

# EdgeFormer Benchmark Report

Generated on 2025-03-29 19:26

## Performance Highlights:

Best performance: 1186.91 tokens/second

Model size: small

Sequence length: 1024

Device: envy

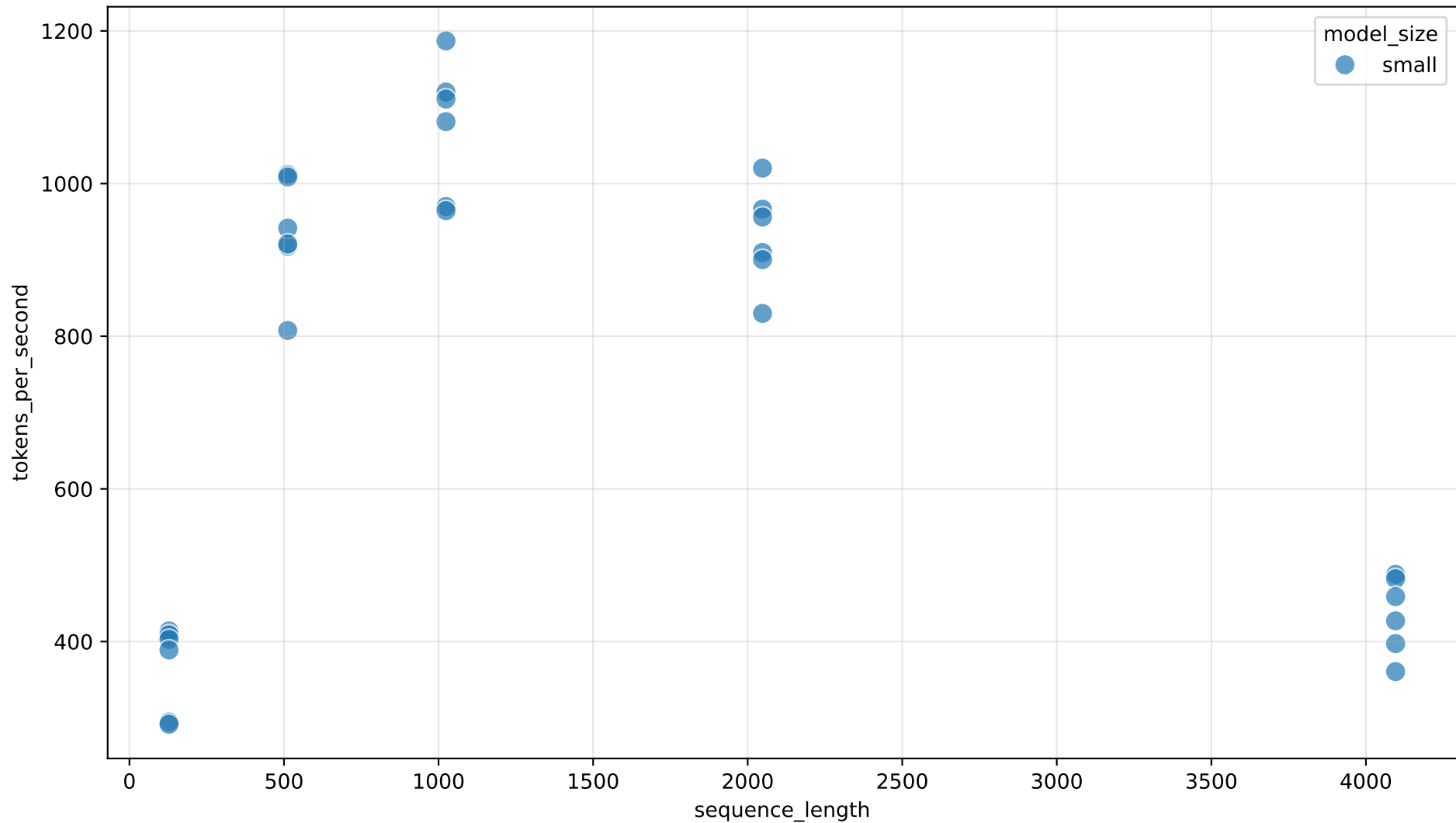
## Recommendations:

Optimal sequence length is around 1024 tokens for best performance

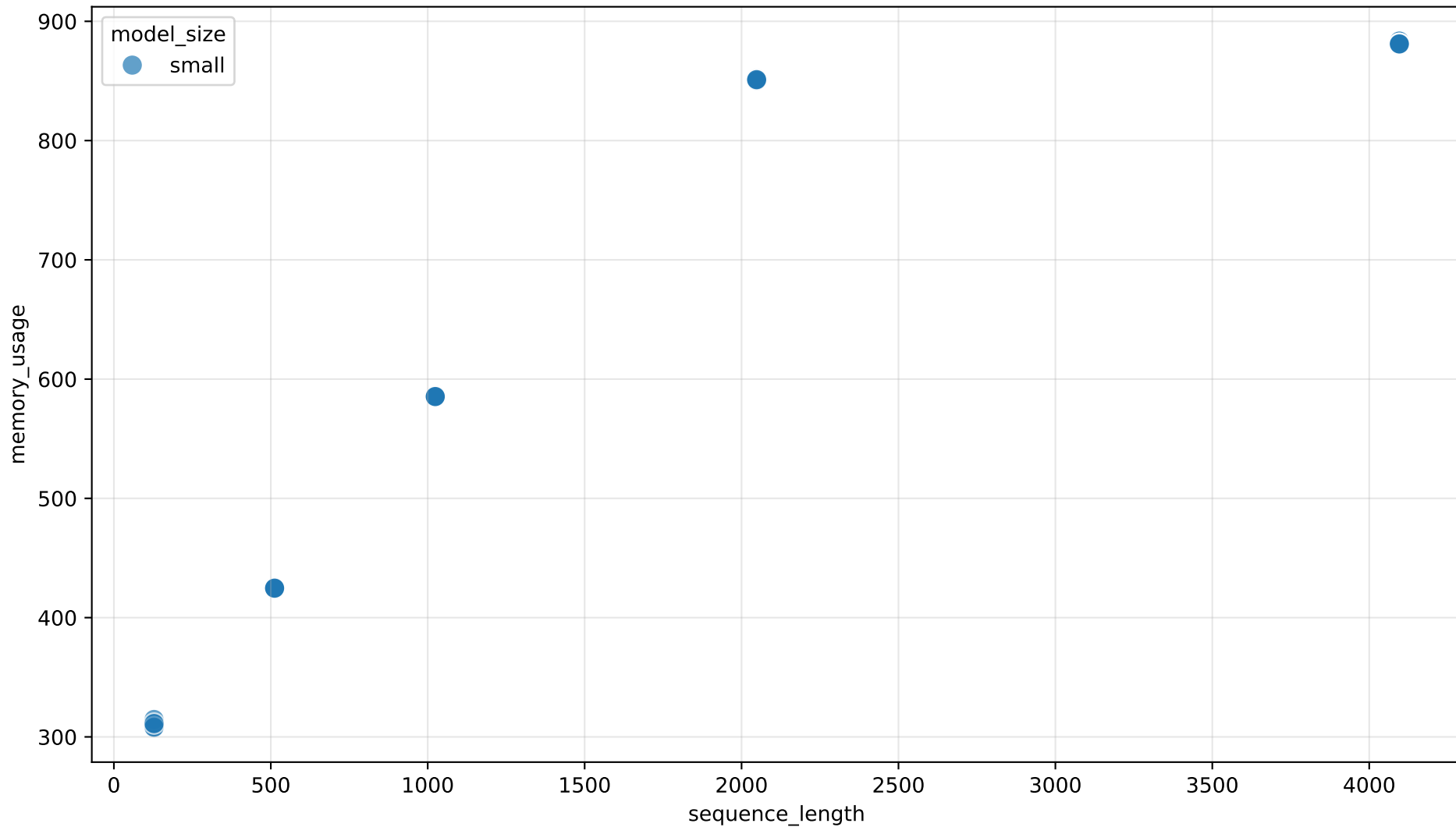
Phase 2 optimizations significantly improve performance on Intel hardware

Memory usage increases linearly with sequence length

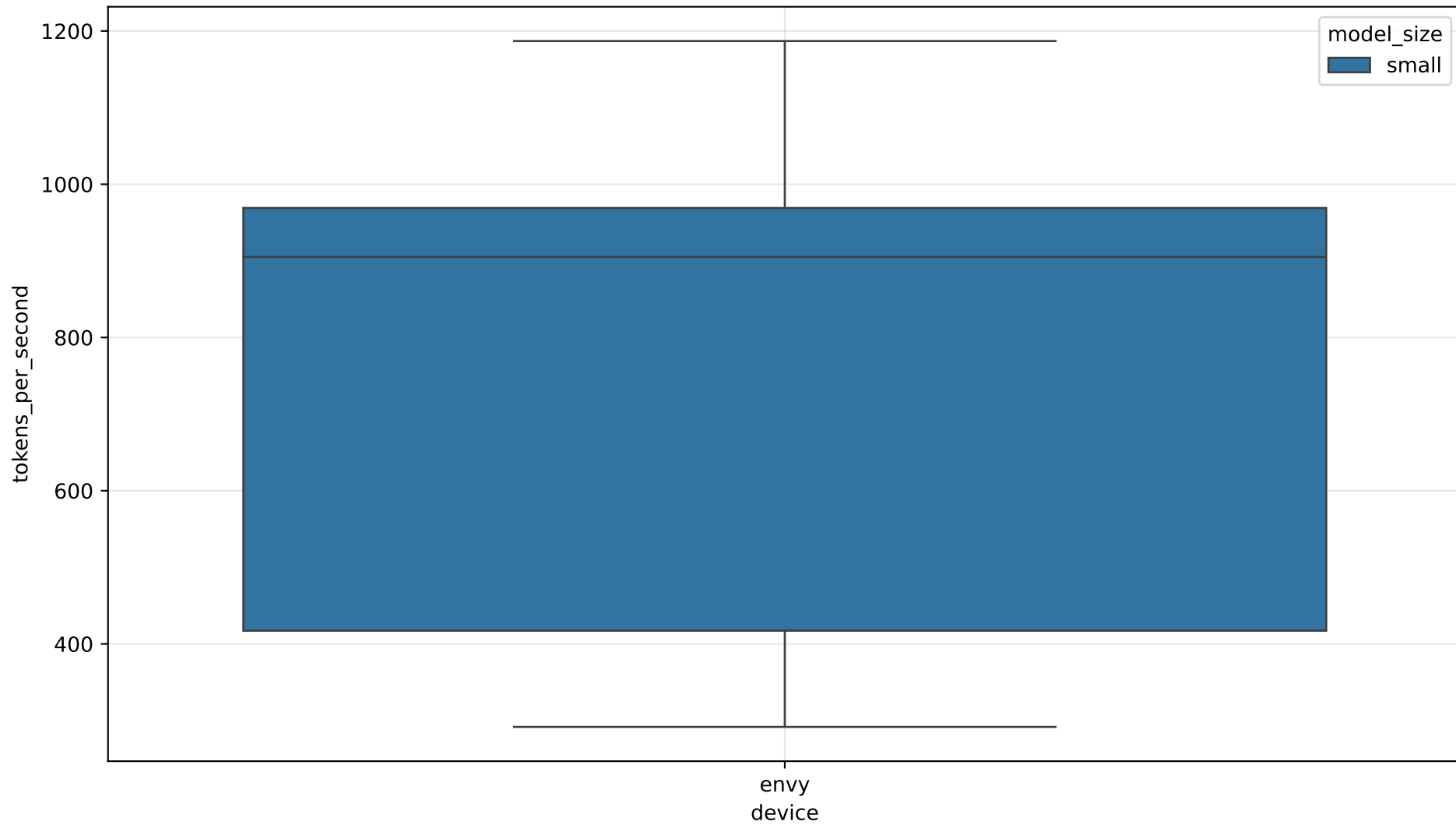
Performance vs Sequence Length



Memory Usage vs Sequence Length



Performance by Device



Performance Summary by Model Size

model_size	mean	min	max	std
small	747.98	291.66	1186.91	300.25

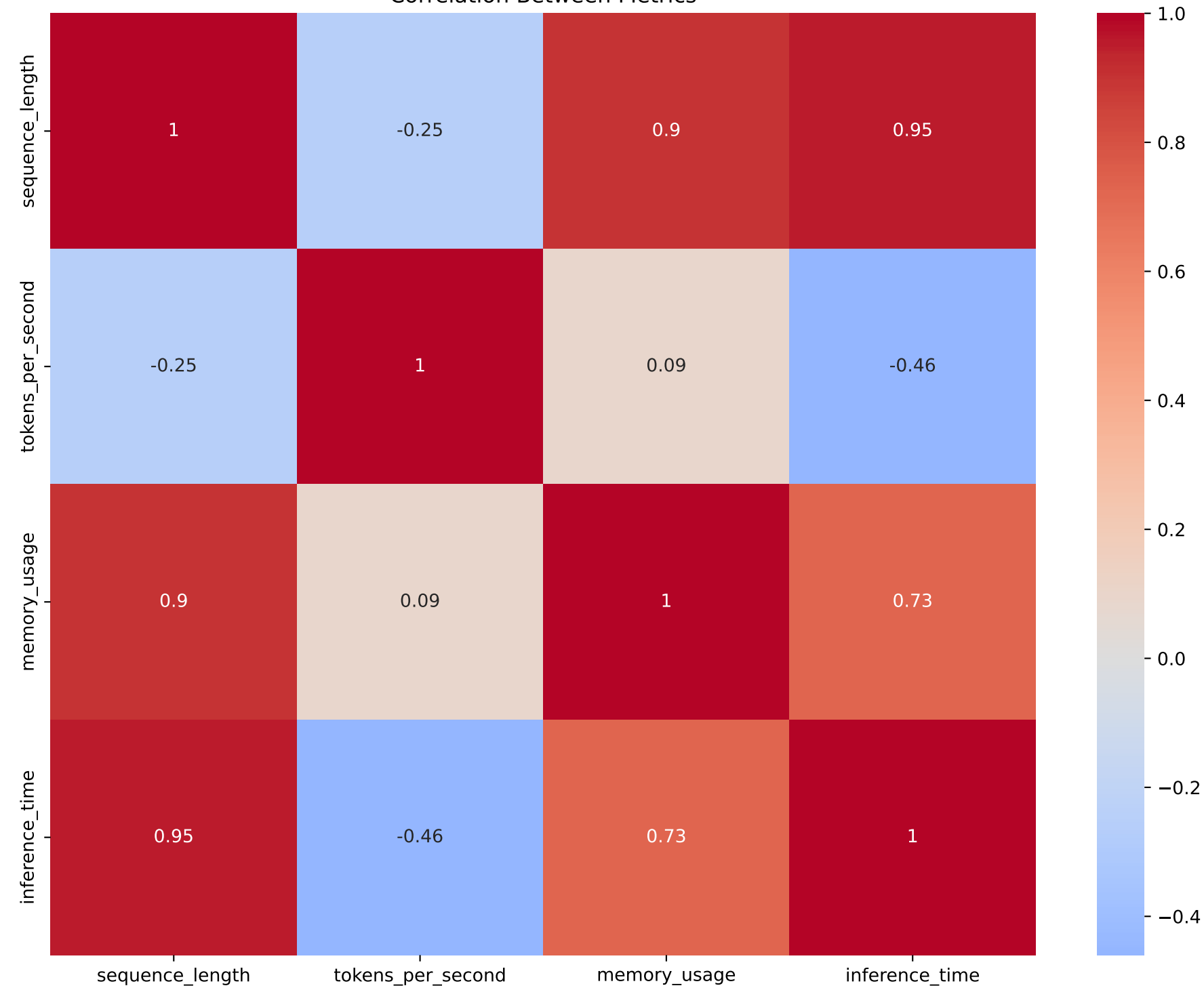
Performance Summary by Sequence Length

sequence_length	mean	min	max	std
128.0	366.79	291.66	413.95	57.65
512.0	934.7	807.62	1011.67	74.92
1024.0	1072.18	964.5	1186.91	88.39
2048.0	930.47	829.89	1020.23	65.52
4096.0	435.73	360.68	488.07	50.25

## Average Performance (tokens/s) by Device and Model Size

envy	small
	747.98

Correlation Between Metrics





Distribution of Performance

