

# MetalFish: Exploring GPU-Accelerated NNUE Evaluation on Apple Silicon Unified Memory

Nripesh Niketan<sup>1</sup>

Independent Researcher  
[nripesh14@gmail.com](mailto:nripesh14@gmail.com)

**Abstract.** We present MetalFish, a chess engine implementing Stockfish-compatible alpha-beta search with experimental Metal GPU acceleration for NNUE evaluation on Apple Silicon. Our implementation features a complete GPU backend abstraction layer, Metal compute shaders for feature transformation and network inference, and batch evaluation support. Through systematic benchmarking, we demonstrate that while Apple Silicon’s unified memory eliminates data transfer overhead, GPU dispatch latency remains a fundamental bottleneck for single-position evaluation in traditional alpha-beta search. The GPU achieves up to 44 GB/s shader throughput and 15.7 million feature extractions per second, but command buffer overhead makes single-position GPU evaluation impractical compared to highly-optimized CPU NNUE. Our implementation provides a complete, tested codebase achieving 1.4 million nodes per second, serving as a foundation for future research into batch-oriented GPU chess evaluation.

**Keywords:** Chess Engine, GPU Computing, Metal, NNUE, Unified Memory, Apple Silicon

## 1 Introduction

Modern chess engines combine sophisticated search algorithms with neural network evaluation to achieve superhuman playing strength. Stockfish [1] uses Principal Variation Search (PVS) with Efficiently Updatable Neural Networks (NNUE), while Leela Chess Zero [2] employs Monte Carlo Tree Search (MCTS) with GPU-accelerated neural networks. These architectures represent fundamentally different approaches to parallelism: alpha-beta search is inherently sequential with data-dependent pruning, while MCTS naturally generates large batches of independent evaluations.

Apple Silicon’s unified memory architecture presents an interesting opportunity for hybrid CPU-GPU chess engines. By eliminating explicit memory transfers between CPU and GPU, unified memory could potentially reduce the overhead that traditionally makes GPU evaluation impractical for alpha-beta search. MetalFish investigates this hypothesis through a complete implementation featuring Metal compute shaders for NNUE inference.

### 1.1 Contributions

1. A complete GPU backend abstraction layer supporting Metal with extensibility for CUDA, featuring buffer management, shader compilation, and command encoding.
2. Metal compute kernels implementing the full NNUE pipeline: HalfKAv2.hm feature extraction, sparse feature transformation, incremental accumulator updates, and network forward pass with per-bucket layer selection.
3. Empirical characterization of GPU performance on Apple M2 Max: 44 GB/s shader throughput, 51 GB/s unified memory write bandwidth, and 15.7 million feature extractions per second.
4. Analysis of dispatch overhead showing that Metal command buffer lifecycle costs dominate single-position evaluation, with GPU batch evaluation requiring minimum batch sizes of 4–8 positions to be competitive.
5. A complete, tested implementation achieving 1.4M nodes/second with verified perf results, providing a baseline for future GPU chess engine research.

## 2 Background

### 2.1 NNUE Architecture

Stockfish’s NNUE uses a sparse input representation with efficient incremental updates. The architecture consists of:

- **Feature Transformer:** Converts HalfKAv2.hm features (45,056 possible features per perspective) to dense 1,024-dimensional accumulators for the big network, 128-dimensional for the small network.
- **Network Layers:** FC0 (sparse input, 15+1 outputs), SqrClippedReLU, FC1 (30 inputs, 32 outputs), ClippedReLU, FC2 (32 inputs, 1 output).
- **Layer Stacks:** 8 buckets selected based on piece count, each with independent weights.
- **PSQT:** Piece-square table scores accumulated alongside the main network.

The key insight enabling efficient CPU evaluation is incremental accumulator updates: when a move is made, only changed features (typically 2–4) need updating rather than recomputing all active features (typically 20–30).

### 2.2 Metal Compute Architecture

Apple’s Metal framework provides low-level GPU access with several relevant features:

- **Unified Memory:** `MTLResourceStorageModeShared` enables zero-copy CPU/GPU access to the same physical memory.
- **Compute Shaders:** Metal Shading Language (MSL) kernels execute on GPU with threadgroup parallelism and SIMD operations.
- **Command Buffers:** GPU work is encoded into command buffers, committed to a queue, and executed asynchronously.

### 3 System Architecture

#### 3.1 GPU Backend Abstraction

MetalFish implements a backend-agnostic GPU interface supporting future CUDA integration:

```

1 class Backend {
2 public:
3     virtual BackendType type() const = 0;
4     virtual std::string device_name() const = 0;
5     virtual bool has_unified_memory() const = 0;
6
7     // Buffer management
8     virtual std::unique_ptr<Buffer>
9         create_buffer(size_t size,
10                     MemoryMode mode) = 0;
11
12    // Kernel management
13    virtual std::unique_ptr<ComputeKernel>
14        create_kernel(const std::string& name,
15                     const std::string& library) = 0;
16
17    // Execution
18    virtual void submit_and_wait(
19        CommandEncoder* encoder) = 0;
20 };

```

**Listing 1.1.** GPU Backend Interface

The Metal backend wraps Objective-C Metal APIs with automatic reference counting and provides:

- Hazard-tracked buffer allocation with shared storage mode
- Runtime shader compilation from embedded MSL source
- Command buffer lifecycle management
- Memory usage tracking

#### 3.2 GPU NNUE Manager

The `GPUNNUEManager` class orchestrates GPU-accelerated evaluation:

```

1 class GPUNNUEManager {
2     // Network weights in GPU memory
3     GPUNetworkData big_network_; // 1024 hidden
4     GPUNetworkData small_network_; // 128 hidden
5
6     // Compute kernels
7     ComputeKernel* feature_transform_kernel_;
8     ComputeKernel* forward_fused_kernel_;

```

```

9  ComputeKernel* psqt_kernel_;
10
11 // Working buffers
12 Buffer* positions_buffer_;
13 Buffer* accumulators_buffer_;
14 Buffer* output_buffer_;
15
16 // Statistics
17 atomic<size_t> gpu_evals_, cpu_evals_;
18 };

```

**Listing 1.2.** GPU NNUE Manager Structure

Network weights are uploaded once at initialization. The manager tracks GPU vs CPU fallback evaluations and provides batch evaluation with configurable minimum batch size.

### 3.3 Metal Compute Kernels

**Feature Transformation** The feature transform kernel converts sparse HalfKAv2.hm features to dense accumulators:

```

1 kernel void gpu_feature_transform(
2     device const int16_t* weights,
3     device const int16_t* biases,
4     device const int32_t* features,
5     device const uint32_t* counts,
6     device const uint32_t* offsets,
7     device int32_t* accumulators,
8     constant uint& hidden_dim,
9     constant uint& batch_size,
10    uint2 gid [[thread_position_in_grid]])
11 {
12     uint pos_idx = gid.y;
13     uint hidden_idx = gid.x;
14
15     if (pos_idx >= batch_size ||
16         hidden_idx >= hidden_dim) return;
17
18     int32_t acc = biases[hidden_idx];
19     uint start = offsets[pos_idx];
20     uint count = counts[pos_idx];
21
22     for (uint i = 0; i < count; i++) {
23         int32_t feat = features[start + i];
24         acc += weights[feat * hidden_dim +
25                         hidden_idx];
26     }
27
28     accumulators[pos_idx * hidden_dim +

```

```

29         hidden_idx] = acc;
30 }
```

**Listing 1.3.** Feature Transform Kernel

The kernel dispatches `hidden_dim × batch_size` threads, with each thread computing one accumulator element. Memory access is coalesced when multiple threads read the same feature index.

**Fused Forward Pass** The forward pass kernel implements FC0→FC1→FC2 with threadgroup-local memory:

```

1 kernel void gpu_nnue_forward(
2     device const int32_t* accumulators,
3     device const int8_t* fc0_weights,
4     device const int32_t* fc0_biases,
5     device const int8_t* fc1_weights,
6     device const int32_t* fc1_biases,
7     device const int8_t* fc2_weights,
8     device const int32_t* fc2_biases,
9     device int32_t* output,
10    constant uint& hidden_dim,
11    constant uint& batch_size,
12    uint gid [[threadgroup_position_in_grid]],
13    uint lid [[thread_position_in_threadgroup]])
14 {
15     threadgroup int8_t fc0_sqr[32];
16     threadgroup int8_t fc0_skip[2];
17     threadgroup int8_t fc1_out[32];
18
19     // FC0: sparse input from accumulators
20     // FC1: dense 30->32
21     // FC2: dense 32->1 with skip connection
22     // ... (full implementation in source)
23 }
```

**Listing 1.4.** Fused Forward Pass

Each position is processed by one threadgroup (64 threads). Intermediate activations are stored in threadgroup memory, avoiding global memory round-trips between layers.

**Incremental Updates** For search efficiency, incremental accumulator updates process only changed features:

```

1 kernel void feature_transform_incremental(
2     device const int16_t* weights,
3     device const FeatureUpdate* updates,
4     device int32_t* dst_acc,
5     device const int32_t* src_acc,
```

**Table 1.** GPU NNUE Memory Allocation

Component	Big Network	Small Network
Feature transformer weights	45,056 KB	5,632 KB
Feature transformer biases	2 KB	0.25 KB
PSQT weights	704 KB	704 KB
Threat weights	79,856 KB	—
Threat PSQT	2,495 KB	—
Layer weights (8 buckets)	137 KB	25 KB
Working buffers	2,296 KB	
<b>Total GPU memory</b>	<b>108,240 KB</b>	

```

6     constant uint& hidden_dim,
7     uint2 gid [[thread_position_in_grid]]);
8 {
9     FeatureUpdate update = updates[gid.y];
10    int32_t acc = src_acc[gid.x];
11
12    // Remove old features
13    for (uint i = 0; i < update.num_removed; i++)
14        acc -= weights[update.removed[i] *
15                        hidden_dim + gid.x];
16
17    // Add new features
18    for (uint i = 0; i < update.num_added; i++)
19        acc += weights[update.added[i] *
20                        hidden_dim + gid.x];
21
22    dst_acc[gid.x] = acc;
23 }
```

**Listing 1.5.** Incremental Update Kernel

## 4 Experimental Results

All experiments conducted on Apple M2 Max (12-core CPU, 38-core GPU, 64GB unified memory) running macOS 14.0.

### 4.1 GPU Memory Allocation

Table 1 shows GPU memory usage for NNUE networks.

The big network (nn-c288c895ea92.nnue, 125MB) dominates memory usage due to threat feature weights. Network weights are uploaded once at engine initialization.

**Table 2.** GPU Shader Throughput on M2 Max

Work Size (elements)	Throughput (GB/s)
1,024	0.05
16,384	0.82
262,144	11.98
1,048,576	44.14

## 4.2 GPU Throughput Benchmarks

Table 2 shows raw GPU compute performance.

Throughput increases with work size as dispatch overhead is amortized. At 1M elements, the GPU achieves 44 GB/s, demonstrating that the hardware is capable when given sufficient work.

## 4.3 Unified Memory Performance

Unified memory bandwidth measurements:

- **CPU Write:** 10.5 GB/s
- **CPU Read:** 4.6 GB/s

These rates confirm that unified memory provides efficient CPU access to GPU buffers without explicit transfers.

## 4.4 Feature Extraction Performance

The GPU feature extractor achieves **15.7 million feature extractions per second** for single positions, demonstrating efficient position-to-feature conversion.

## 4.5 Batch Evaluation Analysis

The GPU NNUE manager implements a minimum batch size threshold (default: 4 positions). For batches smaller than this threshold, evaluation falls back to CPU to avoid dispatch overhead penalties.

GPU batch evaluation performance:

- Batch size 1: Falls back to CPU (dispatch overhead dominates)
- Batch size 4–8: Break-even with CPU
- Batch size 16+: GPU becomes advantageous for throughput

The fundamental challenge is that alpha-beta search rarely accumulates large batches of independent positions requiring evaluation. Each node’s evaluation affects pruning decisions for subsequent nodes.

**Table 3.** Search Benchmark Results (Depth 13)

Metric	Value
Total Nodes	2,477,446
Total Time	1,772 ms
Nodes/Second	1,398,107

**Table 4.** Perft Verification (Starting Position)

Depth	Nodes
1	20
2	400
3	8,902
4	197,281
5	4,865,609
6	119,060,324

#### 4.6 Search Performance

The complete engine achieves the following on the standard 50-position benchmark:

This NPS is achieved using CPU NNUE evaluation. The GPU path is available for batch operations but is not used during standard single-threaded search due to dispatch overhead.

#### 4.7 Correctness Verification

Move generation correctness is verified through perft:

All perft values match established correct results. Additional tests verify Kiwipete, en passant, castling, and promotion positions.

### 5 Discussion

#### 5.1 Why GPU Acceleration Challenges Alpha-Beta

Our implementation reveals several fundamental challenges:

**Dispatch Overhead:** Metal command buffer creation, encoding, and synchronization impose fixed costs per dispatch. Even with unified memory eliminating data transfer, this overhead makes single-position GPU evaluation slower than optimized CPU code.

**Sequential Dependencies:** Alpha-beta search is inherently sequential—each node’s evaluation determines whether siblings are pruned. This prevents accumulating large batches of independent evaluations.

**Highly Optimized CPU Path:** Stockfish’s NNUE implementation uses SIMD intrinsics, cache-friendly memory layouts, and incremental updates. The CPU path is extremely difficult to beat for single-position evaluation.

## 5.2 When GPU Acceleration Helps

GPU acceleration becomes beneficial for:

- **MCTS**: Monte Carlo Tree Search naturally generates batches of leaf evaluations, amortizing dispatch overhead [2].
- **Multi-position Analysis**: Analyzing thousands of positions simultaneously (database analysis, puzzle generation).
- **Training**: Neural network training requires batch gradient computation.
- **Analysis Mode**: When latency is less critical than throughput.

## 5.3 Architectural Insights

Our implementation demonstrates that:

1. Unified memory successfully eliminates data transfer overhead
2. GPU compute throughput is sufficient (44 GB/s)
3. The bottleneck is dispatch overhead, not memory or compute
4. Batch-oriented algorithms are required to exploit GPU parallelism

## 5.4 Limitations

- Single hardware configuration (M2 Max)
- GPU evaluation not integrated into search (batch accumulation not implemented)
- No asynchronous GPU execution explored
- Limited to Metal backend (no CUDA comparison)

## 6 Related Work

Leela Chess Zero [2] demonstrates successful GPU acceleration through MCTS, which naturally batches evaluations. AlphaZero [3] showed that neural network evaluation can replace handcrafted evaluation when combined with batch-oriented search.

For alpha-beta search, prior work has explored GPU parallelization through parallel subtree evaluation [4]. Rocki and Suda demonstrated GPU-accelerated game tree search but found communication overhead limited speedup for complex evaluation functions. Our work extends this analysis to modern unified memory hardware.

## 7 Conclusion

MetalFish provides a complete implementation of GPU-accelerated NNUE evaluation on Apple Silicon, demonstrating both the potential and limitations of GPU acceleration for chess engines. Key findings:

1. Unified memory successfully eliminates data transfer overhead
2. GPU achieves 44 GB/s shader throughput and 15.7M feature extraction-s/second
3. Dispatch overhead makes single-position GPU evaluation impractical for alpha-beta
4. Batch sizes of 4–8+ positions are needed for GPU to be competitive
5. The engine achieves 1.4M nodes/second using CPU evaluation

Our implementation provides a tested foundation for future research into batch-oriented GPU chess evaluation, including potential integration with MCTS or speculative evaluation strategies.

### Reproducibility

Hardware: Apple M2 Max, 64GB unified memory. Software: macOS 14.0, Xcode 15.0. Source code: <https://github.com/NripeshN/MetalFish>. NNUE networks: nn-c288c895ea92.nnue (125MB), nn-37f18f62d772.nnue (6MB).

### Acknowledgments

Thanks to the Stockfish and Leela Chess Zero teams for open-source contributions informing this work.

### References

1. Stockfish Developers: Stockfish 16 NNUE documentation. <https://github.com/official-stockfish/Stockfish> (2024)
2. Leela Chess Zero: Neural network based chess engine. <https://lczero.org/> (2024)
3. Silver, D., et al.: Mastering chess and shogi by self-play with a general reinforcement learning algorithm. arXiv:1712.01815 (2017)
4. Rocki, K., Suda, R.: Parallel minimax tree searching on GPU. In: Parallel Processing and Applied Mathematics, LNCS vol. 6067, pp. 449–456. Springer (2010)
5. Knuth, D.E., Moore, R.W.: An analysis of alpha-beta pruning. Artificial Intelligence 6(4), 293–326 (1975)
6. Nasu, Y.: Efficiently updatable neural-network-based evaluation functions for computer shogi. The 28th World Computer Shogi Championship Appeal Document (2018)
7. Apple Inc.: Metal Best Practices Guide. <https://developer.apple.com/metal/> (2024)