

ANE Batch Prefill for On-Device Parallel LLM Inference: Enabling Concurrent ANE Prefill and GPU Decode on Apple Silicon

AtomGradient

github.com/AtomGradient/hybrid-batch-prefill-on-ane

March 2026

Abstract

Our prior work [1] demonstrated that Apple Silicon’s Neural Engine (ANE) can offload LLM prefill from the GPU, reducing GPU power draw by $282\times$ and freeing thermal budget for sustained decode. However, the private `AppleNeuralEngine.framework` API used in that study dispatches only one token per call, yielding a constant latency of ~ 42 ms/token — impractical for interactive inference.

This companion paper addresses that bottleneck by introducing **ANE batch prefill**: processing up to 32 tokens per ANE dispatch (`ANE_SPATIAL=32`) via fused matrix–vector kernels compiled at model-load time. We benchmark three inference modes — pure ANE, pure MLX (GPU), and hybrid (ANE batch prefill + MLX GPU decode) — on Qwen3.5-0.8B and Qwen3.5-2B across short (13 tokens), medium (28 tokens), and long (74 tokens) prompts on an M2 Ultra.

Key results: (i) ANE batch prefill reaches 268 tok/s on 0.8B (vs. 23.7 tok/s for sequential dispatch — an **$11.3\times$ speedup**); (ii) ANE batch prefill *exceeds* GPU prefill throughput on short prompts (149 vs. 138 tok/s for 0.8B at 13 tokens); (iii) hybrid decode throughput reaches 39–54 tok/s on 0.8B (vs. ~ 18 tok/s for pure ANE decode); (iv) state transfer overhead (save + read + cache construction) is < 30 ms total, negligible relative to inference time. These results validate the feasibility of **concurrent ANE prefill + GPU decode pipelines** for multi-request serving on edge devices.

1. Introduction

On-device LLM inference faces a fundamental tension: the prefill phase is compute-bound (matrix multiplications over the full prompt), while the decode phase is memory-bandwidth-bound (sequential token generation). Apple Silicon provides three compute units — CPU, GPU, and Neural Engine (ANE) — sharing a unified memory bus, offering an opportunity to *disaggregate* these phases across hardware.

Our prior work [1] explored four inference pipelines on the same hardware, concluding that:

- CoreML batched ANE prefill reaches GPU parity at ~ 410 tokens,
- Private API sequential dispatch has a constant ~ 42 ms/token latency (the fundamental bottleneck),

- ANE prefill reduces GPU power from 62.05 W to 0.22 W.

The critical limitation was **sequential dispatch**: the private `AppleNeuralEngine.framework` API processes only one token per call, making prefill throughput scale as $O(N)$ regardless of ANE hardware capability. CoreML’s batched dispatch solves this but requires a heavy framework layer and model compilation step.

This paper demonstrates that the private API *can* achieve batched prefill by compiling fused ANE kernels that process up to 32 tokens simultaneously (`ANE_SPATIAL=32`). The key contributions are:

1. **ANE batch prefill implementation**: Fused matrix–vector kernels compiled at model-load time, processing 32 tokens per dispatch while CPU handles per-token state updates (RMSNorm, RoPE, SSM recurrence, GQA attention).
2. **Hybrid pipeline validation**: ANE batch prefill + MLX GPU decode via binary state serialization (64-byte header + per-layer state), demonstrating <30ms transfer overhead.
3. **Throughput characterization**: Benchmarks across two model sizes (0.8B, 2B) and three prompt lengths, showing batch prefill achieves $11.3\times$ speedup over sequential dispatch.
4. **Concurrent pipeline feasibility**: Since ANE batch prefill leaves the GPU idle (0.22 W measured in [1]), GPU decode of a previous request can execute *simultaneously*, enabling pipelined multi-request serving.

2. Background

2.1. Hardware Platform

All experiments were conducted on a Mac Studio with an Apple M2 Ultra (Table 1).

Table 1: Test platform specifications.

Spec	Value
Chip	Apple M2 Ultra
CPU cores	24 (16P + 8E)
GPU cores	76
ANE cores	32
ANE TOPS	31.6
Memory	192 GB unified
Bandwidth	800 GB/s
Metal support	Metal 4

2.2. Model Architecture: Qwen3.5 Hybrid

Qwen3.5 uses a hybrid architecture combining **DeltaNet** linear attention layers with periodic **full attention** layers (every `full_attention_interval` layers). This creates two distinct state types:

- **DeltaNet layers**: SSM state $[H, K, V]$ + circular convolution buffer $[C, \text{kernel} - 1]$
- **Full attention layers**: KV cache [capacity, heads, dim] with circular buffer (capacity = 2,048)

The 0.8B model has 24 layers (18 DeltaNet + 6 FullAttn); the 2B model has 36 layers (27 DeltaNet + 9 FullAttn).

2.3. Sequential vs. Batched ANE Dispatch

Our prior work [1] identified two dispatch modes for ANE inference:

Sequential dispatch (prior work, Pipelines 3–4): Each token requires a separate ANE model descriptor evaluation via `IOSurface` shared memory. Dispatch latency is ~ 42 ms/token regardless of prompt length, yielding $\text{TTFT} \propto N$.

Batched dispatch (this work): ANE kernels are compiled to process `ANE_SPATIAL=32` tokens per evaluation. The ANE performs the matrix multiplications for all 32 tokens in parallel, while the CPU sequentially updates per-token state (RMSNorm, RoPE, SSM recurrence, attention). This reduces ANE dispatch overhead by up to $32\times$.

3. Methods

3.1. ANE Batch Prefill

The inference pipeline splits each transformer layer between ANE and CPU:

1. **ANE**: Executes all linear projections as fused matrix–vector operations (`first_proj`, `o_proj`, `fused_ffn`, `lm_head`). Kernels are compiled to accept spatial dimension = 32, enabling batch processing.
2. **CPU**: Executes non-linear operations sequentially per token: RMSNorm (via `vDSP`), RoPE, SSM recurrence (DeltaNet state update), `Conv1d`, GQA attention with KV cache, and `softmax` (all via Accelerate framework `CBLAS/vDSP`).

Prompt tokens are processed in batches of up to 32. For a 74-token prompt, this requires $\lceil 74/32 \rceil = 3$ ANE dispatches (vs. 74 sequential dispatches in the prior approach). ANE kernels are compiled once at model load and cached persistently to `{model_dir}/ane_weights/`.

3.2. State Serialization

To bridge ANE prefill output to MLX GPU decode, we serialize the full model state to a binary file:

- **Header** (64 bytes): Magic (`ANELMS\0\0`), version, layer count, prompt length, logits dimension.
- **Per-layer state**: DeltaNet layers store `conv_pos` + SSM state + conv state; full attention layers store `kv_len` + `kv_start` + K cache + V cache.
- **Footer**: Final logits vector for next-token sampling.

The MLX decode script reads this binary state, constructs `ArraysCache` (DeltaNet) or `KVCache` (`FullAttn`) objects, and continues autoregressive generation on the GPU.

3.3. Three Inference Modes

We benchmark three modes:

1. **Pure ANE**: Both prefill and decode on ANE+CPU (batch prefill for prompt, sequential decode for generation).
2. **Pure MLX**: Both prefill and decode on GPU via MLX framework (baseline).
3. **Hybrid**: ANE batch prefill \rightarrow state serialization \rightarrow MLX GPU decode. The state transfer includes save (ANE side), read (MLX side), and cache construction.

4. Experiments

4.1. Setup

We tested two Qwen3.5 models:

- **Qwen3.5-0.8B** (FP16, 24 layers, $d = 1,024$)
- **Qwen3.5-2B** (BF16, 36 layers, $d = 1,536$)

Three prompt lengths were used:

- **Short:** “Hello” (13 tokens with chat template)
- **Medium:** “Explain the theory of relativity...” (28 tokens)
- **Long:** “Write a detailed technical comparison...” (74 tokens)

Each configuration was run 10 times; we report averages. Decode was capped at 20 tokens with temperature = 0 (greedy sampling).

4.2. Throughput Results: Qwen3.5-0.8B

Table 2: Throughput comparison for Qwen3.5-0.8B FP16 on M2 Ultra. All values averaged over 10 runs.

Prompt	Mode	Prefill (tok/s)	Decode (tok/s)	Total (ms)
Short (13 tok)	ANE	148.7	17.6	2,931
	MLX	138.2	61.2	3,679
	Hybrid	148.3	39.4	6,050
Medium (28 tok)	ANE	237.5	16.4	3,645
	MLX	290.0	64.5	3,836
	Hybrid	237.9	53.9	6,280
Long (74 tok)	ANE	268.0	17.1	3,774
	MLX	736.2	64.1	3,846
	Hybrid	271.1	54.0	6,335

4.3. Throughput Results: Qwen3.5-2B

4.4. State Transfer Overhead

Total state transfer overhead is ~ 27 ms for both models, confirming that the binary serialization format is efficient and does not bottleneck the hybrid pipeline.

4.5. Batch vs. Sequential Prefill Comparison

Table 5 compares the batch prefill results from this work with the sequential dispatch measurements from our prior study [1].

For 0.8B, batch prefill achieves an $11.3\times$ speedup over sequential dispatch, reducing the gap with GPU prefill from $31.1\times$ to $2.75\times$ at 74 tokens. For 2B, batch prefill yields a $7.3\times$ speedup, with GPU maintaining a larger lead ($4.8\times$ over batch) due to the 2B model’s higher compute requirements.

Table 3: Throughput comparison for Qwen3.5-2B BF16 on M2 Ultra. All values averaged over 10 runs.

Prompt	Mode	Prefill (tok/s)	Decode (tok/s)	Total (ms)
Short (13 tok)	ANE	91.7	12.4	5,271
	MLX	179.6	90.8	4,076
	Hybrid	92.8	22.6	8,878
Medium (28 tok)	ANE	160.4	12.5	6,189
	MLX	358.0	94.1	4,212
	Hybrid	159.3	29.0	9,281
Long (74 tok)	ANE	172.7	12.7	6,455
	MLX	828.6	92.4	4,274
	Hybrid	173.1	29.2	9,518

Table 4: State transfer overhead for hybrid pipeline (averaged over all prompt lengths and runs).

Model	Save (ms)	Read (ms)	Cache (ms)
0.8B	12	9	6
2B	11	9	6

5. Discussion

5.1. Prefill: ANE Batch vs. GPU

ANE batch prefill *outperforms* GPU prefill on short prompts: 149 vs. 138 tok/s for 0.8B at 13 tokens. This advantage diminishes as prompt length increases — GPU prefill scales more efficiently with batch size, reaching 736 tok/s at 74 tokens while ANE saturates at 268 tok/s.

The crossover point lies between 13 and 28 tokens for 0.8B. Below this threshold, ANE batch prefill is strictly faster than GPU. For the 2B model, GPU prefill already leads at 13 tokens (180 vs. 92 tok/s), with the gap widening at longer prompts.

This **short-prompt advantage** is significant for interactive applications (chatbots, code completion) where prompts are typically 10–30 tokens.

5.2. Decode Throughput Gap

Pure ANE decode (~ 18 tok/s for 0.8B, ~ 12 tok/s for 2B) is significantly slower than GPU decode (~ 61 – 65 tok/s for 0.8B, ~ 91 – 94 tok/s for 2B). This confirms the finding from our prior work that decode is memory-bandwidth-bound, and the GPU’s 800 GB/s bandwidth advantage is decisive.

Hybrid decode (39–54 tok/s for 0.8B) falls between pure ANE and pure MLX. The gap from pure MLX is due to the state reconstruction overhead in the MLX cache objects, not the serialization I/O (which is < 30 ms).

Table 5: Batch prefill (this work) vs. sequential dispatch (prior work [1]) for both models at 74 tokens on M2 Ultra.

Model	Dispatch mode	Prefill (tok/s)	Speedup	TTFT
0.8B	Sequential (prior)	23.7	1.0×	~3,122 ms
	Batch (this work)	268.0	11.3 ×	~276 ms
	GPU baseline (MLX)	736.2	31.1×	~101 ms
2B	Sequential (prior)	23.7	1.0×	~3,122 ms
	Batch (this work)	172.7	7.3 ×	~429 ms
	GPU baseline (MLX)	828.6	35.0×	~89 ms

5.3. The Concurrent Pipeline Argument

The hybrid pipeline’s total wall-clock time (6,050–6,335 ms for 0.8B) is higher than pure MLX (3,679–3,846 ms) when measured as a single request. This is expected — the sequential execution adds the state transfer overhead.

However, the **real value** of ANE batch prefill is concurrent execution. Since ANE prefill consumes only ~ 0.22 W of GPU power (measured in [1]), the GPU is available for *simultaneously* decoding a previous request at full speed (~ 90 tok/s). In a multi-request scenario:

- **Request k** : GPU decoding at ~ 90 tok/s
- **Request $k + 1$** : ANE batch prefilling at 149–268 tok/s
- **Net effect**: Zero idle time between requests; no GPU contention during prefill

This pipelining is impossible with GPU-only inference, where prefill and decode must be serialized on the same hardware.

5.4. Power and Thermal Implications

Combining the power measurements from our prior work with the throughput results here:

- **GPU prefill**: ~ 62 W GPU power (compute-bound)
- **ANE batch prefill**: ~ 0.22 W GPU + ~ 1.58 W ANE = ~ 1.8 W total
- **GPU decode**: ~ 14 W GPU power (bandwidth-bound)

During concurrent ANE prefill + GPU decode, total power is approximately $1.8 + 14 = 15.8$ W, compared to $62 + 14 = 76$ W if both prefill and decode used the GPU (hypothetical sequential scenario with two GPU requests). This represents a **79% power reduction** for the prefill component.

On mobile devices (A-series, TDP ~ 3 – 8 W), GPU prefill at 62 W would immediately trigger thermal throttling. ANE prefill at ~ 1.8 W fits within the mobile thermal envelope, enabling sustained inference without clock-speed degradation.

5.5. Limitations

1. **Private API dependency**: The `AppleNeuralEngine.framework` API is undocumented and may change across macOS versions.
2. **Hybrid decode gap**: Hybrid decode (39–54 tok/s) does not match pure MLX decode (~ 61 – 94 tok/s), likely due to cache reconstruction overhead in MLX’s Python layer.

3. **Single-machine results:** All measurements are from a single M2 Ultra; scaling behavior on other Apple Silicon variants (M1, M3, M4, A-series) remains to be validated.
4. **No concurrent measurement:** We measured ANE prefill and GPU decode separately; true concurrent execution was not yet benchmarked end-to-end.

6. Conclusion

We demonstrated that ANE batch prefill via the private `AppleNeuralEngine.framework` API achieves an $11.3\times$ speedup over sequential dispatch, reaching 268 tok/s for Qwen3.5-0.8B on M2 Ultra. On short prompts (≤ 13 tokens), ANE batch prefill *outperforms* GPU prefill (149 vs. 138 tok/s).

Combined with the power measurements from our prior work (0.22 W GPU during ANE prefill vs. 62 W during GPU prefill), these results validate the feasibility of **concurrent ANE prefill + GPU decode** for on-device multi-request serving. The state transfer overhead (< 30 ms) is negligible, and the ANE’s near-zero GPU power consumption during prefill enables true hardware-level parallelism without thermal contention.

Future work includes end-to-end concurrent pipeline benchmarking, extension to larger models (9B+), and investigation of the ANE batch size ceiling beyond `ANE_SPATIAL=32`.

Acknowledgments

Experiments were conducted on a Mac Studio (M2 Ultra, 192 GB). Power measurements referenced from our prior study used `powermetrics` on the same machine.

This work builds upon the foundational ANE reverse-engineering efforts by `maderix` [2] and `johnmai-dev` [3]. Without their pioneering work on the private `AppleNeuralEngine.framework` API, this research would not have been possible.

References

- [1] AtomGradient. *Disaggregated LLM Inference on Apple Silicon: CoreML ANE Prefill + MLX GPU Decode*. 2026. <https://github.com/AtomGradient/hybrid-ane-mlx-bench>
- [2] maderix. *ANE: Apple Neural Engine reverse-engineering and experiments*. <https://github.com/maderix/ANE>
- [3] johnmai-dev. *ANE-LM: LLM inference on Apple Neural Engine*. <https://github.com/johnmai-dev/ANE-LM>