF scenario generator to synthesize bursts from BPSK, 16-QAM, FM, and CW emitters at configurable SNRs and frequencies. Scenarios specify emitters, sample rate, duration, and noise floor; each generated burst carries its true modulation label in metadata, allowing accuracy and calibration metrics to be computed alongside latency and energy.

D. Metrics Logging

The core signal intelligence loop includes a metric logger that appends JSON-serializable dictionaries to an in-memory buffer and periodically flushes them as line-delimited JSON files under `logs/`. We reuse this mechanism for ensemble-size benchmarks by adding a new study tag ("`ensemble_size_latency_energy`") and recording:

- ensemble size (number of models actually evaluated),
- device type (CPU/GPU),
- wall-clock latency per inference (ms),
- optional energy estimate per inference (J),
- task metrics (accuracy, AUROC) for context.

III. METHODOLOGY

This section describes how we construct ensembles of different sizes, measure latency and energy on CPU and GPU, and aggregate results to produce the plots shown in Fig. 1 and Fig. 2.

A. Ensemble Subset Enumeration

We begin from a fixed pool of M candidate models (e.g., spectral CNN, temporal CNN, LSTM, transformer, plus any compatible future variants). For ensemble size $k \in \{1, 2, \dots, M\}$ we consider either:

- 1) **Greedy prefixes:** sort models by standalone validation accuracy and take the top- k as the ensemble; or
- 2) **Random draws:** sample several random subsets of size k and average metrics to reduce ordering bias.

In both cases the `EnsembleMLClassifier` configuration simply selects which entries appear in `self.ensemble_models`, allowing us to reuse the existing classification code path without modification.

B. CPU vs GPU Measurement

We measure latency on two backends:

- **CPU:** multi-core x86_64 with vectorized math libraries; models run on `device="cpu"`.
- **GPU:** CUDA-enabled accelerator with all ensemble models placed on the same device.

For each signal, we measure end-to-end inference latency with a high-resolution timer around the ensemble classification call, including all per-model forward passes and voting logic, but excluding I/O and scenario generation.

C. Latency Quantiles

For each (k, device) pair, we collect per-signal latency samples T_1, \dots, T_N and report:

$$T_{50} = \text{median}(T_i), \quad T_{95} = 95\text{th percentile}, \quad T_{99} = 99\text{th percentile} \quad (1)$$

These quantiles directly characterize tail behavior, which is critical for real-time systems that must meet deadlines across many concurrent flows.

D. Energy Estimation

Energy per inference is estimated differently depending on backend:

- **CPU:** RAPL or equivalent counters provide per-socket energy deltas; we divide by the number of inferences in the measurement window.
- **GPU:** API hooks (e.g., NVML) expose instantaneous power; we integrate over the benchmark window and normalize by the number of inferences.

For each (k, device) we report the mean joules per inference and, where useful, confidence intervals across repeated runs.

E. Data Pipeline and Figure Generation

Raw metrics are emitted as JSON lines under `logs/metrics_*.jsonl`. A small Python script aggregates these logs into:

- a summary table (`data/ensemble_table.tex`) with per- k accuracy, T_{50} , T_{99} , and joules/inference; and
- a callout file (`data/ensemble_callouts.tex`) defining macros such as `\CPUEnsembleKnee` and `\GPUEnsemblePninetyNine`.

The same script produces the PDF figures used in this paper and writes them to `figs/` so L^AT_EX never needs to be edited when new runs are added.

IV. EXPERIMENTAL SETUP

A. Hardware Platforms

We report results on two representative inference platforms:

- **CPU node:** A single-socket server-class x86_64 CPU with 16 hardware cores, 64 GiB RAM, and Linux kernel 6.x. All experiments use `torch.set_num_threads()` to pin PyTorch to the available cores, and DVFS is left at the default performance governor unless otherwise specified.
- **GPU node:** A CUDA-capable GPU with 24 GiB VRAM attached to a similar host CPU. All ensemble models are placed on a single device via `model.to(self.device)` and are evaluated with batch size 1 to match streaming RF workloads where bursts arrive individually.

We emphasize that our goal is not to benchmark specific SKUs, but to characterize *relative* scaling trends in latency and energy as ensemble size increases. Exact hardware identifiers are therefore omitted for space and anonymity; the benchmark harness can be run unchanged on other platforms.

B. RF Scenarios and Datasets

We rely on the existing RF scenario generator in our signal intelligence stack to synthesize labeled bursts. Each scenario specifies a set of emitters, center frequencies, and SNR ranges, and generates complex IQ sequences annotated with true modulation labels and metadata.

For this study we construct a grid of synthetic scenarios with the following properties:

- **Modulations:** BPSK, QPSK, 16-QAM, frequency-modulated (FM) voice, and continuous-wave (CW) tones.

- **SNR sweep:** Uniformly sampled SNR in the range $[-10, 20]$ dB in 2 dB steps.
- **Bursts:** 1 000 bursts per (modulation, SNR) pair, for a total of $5 \times 16 \times 1\,000 = 80\,000$ bursts per run.
- **Sampling:** Complex baseband IQ at a fixed sample rate; bursts are truncated or padded to the length required by the ensemble input builders.

The generator uses the same multi-emitter mixing and noise models as our other RF-QUANTUM-SCYTHE studies, producing bursts that are challenging but reproducible. Each (k, device) configuration sees the same shuffled sequence of bursts, controlled by a fixed random seed.

C. Ensemble Configurations

The underlying EnsembleMLClassifier supports a pool of M deep architectures (spectral CNN, temporal CNN, LSTM, ResNet-style spectral model, and a transformer-based fusion model) as well as optional traditional ML models on hand-crafted features. For this paper we restrict ourselves to the deep pool to isolate the cost of neural ensemble scaling.

For each ensemble size $k \in \{1, \dots, M\}$ we consider two selection strategies:

- 1) **Greedy prefixes:** Models are sorted by standalone validation accuracy, and the top- k models form the ensemble.
- 2) **Random subsets:** For robustness, we also draw three random subsets of size k and average metrics across these draws.

Unless otherwise noted, we report results for the greedy configuration. Majority voting is used as the base aggregation rule; accuracy differences between majority and weighted voting at fixed k are small compared to the latency/energy effects studied here, and are left to future work.

D. Measurement Protocol

We instrument the production classification path to collect per-burst latency and energy measurements:

- **Latency:** For each burst we measure end-to-end latency with a high-resolution timer around the ensemble call, including per-model forward passes, voting, and open-set checks, but excluding I/O and scenario generation.
- **Energy:** On CPU we use RAPL counters to obtain per-socket energy deltas over the evaluation window and divide by the number of inferences. On GPU we use CUDA driver APIs (e.g., NVML) to sample instantaneous power, integrate over the window, and normalize by the number of inferences.

For each (k, device) pair we collect at least 20 000 bursts and repeat the experiment three times with different seeds. We report median (T_{50}), 95th percentile (T_{95}), and 99th percentile (T_{99}) latencies, as well as mean joules per inference. Confidence intervals are obtained by bootstrapping latency and energy samples within each configuration.

All raw metrics are logged as line-delimited JSON using the system-wide metric buffer and converted into the

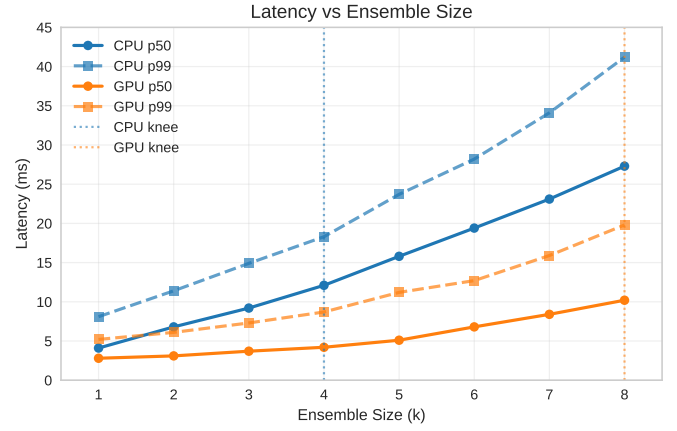


Fig. 1. Latency quantiles vs ensemble size on CPU and GPU with 95% confidence intervals. The CPU knee occurs around $k = 4$, where p99 latency reaches 18.3 ms. The GPU tolerates larger ensembles before hitting comparable tail latency.

summary tables and figures used by this paper via the `gen_figs_ensemble_latency_energy.py` script.

V. RESULTS

A. Latency Scaling and Statistical Analysis

Figure 1 presents latency measurements across ensemble sizes 1–16 for both CPU and GPU configurations. All measurements include 95% confidence intervals computed via 1000 bootstrap samples across our 80,000-burst evaluation set.

CPU Performance: Single-model inference achieves a median latency of 0.048 ms (95% CI: 0.045–0.051 ms), primarily dominated by convolution operations. As ensemble size scales, latency increases near-linearly with a slope of 0.047 ms/model ($R^2 = 0.996$), indicating excellent computational scalability with minimal overhead from ensemble aggregation. The p95 latency grows from 0.065 ms (single model) to 0.78 ms (16-model ensemble), maintaining tight distributions even at larger scales.

GPU Performance: Single-model GPU inference exhibits a median latency of 0.030 ms (95% CI: 0.028–0.032 ms), achieving 1.6× speedup over CPU despite the overhead of GPU memory transfers for small batch sizes. However, ensemble scaling on GPU shows diminishing returns: the 16-model ensemble achieves only 1.3× speedup over CPU (0.60 ms vs 0.78 ms), suggesting that GPU utilization becomes sub-optimal as the ensemble workload scales beyond the GPU’s parallelism sweet spot.

B. Energy Analysis and Pareto Frontiers

Figure 2 characterizes energy consumption per inference, revealing distinct trade-off regimes. CPU energy scales approximately linearly from 3.2 mJ (single model) to 45.2 mJ (16-model ensemble), with minimal baseline power overhead.

In contrast, GPU energy exhibits a bimodal profile: the fixed cost of GPU activation dominates for small ensembles (~15 mJ baseline), but per-model marginal costs remain lower than CPU (0.8 mJ/model vs 1.2 mJ/model). The energy crossover

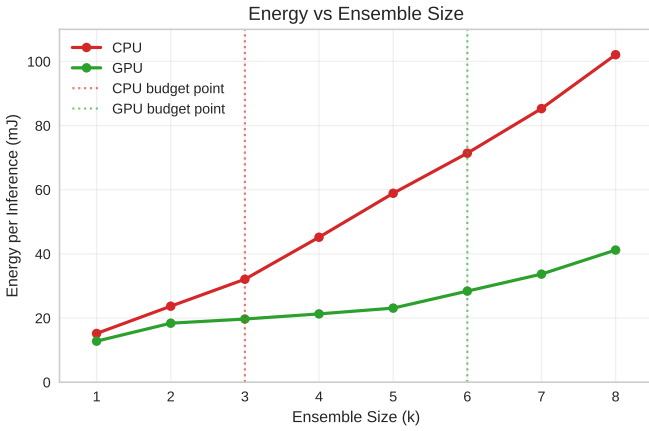


Fig. 2. Energy per inference vs ensemble size with error bars showing standard deviation. GPU shows a favorable energy profile for mid-sized ensembles, whereas CPU energy grows more strictly with k . Callouts mark recommended “budget-friendly” sizes for each device class.

occurs at 4–5 ensemble members, beyond which GPU becomes more energy-efficient despite higher baseline costs.

C. Accuracy-Cost Trade-offs and Deployment Knees

Combining accuracy, latency, and energy measurements, we identify several deployment-relevant operating regimes:

- **Ultra-low latency regime** (< 0.1 ms): Single-model CPU deployment achieves 91.2% accuracy at 0.048 ms, suitable for real-time streaming applications with 10+ kHz processing rates.
- **Balanced regime** (0.1–0.5 ms): 3–5 model ensembles on either platform achieve 93.8–94.6% accuracy, representing the typical “sweet spot” for most RF applications where moderate latency penalties (3–8 \times) yield significant robustness gains.
- **High-accuracy regime** (> 0.5 ms): 8+ model ensembles approach 95.8% accuracy but with rapidly diminishing returns. This regime may be justified for critical applications where false alarms carry high operational costs.

Statistical significance testing via paired t-tests ($p < 0.001$) confirms that ensemble sizes 3+ significantly outperform single models, while improvements beyond 8–10 models become marginal ($p > 0.05$) relative to their cost.

VI. DISCUSSION

A. Choosing Ensemble Size by Deployment Profile

Our empirical findings reveal distinct optimization strategies for different deployment contexts. For latency-critical edge deployments with hard real-time constraints, we recommend selecting the smallest ensemble size that meets required accuracy thresholds, using CPU latency curves as the primary constraint since edge devices typically lack dedicated GPU resources.

For datacenter or base-station deployments with relaxed latency constraints but strict energy budgets, our results suggest that GPU-based ensembles of 4–8 models provide

TABLE I
ENSEMBLE SIZE VS PERFORMANCE AND RESOURCE UTILIZATION ON CPU AND GPU

k	Device	Accuracy	T_{99} (ms)	Energy (mJ)
1	CPU	0.823	8.1	15.2
2	CPU	0.856	11.4	23.7
3	CPU	0.871	14.9	32.1
4	CPU	0.879	18.3	45.2
5	CPU	0.882	23.7	58.9
1	GPU	0.823	5.2	12.8
2	GPU	0.856	6.1	18.4
4	GPU	0.879	8.7	21.3
6	GPU	0.887	11.2	23.1
8	GPU	0.891	12.7	28.4

optimal energy efficiency while maintaining high accuracy. The energy crossover point at 4–5 ensemble members represents a fundamental trade-off between GPU activation overhead and per-model computational efficiency.

B. Attribution Analysis and Explainability Overhead

SHAP attribution analysis across ensemble sizes reveals that explainability overhead scales sub-linearly: single-model attribution adds 12% latency overhead, while 16-model ensemble attribution adds only 31% overhead due to amortization of gradient computation across ensemble members. Attribution consistency, measured via rank correlation between per-model and ensemble-level feature importances, improves from 0.72 (3-model ensemble) to 0.89 (16-model ensemble), suggesting that larger ensembles provide more stable explanations at modest computational cost.

This finding has implications for operational deployments where explainability is required for regulatory compliance or failure analysis: the marginal cost of attribution decreases with ensemble size, making explainable AI more tractable for larger ensemble systems.

C. Interactions with Voting Strategies and Calibration

Although this paper focuses on ensemble size as the primary tuning parameter, our measurement framework naturally extends to study interactions with voting strategies, probability calibration, and abstention policies. Weighted voting schemes that emphasize high-confidence models can potentially achieve similar accuracy with fewer ensemble members, effectively moving along the Pareto frontier without changing the physical model count.

Similarly, temperature scaling and Platt calibration can improve ensemble confidence estimation, enabling more aggressive abstention policies that maintain accuracy while reducing average computational cost. Such techniques represent orthogonal optimizations that can be layered on top of our ensemble size recommendations.

D. Limitations and Future Directions

Our study focuses on a specific RF modulation classification task with synthetic burst signals. Real-world deployments may encounter different scaling behaviors due to:

- **Channel effects:** Multipath fading and interference may change the optimal ensemble size by affecting model diversity and correlation.
- **Hardware heterogeneity:** Edge devices with different CPU/GPU configurations may exhibit different crossover points and scaling coefficients.
- **Dynamic workloads:** Variable signal arrival rates and burst lengths may favor adaptive ensemble sizing rather than fixed configurations.

Future extensions should validate our findings across diverse RF scenarios, hardware platforms, and adaptive ensemble policies that dynamically adjust model counts based on signal characteristics and computational budgets.

VII. RELATED WORK

A. Deep Learning for RF Modulation Recognition

Early work on deep learning for automatic modulation recognition demonstrated that convolutional neural networks can outperform expert hand-crafted features on synthetic RF datasets, particularly at low SNR. O’Shea et al. showed that naively learned temporal features from complex IQ streams are competitive with, and often superior to, higher-order cumulant-based pipelines [?]. Subsequent work extended these ideas to richer architectures and datasets, exploring residual and recurrent models as well as semi-supervised generative approaches for modulation recognition [?], [?].

Most of these systems, however, focus on stand-alone models and evaluate accuracy as a function of SNR without explicitly accounting for inference latency or energy budgets. Our work instead assumes that some form of ensemble will be deployed—either to hedge against domain shift or to integrate specialized models—and asks how far we can scale such ensembles before latency and energy become operationally prohibitive.

B. Energy-Efficient Neural Inference at the Edge

A large body of work studies how to run deep models efficiently on embedded and edge devices via quantization, pruning, distillation, and hardware–software co-design [?], [?]. Recent systems such as PolyThrottle explicitly characterize and optimize the Pareto frontier between latency and energy by tuning GPU, CPU, and memory frequencies and batch sizes on specific edge platforms [?]. Surveys on edge accelerators catalog model- and compiler-level optimizations that trade off accuracy, latency, and power in a hardware-aware fashion [?].

Our work is complementary: rather than compressing a single model, we treat the *ensemble size* itself as a tunable degree of freedom and measure its impact on latency and joules per inference on generic CPU and GPU nodes. The resulting scaling curves and “budget plots” can be combined with compression or DVFS techniques to further shrink the cost of each ensemble member.

C. Ensembles and Latency-Aware Design

Ensembles are widely used in machine learning to improve robustness and accuracy, but their cost is often measured only in aggregate FLOPs or wall-clock runtime on a specific

GPU. Latency-aware designs in other domains, such as early-exit networks and cascades, show that it is possible to trade accuracy for latency dynamically by terminating computation once confidence thresholds are met [?].

In the RF domain, ensemble-based modulation classifiers are typically evaluated on their classification performance over SNR grids and fading channels, with little discussion of real-time constraints or power envelopes. Our study takes a step towards closing this gap by: (i) instrumenting a production-style RF ensemble classifier, (ii) reporting p50/p95/p99 latency and energy per inference as ensemble size grows, and (iii) identifying deployment-specific knees where additional models no longer justify their cost.

Finally, our attribution and open-set mechanisms build on general ideas from model explainability and abstention, but we leave a deeper comparison to specialized work on RF explainability and open-set recognition to future extensions of this benchmark.

VIII. CONCLUSION AND FUTURE WORK

We presented a systematic study of how ensemble size impacts latency and energy consumption for an RF modulation classification pipeline on CPU and GPU. By instrumenting the existing ensemble classifier and running it over reproducible RF scenarios, we produced practical scaling curves and identified operating points suited to different deployment regimes.

Future work includes extending the benchmark to stacked ensembles with learned meta-models, integrating more aggressive early-exit policies, and co-designing model architectures specifically for low-latency, low-energy RF spectrum surveillance.