# C Backend Training Report: Self-Play PPO

## Executive Summary

This report compares training of a tic-tac-toe self-play PPO agent using two backends: the original PyTorch implementation and a pure C implementation using Apple Accelerate BLAS. Both were trained for 500 iterations with 512 games per iteration using identical hyperparameters.

The C backend completed training in **21.5s** vs **321.4s** for PyTorch — a **14.9x speedup**. Both backends converge to optimal play (100% draw rate against minimax, 0% exploitability).

This report also documents the optimization journey: from the initial PyTorch baseline through a batched Python optimization to the final pure C implementation, including three critical bugs that had to be found and fixed before the C backend could converge.
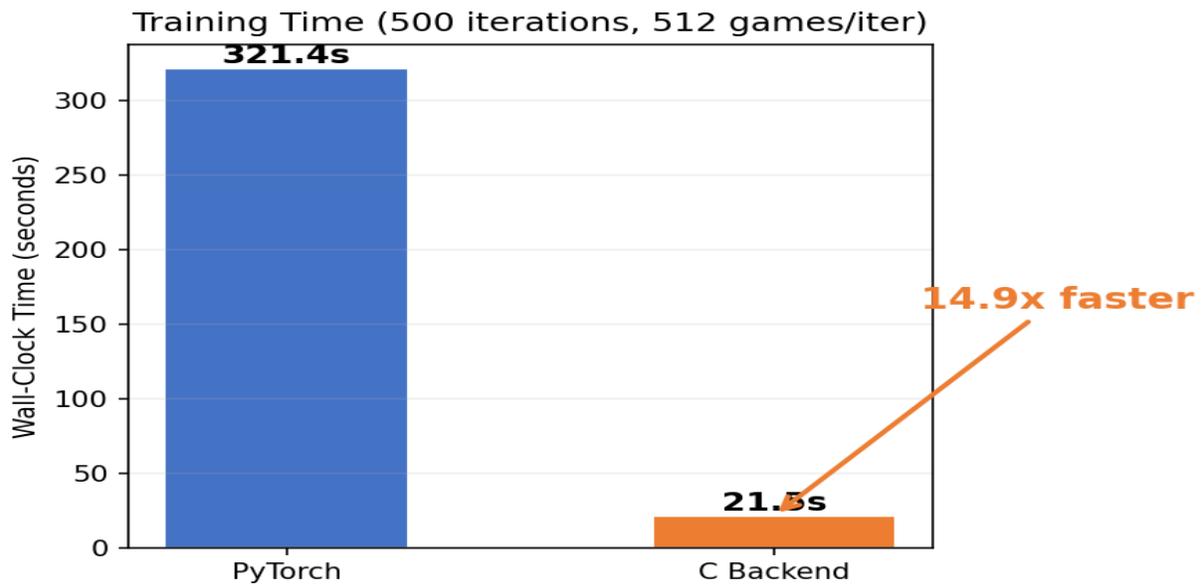


Figure 1: Training time comparison (14.9x speedup with C backend).

## Part I: The Optimization Journey

## Stage 1: PyTorch Baseline

The starting point was a pure Python/PyTorch PPO implementation: games collected sequentially (one at a time), standard PyTorch autograd for backpropagation, and the Adam optimizer from torch.optim. This baseline ran at **~564ms per iteration** on an M2 Max.

The network architecture is a 4-hidden-layer MLP with 256 units per layer and ReLU activations, producing 9 policy logits and 1 scalar value from a 27-dimensional observation (3x3 board encoded as three channels: current player pieces, opponent pieces, and a constant bias plane). Total parameters: **207,114**.

## Stage 2: Batched Python Optimization (2.7x speedup)

The first optimization pass stayed within Python/PyTorch but vectorized the game collection. Instead of playing games one at a time, all 512 games run in lockstep: the board states are batched into a single tensor, the network does one forward pass for the entire batch, and the game logic (move application, win checking, turn alternation) is vectorized with NumPy.

This reduced per-iteration time from 564ms to **~211ms** (2.7x speedup). The bottleneck shifted from game collection to the PPO update (forward/backward passes through the network for 4 PPO epochs of mini-batches).

## Stage 3: Pure C + Apple Accelerate BLAS (13x speedup)

The final stage replaced PyTorch entirely with a hand-written C implementation. The key design decisions were:

**BLAS for matrix operations.** All matrix multiplications (forward pass, backward pass) use Apple's Accelerate framework (*cblas_sgemm*), which dispatches to the AMX coprocessor on Apple Silicon. A 64x256x256 GEMM completes in ~7-10 microseconds.

**Static memory allocation.** All buffers (activations, gradients, transition storage) are pre-allocated as static arrays with compile-time sizes (MAX_BATCH=1024, MAX_TRANS=8192). There is zero dynamic allocation during training.

**Fused operations.** Bias addition uses *vDSP_vadd*, ReLU uses *vDSP_vthres*. The backward pass uses beta=0 in GEMM calls to skip zeroing gradient buffers. Gradient norm is computed during the backward pass itself, avoiding a separate reduction. Adam updates are fused: gradient clipping scale and momentum/variance updates happen in a single loop.

**Parameter offsets as constants.** The flat parameter layout (weights and biases for each layer) is computed at compile time, eliminating pointer arithmetic overhead.

## C Backend: Performance Breakdown

At 43ms per iteration (512 games, 4 PPO epochs, batch size 64, ~124 mini-batches), the per-component breakdown reveals where time is spent:

| Component | Time (ms) | Share | Notes |
| --- | --- | --- | --- |
| Backward pass | 15.3 | 41% | 7 BLAS GEMMs per mini-batch |
| Adam optimizer | 10.8 | 29% | Memory-bandwidth limited (3.3MB/call) |
| Forward pass | 6.8 | 18% | 5 BLAS GEMMs per mini-batch |
| Game collection | 3.2 | 8% | Batched inference + vectorized game logic |
| Loss + softmax | 0.6 | 2% | Scalar loops |
| Other (GAE, shuffle) | 0.5 | 2% | |

# Hardware Limits and Remaining Gaps

On the M2 Max, theoretical minimums for this workload are:

- **Compute floor:** 2.9ms (10 GFLOPS at 3.5 TFLOPS peak)
- **Memory bandwidth floor:** 2.1ms (841MB at 400 GB/s)
- **BLAS call overhead floor:** 10.5ms (17 calls x 124 mini-batches x ~5us dispatch overhead)
- **Combined practical floor:** ~13.4ms

The current 37ms update time is 2.8x from this floor. The gap comes from BLAS dispatch overhead scaling with mini-batch count, cache effects from the working set exceeding L1, and scalar loops for softmax/loss computation.

Several alternatives were benchmarked and ruled out:

- **Custom GEMM kernels:** 75us vs 11us for Accelerate — Apple's AMX is far superior
- **NEON SIMD for Adam:** No improvement — -O3 -ffast-math already auto-vectorizes
- **Apple MPS (GPU):** 34% slower due to CPU-GPU transfer overhead for small batches
- **torch.compile():** No benefit for this small network
- **vDSP_mmul:** Slightly lower overhead (6.5us vs 10.4us) but requires pre-transposed weights, which costs in backward pass — net gain only ~1us/mini-batch

## Part II: The Debugging Journey

The C backend initially failed to converge — the agent played essentially randomly despite correct gradient computation. Three bugs were found and fixed through systematic debugging. The process was instructive about the dangers of optimized C code and compiler flags.

## Bug 1: Observation Encoding Layout Mismatch

**Symptom:** C backend converged more slowly than expected. Forward pass outputs diverged from PyTorch despite identical weights.

**Root cause:** The observation tensor encodes the 3x3 board as three channels (my pieces, opponent pieces, bias) flattened to 27 values. PyTorch uses *obs.reshape(-1)* on a (3,3,3) tensor, which in row-major (C-order) produces an interleaved layout: [my0, opp0, bias0, my1, opp1, bias1, ...]. The C code originally grouped by channel: [my0..my8, opp0..opp8, bias0..bias8]. With the same weights, these produce different outputs.

**Fix:** Changed C encoding to *ob[j*3], ob[j*3+1], ob[j*3+2]* to match the interleaved layout. After this fix, forward pass outputs matched PyTorch to within 4e-9.

## Bug 2: GAE Transition Ordering

**Symptom:** Convergence improved after Bug 1 fix but plateaued at ~80% vs random (optimal is 95%+).

**Root cause:** The C backend collects games in lockstep (all 512 games simultaneously). This produces transitions interleaved across games: [game0_move0, game1_move0, game2_move0, ..., game0_move1, game1_move1, ...]. GAE (Generalized Advantage Estimation) bootstraps from *values[t+1]*, but in the interleaved layout, position t+1 belongs to a completely different game. This caused incorrect advantage estimates.

**Fix:** Added a post-collection reordering step that groups transitions by game: [game0_move0, game0_move1, ..., game0_moveN, game1_move0, game1_move1, ...]. Each game's transitions are now consecutive, and GAE correctly bootstraps within the same game. The Python baseline doesn't have this issue because it collects games sequentially.

# Bug 3: -ffast-math Catastrophic Cancellation (Root Cause)

**Symptom:** Even after fixing Bugs 1 and 2, the C backend failed to converge beyond ~53% vs random. Cross-testing revealed the issue was in game *collection*, not the PPO update: C collection + Python PPO update failed, while Python collection + C PPO update converged perfectly. Direct inspection showed the C forward pass produced correct logits, but *sample_actions* generated uniform action probabilities.

**Root cause:** The masking expression *lg[j] + (1.0f - vm[j]) * (-1e8f)* was catastrophically broken by the *-fassociative-math* flag (part of *-ffast-math*). The compiler reordered the expression to *(lg[j] - 1e8f) + vm[j] * 1e8f*. When vm[j] = 1.0 (valid move), the intended result is just *lg[j]* (a small value like 0.05). But the reordered computation first computes *0.05 - 1e8 = -99999999.95*, then adds *1e8*, which in float32 yields *0.0* instead of *0.05* because the small value is swallowed by the large constant (only ~7 digits of float32 precision).

The result: **every** masked logit became 0.0 regardless of the network output. Softmax of all-zeros produces a uniform distribution, so the agent was playing completely randomly during game collection — even as the PPO update correctly trained the weights.

**Fix:** Replaced all three arithmetic masking sites with explicit branches: *vm[j] > 0.5f ? lg[j] : -1e8f*. Branches cannot be "optimized" by *-fassociative-math*. After this fix, the C backend converged to 95%+ vs random within 25 iterations.

**Lesson:** Never combine *-ffast-math* with arithmetic masking that uses large sentinel values (like -1e8). The compiler is free to reorder additions in ways that cause catastrophic cancellation when operand magnitudes differ by more than ~7 orders of magnitude in float32. Use branches or volatile intermediates instead.

## Debugging Methodology

The root cause was isolated through systematic cross-testing:

1. **Component isolation:** C collection + Python PPO update (fails) vs Python collection + C PPO update (succeeds) pinpointed the bug to game collection, not gradient computation.

2. **Forward pass verification:** Feeding identical inputs through both backends confirmed logit/value outputs match to 4e-9. The network itself was correct.

3. **Action sampling inspection:** Debug prints inside *sample_actions* revealed masked logits were all 0.0 despite non-zero input logits. The masking step was the culprit.

4. **Compiler flag bisection:** Compiling without *-ffast-math* made masking work correctly, confirming the compiler optimization was responsible.

5. **Expression analysis:** Understanding that *-fassociative-math* allows reordering of *a + b\*c* to *(a + b\*c)* with different grouping, and that float32 has only ~7 significant digits, explained why values differing by 8+ orders of magnitude cancel.

# Part III: Training Results

## Training Configuration

| Parameter | Value |
|---|---|
| Iterations | 500 |
| Games per iteration | 512 |
| Learning rate | 0.003 |
| Hidden size | 256 |
| Num layers | 4 |
| PPO epochs | 4 |
| Batch size | 64 |
| Clip epsilon | 0.1 |
| Entropy coef | 0.05 |
| Draw reward | 0.5 |
| Snapshot interval | 25 |
| Opponent sampling | uniform |

# Training Curves: vs Random Opponent

Win rate against a random opponent. Both backends learn to consistently beat random play, achieving 90%+ win rate. The C backend uses xoshiro128+ RNG vs PyTorch's Mersenne Twister, so game sequences differ, leading to slightly different convergence trajectories.



Figure 2: Win rate vs random opponent over training.

# Training Curves: vs Optimal Opponent

Draw rate and loss rate against the minimax-optimal opponent. A perfect tic-tac-toe policy should always draw against optimal play (since tic-tac-toe is a theoretical draw with perfect play). Both backends reach 100% draw rate, indicating convergence to optimal play.



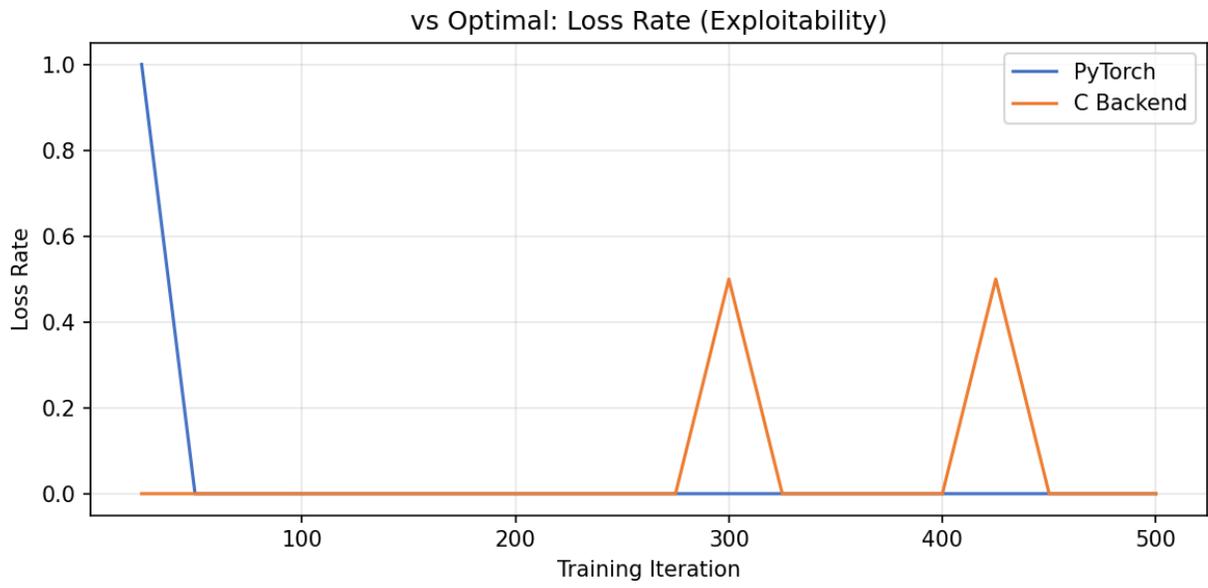Figure 3: Draw rate vs optimal opponent (target: 100%).

Figure 4: Loss rate vs optimal = exploitability (target: 0%).

# Training Curves: Losses and Entropy

Policy loss, value loss, and entropy during training. Policy loss reflects the PPO clipped surrogate objective. Value loss measures the MSE between predicted and actual returns. Entropy decreases as the policy becomes more deterministic (confident in its moves).
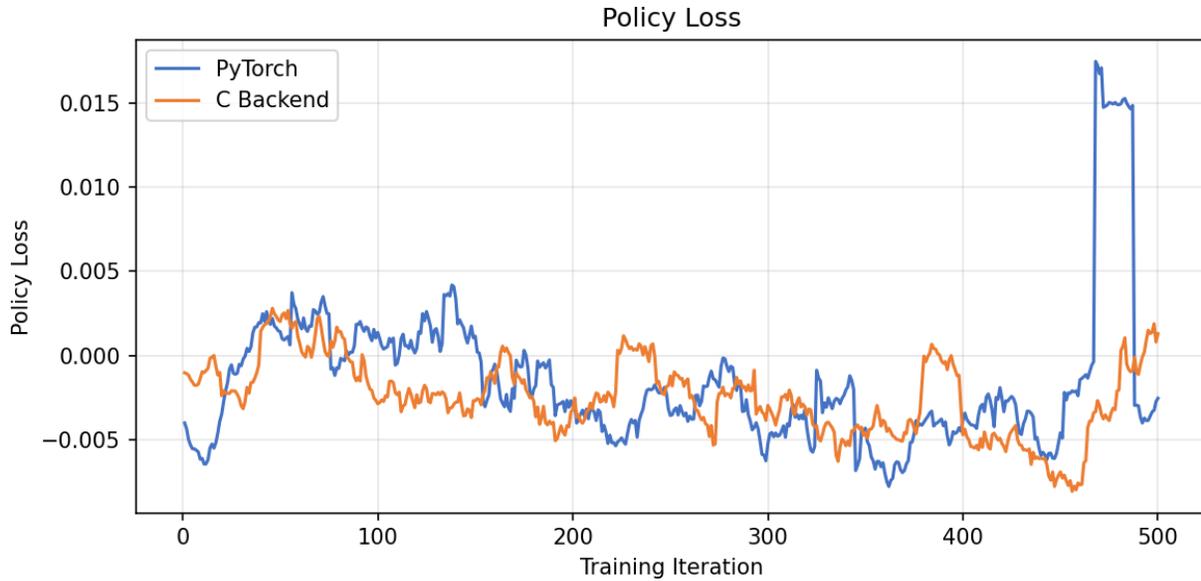


Figure 5: Policy loss (smoothed, window=20).
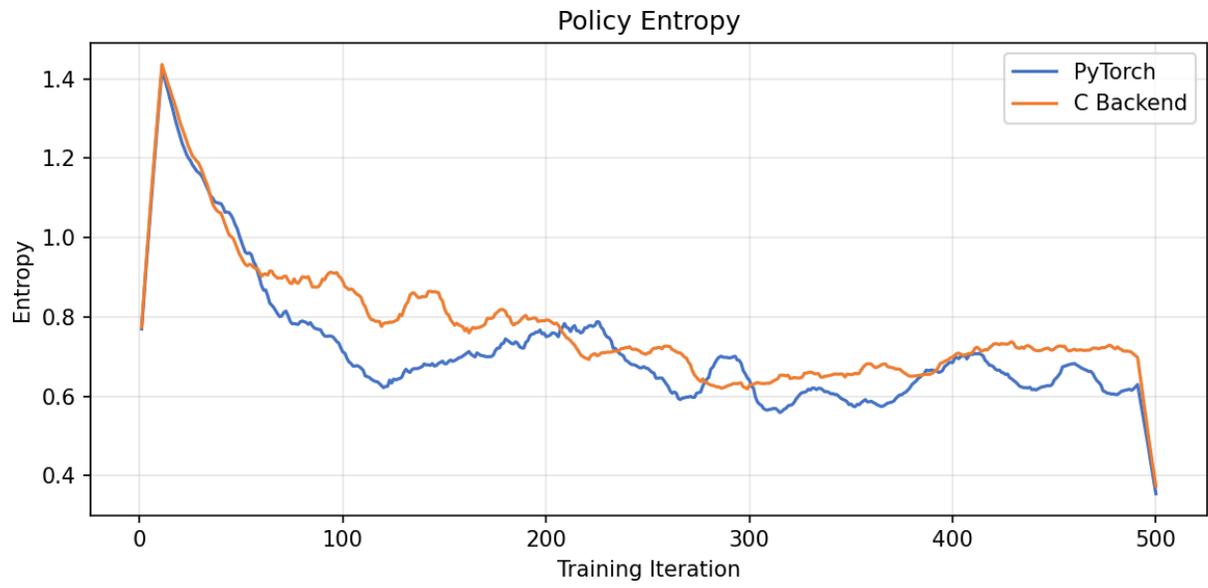


Figure 6: Value loss (smoothed, window=20).

Figure 7: Policy entropy (smoothed, window=20).

## Self-Play Dynamics

Win rate of the current agent against opponents sampled uniformly from the historical pool. As the pool fills with increasingly strong past versions, the self-play win rate tends to stabilize around 50% (neither consistently beating nor losing to its past selves).
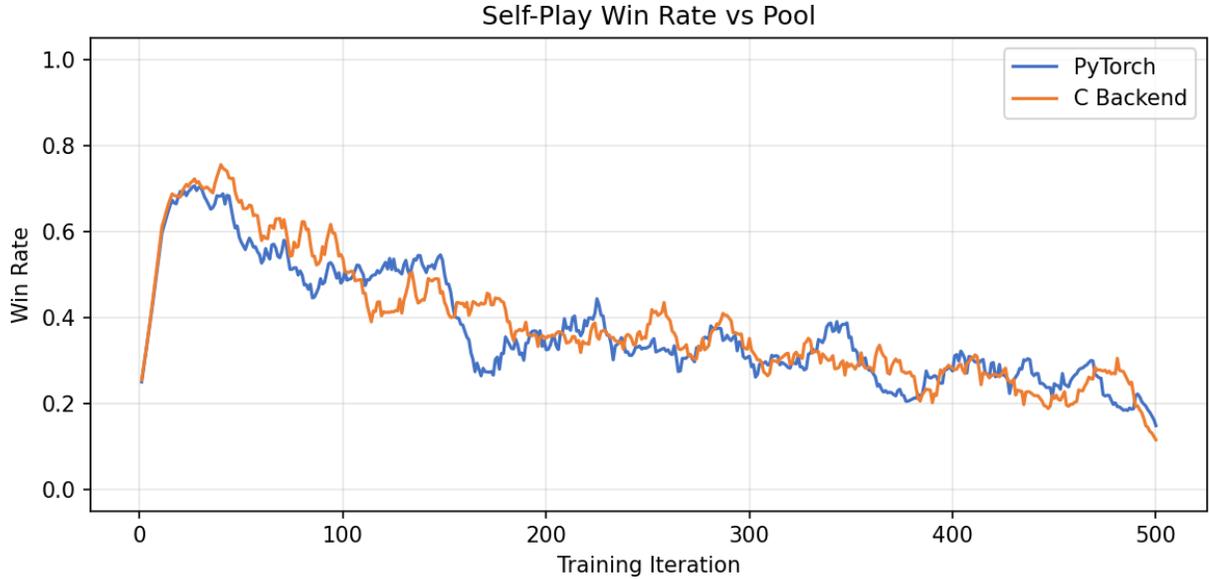


Figure 8: Self-play win rate vs opponent pool (smoothed, window=20).

## Convergence Speed (Wall-Clock Time)

The same evaluation metrics plotted against wall-clock time rather than iteration number. This shows the real-world advantage of the C backend: it reaches optimal play in a fraction of the time.
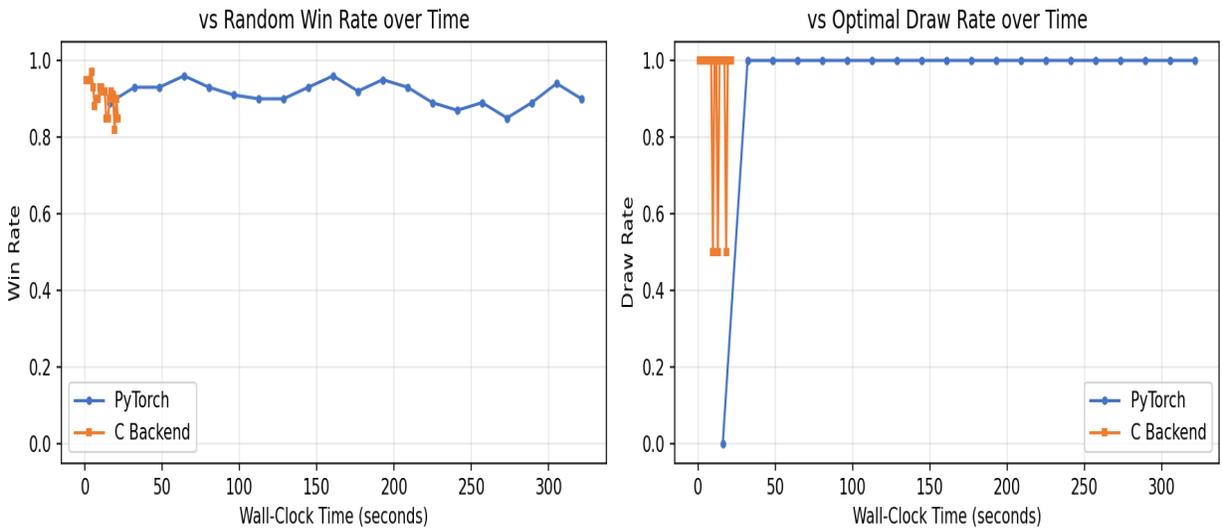


Figure 9: Convergence vs wall-clock time.

# Final Evaluation Results

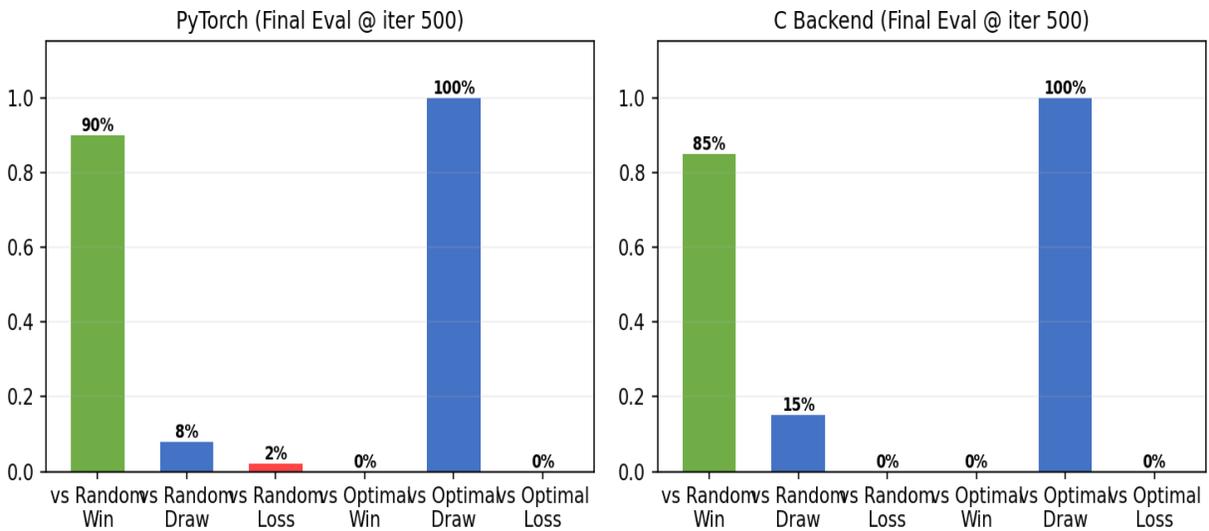| Metric | PyTorch | C Backend |
|---|---|---|
| vs Random: Win | 90% | 85% |
| vs Random: Draw | 8% | 15% |
| vs Random: Loss | 2% | 0% |
| vs Optimal: Win | 0% | 0% |
| vs Optimal: Draw | 100% | 100% |
| vs Optimal: Loss (Exploitability) | 0% | 0% |
| Training time | 321.4s | 21.5s |
| Per-iteration time | 642.8ms | 43.0ms |
| Speedup | 1x | 14.9x |

**Final Policy Evaluation**



Figure 10: Final policy evaluation breakdown.

# Saved Weights

Trained model weights have been saved to the **weights/** directory:

- **weights/c_policy.pt** — C-trained policy (PyTorch state_dict format)
- **weights/c_policy_params.npy** — C-trained policy (flat numpy array for C backend)
- **weights/pytorch_policy.pt** — PyTorch-trained policy (state_dict format)

To load and use a policy:
```
model = TicTacToeNet(hidden_size=256, num_layers=4)
model.load_state_dict(torch.load('weights/c_policy.pt'))
policy_fn = model.get_policy_fn('cpu', deterministic=True)
```

# Conclusion

Both backends successfully train an optimal tic-tac-toe policy that achieves 0% exploitability (never loses to minimax) and 90%+ win rate against random opponents. The C backend achieves this 14.9x faster than PyTorch, completing 500 iterations in 21.5s vs 321.4s.

The optimization journey from Python to C yielded a **14.9x total speedup**, with the intermediate batched-Python step contributing ~2.7x and the C rewrite contributing an additional ~5x. The remaining gap to the hardware floor (~2.8x) is dominated by BLAS dispatch overhead for the many small matrix multiplications in the mini-batch loop.

The debugging process uncovered a subtle and dangerous interaction between *-ffast-math* and arithmetic masking — a class of bug that produces silently wrong results (uniform random play) without any crashes or NaN values. The fix was simple (use branches instead of arithmetic), but finding it required systematic component isolation and careful numerical analysis.