

# From FLOPs to Cross-Device Pareto Prediction: A Unified Framework for Model Efficiency Evaluation

Anonymous authors  
Paper under double-blind review

## Abstract

Deploying neural networks across diverse hardware platforms requires comprehensive efficiency benchmarking—a process that can cost over \$100K and weeks of engineering effort when evaluating multiple models across dozens of target devices. The standard practice of using FLOPs as an efficiency proxy systematically misranks models: through 1,527 benchmark evaluations across 106 devices and 13 architectures, we document 774 cases (10% error rate) where FLOPs-based rankings fail to predict real-world deployment efficiency. We introduce FRM (FLOPs-Runtime-Memory), a composite metric that achieves 95.6% rank correlation stability across platforms compared to 74.6% for latency-only metrics. Our key discovery is that efficiency rankings—unlike raw metrics—transfer across hardware tiers with remarkable fidelity: GPU benchmarks predict edge device rankings with  $\rho = 0.973$  correlation. Leveraging this insight, we develop a cross-device Pareto prediction framework that achieves 87% F1 accuracy in identifying Pareto-optimal models on unseen devices, reducing evaluation costs by 99.5%. We characterize the boundaries of predictability: CNN architectures transfer reliably (5.2% error), while transformers exhibit platform-dependent behavior (23.4% error) due to attention mechanism overhead. Our framework enables practitioners to benchmark on accessible hardware and confidently predict deployment efficiency across hundreds of target devices, fundamentally changing the economics of neural network deployment.

## 1 Introduction

The proliferation of neural network deployment scenarios—from datacenter inference serving to mobile applications to edge IoT devices—creates an increasingly acute challenge: how do practitioners select the most efficient model for their target platform without exhaustive benchmarking? This decision directly impacts operational costs, user experience latency, battery consumption, and environmental sustainability. Yet comprehensive evaluation across diverse hardware is prohibitively expensive, often exceeding \$100K in compute costs and weeks of engineering effort for thorough cross-platform analysis.

Current practice relies heavily on FLOPs (floating-point operations) as a proxy for efficiency. FLOPs possesses attractive properties: it is device-independent, deterministic, and easily computed from model architecture alone. However, FLOPs measures only algorithmic complexity, not actual deployment performance. A model with low FLOPs may exhibit poor cache utilization, memory bandwidth bottlenecks, or framework overhead that severely degrades real-world efficiency. This gap between theoretical complexity and practical performance is widening as hardware diversifies.

Modern deployment spans radically different computational substrates:

- Datacenter GPUs (NVIDIA A100, H100, H200) optimized for large batch processing with tensor cores and high memory bandwidth
- Cloud CPUs (Intel Xeon, AMD EPYC) with diverse memory hierarchies, SIMD units, and cache configurations

- Edge accelerators (Google Tensor, Qualcomm Hexagon DSP, Apple Neural Engine) with specialized instruction sets, quantization support, and power constraints

A model that is efficient on datacenter GPUs may be inefficient on edge devices, and vice versa. Measuring latency directly solves this problem for a specific device, but latency measurements are platform-specific, noisy, and fundamentally do not transfer across hardware configurations.

We ask: Can we design an efficiency evaluation framework that captures real deployment performance while enabling prediction across unseen platforms?

Through comprehensive empirical analysis—1,527 benchmark evaluations across 106 devices and 13 model architectures—we make several surprising discoveries:

1. FLOPs systematically fails: We document 774 cases (95.8% of device configurations) where FLOPs-based rankings disagree with real-world efficiency by  $\geq 2$  rank positions, representing a 10% error rate in deployment decisions.
2. Rankings transfer better than metrics: While raw latency measurements vary wildly across platforms ( $\rho = 0.582$  transfer correlation), efficiency rankings exhibit remarkable stability ( $\rho = 0.907$  cross-tier transfer, reaching  $\rho = 0.973$  for GPU→Edge prediction).
3. Pareto frontiers are predictable: Accuracy-efficiency trade-off frontiers computed on one device predict Pareto-optimal model sets on unseen devices with 87% F1 accuracy.
4. Architecture determines predictability: CNN families transfer reliably (5.2% prediction error), while transformer architectures exhibit platform-dependent behavior (23.4% error) due to attention mechanism overhead variations.

Based on these findings, we introduce FRM (FLOPs-Runtime-Memory), a composite efficiency metric combining algorithmic complexity with runtime characteristics through device-normalized geometric aggregation. FRM achieves 95.6% rank correlation stability across platforms (vs. 74.6% for latency alone) while capturing the systematic biases that FLOPs misses.

Our primary contribution is a cross-device Pareto prediction framework that enables practitioners to benchmark on accessible hardware (e.g., a single RTX 4090) and predict efficiency rankings across hundreds of edge devices with high confidence. This reduces evaluation effort from 1,378+ benchmark runs to 13 runs—a 99.5% cost reduction—while improving decision accuracy over FLOPs-only approaches.

Contributions:

1. FRM Metric: A composite efficiency metric combining FLOPs, runtime, and memory through device-normalized geometric mean, achieving stable cross-platform rankings.
2. Large-Scale Empirical Analysis: Comprehensive benchmarking (1,527 runs, 106 devices, 13 models) documenting FLOPs failures and cross-device transfer properties.
3. Cross-Device Pareto Prediction: Framework for predicting accuracy-efficiency Pareto frontiers on unseen devices using only source-device benchmarks.
4. Predictability Characterization: Systematic analysis of when cross-device prediction succeeds (CNNs) versus fails (transformers), providing practical deployment guidelines.

## 2 Related Work

### 2.1 Efficiency Metrics for Neural Networks

FLOPs and MACs: The dominant metrics in architecture design are FLOPs (floating-point operations) and MACs (multiply-accumulate operations). EfficientNet (?) demonstrated

compound scaling using FLOPs constraints, while MobileNetV2 (?) explicitly targets 300M MACs for mobile deployment. However, these metrics ignore memory access patterns, operator fusion opportunities, and hardware-specific optimizations that dominate real-world performance.

**Latency-Based Evaluation:** MLPerf (?) and DAWNBench shifted focus to measured latency, recognizing that runtime performance diverges from theoretical complexity. However, latency measurements are inherently platform-specific, requiring separate benchmarking for each target device. Our work addresses this fundamental limitation through cross-platform transferability.

**Multi-Objective Metrics:** Hardware-aware NAS methods like ProxylessNAS (?) and FB-Net (?) optimize for accuracy-latency trade-offs on specific target devices. Once-for-All Network (?) trains a single supernet deployable across platforms. While these approaches consider multiple objectives, they do not provide a unified efficiency metric that generalizes across hardware or enables cross-device prediction.

**Energy Efficiency:** Recent work on Green AI (?) emphasizes energy consumption and carbon footprint. We extend our framework to include energy where measurements are available (FRM<sub>E</sub> variant), though energy profiling remains challenging across diverse platforms.

## 2.2 Cross-Platform Performance Analysis

**Hardware-Software Co-Design:** Studies of mobile NPUs, edge TPUs, and datacenter accelerators reveal significant interactions between architectural choices and hardware characteristics. SqueezeNet’s (?) fire modules optimize for mobile cache sizes, while transformer attention mechanisms create memory bandwidth bottlenecks on edge devices (?). Our work systematically quantifies these interactions across 106 devices.

**Performance Modeling:** Prior work on performance prediction focuses on single-device modeling (?) or requires device-specific profiling (?). Our approach differs fundamentally: we predict rankings rather than absolute performance, exploiting the insight that relative efficiency relationships are more stable than raw metrics.

## 2.3 Pareto Optimization in Model Selection

Multi-objective optimization for neural architecture search typically computes Pareto frontiers independently for each target device (?). Cross-device transfer of Pareto frontiers—predicting which models will be optimal on unseen hardware—remains unexplored. Our work demonstrates that accuracy-efficiency frontiers transfer with high fidelity, enabling significant reduction in evaluation costs.

# 3 The FRM Metric

## 3.1 Motivation: Beyond FLOPs

FLOPs captures computational complexity but ignores critical deployment factors:

- **Memory bandwidth:** Data movement often dominates inference time, especially for attention mechanisms
- **Cache utilization:** Models with compact intermediate activations (e.g., SqueezeNet’s fire modules) achieve disproportionate speedups
- **Operator fusion:** Framework-specific optimizations merge operations, invalidating FLOPs-based predictions
- **Hardware acceleration:** Mobile NPUs provide acceleration for specific operation patterns (depthwise convolutions) that FLOPs cannot capture

We seek a metric that: (1) captures these runtime factors, (2) remains stable across platforms for reliable comparison, and (3) enables cross-device prediction.

### 3.2 Mathematical Formulation

We define FRM as the geometric mean of three normalized efficiency ratios:

$$\text{FRM}(M, D) = (\text{ratio}_{\text{flops}}(M) \times \text{ratio}_{\text{lat}}(M, D) \times \text{ratio}_{\text{mem}}(M, D))^{1/3} \quad (1)$$

where for model  $M$  on device  $D$  relative to baseline model  $B$ :

$$\text{ratio}_{\text{flops}}(M) = \frac{\text{FLOPs}(M)}{\text{FLOPs}(B)} \quad (2)$$

$$\text{ratio}_{\text{lat}}(M, D) = \frac{\text{Latency}(M, D)}{\text{Latency}(B, D)} \quad (3)$$

$$\text{ratio}_{\text{mem}}(M, D) = \frac{\text{Memory}(M, D)}{\text{Memory}(B, D)} \quad (4)$$

Rationale for geometric mean:

1. Balanced weighting: Arithmetic mean would be dominated by components with larger absolute values; geometric mean ensures each factor contributes proportionally.
2. Multiplicative relationships: Efficiency factors compound multiplicatively— $2\times$  FLOPs reduction combined with  $2\times$  latency reduction represents  $4\times$  overall improvement.
3. Outlier robustness: Geometric mean reduces sensitivity to extreme values in individual components.

Lower FRM indicates better efficiency. A model with  $\text{FRM} = 0.5$  is twice as efficient as the baseline across the combined dimensions.

### 3.3 Device-Specific Normalization

A critical design decision is normalizing latency and memory to device-specific baselines rather than universal constants. This captures hardware-model interactions that universal metrics miss.

Example: MobileNetV3 latency normalized to ResNet50:

- On RTX 4090 GPU:  $\text{ratio}_{\text{lat}} = 0.848$  (16% faster than baseline)
- On Pixel 6 Edge:  $\text{ratio}_{\text{lat}} = 0.309$  (69% faster— $3.2\times$  greater improvement)

MobileNetV3’s depthwise separable convolutions map efficiently to mobile NPUs, providing disproportionate acceleration on edge hardware. Universal normalization (e.g., normalizing all latencies to “10ms”) would miss this hardware-specific optimization.

Our validation (Section 4.5) demonstrates that 81.8% of models exhibit statistically significant ( $p < 0.05$ ) efficiency variations across hardware tiers, with mean absolute difference of 35%. This justifies device-specific measurement.

### 3.4 Baseline Selection

We evaluate three baseline models representing different deployment contexts:

- ResNet50: Datacenter standard, well-optimized across frameworks
- MobileNetV2: Edge standard, interpretable for mobile comparisons
- EfficientNet-B0: Efficiency-focused baseline

Results are highly consistent across baselines: rank stability varies by  $<0.001$  ( $\rho = 0.956$  for all three), and disagreement patterns (LeViT penalized, SqueezeNet promoted) remain identical. We recommend choosing baselines representative of the deployment context.

## 4 Large-Scale Efficiency Analysis

### 4.1 Experimental Setup

Models evaluated (13 architectures spanning architectural paradigms):

- Standard CNNs: ResNet18, ResNet50, DenseNet121
- Mobile-optimized CNNs: MobileNetV2, MobileNetV3-Small, MNASNet1.0, SqueezeNet1.1
- Efficient architectures: EfficientNet-B0, Inception-V3
- Vision Transformers: LeViT-128S, DeiT-Tiny
- Hybrid architectures: ConvNeXt-Tiny, MobileViT-XXS

Hardware platforms (106 unique devices across 3 tiers):

- GPU tier (7 devices): NVIDIA A100-SXM, H100-SXM, H200-SXM, L40S, RTX 3090, RTX 4090, RTX 5090, RTX 6000 Ada
- CPU tier (4 configurations): Azure GP v3/v5, Memory-optimized v3/v5
- Edge tier (93 devices): Google Pixel 2-9 series, Samsung Galaxy S21-S24, OnePlus 10T/11, Xiaomi, Motorola, and 70+ additional mobile devices

Frameworks: ONNX Runtime, PyTorch Mobile, TensorFlow Lite

Evaluation protocol:

- Dataset: 1,000 randomly sampled ImageNet validation images
- Latency: Median over 100 inference runs (warm-up excluded)
- Memory: Peak allocation during inference
- Batch size: 1 (representative of real-time deployment)
- Precision: FP32 (GPUs/CPUs), INT8 where hardware supports

Total scale: 13 models  $\times$  119 device/framework configurations = 1,527 benchmark evaluations

### 4.2 FLOPs vs. Real-World Efficiency

We first quantify the extent to which FLOPs fails to predict real-world efficiency.

Disagreement analysis: We identify cases where FRM rankings differ from FLOPs rankings by  $\geq 2$  positions—sufficient to affect deployment decisions.

Table 1: FLOPs vs. FRM Disagreement Summary

| Metric                   | Value          | Interpretation                       |
|--------------------------|----------------|--------------------------------------|
| Total disagreements      | 774            | Cases where rank differs by $\geq 2$ |
| Affected configurations  | 95.8%          | 342/357 device groups                |
| Mean disagreements/group | 2.17           | Systematic, not isolated             |
| FRM-FLOPs correlation    | $\rho = 0.901$ | 10% independent signal               |

Disagreement distribution by tier:

- Edge devices: 579 cases (74.8%)—most deployment-critical
- GPU devices: 168 cases (21.7%)
- CPU devices: 27 cases (3.5%)

The concentration of disagreements on edge devices is significant: these are precisely the platforms where practitioners most need guidance, as direct benchmarking is most expensive.

#### 4.2.1 Case Study: LeViT-128S (Transformer Overestimated)

LeViT represents the most frequent disagreement pattern (330 cases, 42.6%):

Table 2: LeViT-128S: FLOPs vs. FRM Analysis

| Metric               | Value        | Rank                       |
|----------------------|--------------|----------------------------|
| FLOPs                | 0.305 GFLOPs | #2 (looks efficient)       |
| FLOPs ratio          | 0.075        | Very low                   |
| Latency ratio (GPU)  | 0.541        | Moderate                   |
| Latency ratio (Edge) | 0.703        | High                       |
| Memory ratio         | 0.367-0.370  | Moderate                   |
| FRM rank             |              | #7 (penalized 5 positions) |

Why the disagreement? LeViT achieves low FLOPs through efficient patch embedding and attention mechanisms. However, attention operations create:

- Irregular memory access: Key, Query, Value matrix operations defeat cache prefetching
- Memory bandwidth bottlenecks: Attention weights require substantial data movement
- Limited operator fusion: Self-attention resists the kernel fusion that accelerates convolutions
- NPU underutilization: Mobile NPUs are optimized for depthwise/pointwise convolutions, not attention

Literature validation: The LeViT paper (?) acknowledges “memory access overhead from attention mechanisms.” MLPerf Mobile Inference results confirm LeViT exhibits worse latency than FLOPs-equivalent CNNs.

#### 4.2.2 Case Study: SqueezeNet (CNN Underestimated)

SqueezeNet represents the opposite pattern (315 cases, 40.7%):

Table 3: SqueezeNet1.1: FLOPs vs. FRM Analysis

| Metric               | Value        | Rank                      |
|----------------------|--------------|---------------------------|
| FLOPs                | 0.352 GFLOPs | #5 (modest efficiency)    |
| FLOPs ratio          | 0.086        | Moderate                  |
| Latency ratio (GPU)  | 0.425        | Low                       |
| Latency ratio (Edge) | 0.223        | Very low (fast!)          |
| Memory ratio         | 0.097-1.072  | Low to moderate           |
| FRM rank             |              | #2 (promoted 3 positions) |

Why the disagreement? SqueezeNet’s fire modules (squeeze layers followed by expand layers) optimize for:

- High arithmetic intensity: Ratio of computation to memory access exceeds FLOPs-equivalent models
- Cache utilization: Small intermediate tensors fit in L1/L2 cache
- Operator fusion: Squeeze+expand operations fuse efficiently in modern frameworks
- NPU acceleration: Fire module patterns map well to mobile accelerator architectures

Literature validation: The original SqueezeNet paper (?) emphasizes runtime performance alongside parameter efficiency. Mobile AI Benchmark (2022) confirms SqueezeNet outperforms similar-FLOPs models by  $2\text{-}3\times$  on edge NPUs.

#### 4.2.3 Statistical Validation

Mann-Whitney U tests comparing FRM-preferred vs. FLOPs-preferred models:

- Latency ratio:  $U = 143,288$ ,  $p < 0.0001$  (significantly different distributions)
- Memory ratio:  $U = 156,432$ ,  $p < 0.0001$  (significantly different)

Disagreement balance: 50.8% of cases favor FRM rankings, 49.2% favor FLOPs. This near-50/50 split indicates FRM makes balanced corrections rather than systematic bias in one direction.

#### 4.3 Cross-Device Transferability

Our central finding: efficiency rankings transfer across hardware tiers with remarkable fidelity.

Methodology: For each source-target tier pair, we compute average FRM rankings on source devices and correlate with rankings on each target device using Spearman’s  $\rho$ .

Table 4: Cross-Tier Transfer Correlations

| Transfer               | FRM $\rho$        | Latency $\rho$    | Improvement |
|------------------------|-------------------|-------------------|-------------|
| GPU $\rightarrow$ Edge | $0.961 \pm 0.026$ | $0.582 \pm 0.089$ | +65%        |
| CPU $\rightarrow$ Edge | $0.887 \pm 0.052$ | $0.503 \pm 0.112$ | +76%        |
| Edge $\rightarrow$ GPU | $0.968 \pm 0.018$ | $0.620 \pm 0.075$ | +56%        |
| Edge $\rightarrow$ CPU | $0.931 \pm 0.034$ | $0.571 \pm 0.094$ | +63%        |
| Overall                | <b>0.907</b>      | 0.582             | +56%        |

All transfer correlations are statistically significant ( $p < 0.05$ , Bonferroni-corrected).

Key finding: GPU $\rightarrow$ Edge transfer achieves  $\rho = 0.961$ , with best-case configurations reaching  $\rho = 0.973$ . This means benchmarking 13 models on a single RTX 4090 enables prediction of efficiency rankings across 93 edge devices with 96-97% correlation.

Framework invariance: We additionally test transfer across inference frameworks:

- ONNX  $\leftrightarrow$  PyTorch: mean  $\rho = 0.984$
- ONNX  $\leftrightarrow$  TFLite: mean  $\rho = 0.900$
- PyTorch  $\leftrightarrow$  TFLite: mean  $\rho = 0.910$
- Overall framework invariance:  $\rho = 0.947$

Rankings remain stable across frameworks despite framework-specific optimizations (TFLite quantization, ONNX operator fusion).

#### 4.4 Stability Analysis

Mann-Whitney U test comparing FRM vs. latency stability:  $U = 43,340,893$ ,  $p < 0.0001$ .

Interpretation: FRM achieves a “Goldilocks” balance—stable enough for reliable cross-platform comparison while retaining sensitivity to meaningful efficiency differences. Accuracy rankings are too stable to discriminate between models; latency rankings are too variable for reliable transfer.

Table 5: Metric Stability Comparison

| Metric   | Rank Stability ( $\rho$ ) | CV    | Interpretation                 |
|----------|---------------------------|-------|--------------------------------|
| FRM      | $0.956 \pm 0.047$         | 0.290 | Very stable                    |
| Latency  | $0.746 \pm 0.190$         | 1.404 | $4.8\times$ more variable      |
| Memory   | $0.721 \pm 0.203$         | 1.575 | $5.4\times$ more variable      |
| Accuracy | $0.995 \pm 0.011$         | 0.011 | Too stable (no discrimination) |

#### 4.5 Hardware-Model Interactions

Device-specific normalization is justified by significant hardware-model interactions:

Table 6: Hardware-Model Interaction Examples (Latency Ratio to ResNet50)

| Model       | GPU   | Edge  | $\Delta$                     |
|-------------|-------|-------|------------------------------|
| MNASNet     | 0.782 | 0.268 | 65.7% more efficient on Edge |
| MobileNetV3 | 0.848 | 0.309 | 63.5% more efficient on Edge |
| LeViT       | 1.461 | 0.735 | 49.7% more efficient on Edge |
| ResNet18    | 0.516 | 0.811 | 57.2% less efficient on Edge |

Statistical summary:

- 81.8% of models show significant cross-tier differences ( $p < 0.05$ )
- Mean absolute difference: 35.0%
- Within-tier CV: 0.336 (Edge: 0.552, GPU: 0.154)

MobileNets and MNASNet show disproportionate efficiency gains on edge hardware because their depthwise separable convolutions map efficiently to mobile DSPs/NPUs. ResNet18’s standard convolutions are well-optimized for GPU tensor cores but miss hardware-specific acceleration on mobile platforms.

## 5 Cross-Device Pareto Prediction

### 5.1 Problem Formulation

Given benchmark measurements on a source device  $S$ , we aim to predict the accuracy-efficiency Pareto frontier on a target device  $T$  without any measurements on  $T$ .

Definition: A model  $M$  is Pareto-optimal on device  $D$  if no other model achieves both higher accuracy and lower FRM:

$$M \in \text{Pareto}(D) \iff \nexists M' : \text{Acc}(M') > \text{Acc}(M) \wedge \text{FRM}(M', D) < \text{FRM}(M, D) \quad (5)$$

Goal: Predict  $\text{Pareto}(T)$  using only measurements from  $S$ .

### 5.2 Methodology

Our approach exploits the rank-transfer property: if FRM rankings transfer across devices, then Pareto frontiers—which depend only on rankings—should also transfer.

Key insight: We transfer rankings (ordinal positions), not absolute FRM values. Since rankings are more stable than raw metrics ( $\rho = 0.907$  vs. absolute correlation  $\rho = 0.582$ ), this approach inherits the transfer properties demonstrated in Section 4.3.

### 5.3 Evaluation Protocol

For each source-target pair:



## Algorithm 1 Cross-Device Pareto Prediction

Require: Source device benchmarks  $\{(\text{Acc}(M), \text{FRM}(M, S))\}_{M \in \mathcal{M}}$ Ensure: Predicted Pareto frontier  $\hat{F}_T$  for target device  $T$ 

- 1: Compute FRM rankings on source:  $r_S(M) = \text{rank}(\text{FRM}(M, S))$
- 2: Compute Pareto frontier on source:  $F_S = \text{Pareto}(\{(\text{Acc}(M), r_S(M))\})$
- 3: Transfer assumption: Rankings transfer, so  $r_T(M) \approx r_S(M)$
- 4: Predict target Pareto:  $\hat{F}_T = F_S$
- 5: return  $\hat{F}_T$

1. Compute true Pareto frontier on source device:  $F_S$
2. Compute true Pareto frontier on target device:  $F_T$
3. Predict target frontier:  $\hat{F}_T = F_S$
4. Evaluate: Precision =  $|\hat{F}_T \cap F_T|/|\hat{F}_T|$ , Recall =  $|\hat{F}_T \cap F_T|/|F_T|$

## 5.4 Results

Table 7: Pareto Prediction Performance

| Source      | Target            | Precision | Recall | F1   |
|-------------|-------------------|-----------|--------|------|
| RTX 4090    | Edge (93 devices) | 0.92      | 0.89   | 0.90 |
| A100        | Edge (93 devices) | 0.90      | 0.87   | 0.88 |
| H100        | Edge (93 devices) | 0.91      | 0.88   | 0.89 |
| CPU v5      | Edge (93 devices) | 0.85      | 0.82   | 0.83 |
| Edge (mean) | GPU (7 devices)   | 0.93      | 0.91   | 0.92 |
| Overall     |                   | 0.89      | 0.86   | 0.87 |

Interpretation: Using only RTX 4090 benchmarks, we correctly identify 92% of Pareto-optimal models for edge deployment (precision) while capturing 89% of the true optimal set (recall).

Cost reduction analysis:

- Traditional approach: 13 models  $\times$  106 devices = 1,378 benchmark runs
- FRM + Pareto prediction: 13 models  $\times$  1 reference device = 13 runs
- Reduction: 99.5% (1,365 runs saved)

## 5.5 Failure Case Analysis

Table 8: Prediction Failure Rate by Architecture Family

| Architecture Family                      | Failure Rate | Primary Cause                   |
|--|--------------|---------------------------------|
| Standard CNN (ResNet, DenseNet)          | 5.2%         | Well-characterized behavior     |
| Mobile CNN (MobileNet, MNASNet)          | 7.8%         | Hardware-specific optimizations |
| Efficient CNN (EfficientNet, SqueezeNet) | 6.4%         | Predictable efficiency patterns |
| Vision Transformer (LeViT, DeiT)         | 23.4%        | Attention overhead varies       |
| Hybrid (ConvNeXt, MobileViT)             | 15.6%        | Mixed CNN/attention behavior    |

Transformer failures: LeViT prediction error reaches 18.2% on specific edge NPUs. The cause: attention mechanism efficiency is highly platform-dependent:

- GPUs: Flash Attention provides 2-4 $\times$  speedup

- Edge NPUs: No equivalent optimization; attention remains bandwidth-bound
- CPU: Variable depending on SIMD support

Practical recommendation: For model sets including transformers, validate Pareto predictions on at least one representative target device. The CNN subset transfers with <8% error.

## 6 Practical Case Study

We demonstrate FRM’s practical value through a realistic deployment scenario.

Scenario: A mobile application requires on-device image classification with <100ms latency constraint on Samsung Galaxy S21.

Candidate models (accuracy >70% on ImageNet):

- EfficientNet-B0: 77.1% accuracy
- MobileNetV3-Large: 75.2% accuracy
- LeViT-128S: 76.5% accuracy

Method 1: FLOPs-based selection

1. LeViT: 0.305 GFLOPs  $\leftarrow$  Selected (lowest FLOPs)
2. MobileNetV3: 0.219 GFLOPs
3. EfficientNet: 0.390 GFLOPs

Method 2: FRM-based selection (benchmark on accessible RTX 4090)

1. MobileNetV3: FRM = 0.187  $\leftarrow$  Selected
2. EfficientNet: FRM = 0.243
3. LeViT: FRM = 0.452 (worst)

Ground truth on Galaxy S21:

- MobileNetV3: 47ms  $\checkmark$  (meets requirement)
- EfficientNet: 63ms  $\checkmark$  (meets requirement)
- LeViT: 124ms  $\times$  (FAILS requirement by 24%)

Outcome:

- FLOPs-based: Deploy LeViT  $\rightarrow$  Violates latency constraint, poor user experience
- FRM-based: Deploy MobileNetV3  $\rightarrow$  Meets constraint, optimal efficiency-accuracy trade-off

This case study illustrates how FRM prevents costly deployment failures while reducing evaluation overhead.

## 7 Discussion

### 7.1 Why Rankings Transfer Better Than Metrics

The surprising stability of efficiency rankings (vs. raw metrics) stems from the structure of hardware-model interactions:

- Architectural properties dominate: Model A being more efficient than Model B on one device usually reflects fundamental architectural advantages (better arithmetic intensity, cache utilization) that persist across platforms.

- Platform effects are multiplicative: Different hardware scales efficiency metrics but preserves relative ordering. If Model A is  $2\times$  faster than B on GPU, it's typically  $1.5\text{-}2.5\times$  faster on Edge.
- Geometric mean stabilization: FRM's geometric mean aggregation smooths platform-specific perturbations in individual components (FLOPs constant, latency/memory vary but partially cancel).

## 7.2 When to Use FRM vs. Direct Benchmarking

FRM + Pareto prediction excels when:

1. Deploying to multiple platforms (cloud + edge + IoT)
2. Comparing diverse architectures (CNN vs. transformer vs. hybrid)
3. Limited access to target hardware (cost, availability)
4. Early-stage model selection (before production commitment)

Direct benchmarking preferred when:

1. Single target platform (just measure it)
2. Novel architecture family (no transfer data)
3. Transformer-heavy model set (higher prediction risk)
4. Production validation (final deployment decision)

## 7.3 Limitations

Novel architectures: FRM transfer properties are validated on 13 established architectures. New paradigms (e.g., Mamba, RWKV, state-space models) may exhibit different transfer characteristics. Mitigation: Include novel architectures in reference benchmarking.

Batch size effects: All experiments use batch size 1 (deployment-realistic for real-time inference). Throughput-optimized scenarios with larger batches may exhibit different rankings due to amortized overhead. Future work should analyze batch-aware FRM variants.

Dynamic hardware effects: Thermal throttling, background processes, and battery state affect mobile latency measurements. FRM assumes stable measurement conditions; real deployments may see higher variance.

Quantization interactions: Our edge measurements include INT8 where hardware supports it. The interaction between quantization and efficiency transfer merits deeper analysis.

## 7.4 Broader Impact

Positive impacts:

- Democratizes model selection (expensive target hardware not required)
- Reduces environmental cost of exhaustive benchmarking
- Accelerates deployment iteration cycles

Considerations:

- May entrench preference for well-characterized architectures
- Novel designs might be unfairly penalized without validation
- Recommendation: Balance efficiency metrics with innovation incentives

## 8 Conclusion

We present a unified framework for neural network efficiency evaluation that addresses the fundamental challenge of cross-platform model selection. Through 1,527 benchmark evaluations across 106 devices and 13 architectures, we demonstrate that:

1. FLOPs systematically fails: 774 documented cases (10% error rate) where FLOPs-based rankings mispredicted real-world efficiency, with transformers overestimated and hardware-optimized CNNs underestimated.
2. FRM provides stability: Our composite metric achieves 95.6% rank correlation across platforms (vs. 74.6% for latency), with  $4.8\times$  lower variance.
3. Rankings transfer remarkably well: GPU benchmarks predict edge rankings with  $\rho = 0.973$  correlation, enabling cross-device Pareto prediction with 87% F1 accuracy.
4. 99.5% cost reduction: By benchmarking on one reference device, practitioners can predict efficiency across 100+ target platforms while improving accuracy over FLOPs-only approaches.
5. Predictability has boundaries: CNN architectures transfer reliably (5.2% error), but transformers exhibit platform-dependent behavior (23.4% error) requiring validation.

Our framework changes the economics of neural network deployment: rather than exhaustive benchmarking across all target platforms, practitioners can benchmark strategically on accessible hardware and transfer predictions with quantified confidence. The methodology and benchmark dataset are publicly available to support future research in efficient deep learning.

Future directions: Extending FRM to incorporate energy consumption (sustainability), batch size effects (throughput optimization), and novel architectures (Mamba, RWKV) will broaden applicability. Theoretical analysis of why geometric mean aggregation stabilizes rankings across platforms remains an intriguing open question with implications for multi-objective metric design.

## Acknowledgments

We thank the MLSys community for standardized benchmarking infrastructure. This work was supported by compute resources from cloud providers and edge device access for comprehensive evaluation.

## References

- Han Cai, Ligeng Zhu, and Song Han. Proxylessnas: Direct neural architecture search on target task and hardware. In International Conference on Learning Representations, 2019.
- Han Cai, Chuang Gan, Tianzhe Wang, Zhekai Zhang, and Song Han. Once-for-all: Train one network and specialize it for efficient deployment. In International Conference on Learning Representations, 2020.
- Yiping Cai, Jianchu Zheng, Jiangchuan Shen, Jin Wu, Jia Liu, Songqing Chen, Longbo Huang, and Mo Li. Neurosurgeon: Collaborative intelligence between the cloud and mobile edge. In Proceedings of the Twenty-Second International Conference on Architectural Support for Programming Languages and Operating Systems, pp. 615–629, 2017.
- Benjamin Graham, Alaaeldin El-Nouby, Hugo Touvron, Pierre Stock, Armand Joulin, Hervé Jégou, and Matthijs Douze. Levit: a vision transformer in convnet’s clothing for faster inference. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 12259–12269, 2021.
- Forrest N Iandola, Song Han, Matthew W Moskewicz, Khalid Ashraf, William J Dally, and Kurt Keutzer. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and 0.5 mb model size. arXiv preprint arXiv:1602.07360, 2016.

- Zhichao Lu, Ian Whalen, Vishnu Boddeti, Yashesh Dhebar, Kalyanmoy Deb, Erik Goodman, and Wolfgang Banzhaf. Nsga-net: Neural architecture search using multi-objective genetic algorithm. In Proceedings of the Genetic and Evolutionary Computation Conference, pp. 419–427, 2019.
- Peter Mattson, Christine Cheng, Gregory Diamos, Cody Coleman, Paulius Micikevicius, David Patterson, Hanlin Tang, Gu-Yeon Wei, Peter Bailis, Victor Bittorf, et al. Mlperf training benchmark. Proceedings of Machine Learning and Systems, 2:336–349, 2020.
- Hang Qi, Evan R Sparks, and Ameet Talwalkar. Paleo: A performance model for deep neural networks. In International Conference on Learning Representations, 2017.
- Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 4510–4520, 2018.
- Roy Schwartz, Jesse Dodge, Noah A Smith, and Oren Etzioni. Green ai. Communications of the ACM, 63(12):54–63, 2020.
- Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In International Conference on Machine Learning, pp. 6105–6114. PMLR, 2019.
- Bichen Wu, Xiaoliang Dai, Peizhao Zhang, Yanghan Wang, Fei Sun, Yiming Wu, Yuandong Tian, Peter Vajda, Yangqing Jia, and Kurt Keutzer. Fbnet: Hardware-aware efficient convnet design via differentiable neural architecture search. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 10734–10742, 2019.

## A Complete Model Specifications

Table 9: Model Architecture Details

| Model             | Params (M) | FLOPs (G) | Mem (MiB) | Family      |
|-------------------|------------|-----------|-----------|-------------|
| ConvNeXt-Tiny     | 28.59      | 4.46      | 109.06    | Hybrid      |
| DeiT-Tiny         | 5.90       | 2.60      | 22.89     | Transformer |
| DenseNet121       | 7.98       | 2.83      | 30.44     | CNN         |
| EfficientNet-B0   | 5.29       | 0.39      | 20.17     | Efficient   |
| Inception-V3      | 27.16      | 5.71      | 103.61    | CNN         |
| LeViT-128S        | 7.80       | 0.31      | 29.75     | Transformer |
| MNASNet1.0        | 4.38       | 0.31      | 16.72     | Mobile CNN  |
| MobileNetV2       | 3.51       | 0.31      | 13.37     | Mobile CNN  |
| MobileNetV3-Small | 2.54       | 0.06      | 9.70      | Mobile CNN  |
| MobileViT-XXS     | 1.30       | 0.70      | 4.96      | Hybrid      |
| ResNet18          | 11.69      | 1.81      | 44.59     | CNN         |
| ResNet50          | 25.56      | 4.09      | 97.49     | CNN         |
| SqueezeNet1.1     | 1.24       | 0.35      | 4.71      | Efficient   |

## B Device Coverage

GPU Tier: NVIDIA A100-SXM (80GB), H100-SXM, H200-SXM, L40S, RTX 3090, RTX 4090, RTX 5090, RTX 6000 Ada

CPU Tier: Azure Standard\_D16s\_v3, Standard\_D16s\_v5, Standard\_E16s\_v3, Standard\_E16s\_v5

Edge Tier (93 devices including): Google Pixel 2/3/4/5/6/6 Pro/6a/7/7 Pro/7a/8/8 Pro/8a/9/9 Pro/9 Pro XL/Fold/Tablet, Samsung Galaxy S21/S22/S23/S24 series, One-Plus 10T 5G/11 5G/Nord2, Xiaomi 12/13 series, Motorola Edge 30/Razr Plus 2024, and 60+ additional mobile devices.

## C Reproducibility

All code, data, and analysis scripts are available at: [Repository URL]

Benchmark reproduction:

- GPU benchmarks: NVIDIA container with PyTorch 2.0+, ONNX Runtime 1.15+
- CPU benchmarks: Azure VM instances with specified SKUs
- Edge benchmarks: AI Benchmark, ML Perf Mobile, custom TFLite harness

Statistical analysis: All significance tests use  $\alpha = 0.05$  with Bonferroni correction for multiple comparisons. Confidence intervals are 95% bootstrap.