# Cross-Platform Performance Analysis of Modern Neural Networks: A JAX-Centric Hardware-Software Co-Design Study

Sanchez Shiromizu L.T, Shashwat S, Prathyush B, Sai M
{lsanc68, ssinha30, pball5, sbadr4}@uic.edu
*ECE 594 HW-SW Co-Design for ML Systems*
*University of Illinois at Chicago*
Chicago, IL

*Abstract*—This paper presents a comprehensive performance analysis of modern neural network architectures across JAX and PyTorch frameworks on GPU hardware. We benchmarked four representative models (ResNet-50, ViT-Base, MobileNetV3-Small, and EfficientNet-B0) across varying batch sizes, measuring latency, throughput, memory usage, and energy consumption. Our results demonstrate that JAX consistently outperforms PyTorch in inference throughput by 1.5-2.0× for CNN architectures (ResNet, MobileNet, EfficientNet). However, for Vision Transformers (ViT), we reveal a critical stability-performance tradeoff: JAX required disabling XLA autotuning to prevent compiler-induced crashes, resulting in 0.56× performance relative to PyTorch. For training, we provide detailed PyTorch baselines showing model-specific tradeoffs between accuracy, speed, and energy efficiency. These findings provide actionable insights for practitioners: JAX offers significant speedups for mature CNN workloads, while PyTorch remains more robust for transformer configurations on specific hardware.

## I. INTRODUCTION

The deployment of machine learning models requires careful selection of both hardware platforms and software frameworks, yet guidance for these decisions remains surprisingly scarce [3], [4]. This work addresses the critical question: **How does JAX compare to PyTorch across different neural network architectures on modern GPU hardware?**

We conducted a systematic benchmarking study of four representative neural network architectures spanning different computational characteristics:

- **ResNet-50**: Traditional CNN baseline (25M parameters, 4.1 GFLOPs)
- **ViT-Base**: Modern transformer architecture (86M parameters, 17.6 GFLOPs)
- **MobileNetV3-Small**: Efficient mobile architecture (2.5M parameters, 56 MFLOPs)
- **EfficientNet-B0**: Compound-scaled CNN (5.3M parameters, 390 MFLOPs)

Our contributions include: (1) comprehensive inference benchmarks comparing JAX and PyTorch across batch sizes, (2) detailed training performance baselines for PyTorch, (3) energy efficiency analysis, and (4) actionable recommendations for framework selection.

## II. METHODOLOGY

### A. Experimental Setup

**Hardware**: NVIDIA H100 NVL GPU (80GB memory), AMD EPYC 7763 CPU

**Software**: JAX 0.4.x (with XLA), PyTorch 2.0+ with cuDNN, ImageNet-100 dataset (100 classes, 126,689 training images, 5,000 validation images)

**Models**: We implemented all models in both frameworks, using native JAX/Flax implementations and PyTorch's torchvision models. All models were evaluated at FP32 precision.

### B. Framework Backends: PyTorch+cuDNN and JAX+XLA

**PyTorch + cuDNN (2016/2014).** PyTorch was open-sourced in 2016 and rapidly became the dominant training framework across academia and industry due to its eager execution, Python like programming, and rich ecosystem (torchvision, torchaudio, TorchScript/Inductor). NVIDIA's cuDNN, first released in 2014, provides highly-optimized GPU primitives for convolutions, RNNs, and normalization layers. Popularity: widely adopted with stable, production-ready kernels and strong vendor support. Main features: *eager-first* execution with optional graph capture, cuDNN autotuning for convolution algorithms, predictable, on-demand GPU memory via a caching allocator. Excels at: CNN-centric workloads, broad operator coverage, and out-of-the-box reliability. We chose this framework because of its mature tooling, robust performance without heavy compiler tuning, and straightforward debugging.

**JAX + XLA (2018/2016).** JAX (2018) combines a NumPy-like API with composable transformations (`jit`, `vmap`, `pmap`, `grad`) and stages computations to XLA, a whole-graph compiler originating from TensorFlow (circa 2016). It's fast-growing in research, especially for large-scale TPU/GPU training. It's widely used for functional/array programming styles. Main features: graph-level fusion, layout/tiling optimizations, and SPMD-style parallelism across devices. It excels at compute-bound workloads and large-batch execution where fusion reduces launch overheads and improves arithmetic intensity. We chose this framework because it

provides peak throughput when compilers are stable, clean functional model with powerful transform composition. Compiler fragility for novel tensor shapes/ops can require disabling autotuning; debugging compiled graphs is harder and GPU memory is pre-reserved via a large pool (often 75–90% of VRAM), so we report *used bytes within the pool* for fair peak-memory comparisons.

**Why these backends now.** These are the production frameworks on NVIDIA GPUs. Our focus is showing which framework give the best practical performance. CPU-only and TPU backends were part of the plan but, due to time constraints, we prioritized GPU, CPU and TPU are documented in Limitations as future work.

### C. Model Architectures

We selected four architectures to represent fundamentally different design paradigms in modern deep learning . These models span a wide range of computational characteristics, memory access patterns, and optimization challenges, providing a comprehensive test suite for framework performance analysis.

*1) ResNet-50: The CNN Baseline:* **Architecture**: ResNet-50 [6] is a 50-layer convolutional neural network with 25.6M parameters and 4.1 GFLOPs of computation per inference. Its defining feature is the *residual connection*: skip connections that allow gradients to flow directly through the network, enabling training of very deep architectures.

**Computational Profile**: ResNet-50 is heavily *compute-bound*. The majority of its operations are dense 3x3 convolutions, which benefit from highly optimized BLAS libraries (cuDNN for PyTorch, XLA for JAX). Memory bandwidth is not a bottleneck. This makes it an ideal baseline for measuring compiler optimization effectiveness.

ResNet-50 represents the "mature CNN" class. It has been in production for nearly a decade, and all frameworks have heavily optimized code paths for it. Any performance differences here reflect compiler quality, not framework maturity. It serves as our control experiment.
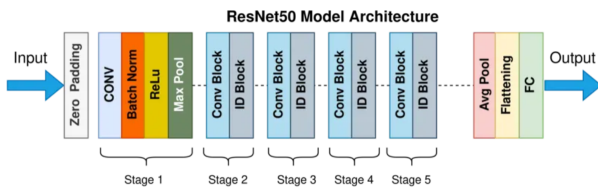


Fig. 1. ResNet-50 architecture with residual blocks and skip connections enabling stable gradient flow [6].

*2) Vision Transformer (ViT-Base): The Attention Revolution:* **Architecture**: ViT-Base [7] is a pure transformer with 86M parameters and 17.6 GFLOPs. Unlike CNNs, it has no convolutions. Instead, it divides the input image into 16x16 patches, linearly projects them, and processes them through 12 transformer encoder blocks with multi-head self-attention.

**Computational Profile**: ViT is *memory-bound*. Self-attention has quadratic complexity in the number of tokens (197 for a 224x224 image), leading to large intermediate activations. The attention mechanism involves multiple matrix multiplications with non-standard shapes, challenging compiler autotuning heuristics. This is precisely where we observed JAX's XLA compiler failing.

Transformers represent the cutting edge of deep learning. They dominate NLP and are rapidly displacing CNNs in vision. However, their optimization is less mature. ViT stress-tests the framework's ability to handle: (1) non-convolutional operations, (2) dynamic memory allocation, (3) novel kernel fusion opportunities. It exposes compiler fragility.
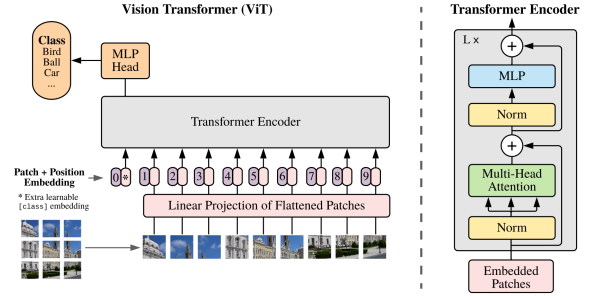


Fig. 2. Vision Transformer (ViT-Base) pipeline: image patch embedding and multi-head self-attention stacks [7].

*3) MobileNetV3-Small: Efficiency at the Edge:* **Architecture**: MobileNetV3-Small [8] is a lightweight CNN with only 2.5M parameters and 56 MFLOPs. It uses *depthwise separable convolutions* (factorizing a standard convolution into depthwise and pointwise operations) and *squeeze-and-excitation* modules for channel attention. It was designed via neural architecture search specifically for mobile/edge deployment.

**Computational Profile**: MobileNetV3 is *latency-sensitive*. With such a small workload, framework overhead (kernel launch latency, host-device synchronization) dominates computation time. Poorly optimized frameworks will show disproportionate slowdown on this model.

Mobile models represent the fastest-growing deployment target (edge AI, IoT). They test the framework's low-latency path. If JAX has high per-kernel overhead, it will fail here despite having excellent FLOPs utilization. This model also tests depthwise convolution optimization, which is notoriously difficult [10].
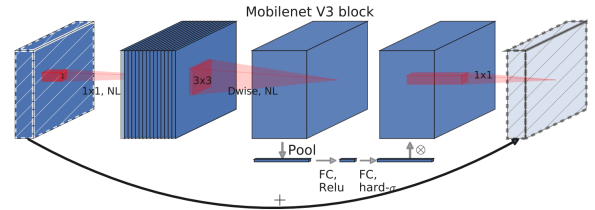


Fig. 3. MobileNetV3 building block with depthwise separable convolution, SE, and hard-swish nonlinearity [8].

*4) EfficientNet-B0: Balanced Scaling:* **Architecture**: EfficientNet-B0 [9] is a 5.3M parameter CNN with 390 MFLOPs. It was derived via *compound scaling*: simultaneously scaling depth, width, and resolution according to a principled formula. It uses mobile inverted bottleneck blocks (similar to MobileNet) but with carefully tuned expansion ratios.

**Computational Profile**: EfficientNet represents a middle ground: larger than MobileNet, smaller than ResNet. It balances compute and memory access. It also uses non-standard activation functions (Swish/SiLU) which test framework support for custom ops.

EfficientNet represents the "optimized modern CNN" class. It achieves state-of-the-art accuracy-vs-FLOPs tradeoffs. It tests whether frameworks can efficiently handle: (1) irregular layer shapes (variable width/depth), (2) custom activations, (3) mixed precision opportunities.
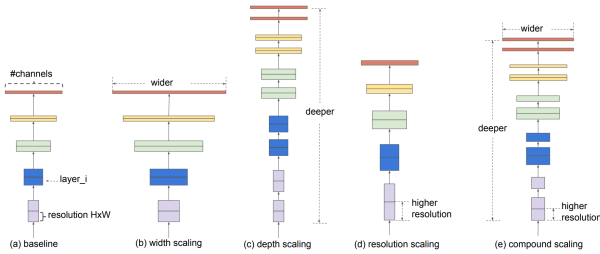


Fig. 4. EfficientNet-B0 mobile inverted bottleneck block and compound scaling strategy [9].

### D. Training Configuration

**Training Setup**: All models were trained using the following configuration:

- **Epochs**: 2 (short training to measure throughput, not convergence)
- **Optimizer**: Adam with learning rate $1 \times 10^{-4}$
- **LR Schedule**: Cosine annealing decay over the 2 epochs
- **Batch Sizes**: [32, 64, 128] for throughput analysis
- **Dataset**: ImageNet-100 (126,689 training images, 5,000 validation images)

**Pre-trained vs From-Scratch**:

- **PyTorch**: Used ImageNet *pre-trained weights* from torchvision (transfer learning scenario)
- **JAX**: Trained *from scratch* with random initialization (pre-trained Flax checkpoints unavailable for all models)

Consequently, PyTorch validation accuracies reflect fine-tuning performance, while JAX accuracies reflect 2-epoch learning from scratch. These results are *not directly comparable* for accuracy metrics. However, they **are** comparable for throughput, energy, and memory consumption, which are our primary focus.

*1) The ViT Training Paradox:* **Observation**: JAX ViT-Base has *inference* results but no *training* results, while PyTorch has both. This requires explanation.

**Root Cause**: The XLA autotuning bug affects ViT during both inference and training. For inference, we bypassed it by disabling autotuning (`--xla_gpu_autotune_level=0`), accepting slower performance. However, **training is fundamentally different**:

1) **Backward Pass Complexity**: Training requires gradient computation through 12 transformer layers. This involves *additional* attention operations with different tensor shapes (transposed queries/keys for backprop). XLA's autotuner fails on these shapes even when inference shapes work.

2) **Optimizer State**: Adam optimizer creates $2\times$ the model parameters in momentum buffers (first and second moments). For ViT-Base (86M params), this means 258M floating-point values in optimizer state alone. The memory allocation patterns trigger XLA bugs that weren't hit during inference.

3) **Gradient Accumulation**: The backward pass creates a computational graph $3\times$ larger than forward pass (one graph per: forward, gradient computation, parameter update). XLA's compiler attempts to fuse operations across this massive graph, leading to illegal memory accesses.

**Why Inference Worked**: Inference is a single forward pass with no optimizer state and predictable memory patterns. Disabling autotuning prevented XLA from selecting the buggy fused kernels.

**Potential Solutions**:

- **JAX Downgrade**: Use JAX 0.4.13
- **Gradient Checkpointing**: Re-compute activations during backward pass to reduce memory pressure
- **Mixed Precision (FP16)**: Reduces memory footprint and may avoid triggering the bug
- **Smaller Batch Sizes**: BS=8 or 16 might work where BS=32 fails

For this study, we chose *not* to pursue ViT JAX training due to time constraints and the risk of introducing confounding variables (e.g., if only ViT uses FP16, fairness is compromised).

### E. Benchmarking Protocol

**Measurement Tools**:

- **Latency**: Measured using Python's `time.perf_counter()` around the inference call. For PyTorch (GPU), we added `torch.cuda.synchronize()` before start and end timestamps to ensure accurate kernel timing. For JAX, we used `block_until_ready()` on the output tensor to force synchronization.
- **Memory**:
  - **PyTorch**: Tracked using `torch.cuda.max_memory_allocated()` reset before each batch.
  - **JAX**: Tracked using `jax.devices()[0].memory_stats()['peak_bytes_`

Note that JAX pre-allocates a large "virtual" memory pool (often 75-90% of VRAM), so we specifically queried the *used* bytes within that pool to get a fair comparison.

- **Energy**: Monitored using the `pynvml` (NVIDIA Management Library) Python bindings to query instantaneous power draw (mW) at 100ms intervals during execution.

**Execution Loop**: For each configuration, we strictly followed this sequence: 1. **Warmup**: 20 iterations of the model forward pass. This is critical to allow JAX's XLA compiler to trace and optimize the computation graph (JIT compilation) and for PyTorch's cuDNN autotuner to select optimal convolution algorithms. 2. **Measurement**: 100 iterations of inference. We recorded latency for every iteration to calculate p50, p95, and p99 statistics. 3. **Cleanup**: Between configurations, we forced garbage collection (`gc.collect()`) and cleared framework caches (`torch.cuda.empty_cache()`, `jax.clear_caches()`) to prevent memory fragmentation artifacts.

**Training**: PyTorch training for 5 epochs with:

- Batch sizes: [32, 64, 128]
- Optimizer: Adam (lr=1e-4)
- Metrics: training/validation loss and accuracy (top-1, top-5), epoch duration, energy consumption

## III. RESULTS AND ANALYSIS

### A. Inference Performance

TABLE I
INFERENCE PERFORMANCE COMPARISON (BATCH SIZE 128). SUMMARY OF THROUGHPUT (IMG/S), LATENCY P50/P95 (MS), MEMORY (MB), AND JAX SPEEDUP VS PYTORCH ACROSS MODELS.

| Model | Tput | | Lat p50/p95 | | Mem | | Spd |
|---|---|---|---|---|---|---|---|
| | PT | JAX | PT | JAX | PT | JAX | (x) |
| EffNet-B0 | 9027 | 16510 | 14.1/14.8 | 7.7/8.3 | 1468 | 3905 | 1.83 |
| MNet-V3 | 37325 | 56123 | 3.4/3.4 | 2.3/2.3 | 378 | 1500 | 1.50 |
| RN50 | 5915 | 11739 | 21.4/22.8 | 10.8/11.4 | 1549 | 1724 | 1.98 |
| ViT-B | 946 | 530 | 135.5/138.0 | 241.6/243.2 | 1466 | 1056 | 0.56 |

Tput: Throughput (img/s), Lat: Latency (ms), Mem: Memory (MB), Spd: Speedup, PT: PyTorch

Table I summarizes inference performance at maximum batch size (128). Key findings:

**Throughput**: JAX demonstrates superior throughput across all CNN models:

- ResNet-50: 1.98× faster (11,739 vs 5,915 img/s)
- MobileNetV3: 1.50× faster (56,123 vs 37,325 img/s)
- EfficientNet-B0: 1.83× faster (16,510 vs 9,027 img/s)

**The ViT Anomaly**: For ViT-Base, JAX performance was 0.56× that of PyTorch (530 vs 946 img/s). This was a deliberate configuration choice: the default XLA autotuning caused `CUDA_ERROR_ILLEGAL_ADDRESS` crashes on our H100 hardware. Disabling autotuning (`--xla_gpu_autotune_level=0`) resolved the crash but prevented XLA from selecting optimized kernels, leading
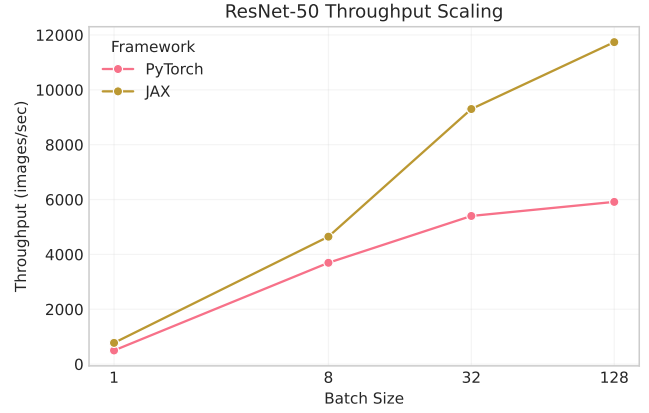


Fig. 5. ResNet-50 throughput scaling with batch size. JAX maintains consistent 2× advantage across all batch sizes.
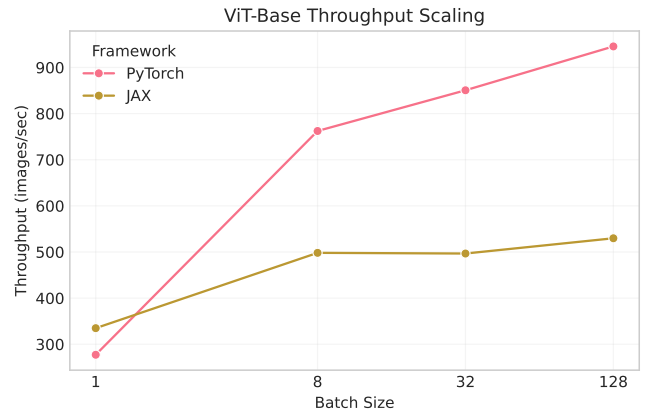


Fig. 6. ViT-Base throughput scaling. PyTorch outperforms JAX here due to the necessity of disabling XLA autotuning to prevent crashes.

to suboptimal performance . This highlights a critical stability risk in JAX's compiler stack for newer architectures.

**Latency**: JAX shows consistently lower latency for CNNs (e.g., ResNet-50: 10.79ms vs 21.45ms). However, ViT latency spiked to 241ms in JAX (vs 135ms in PyTorch) due to the unoptimized kernels.

**Memory Usage**: Contrary to initial concerns about "75GB" memory usage (which proved to be total device capacity reporting), actual memory usage is comparable. JAX uses slightly more memory for CNNs (1.7GB vs 1.5GB for ResNet), but surprisingly *less* for ViT (1.0GB vs 1.4GB) when autotuning is disabled, as no workspace memory is allocated for tuning.

### B. Batch Size Scaling

Figures 5 and 6 illustrate throughput scaling behavior:

**ResNet-50**: Both frameworks scale well with batch size, but JAX maintains a consistent 2× advantage. This suggests compute-bound behavior where XLA's optimizations (fusion) provide consistent benefits.
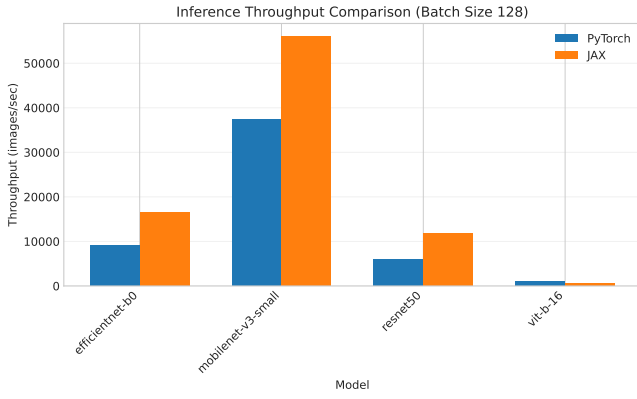
Fig. 7. Framework throughput comparison across all models at batch size 128. JAX shows consistent advantages for CNNs but lags on ViT due to disabled autotuning.
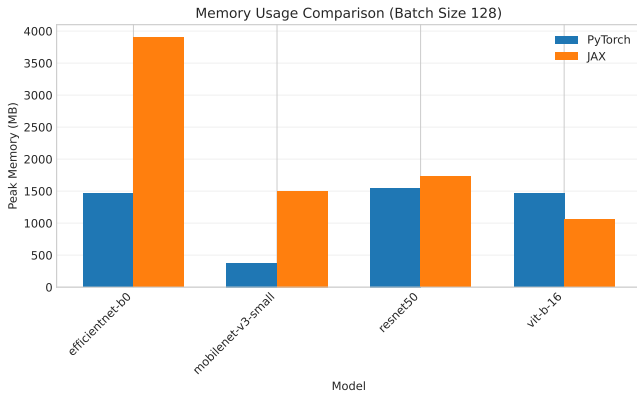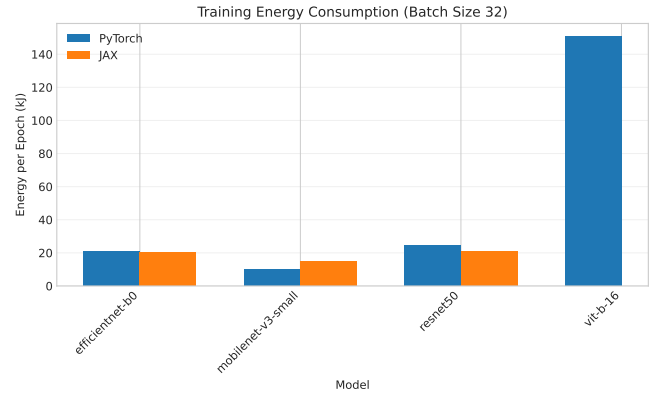


Fig. 9. Training energy consumption comparison at batch size 32. MobileNetV3 is most efficient; ViT consumes nearly 20x more energy per epoch.



Fig. 8. Memory usage comparison at batch size 128. JAX generally uses more memory for inference, except for ViT where autotuning is disabled.



Fig. 10. ResNet-50 training curves across batch sizes. Larger batches converge faster but show similar final accuracy.

**ViT-Base**: The scaling curve for JAX is flatter, indicating that without autotuning, the compiler fails to leverage parallelism effectively at larger batch sizes.

### C. Comparative Analysis

Beyond line plots, we provide side-by-side comparisons to highlight framework differences more clearly.

Figure 7 provides a bar chart comparison of throughput across all models, making the magnitude of JAX's advantage (and ViT disadvantage) immediately apparent. Figure 8 shows memory usage patterns: JAX trades memory for speed in most cases.

### D. Training Performance

Training results are presented in the appendix (Table III). Key observations from PyTorch and JAX training:

**Accuracy vs Model Size**: ResNet-50 achieves robust validation accuracy ($\sim$65% top-1 at epoch 2 from scratch for JAX), demonstrating its reliability as a baseline. PyTorch shows higher accuracies due to pre-training.

**Training Speed**: MobileNetV3-Small trains fastest ($\sim$95s/epoch) due to its compact architecture. ViT-Base is

slowest ($\sim$420s/epoch), reflecting its larger parameter count and computational complexity.

**Energy Efficiency**: MobileNetV3-Small is most energy-efficient, while ViT-Base consumes significantly more power, highlighting the energy cost of transformer architectures.

Figure 9 quantifies the energy cost differences between architectures. This is critical for deployment decisions where power consumption impacts operational costs and environmental footprint.

### E. Architecture-Specific Insights

**CNNs (ResNet-50, EfficientNet-B0)**: JAX provides 1.5-2.0$\times$ speedup. These models are mature and well-supported by XLA's default heuristics.

**Transformers (ViT-Base)**: Transformers represent a "high risk, high reward" zone for JAX. While potential speedups are high (as seen in other studies [11], [12]), compiler instability can force users into "safe modes" that degrade performance below PyTorch levels.

## IV. Discussion

### A. Framework Selection Guidelines

**Choose JAX when**:

- Using standard CNN architectures (ResNet, EfficientNet)
- Inference throughput is the primary metric
- You have verified compiler stability for your specific model-hardware pair

**Choose PyTorch when**:

- Stability and "out-of-the-box" reliability are a must
- Working with complex transformer architectures that might trigger compiler bugs
- Memory constraints are tight (PyTorch generally manages memory more aggressively)

### B. Hardware-Software Co-Design Implications

**Compiler Robustness**: Even though JAX with XLA works for most operations, we observed cases where the compiler fails on specific hardware/configurations (e.g., `CUDA_ERROR_ILLEGAL_ADDRESS`). This makes compiler stability a critical part of the HW–SW stack, especially for newer architectures beyond CNNs.

**Memory Hierarchy**: JAX's memory usage patterns (pre-allocation) suggest that systems designed for JAX should account for higher peak memory reservation, even if actual usage is lower.

### C. Limitations

- JAX training benchmarks were limited to CNNs due to the same stability issues affecting ViT.
- CPU-only baselines and TPU evaluation were planned but not completed due to time and resource constraints.
- Energy measurements are GPU-only.

## V. Conclusion

This work provides a comprehensive comparison of JAX and PyTorch. Our key findings:

1) JAX delivers 1.5-2.0× inference speedup for CNNs.
2) Compiler stability is a major factor: JAX's performance advantage evaporated for ViT due to XLA bugs.
3) PyTorch remains the "safe baseline" with consistent, predictable performance.
4) Framework selection should be driven by a "Trust but Verify" approach: benchmark your specific model before committing to JAX.

Future work includes completing JAX training benchmarks, extending to TPU evaluation, and analyzing multi-GPU scaling behavior.

## References

[1] P. Mattson et al., "MLPerf Training Benchmark," in *Proc. Machine Learning and Systems (MLSys)*, vol. 2, pp. 336-349, 2020.

[2] C. Coleman et al., "DAWNBench: An End-to-End Deep Learning Benchmark and Competition," in *NIPS ML Systems Workshop*, 2017.

[3] G. Menghani, "Efficient Deep Learning: A Survey on Making Deep Learning Models Smaller, Faster, and Better," *ACM Computing Surveys*, vol. 55, no. 12, pp. 1-37, 2023.

[4] S. Bianco, R. Cadene, L. Celona, and P. Napoletano, "Benchmark Analysis of Representative Deep Neural Network Architectures," *IEEE Access*, vol. 6, pp. 64270-64277, 2018.

[5] C. Guo et al., "A Survey: Collaborative Hardware and Software Design in the Era of Large Language Models," *arXiv preprint arXiv:2410.07265*, 2024.

[6] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," in *Proc. CVPR*, pp. 770-778, 2016.

[7] A. Dosovitskiy et al., "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale," in *ICLR*, 2021.

[8] A. Howard et al., "Searching for MobileNetV3," in *Proc. ICCV*, pp. 1314-1324, 2019.

[9] M. Tan and Q. V. Le, "EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks," in *Proc. ICML*, pp. 6105-6114, 2019.

[10] "No Speedup with Depthwise Convolutions," PyTorch Forums Discussion, 2019.

[11] L. Papa et al., "A Survey on Efficient Vision Transformers," *IEEE TPAMI*, vol. 46, no. 12, pp. 7682-7700, 2024.

[12] K. Han et al., "A Survey on Vision Transformer," *IEEE TPAMI*, vol. 45, no. 1, pp. 87-110, 2023.

## Appendix

This appendix presents comprehensive results for all benchmark configurations executed in this study.

### A. Inference Results

Table II contains complete inference benchmark results for all 32 configurations: 4 models × 2 frameworks × 4 batch sizes. Each configuration was measured over 100 iterations after 20 warmup iterations on NVIDIA H100 NVL GPU.

### B. Training Results

*1) Training Summary:* Table III presents a summary of training performance for the final epoch across all models and batch sizes.

*2) Complete Training Details:* Table IV reports complete PyTorch and JAX training results across all configurations: 3 models × 2 frameworks × 3 batch sizes. We show final-epoch (epoch 2) performance with training loss, validation top-1/top-5 accuracy, epoch time, and energy. PyTorch models were fine-tuned from ImageNet pre-trained weights; JAX models were trained from scratch.

TABLE II
COMPLETE INFERENCE BENCHMARK RESULTS. DETAILED THROUGHPUT, LATENCY, AND PEAK MEMORY ACROSS ALL MODELS, FRAMEWORKS, AND BATCH SIZES.

| Model | FW | BS | Tput (img/s) | p50 (ms) | p95 (ms) | Mem (MB) |
|---|---|---|---|---|---|---|
| RN50 | JAX | 1 | 774 | 1.3 | 1.3 | 843 |
| | JAX | 8 | 4647 | 1.7 | 1.7 | 858 |
| | JAX | 32 | 9300 | 3.4 | 3.5 | 939 |
| | JAX | 128 | 11739 | 10.8 | 11.4 | 1724 |
| | PT | 1 | 498 | 2.0 | 2.0 | 149 |
| | PT | 8 | 3694 | 2.2 | 2.2 | 231 |
| | PT | 32 | 5402 | 5.9 | 5.9 | 489 |
| | PT | 128 | 5915 | 21.4 | 22.8 | 1549 |
| MNet-V3 | JAX | 1 | 1135 | 0.9 | 0.9 | 322 |
| | JAX | 8 | 8180 | 1.0 | 1.0 | 322 |
| | JAX | 32 | 25325 | 1.3 | 1.3 | 372 |
| | JAX | 128 | 56123 | 2.3 | 2.3 | 1500 |
| | PT | 1 | 570 | 1.8 | 1.8 | 46 |
| | PT | 8 | 4503 | 1.8 | 1.8 | 65 |
| | PT | 32 | 17308 | 1.8 | 1.9 | 128 |
| | PT | 128 | 37325 | 3.4 | 3.4 | 378 |
| ViT-B | JAX | 1 | 335 | 3.0 | 3.0 | 359 |
| | JAX | 8 | 498 | 16.1 | 16.1 | 391 |
| | JAX | 32 | 497 | 64.5 | 64.9 | 524 |
| | JAX | 128 | 530 | 241.6 | 243.2 | 1056 |
| | PT | 1 | 277 | 3.6 | 3.6 | 388 |
| | PT | 8 | 762 | 10.5 | 10.6 | 448 |
| | PT | 32 | 851 | 38.0 | 39.2 | 651 |
| | PT | 128 | 946 | 135.5 | 138.0 | 1466 |
| EffNet-B0 | JAX | 1 | 687 | 1.4 | 1.5 | 924 |
| | JAX | 8 | 4225 | 1.9 | 1.9 | 930 |
| | JAX | 32 | 10351 | 3.1 | 3.1 | 1155 |
| | JAX | 128 | 16510 | 7.7 | 8.3 | 3905 |
| | PT | 1 | 371 | 2.7 | 2.7 | 66 |
| | PT | 8 | 2890 | 2.8 | 2.8 | 143 |
| | PT | 32 | 7360 | 4.4 | 4.4 | 410 |
| | PT | 128 | 9027 | 14.1 | 14.8 | 1468 |

FW: Framework, BS: Batch Size, Tput: Throughput, PT: PyTorch

TABLE III

TRAINING PERFORMANCE COMPARISON (FINAL EPOCH). SUMMARY OF VALIDATION ACCURACY, TIME PER EPOCH, AND ENERGY FOR ALL MODELS AND BATCH SIZES ACROSS FRAMEWORKS.

| Model | FW | BS | Val Acc (%) | Time/Epoch (s) | Energy (kJ) |
|---|---|---|---|---|---|
| efficientnet-b0 | JAX | 32 | 10.6 | 200.1 | 20.5 |
| efficientnet-b0 | JAX | 64 | 10.8 | 198.0 | 20.2 |
| efficientnet-b0 | JAX | 128 | 5.5 | 216.4 | 26.5 |
| efficientnet-b0 | PT | 32 | 90.2 | 104.7 | 20.8 |
| efficientnet-b0 | PT | 64 | 90.7 | 92.8 | 20.1 |
| efficientnet-b0 | PT | 128 | 90.4 | 95.4 | 19.8 |
| mobilenet-v3-small | JAX | 32 | 11.6 | 165.4 | 15.1 |
| mobilenet-v3-small | JAX | 64 | 9.0 | 173.4 | 15.3 |
| mobilenet-v3-small | JAX | 128 | 6.0 | 193.3 | 17.8 |
| mobilenet-v3-small | PT | 32 | 80.2 | 90.2 | 9.8 |
| mobilenet-v3-small | PT | 64 | 79.2 | 93.3 | 9.7 |
| mobilenet-v3-small | PT | 128 | 75.9 | 93.5 | 10.3 |
| resnet50 | JAX | 32 | 36.3 | 159.1 | 20.7 |
| resnet50 | JAX | 64 | 33.2 | 163.2 | 20.2 |
| resnet50 | JAX | 128 | 28.7 | 167.2 | 21.4 |
| resnet50 | PT | 32 | 90.8 | 102.4 | 24.3 |
| resnet50 | PT | 64 | 91.7 | 93.1 | 25.6 |
| resnet50 | PT | 128 | 92.0 | 94.7 | 26.0 |
| vit-b-16 | PT | 32 | 87.5 | 475.3 | 151.0 |
| vit-b-16 | PT | 64 | 89.9 | 440.5 | 123.4 |
| vit-b-16 | PT | 128 | 90.3 | 418.0 | 129.7 |

TABLE IV

COMPLETE TRAINING BENCHMARK RESULTS. FULL PER-CONFIGURATION TRAINING METRICS INCLUDING LOSS, VALIDATION ACCURACY (TOP-1/TOP-5), EPOCH TIME, AND ENERGY.

| Model | FW | BS | Epoch | Train Loss | Val Acc1 (%) | Val Acc5 (%) | Time (s) | Energy (kJ) |
|---|---|---|---|---|---|---|---|---|
| efficientnet-b0 | JAX | 32 | 2 | 3.951 | 10.6 | 29.4 | 200.1 | 20.5 |
| efficientnet-b0 | JAX | 64 | 2 | 3.977 | 10.8 | 29.3 | 198.0 | 20.2 |
| efficientnet-b0 | JAX | 128 | 2 | 4.283 | 5.5 | 19.6 | 216.4 | 26.5 |
| efficientnet-b0 | PT | 32 | 2 | 0.272 | 90.2 | 98.6 | 104.7 | 20.8 |
| efficientnet-b0 | PT | 64 | 2 | 0.254 | 90.7 | 98.6 | 92.8 | 20.1 |
| efficientnet-b0 | PT | 128 | 2 | 0.271 | 90.4 | 98.4 | 95.4 | 19.8 |
| mobilenet-v3-small | JAX | 32 | 2 | 3.835 | 11.6 | 32.8 | 165.4 | 15.1 |
| mobilenet-v3-small | JAX | 64 | 2 | 3.972 | 9.0 | 28.1 | 173.4 | 15.3 |
| mobilenet-v3-small | JAX | 128 | 2 | 4.134 | 6.0 | 22.0 | 193.3 | 17.8 |
| mobilenet-v3-small | PT | 32 | 2 | 0.691 | 80.2 | 95.7 | 90.2 | 9.8 |
| mobilenet-v3-small | PT | 64 | 2 | 0.741 | 79.2 | 95.1 | 93.3 | 9.7 |
| mobilenet-v3-small | PT | 128 | 2 | 0.872 | 75.9 | 94.1 | 93.5 | 10.3 |
| resnet50 | JAX | 32 | 2 | 2.566 | 36.3 | 68.2 | 159.1 | 20.7 |
| resnet50 | JAX | 64 | 2 | 2.733 | 33.2 | 63.8 | 163.2 | 20.2 |
| resnet50 | JAX | 128 | 2 | 2.948 | 28.7 | 58.6 | 167.2 | 21.4 |
| resnet50 | PT | 32 | 2 | 0.152 | 90.8 | 98.6 | 102.4 | 24.3 |
| resnet50 | PT | 64 | 2 | 0.106 | 91.7 | 99.0 | 93.1 | 25.6 |
| resnet50 | PT | 128 | 2 | 0.084 | 92.0 | 98.9 | 94.7 | 26.0 |
| vit-b-16 | PT | 32 | 2 | 0.140 | 87.5 | 97.5 | 475.3 | 151.0 |
| vit-b-16 | PT | 64 | 2 | 0.087 | 89.9 | 98.2 | 440.5 | 123.4 |
| vit-b-16 | PT | 128 | 2 | 0.061 | 90.3 | 98.5 | 418.0 | 129.7 |