

Boudica Torc: Right-Sized Generative AI for the Enterprise

Executive Summary

Generic, oversized Large Language Models (LLMs) are inefficient, costly, and present significant data privacy risks for the modern enterprise. The Boudica Torc engine shifts the paradigm from "Research AI" to "Engineering AI," empowering organizations to deploy domain-specific models that are precisely tailored to business needs. Built on an ultra-efficient, secure, on-premise C++ platform, Boudica Torc delivers faster, safer, and more accurate AI by eliminating cloud dependencies, grounding responses in factual data, and providing a clear return on investment.

The Challenge: The Inefficiency of Generic, Cloud-Based AI

While powerful, mainstream cloud-based AI solutions introduce critical friction points for enterprises:

- **Data Privacy & Sovereignty Risks:** Sending sensitive enterprise or patient data to third-party cloud APIs increases the attack surface and creates compliance challenges, even with Business Associate Agreements (BAAs).
- **Unpredictable Costs:** Usage-based, per-query, and subscription models create spiraling Operational Expenditures (OpEx) that are difficult to forecast and control, especially at scale.
- **High Latency & Unreliability:** Network-dependent performance results in unpredictable response times (500ms - 2s), and a loss of internet connectivity renders the system useless.

- **Infrastructure Bloat:** Traditional Python-based frameworks require complex virtual environments and massive dependency chains, increasing management overhead and performance bottlenecks.

The Boudica Solution: An Engineered Foundation for Enterprise AI

Boudica Torc addresses these challenges with a native C++17 implementation of the Transformer architecture, designed for bare-metal performance, total data control, and simplified deployment.

Data Sovereignty and "Zero-Bloat" Architecture

The core of Boudica Torc is its **"Zero-Bloat" C++ stack**. Unlike Python frameworks, Boudica compiles to a single binary with minimal dependencies. This allows for an "Edge Native" deployment where data lives and is processed entirely on the hospital's own servers or a local VPC.



Physically Isolated Data

Data is processed on-premise, never traversing the public internet. This "Trust Physics" approach is ideal for high-profile data and regulated industries.



Deterministic Low Latency

Running on the local network bus delivers consistent, sub-100ms response times. The system works even if the internet is down.



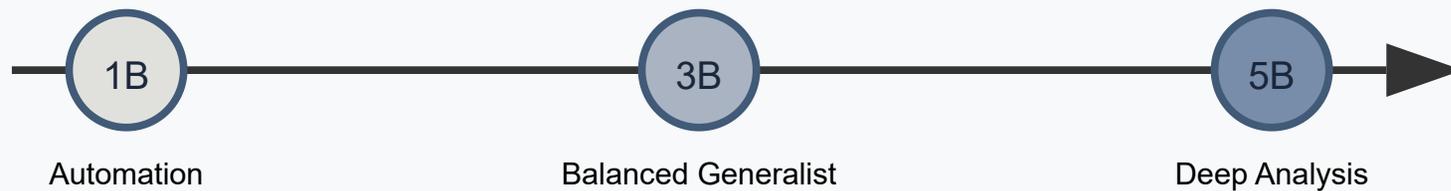
Simplified Deployment

A single C++ binary eliminates complex containerized microservices, reducing operational overhead and dependency management.

Choosing the Right Tool: Model Tiers for Business Needs

Boudica enables a fleet of specialized models where size and capability are matched to the task, maximizing efficiency and value.

Model Capability & Complexity



Model Tier	Primary Use Cases	Enabling Boudica Torc Feature
1 Billion Parameters (The Automation Engine)	Real-time PII sanitization, content classification, structured data extraction, and simple Q&A bots. Optimized for speed and high throughput.	Native C++ PiiSanitizer with validated logic (e.g., Luhn Algorithm) for high-accuracy redaction with deterministic low latency.
3 Billion Parameters (The Balanced Generalist)	Internal knowledge management (HR/IT Helpdesk), advanced customer support conversations, and first-draft generation for reports and emails.	Native LoRA (Low-Rank Adaptation) for parameter-efficient fine-tuning on company data and a Hybrid RAG that combines keyword and semantic search.
5 Billion Parameters (The Deep-Domain Analyst)	Complex legal contract analysis, financial report summarization, scientific research assistance, and cybersecurity threat intelligence. Requires deep reasoning.	The FactualityEnhancer calculates a RAG Grounding Score to flag hallucinations, while Citation Injection provides auditable, click-to-verify evidence for every claim.

Engineered Trust: Minimizing Hallucination for More Accurate Inferences

Boudica Torc is designed with a "Fact-Check Mode" to mitigate the risk of confidently incorrect answers. It provides precise answers with source citations, ensuring clinicians and analysts can trust the output.

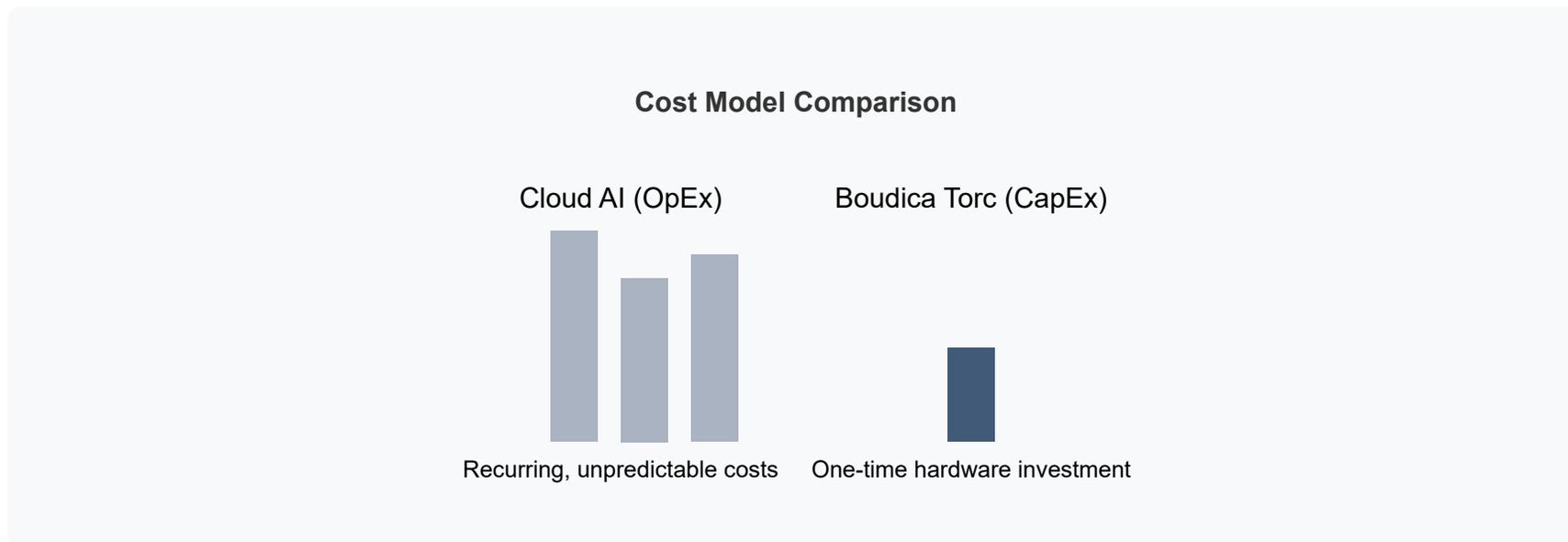
The Factuality Suite

- **Hybrid RAG Search:** Instead of relying on a separate vector database, Boudica leverages PostgreSQL directly. A two-stage retrieval process first uses keyword search for domain-specific terms and then semantic search for contextual meaning. This fetches more relevant information, leading to a more accurate final answer.
- **RAG Grounding Score:** After generating a response, the **FactualityEnhancer** module calculates the similarity between the model's output and the retrieved source documents. A low score flags a potential hallucination for human review.

- **Strict Citation Injection:** In its strictest mode, Boudica must provide a document source for every claim it makes. This creates a transparent and auditable trail, allowing users to "click-to-verify" the information and build trust in the system's output.

A Paradigm Shift in Cost: From Subscription OpEx to Hardware CapEx

Boudica Torc disrupts the expensive SaaS Scribe and AI market by converting a recurring operational expense into a one-time capital expenditure. For a large organization, this represents a fundamental shift in economic value.



For a hospital system with 1,000 providers, a typical cloud scribe solution can cost over \$2.4 Million per year. With Boudica, the cost is a one-time hardware upgrade (e.g., a \$500 GPU per workstation), leading to an ROI in under 3 months. By investing in on-premise hardware once, organizations can run the models forever without ongoing license fees.

Conclusion: The Power of Purpose-Built AI

Boudica Torc provides the essential infrastructure to move beyond the limitations of generic AI. By enabling the efficient, secure, and affordable deployment of domain-specific models, your organization can build an AI ecosystem that is:

- **Faster:** Reduced overhead and better memory locality from a C++ core.
- **Safer:** Total ownership of the data path with on-premise deployment.
- **Smarter:** More accurate, factual, and auditable responses grounded in your data.

With a single, self-contained binary, you gain the power to deploy the right model for the right job, on your own infrastructure, perfectly aligned with your business objectives.

Created by Boudica