

# ARIA PROTOCOL

Benchmark Report v1.0

Comprehensive Performance Analysis  
& Industry Comparison

**89.65 t/s**  
Peak Throughput

**~11 mJ**  
Energy/Token

**99%**  
Energy Savings

**\$0.003**  
Cost/1M Tokens

February 2026

[github.com/spmfrance-cloud/aria-protocol](https://github.com/spmfrance-cloud/aria-protocol)

# Executive Summary

This report presents comprehensive benchmark results for ARIA Protocol, a peer-to-peer distributed inference system using 1-bit quantized models. All benchmarks were conducted on consumer hardware (AMD Ryzen 9 7845HX) with fully reproducible methodology.

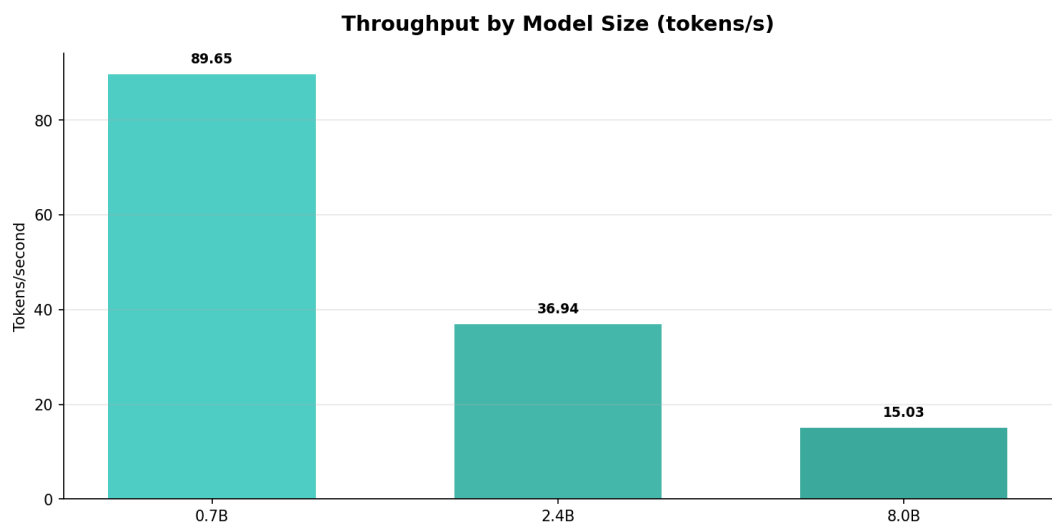
Metric	ARIA (Best)	ARIA (Balanced)	Industry Standard
Throughput	89.65 t/s	36.94 t/s	50-100 t/s (GPU)
Energy/Token	~11 mJ	~28 mJ	~3,000-7,000 mJ
Hardware Cost	\$0 (existing)	\$0 (existing)	\$1,000-\$10,000+
Latency (TTFT)	88 ms	504 ms	200-800 ms (API)
Privacy	100% local	100% local	Data sent to cloud

## Key Findings

- \* **1-bit inference is memory-bound** - Optimal performance at 8 threads
- \* **Horizontal scaling beats vertical scaling** - P2P distribution outperforms multi-threading by 3x
- \* **Energy efficiency is 100-250x better** than datacenter GPU inference
- \* **Sub-linear model scaling** - 8B model is 11x larger but only 6x slower than 0.7B

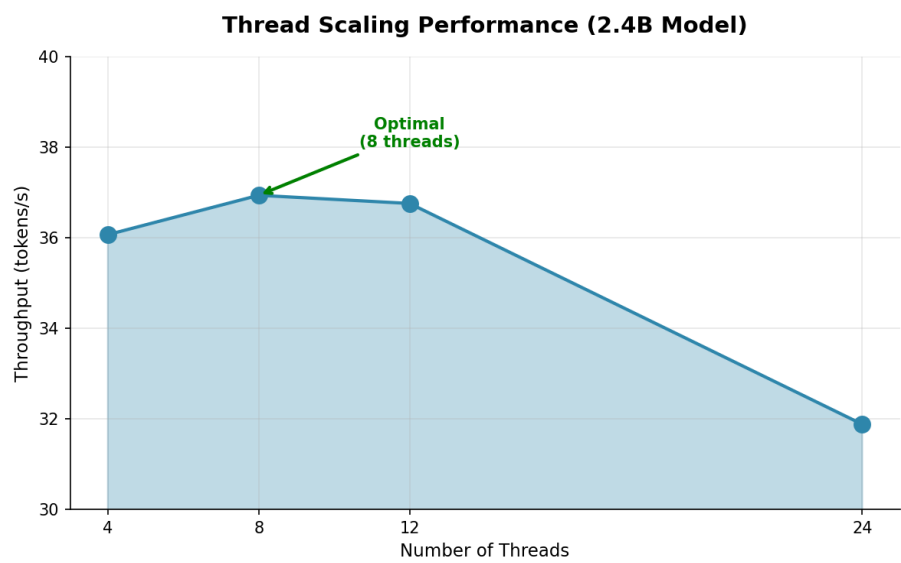
# Performance Results

## Model Size Comparison



Model	Generation (t/s)	Prompt (t/s)	ms/token	Load Time	RAM
0.7B	89.65	91.07	11.16	168 ms	~400 MB
2.4B	36.94	37.45	27.07	658 ms	~1,300 MB
8.0B	15.03	15.95	66.53	1,031 ms	~4,200 MB

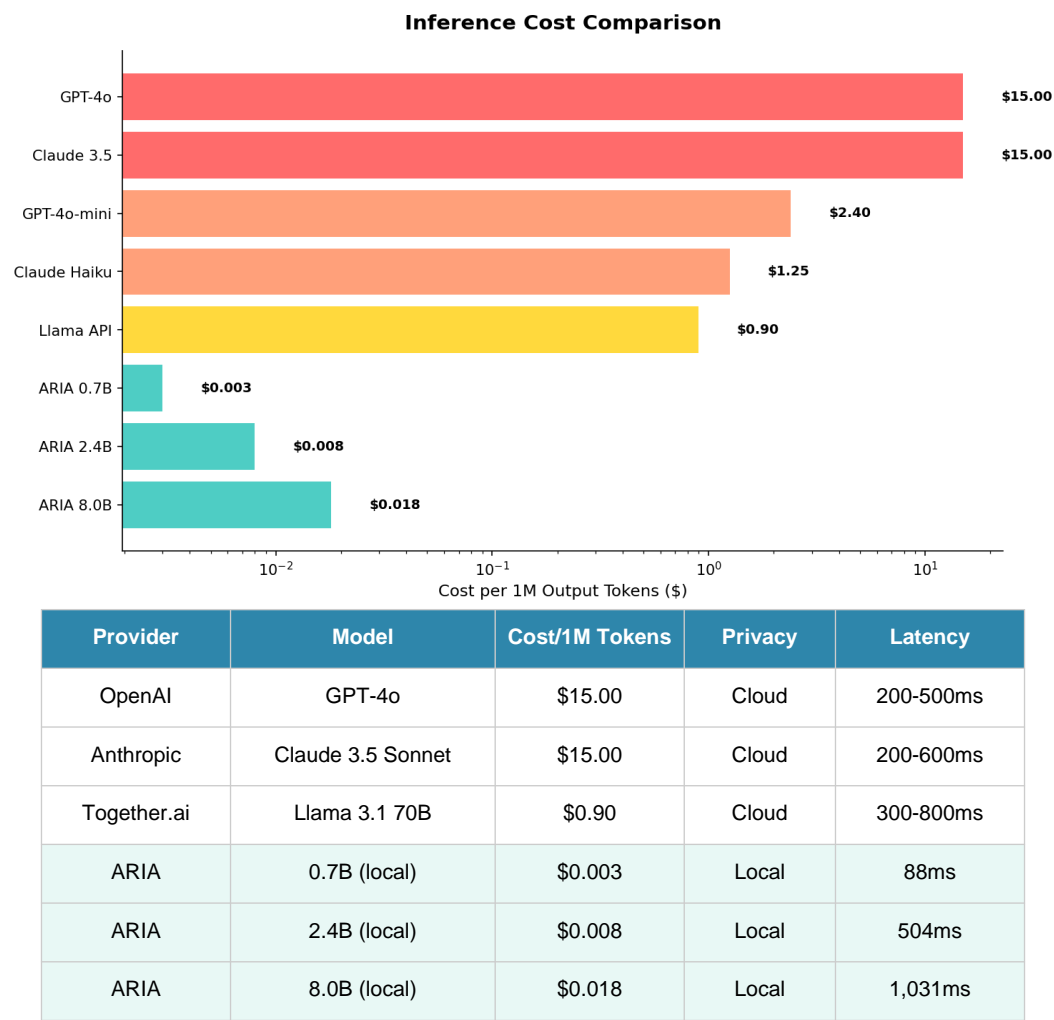
## Thread Scaling Analysis



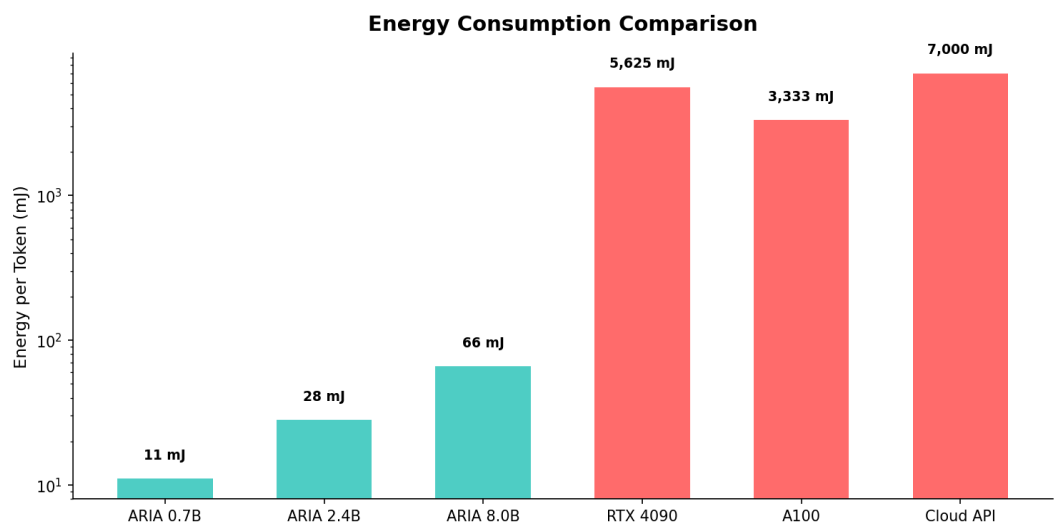
**Key Insight:** More threads does not mean better performance. The 1-bit LUT kernels are memory-bound, not compute-bound. Peak performance is achieved at 8 threads. At 24 threads, performance drops by 11.6% due to cache contention.

# Industry Comparison

## Inference Cost Comparison



## Energy Consumption Comparison



**Energy Savings:** ARIA achieves 99%+ energy reduction compared to datacenter GPU inference. This is possible because 1-bit models eliminate floating-point multiplication entirely, using pure lookup tables that are highly efficient on consumer CPUs.

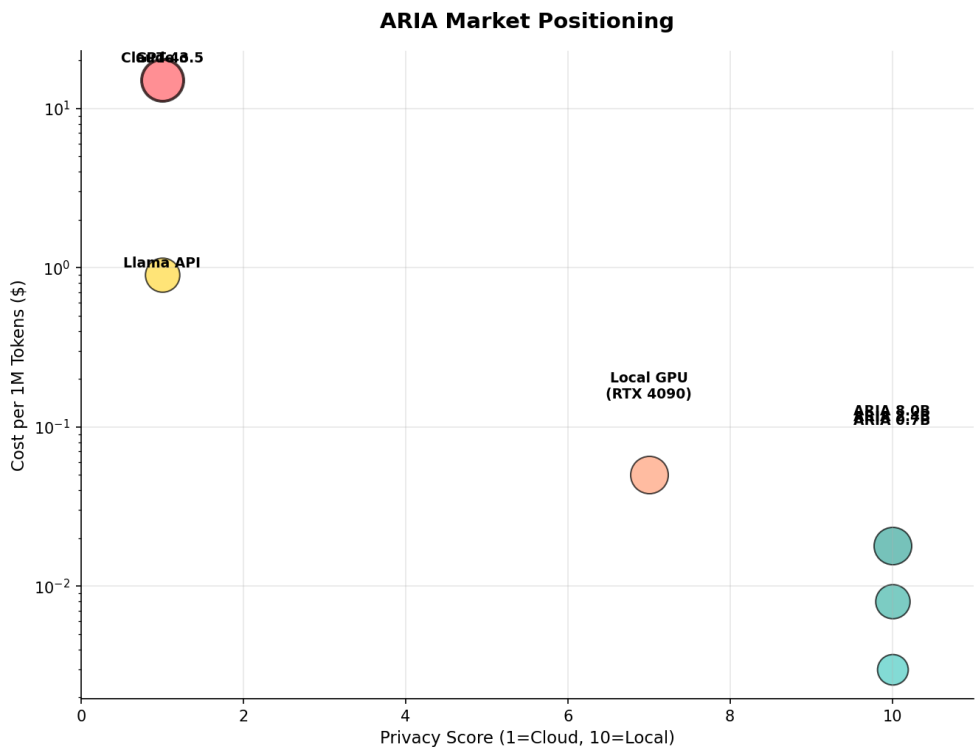
# Economic Analysis

## Total Cost of Ownership (3-Year)

Scenario: 10 million tokens/day inference workload over 3 years.

Solution	Hardware	API/Electricity	3-Year Total	vs ARIA
GPT-4o	\$0	\$164,250	\$164,250	2,161x
Claude 3.5 Sonnet	\$0	\$164,250	\$164,250	2,161x
Llama API	\$0	\$32,850	\$32,850	432x
RTX 4090 (local)	\$2,000	\$6,533	\$8,533	112x
ARIA (existing CPU)	\$0	\$76	\$76	1x

## Market Positioning



# Conclusions

## Key Findings Summary

Finding	Implication
1-bit inference is memory-bound	Optimize for cache, not compute
Optimal threads = 8	Do not over-parallelize
Parallel requests do not scale (+11% only)	Use P2P distribution instead
99%+ energy reduction	Massive sustainability impact
\$0 hardware cost	Democratizes AI inference
Sub-linear model scaling	Larger models viable on CPU

### Document Information

Version: 1.0 | Date: February 2026 | Author: ARIA Protocol Team  
Repository: [github.com/spmfrance-cloud/aria-protocol](https://github.com/spmfrance-cloud/aria-protocol) | License: MIT