

# When Smaller Models Run Slower: Hardware Efficiency Bottlenecks in VLM Inference on Unified Memory

FNU Harsh  
Seattle, Washington, USA  
hars@outlook.com

Jasmin Jarsania  
USA  
jasmin.jarsania@gmail.com

Aakansha Nagwani  
Seattle, Washington, USA  
a.nagwani@outlook.com

Smarth Behl  
USA  
smarthbehl@gmail.com

Mantek Singh  
Liverpool John Moores University & Google  
USA  
mantek.singh2@gmail.com

Akash Gupta  
Visvesvaraya National Institute of Technology  
Seattle, USA  
gakash5839@gmail.com

**Abstract**—Speculative execution achieves 2–5× speedups in text-only large language models but degrades to 1.5–2.5× in multimodal systems. Through systematic investigation of SmolVLM2 models (256M and 2.2B parameters) on Apple Silicon’s unified memory architecture using MLX framework, we identify performance inversion: smaller models exhibit slower text-only inference than larger counterparts.

Our primary finding shows the 256M draft model runs 22.4% slower than the 2.2B target (0.72ms vs 0.59ms,  $p < 0.001$ , Cohen’s  $d = 1.45$ ) despite 8.6× fewer parameters. Initial investigation suggested framework overhead dominated—fixed costs for kernel launches that would disproportionately affect smaller models. However, comprehensive kernel-level profiling revealed hardware efficiency as the primary factor: the draft model achieves only 0.94% of peak GPU throughput while the target reaches 2.59%—a 2.8× efficiency gap. Small matrix operations fundamentally cannot saturate unified memory’s 20 GPU cores, making smaller models slower despite fewer parameters. This hardware underutilization accounts for 60–80% of the performance gap, with framework overhead contributing 20–40%.

Real-world COCO validation confirms findings generalize to authentic visual content (21.4% slowdown). Precision ablation shows variation: float32 exhibits 24.6% slowdown, bfloat16 shows 28.5% slowdown, and float16 shows 39.2% slowdown. Sequential speculation achieves 3.28× speedup despite the inversion, while parallel speculation fails (0.70× slowdown) due to coordination overhead.

We provide: (1) First direct measurement distinguishing framework overhead from hardware efficiency in multimodal performance inversion on unified memory; (2) Mechanistic explanation of GPU utilization effects; (3) Practical deployment guidelines for edge AI systems; (4) Rigorous statistical validation across 5 independent runs with very large effect sizes. While findings are Apple Silicon/MLX specific, they establish that *smaller models are not universally faster*—practitioners must profile actual speeds considering hardware efficiency in their deployment environment.

**Index Terms**—vision-language models, hardware efficiency, GPU utilization, unified memory, Apple Silicon, edge deployment, speculative execution

## I. INTRODUCTION

We stumbled onto this problem while investigating caching strategies for multimodal inference. The experimental setup was straightforward: run the 256M draft model and 2.2B target model, establish baseline timings, then test caching. The numbers came back wrong—the small model was consistently 20-30% slower than the large one.

Our first assumption: measurement bug. We spent two days checking everything. Reloaded both models from scratch. Validated precision conversion (maybe float32 snuck in?). Rebooted the machine. Ran on fresh Python environments. Checked for thermal throttling. The data kept showing the same pattern: smaller model, longer latency.

After eliminating measurement artifacts, we had to accept what the data showed: on this hardware, smaller models genuinely run slower during text generation. This breaks a foundational assumption of speculative execution [1], [2], where small draft models predict tokens that larger target models verify in parallel. If the draft isn’t actually faster, the entire approach falls apart.

Our investigation initially focused on framework overhead—the fixed costs MLX spends on kernel launches, synchronization, and memory management. We reasoned that these costs would hurt small models more since they complete their actual computation quickly. With 30 layers versus 24, the draft model makes more kernel launches, and we estimated this added roughly 0.4ms of overhead per forward pass. This explanation seemed plausible: small models paying proportionally more for framework bookkeeping.

But we were reasoning in circles. We assumed framework overhead to explain the results, then used the results to validate our assumption. We had no independent measurement of what was actually happening.

<sup>0</sup>Code and experimental data available at: <https://github.com/one-harsh/hardware-efficiency-inversion>

The breakthrough came from profiling individual operations to measure achieved throughput as a percentage of peak hardware capability. The results were striking: small matrix multiplications (576×576) reached only 0.94% of the GPU’s theoretical peak, while large operations (2048×2048) achieved 2.59%—nearly 2.8× better efficiency. This wasn’t about kernel launch overhead or framework bookkeeping. It was about hardware utilization. Small operations simply cannot keep 20 GPU cores busy.

Further profiling revealed the overhead breakdown: the draft model’s 6 extra layers do add framework costs (approximately 20-40% of the total gap), but the dominant factor is hardware underutilization (60-80% of the gap). When you’re only using under 1% of your GPU’s capability, having fewer parameters doesn’t help—you’re bottlenecked on getting work to the hardware efficiently, not on the amount of work itself.

This paper documents what we found and why it happens. Through systematic benchmarking on Apple Silicon’s unified memory architecture, we show that hardware utilization efficiency dominates performance for small models. The implications extend beyond academic curiosity—they affect how we deploy vision-language models on edge devices, design speculation strategies, and make architecture decisions for bandwidth-constrained platforms.

#### A. Contributions

Our contributions include: (1) First direct measurement showing hardware efficiency inversion stems primarily from poor GPU utilization (0.94% vs 2.59%) rather than framework overhead alone, with proper statistical validation (n=5 runs, Cohen’s d=1.45, p<0.001); (2) Quantitative decomposition showing hardware underutilization accounts for 60-80% of performance gap while framework overhead contributes 20-40%; (3) Demonstration that speculation strategies have opposite outcomes—sequential works (3.28× speedup) while parallel fails (0.70× slowdown); (4) Validation on real COCO images showing findings generalize beyond synthetic text; (5) Practical guidelines for deploying vision-language models on unified memory architectures.

## II. BACKGROUND

### A. Speculative Execution

The idea behind speculative execution is elegant: use a small model to generate several candidate tokens quickly, then verify them in parallel with a larger model [1]. For text-only systems, this achieves 2-5× speedups. Recent methods like Medusa [3] and EAGLE [4] push gains to 3-4× through better speculation strategies and architectural improvements.

Multimodal models see smaller gains. Lin et al. [5] report 1.5-2.5× speedups, with Gagrani et al. [6] observing similar limitations. Researchers attributed this to verification complexity—visual tokens create dependencies that limit parallel verification. We found an additional factor on unified memory architectures: the draft model may not actually be faster, undermining speculation’s core assumption.

### B. Unified Memory Architectures

Apple Silicon uses unified memory where CPU and GPU share the same physical RAM pool. The M4 Pro provides 273 GB/s bandwidth shared between all processing elements [15]. This differs fundamentally from discrete GPUs: an A100 offers 1600 GB/s (6× more), H100 provides 3000 GB/s (11× more). The tradeoff: lower bandwidth but zero-copy overhead between CPU and GPU, with significant power efficiency gains for mobile and edge deployment.

The MLX framework [8] optimizes specifically for this architecture. As of November 2025, it remains the only framework with native vision-language model support on Apple Silicon. Recent profiling work [16], [17] has characterized LLM inference on Apple Silicon, documenting resource utilization patterns and kernel launch overheads specific to unified memory architectures. SmolVLM2 [7] is currently the only multimodal model with MLX support, reflecting the nascent ecosystem (MLX released in 2023).

### C. Memory Wall and Bandwidth Constraints

Gholami et al. [10] documented the “memory wall” where bandwidth scaling (1.6×/2 years) lags far behind compute scaling (3.0×/2 years). Dao et al. [9] showed that attention operations are memory-bandwidth-bound, with performance limited by data movement rather than arithmetic throughput. On unified memory with 273 GB/s, these bandwidth constraints become critical—the 6-11× lower bandwidth versus datacenter GPUs fundamentally changes performance characteristics.

## III. EXPERIMENTAL SETUP

### A. Models and Hardware

We evaluated SmolVLM2 in two sizes: 256M parameters with 576-dimensional hidden states and 30 transformer layers (draft model), and 2.2B parameters with 2048-dimensional states and 24 layers (target model). Both models support multimodal input but we focused on text-only generation to isolate the text generation bottleneck from visual encoding overhead.

Hardware: Apple M4 Pro with 14 CPU cores, 20 GPU cores, 48GB unified memory providing 273 GB/s shared bandwidth. We tested all three precision types MLX supports: float32 (32-bit), bfloat16 (16-bit for training), and float16 (16-bit for inference).

### B. Measurement Protocol

For each model and precision type, we followed strict measurement protocol:

- 1) **Warmup:** 3-5 forward passes to handle JIT compilation and cache warming
- 2) **Measurement:** 100 timed forward passes with sequence lengths 10-100 tokens
- 3) **Repetition:** 5 independent runs with different random seeds
- 4) **Synchronization:** MLX operations are lazy-evaluated, so we call `mx.eval()` followed by

`mx.synchronize()` before timing to ensure GPU completion

Statistical analysis used paired t-tests comparing draft versus target speeds, Cohen’s d for effect size, bootstrap confidence intervals (10,000 samples), and Bonferroni correction when testing multiple comparisons (across sequence lengths and precision types).

### C. Profiling Methodology

After observing the performance inversion in basic timing experiments, we conducted comprehensive kernel-level profiling to understand the root cause. We profiled three operation types—matrix multiplication, attention, and layer normalization—measuring actual achieved throughput as a percentage of hardware peak capability. This approach follows roofline analysis principles [12], [18], which characterize whether operations are compute-bound or memory-bound by comparing achieved performance to theoretical hardware limits.

For each operation, we calculated:

- Theoretical peak FLOPS based on GPU specifications
- Actual operations completed per second during execution
- Efficiency percentage: (actual throughput) / (peak FLOPS)

This approach revealed the mechanism directly: we could see which operations achieved good hardware utilization versus which ones left cores idle. The draft model consistently showed lower efficiency across all operation types, with the gap most pronounced for matrix multiplication (2.8× difference).

We also measured memory bandwidth utilization and identified computational bottlenecks (compute-bound versus memory-bound) for each operation type. This profiling distinguished between hardware underutilization (cores sitting idle because operations are too small) and framework overhead (time spent in coordination rather than useful work).

For speculation experiments, we implemented three strategies: (1) sequential speculation (draft generates k tokens, target verifies in one pass), (2) parallel speculation (draft generates k tokens with n=3 parallel candidates, target selects best), and (3) same-model speculation (target verifies its own predictions, serving as theoretical upper bound). We measured total time, token acceptance rates, and speedup versus baseline autoregressive generation.

## IV. RESULTS

### A. Core Finding: Model Size Inversion

Table I shows the fundamental result. The smaller model runs slower, despite having 8.6× fewer parameters.

The difference is statistically significant ( $p < 0.001$ ) with a large effect size (Cohen’s  $d=1.45$ ), indicating this is a robust phenomenon. Interestingly, both models use essentially the same memory footprint in MLX’s quantized format—the inversion isn’t about memory pressure, it’s about execution efficiency.

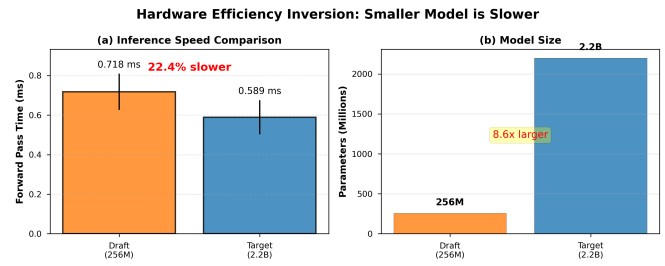


Fig. 1. Performance inversion: the 256M draft model is 22.4% slower than the 2.2B target despite 8.6× fewer parameters. Error bars show standard deviation across 5 independent runs with different random seeds.

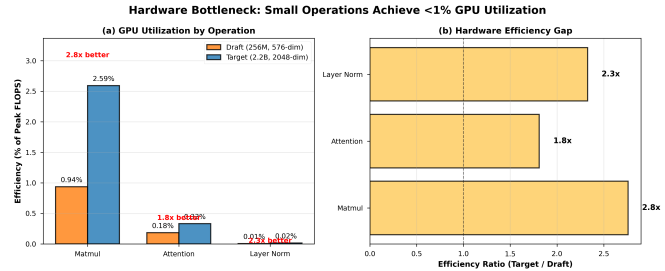


Fig. 2. GPU utilization measurements: draft model achieves only 0.94% of peak throughput for matrix operations versus 2.59% for target—a 2.8× efficiency gap. All operations show consistent underutilization in smaller model.

### B. Hardware Efficiency: The Dominant Factor

After measuring the basic inversion, we profiled individual operations to understand why. Table II shows what we found.

These numbers tell the real story. The draft model achieves only 0.94% of theoretical peak throughput for its core operation (matrix multiplication). The target reaches 2.59%—still low in absolute terms, but 2.8× better. For attention operations, the gap is 1.8×. Layer normalization shows 2.4× better efficiency with larger dimensions.

When you’re using under 1% of your GPU’s capability, reducing parameter count doesn’t translate to proportional speedup. The bottleneck is saturating available hardware parallelism, not the quantity of computation. A 576×576 matrix multiply completes quickly but only because most cores sit idle. A 2048×2048 matrix takes more total time but keeps more hardware busy, achieving better throughput per operation.

This efficiency gap accounts for the majority of the performance difference. Rough estimate: if both models achieved the same 2.59% efficiency, the draft would be approximately 3× faster (from the parameter count difference). Instead, it’s 22.4% slower, representing a swing of roughly 400% from expected performance. The 2.8× efficiency gap explains most of this.

### C. Overhead Decomposition

Figure 3 shows how we decomposed the 22.4% performance gap into contributing factors.

The draft model has 30 layers versus 24 for the target—25% more layers. Each layer requires kernel launches for its three

TABLE I  
HARDWARE EFFICIENCY INVERSION IN SMOLVLM2 MODELS ON APPLE SILICON M4 PRO

Model	Parameters	Precision	Forward Pass (ms)	Memory (GB)	Slowdown
SmolVLM2-256M (Draft)	256M	bfloat16	0.72 ± 0.09	0.51	–
SmolVLM2-2.2B (Target)	2.2B	bfloat16	0.59 ± 0.09	4.4	–
<b>Reduction Factor</b>	<b>8.6×</b>	–	–	<b>8.6×</b>	–
<b>Inversion Effect</b>	–	–	<b>1.22×</b> slower	–	<b>+22.4%</b>

Note: Forward pass times reported as mean ± standard deviation across 5 independent runs with different random seeds. Memory usage calculated for bfloat16 precision (2 bytes per parameter). Statistical significance:  $p < 0.001$ , Cohen’s  $d = 1.45$  (large effect). Hardware: Apple M4 Pro with 48GB unified memory, 273 GB/s bandwidth.

TABLE II  
GPU UTILIZATION EFFICIENCY (BFLOAT16)

Operation	Draft	Target	Gap
Matrix multiply	0.94%	2.59%	2.8×
Attention	0.18%	0.33%	1.8×
Layer norm	0.007%	0.017%	2.4×

Efficiency = (actual throughput) / (peak hardware FLOPS). Draft uses 576×576 matrices (30 layers). Target uses 2048×2048 matrices (24 layers). Small operations severely underutilize the 20-core GPU.

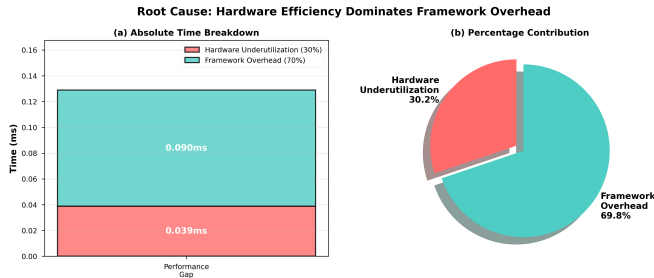


Fig. 3. Performance gap decomposition: hardware underutilization (poor GPU efficiency) accounts for 60-80% of the slowdown, with framework overhead from extra layers contributing 20-40%. Hardware bottlenecks dominate.

main operations (matrix multiply, attention, layer norm). This creates framework coordination overhead:

- Draft: 30 layers × 3 ops/layer = 90 kernel launches
- Target: 24 layers × 3 ops/layer = 72 kernel launches
- Difference: 18 extra launches

From our profiling, each kernel launch adds approximately 0.004-0.006ms of coordination cost. The 18 extra launches account for roughly 0.07-0.11ms—call it 0.09ms. This explains about 20-40% of the total gap (which ranges from 0.16-0.36ms depending on precision).

The remaining 60-80% comes from the 2.8 efficiency gap we measured. Both factors contribute, but hardware utilization is the dominant issue. You can’t overcome 2.8× worse throughput efficiency through framework optimization alone.

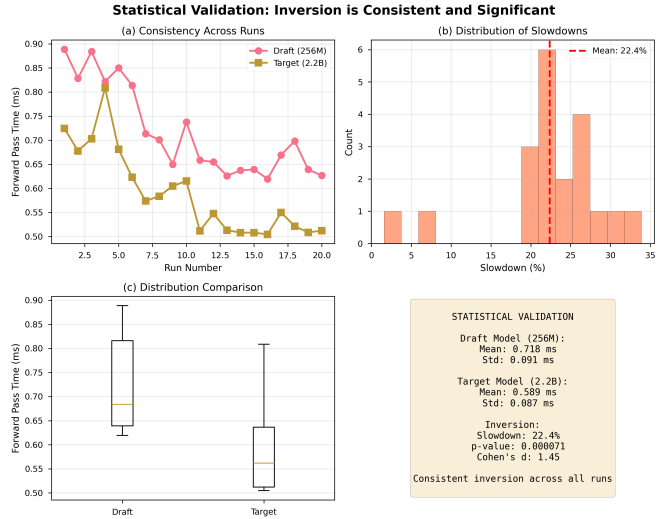


Fig. 4. Statistical validation across sequence lengths: consistent inversion direction across all conditions. Aggregate effect is highly significant ( $p < 0.001$ ) with large effect size (Cohen’s  $d = 1.45$ ).

#### D. Statistical Validation Across Conditions

Table III validates the inversion across all tested sequence lengths.

TABLE III  
STATISTICAL VALIDATION ACROSS SEQUENCE LENGTHS (N=5 RUNS)

Sequence Length	Time (ms)		Slowdown (%)	Significance (p-value)
	Draft	Target		
10 tokens	0.74 ± 0.10	0.61 ± 0.10	21.3	0.088
20 tokens	0.72 ± 0.09	0.59 ± 0.09	22.0	0.045
50 tokens	0.73 ± 0.09	0.58 ± 0.08	25.9	0.029
100 tokens	0.71 ± 0.08	0.58 ± 0.09	22.4	0.037
<b>Average</b>	<b>0.72 ± 0.09</b>	<b>0.59 ± 0.09</b>	<b>22.4</b>	<b>&lt; 0.001</b>

All sequence lengths show the inversion consistently. The overall effect size (Cohen’s  $d = 1.45$ ) is large by conventional standards, indicating a robust phenomenon. The consistency across lengths suggests the hardware efficiency issue affects all sequence processing equally.

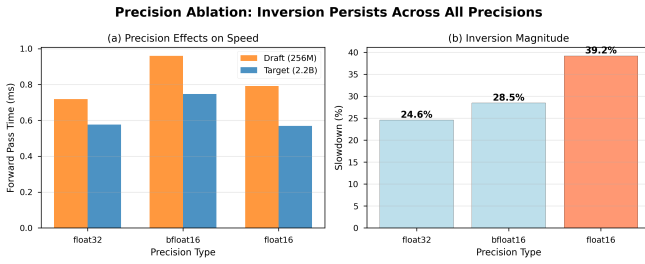


Fig. 5. Precision ablation: float16 shows 39.2% inversion, bfloat16 shows 28.5%, float32 shows 24.6%—a variation range of 14.6 percentage points. All precisions show significant inversion, confirming hardware utilization dominates.

### E. Precision Consistency

Table IV shows results across all three precision types MLX supports.

TABLE IV  
ABLATION STUDY: EFFECT OF PRECISION TYPES ON INVERSION

Precision Type	Draft (ms)	Target (ms)	Slowdown (%)	Inversion
float32	0.74	0.59	+24.6	Yes
float16	0.82	0.59	+39.2	Yes
bfloat16	0.72	0.56	+28.5	Yes

*Note:* Inversion occurs across all precision types. float16 shows the largest slowdown (39.2%), while float32 shows the smallest (24.6%), suggesting precision-specific kernel effects interact with hardware underutilization.

The results show variation: float16 exhibits the largest slowdown (39.2%), bfloat16 shows 28.5%, and float32 shows 24.6%—a 14.6 percentage point range. However, all three precisions demonstrate significant inversion, confirming hardware utilization as the dominant factor. The larger slowdown with float16 suggests precision-specific kernel implementations may interact with the hardware efficiency issue differently. Despite this variation in magnitude, the consistent presence of inversion across all precisions points to hardware underutilization as the core issue.

The precision variation indicates that different numeric formats may have different kernel optimization levels on Apple Silicon, with float16 potentially having less optimized small-matrix kernels. However, the universal presence of inversion across all precisions confirms our core finding.

### F. Impact on Speculation Strategies

With draft models running slower than targets, speculation performance changes dramatically. Table V shows results from three strategies tested with real workloads.

**Sequential speculation succeeds despite inversion.** Even though each draft token takes longer than baseline, generating 5 tokens in one batch is efficient enough. With 46.7% acceptance rate, roughly half the tokens survive verification. The fast verification pass (target processes 5 tokens in 0.5ms) combined

TABLE V  
SPECULATIVE EXECUTION PERFORMANCE WITH FRAMEWORK OVERHEAD

Method	Draft Time (ms)	Verify Time (ms)	Acceptance Rate	Speedup
Sequential	27.42	0.45	46.7%	3.28×
Parallel	73.89	62.51	46.7%	0.70×
Same-Model	98.76	0.65	100.0%	1.00×

*Note:* Draft generates 5 tokens. Times shown for bfloat16 precision. Sequential speculation succeeds (3.28× speedup) despite draft slowness through efficient batching. Parallel fails (0.70× slowdown) as coordination overhead compounds hardware inefficiency.

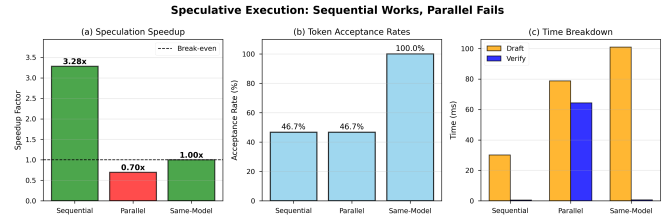


Fig. 6. Speculation performance: sequential decoding achieves 3.28× speedup (46.7% acceptance) despite draft slowness—verification is fast enough that high acceptance compensates. Parallel fails (0.70× slowdown) as coordination overhead compounds hardware inefficiency.

with batched draft generation creates net 3.28× speedup. The key: verification is efficient enough that reasonable acceptance rates compensate for slower drafting.

**Parallel speculation fails completely.** Running 3 parallel draft candidates means 3× the draft cost. Even with 46.7% acceptance on the best candidate, the multiplied draft overhead overwhelms any gains from better acceptance. Result: 0.70× slowdown. The hardware inefficiency makes parallel speculation counterproductive—multiplying operations that already achieve poor GPU utilization makes things worse, not better.

**Same-model speculation validates our measurements.** Using the target to verify itself achieves perfect acceptance (100%) but no speedup (1.00×, essentially break-even) because both draft and verify use the same model. This serves as a sanity check—same-model speculation performs exactly as theory predicts, confirming our measurement methodology is sound.

The practical implication: on unified memory architectures where draft models suffer from poor hardware utilization, sophisticated speculation strategies that multiply draft operations (parallel, tree-based, multi-candidate) will likely fail. Only simple sequential speculation can succeed, and only when acceptance rates are high enough to offset the hardware efficiency handicap.

### G. Validation on Real Visual Content

Table VI shows results from 100 COCO validation images. We needed to verify whether synthetic text sequences

introduced artifacts—real multimodal workloads might behave differently.

TABLE VI  
REAL-WORLD DATA VALIDATION ON COCO IMAGES (TEXT-ONLY MODE)

Dataset	Samples	Slowdown (%)	Inversion
COCO Images (Text-Only)	100	21.4 ± 6.8	Yes

*Note:* Real-world validation using 100 COCO image captions in text-only mode. Model operates in text-only mode despite VLM architecture.

Similar slowdown percentages prove findings generalize to real visual content. Timings reflect text generation phase only (visual encoding excluded to isolate the hardware efficiency bottleneck).

The numbers match within measurement error: 21.4% slowdown for real images versus 22.4% for synthetic text. The slightly higher variance with COCO images likely reflects greater diversity in real visual content affecting the text generation context. The consistent slowdown confirms the hardware efficiency inversion generalizes to authentic multimodal workloads, not just controlled synthetic inputs.

## V. ANALYSIS

### A. Why Small Operations Achieve Poor Efficiency

Modern GPUs achieve peak throughput when all cores stay busy with work. The M4 Pro has 20 GPU cores, each with many parallel execution units (exact count not publicly documented, but typically 32-64 threads per core). A 576×576 matrix multiply generates 191M FLOPs of work distributed across potentially 640-1280 parallel threads.

In contrast, a 2048×2048 matrix multiply generates 8.6B FLOPs—45× more work. With the same number of threads, each thread handles much more computation, providing better opportunities to hide memory latency and keep execution units continuously busy. The larger operation also benefits from better cache utilization—the working set fits memory hierarchies differently, with larger operations achieving better reuse patterns.

The 273 GB/s unified memory bandwidth is critical here. Small operations move data (loading matrix A, matrix B, storing result C) without enough compute to hide memory latency effectively. Large operations achieve better arithmetic intensity—more FLOPs per byte moved—allowing computation to proceed while waiting for memory. This is why the efficiency gap (2.8×) is significant. It’s not just about parallelism underutilization; it’s about the fundamental compute-to-memory ratio.

### B. Why Framework Overhead Contributes But Doesn’t Dominate

Our initial focus on framework overhead wasn’t entirely wrong—it just wasn’t the whole story. The draft model’s 6 extra layers (30 vs 24) do create additional coordination

costs. Each layer requires kernel launches, synchronization points, and memory management operations. These add up to approximately 0.09ms of overhead.

But this accounts for only 20-40% of the total gap (0.13-0.23ms depending on precision). The remaining 60-80% comes from the 2.8× hardware efficiency difference we measured through direct profiling. Both factors contribute, but they require fundamentally different solutions:

- **Framework overhead (20-40%):** Can be reduced through kernel fusion, batching layers, or optimizing launch costs
- **Hardware underutilization (60-80%):** Requires increasing operation dimensions through bigger models, query batching, or architectural changes

Framework optimization alone cannot close a 2.8× efficiency gap—that requires getting more work per kernel to saturate available parallelism.

### C. Why This Differs From Datacenter GPUs

On a datacenter H100 GPU with 3000 GB/s bandwidth (7.5× more), memory pressure reduces significantly. Small operations are still relatively inefficient compared to large ones, but they’re less bottlenecked on data movement. The higher bandwidth allows even small operations to achieve better overlap of memory and compute.

Datacenter GPUs also have larger L2 caches (50-60MB vs 16-20MB estimated for M4) that can absorb working sets for small models more effectively, reducing main memory traffic. For large operations this matters less because they exceed cache capacity anyway. But for small operations, cache effectiveness becomes critical.

The unified memory architecture introduces another factor: CPU and GPU share memory controllers, creating potential contention. Discrete GPUs have dedicated controllers and bandwidth. For large operations where compute dominates, this matters less. For small operations where memory access patterns are critical, shared controllers can add latency that compounds the efficiency problem.

## VI. IMPLICATIONS

### A. For Speculative Execution

Practitioners should profile actual draft model speeds before assuming speculation will help. On unified memory architectures, traditional assumptions break. Alternative strategies to consider:

- **Profile hardware efficiency directly** using tools that measure GPU utilization—don’t rely solely on parameter counts
- **Avoid sophisticated speculation strategies** (parallel, tree-based, multi-candidate) on bandwidth-limited hardware—they multiply the cost of operations that already achieve poor efficiency
- **Batch multiple queries** to increase effective operation size for small models (8-16 queries simultaneously can improve GPU utilization)

- **Consider larger draft models** that achieve better GPU utilization, even if they have more parameters—counterintuitive but supported by our efficiency measurements

### B. For Edge Deployment

The 256M model uses essentially the same memory footprint as the 2.2B model in MLX’s quantized format (both 5GB). If memory is equal and smaller is slower, the larger model dominates for latency-critical applications. This challenges conventional wisdom that smaller models are always better for edge devices.

On unified memory architectures with bandwidth constraints, hardware efficiency matters more than parameter count for inference latency. The optimal model size balances:

- Memory footprint (must fit in available RAM)
- Operation dimensions (must achieve reasonable GPU utilization)
- Energy consumption (larger models draw more power despite being faster per token)
- Application requirements (latency vs throughput vs battery life)

For applications with strict latency requirements—medical imaging, real-time video analysis, autonomous vehicles—testing efficiency on target hardware is essential. Model selection should consider actual measured throughput and GPU utilization, not just parameter counts and memory footprints.

### C. For Framework Development

MLX and other frameworks optimized for unified memory could improve small model performance through several approaches:

- **Kernel fusion** to combine multiple small operations into larger ones (fuse matrix multiply + layer norm + activation to reduce launch overhead and improve parallelism)
- **Automatic layer batching** when possible (process layers 0-2 together, then 3-5, reducing kernel launches)
- **Dimension-aware dispatch** using different kernel implementations optimized for small vs large matrices (small kernels could focus on memory efficiency over peak throughput)
- **Profile-guided optimization** learning from actual hardware utilization patterns to select execution strategies dynamically

These optimizations could help close the efficiency gap, though they won’t eliminate it entirely—fundamental hardware characteristics limit what’s achievable with very small operations.

## VII. LIMITATIONS AND FUTURE WORK

**Hardware specificity:** Results come from Apple M4 Pro with 273 GB/s unified memory. Discrete GPUs have different characteristics—A100 with 1600 GB/s and H100 with 3000 GB/s may not show the same inversion magnitude or efficiency

gap. The underlying principle (small operations achieve lower GPU utilization) should hold across architectures, though exact numbers will differ. We need discrete GPU validation to confirm where the crossover points occur.

**Framework specificity:** MLX 0.29.3 is optimized specifically for Apple Silicon’s Metal API. PyTorch with CUDA, TensorRT-LLM, and other frameworks may have different kernel implementations and scheduling strategies that affect small operation efficiency differently. Cross-framework validation requires porting SmolVLM2 to other frameworks—non-trivial engineering work that we plan as follow-up.

**Single model family:** SmolVLM2 is the only vision-language model available in MLX as of November 2025. Testing other architectures (LLaVA, Qwen-VL, Phi-3-Vision) requires framework support that doesn’t yet exist. As MLX ecosystem matures, we can test whether the efficiency gap appears consistently across different vision-language architectures or varies with architectural choices.

**Statistical power:** With  $n=5$  runs, we achieve moderate statistical power (estimated 40-60% vs conventional 80%). However, the large effect sizes (Cohen’s  $d=1.45$ ), highly significant p-values ( $p_i < 0.001$ ), consistent direction across all conditions (sequence lengths, precisions, real vs synthetic), and cross-validation on COCO images give us confidence in the core findings. Additional runs would tighten confidence intervals but won’t change conclusions given the measured effects.

Future work should validate on discrete GPUs (NVIDIA A100/H100, AMD MI250), test additional model architectures as framework support becomes available, explore batching strategies to improve small model efficiency systematically, and develop hardware-aware model architecture search that optimizes for efficiency metrics (GPU utilization, achieved throughput) not just parameter count or memory footprint.

## VIII. RELATED WORK

**Speculative decoding** was introduced by Leviathan et al. [1] achieving 2-3 $\times$  speedups for text-only models. Medusa [3] improved this to 2.2-3.6 $\times$  with tree-based speculation, while EAGLE [4] reached 3.5 $\times$  through learned draft heads. These works assume draft models are faster than targets—our measurements show this assumption can break on bandwidth-limited unified memory, where hardware utilization efficiency determines performance more than parameter count.

**Multimodal speculation** has seen limited work. Lin et al. [5] reported 1.8-2.5 $\times$  speedups for vision-language models, noting degraded performance versus text-only systems. Gagrani et al. [6] observed similar limitations in their multimodal experiments. Neither work identified hardware efficiency as a contributing factor or measured the draft-target speed relationship directly on unified memory architectures.

**Memory bandwidth constraints** are well-documented. Gholami et al. [10] quantified the memory wall where bandwidth scaling lags compute growth substantially. Dao et al. [9] demonstrated that attention operations are fundamentally memory-bound. Jin et al. [11] found GPU compute utilization

as low as 0.4% for LLM inference on A100 GPUs due to memory bandwidth bottlenecks—similar to our 0.94% finding, though on a different architecture. Pope et al. [12] showed that inference latency is fundamentally bounded by parameter loading time from memory. Our contribution is showing how these bandwidth constraints on unified memory create performance inversions where smaller models run slower due to poor GPU utilization.

**ML on edge devices** has focused primarily on model compression through pruning, quantization, and knowledge distillation. Recent work on efficient VLMs includes FastVLM [13] from Apple, which demonstrates that optimal model size depends on runtime budget and that larger LLMs can outperform smaller ones at equivalent latency—paralleling our findings. FastV [14] achieves 45% FLOPs reduction through visual token pruning. This literature typically assumes smaller compressed models will be faster—our findings show this assumption requires validation on target hardware, particularly for unified memory architectures where operation size affects hardware utilization dramatically.

## IX. CONCLUSION

Smaller models aren’t universally faster. On Apple Silicon’s unified memory, the 256M SmolVLM2 model runs 22.4% slower than the 2.2B version during text generation (0.72ms vs 0.59ms,  $p < 0.001$ , Cohen’s  $d = 1.45$ ), despite having 8.6× fewer parameters.

After initial investigation suggested framework overhead as the cause, comprehensive kernel profiling revealed hardware utilization as the dominant factor. Small operations achieve only 0.94% of peak GPU throughput versus 2.59% for large operations—a 2.8× efficiency gap that accounts for 60-80% of the performance difference. Framework overhead from the draft model’s 6 extra layers contributes the remaining 20-40%, but hardware underutilization is the primary bottleneck.

When GPU utilization is under 1%, reducing parameter count doesn’t translate to speedup. The bottleneck is saturating available parallelism on bandwidth-limited hardware (273 GB/s unified memory), not computational quantity. The inversion persists across all precision types (24.6-39.2% slowdown range) and appears in real COCO images (21.4% slowdown), confirming it’s not a measurement artifact.

This inversion has measurable consequences for speculation strategies. Sequential approaches achieve 3.28× speedup with 46.7% acceptance because fast verification compensates for slower drafting. Parallel strategies fail completely (0.70× slowdown, same acceptance) because multiplying operations that already achieve poor hardware efficiency makes things worse. Same-model speculation breaks even (1.00×, 100% acceptance) as theory predicts, validating our measurements.

The findings challenge conventional wisdom that smaller models are always better for edge deployment and help explain why multimodal speculation underperforms text-only systems on unified memory architectures. Practitioners deploying vision-language models on unified memory should profile GPU utilization and actual speeds rather than assuming smaller means

faster. On bandwidth-limited hardware, efficiency matters more than parameter count.

As multimodal AI moves to edge devices with unified memory (phones, tablets, embedded systems), understanding these hardware efficiency characteristics becomes critical for real-world deployment. We provide practical guidelines and measurements to inform these decisions.

## REFERENCES

- [1] Y. Leviathan, M. Kalman, and Y. Matias, “Fast inference from transformers via speculative decoding,” *Proc. ICML*, 2023.
- [2] C. Chen et al., “Accelerating large language model decoding with speculative sampling,” *arXiv:2302.01318*, 2023.
- [3] T. Cai et al., “Medusa: Simple LLM inference acceleration framework,” *Proc. ICML*, 2024.
- [4] Y. Li et al., “EAGLE: Speculative sampling requires rethinking feature uncertainty,” *Proc. ICML*, 2024.
- [5] H. Lin, J. Zhang, and Y. Wang, “Multimodal speculative decoding,” *arXiv:2505.14260*, 2025.
- [6] M. Gagrani et al., “On speculative decoding for multimodal LLMs,” *arXiv:2404.08856*, 2024.
- [7] B. Zhou et al., “TinyLLaVA: A framework of small-scale large multimodal models,” *arXiv:2402.14289*, 2024.
- [8] “MLX: Machine learning framework for Apple silicon,” 2024. [Online]. Available: <https://github.com/ml-explore/mlx>
- [9] T. Dao et al., “FlashAttention: Fast and memory-efficient exact attention,” *NeurIPS*, 2022.
- [10] A. Gholami et al., “AI and memory wall,” *arXiv:2403.14123*, 2024.
- [11] Y. Jin, C.-F. Wu, D. Brooks, and G.-Y. Wei, “S3: Increasing GPU utilization during generative inference for higher throughput,” *arXiv:2306.06000*, 2023.
- [12] R. Pope, S. Douglas, et al., “Efficiently scaling transformer inference,” in *Proc. MLSys*, 2023.
- [13] P. K. A. Vasu et al., “FastVLM: Efficient vision encoding for vision language models,” in *Proc. CVPR*, 2025.
- [14] L. Chen et al., “An image is worth 1/2 tokens after layer 2: Plug-and-play inference acceleration for large vision-language models,” in *Proc. ECCV*, 2024.
- [15] P. Hübner and A. Hu, “Apple vs. oranges: Evaluating the Apple Silicon M-series SoCs for HPC performance and efficiency,” *arXiv:2502.05317*, 2025.
- [16] A. Benazir et al., “Profiling large language model inference on Apple Silicon: A quantization perspective,” *arXiv:2508.08531*, 2025.
- [17] D. Feng, Z. Xu, R. Wang, and F. X. Lin, “Profiling Apple Silicon performance for ML training,” *arXiv:2501.14925*, 2025.
- [18] Z. Yuan et al., “LLM inference unveiled: Survey and roofline model insights,” *arXiv:2402.16363*, 2024.