

Whitepaper #01: Sovereign AI Infrastructure

Subtitle: *Architecting Intelligent Systems for Disconnected Environments*

1. Executive Summary

The Bottom Line Up Front (BLUF):

While the commercial AI revolution is driven by cloud-hosted APIs (e.g., OpenAI, Anthropic), organizations in defense, healthcare, and critical infrastructure face a stark reality: they cannot send their data to the cloud. For these sectors, the future of AI is not "Cloud First"—it is "Sovereign First."

The Core Problem:

Strict regulatory frameworks (ITAR, HIPAA, CUI) and Zero Trust security mandates often prohibit the use of external inference endpoints. Data leakage risks, however small, are unacceptable when the data involves national security or sensitive intellectual property. Furthermore, mission-critical systems often operate in "air-gapped" or disconnected environments where internet connectivity is physically impossible.

The Solution:

This brief outlines the reference architecture for a Sovereign Inference Unit (SIU)—a fully offline, self-contained AI pipeline. By repatriating workloads to local infrastructure, organizations achieve 100% data sovereignty, eliminate external dependencies, and ensure continuity of operations even during network blackouts. This document serves as a technical guide for architects sizing the hardware (VRAM), network topology, and software stack required to deploy Large Language Models (LLMs) behind the firewall.

2. The Operational Challenge: Why "Local" is Hard

Deploying AI in a connected enterprise is relatively simple: provision an API key, install a Python library, and start sending JSON requests. Deploying AI in a disconnected environment is a fundamental systems engineering challenge. It forces architects to solve problems that cloud providers usually abstract away.

2.1 The Cloud Disconnect

Modern AI development assumes ubiquitous connectivity. Toolchains rely on "lazy loading" resources from the internet:

- `pip install` fetches libraries from PyPI.

- `docker pull` fetches containers from Docker Hub.
- `AutoTokenizer.from_pretrained()` fetches weights from Hugging Face.

In a sovereign environment, all of these commands fail immediately. The "Operational Challenge" is not just running the model; it is rebuilding the entire supply chain of dependencies that allows the model to exist.

2.2 The "No-Egress" Mandate

The defining constraint of a secure facility is the **No-Egress Policy**. Data cannot leave the network. This creates two specific friction points:

1. **Ingress Difficulty:** Getting open-source innovation *into* the environment requires a rigorous "sneaker-net" or "diode" transfer process, involving scanning, hashing, and manual approval.
2. **No Telemetry:** You cannot use monitoring tools like Datadog or LangSmith. You must build your own observability stack (e.g., Prometheus/Grafana) to know if your model is hallucinating or failing.

2.3 The Throughput Reality Check

Stakeholders often expect "ChatGPT-level speed." This expectation must be managed.

- **Cloud Reality:** A massive cluster of H100 GPUs serving thousands of users with auto-scaling.
- **Local Reality:** A single workstation or small rack serving a specific team.

While local hardware cannot match the raw token-per-second generation of a massive cloud cluster, it offers a superior trade-off: **Consistency**. Your local latency is predictable, your queue is yours alone, and your service never goes down because a cloud provider has an outage.

3.1 The VRAM Equation: Calculating Hardware Requirements

In a disconnected environment, you cannot simply "scale up" an instance class when a model crashes. Hardware is fixed, procurement cycles are long, and the cost of under-provisioning is a failed deployment. Therefore, calculating Video Random Access Memory (VRAM) requirements is the single most critical step in architecting a sovereign inference unit.

Many developers rely on rough estimates, but for production-grade systems, we must use a precise formula to account for both the model weights and the dynamic memory required during inference.

The Foundational Formula

The total VRAM required (M_{total}) is the sum of three distinct components: the static model weights, the dynamic Key-Value (KV) cache (the "context window"), and the temporary activation buffers used during computation.

$$M_{\text{total}} \approx (P \times B_p) + (C_{\text{ctx}} \times L \times H \times B_{\text{kv}}) + O_{\text{sys}}$$

Where:

- P = Parameter Count (in billions)
- B_p = Bytes per parameter (determined by Quantization)
- C_{ctx} = Context Window Size (number of tokens)
- L = Number of Layers
- H = Hidden Size (dimension of the model)
- O_{sys} = System Overhead (display buffers, CUDA kernels)

1. The Weight Cost (Static Memory)

This is the baseline memory required just to load the model. It is determined strictly by the model's size and its quantization level.

- **FP16 (Half Precision):** The standard for uncompressed inference. Requires **2 bytes** per parameter.
 - *Example:* Llama-3-8B requires ≈ 16 GB of VRAM.
- **INT8 (8-bit Quantization):** Reduces size by 50% with negligible accuracy loss. Requires **1 byte** per parameter.
 - *Example:* Llama-3-8B requires ≈ 8 GB of VRAM.
- **INT4 (4-bit Quantization - Recommended):** The industry standard for local deployment (GGUF). Requires **0.5 bytes** per parameter (plus metadata overhead).
 - *Example:* Llama-3-8B requires $\approx 5-6$ GB of VRAM.

2. The Context Cost (Dynamic Memory)

This is where most RAG architectures fail. Loading the model is only half the battle; you must leave room for the "KV Cache"—the memory used to store the conversation history and ingested documents.

As the context window fills (e.g., loading a 50-page technical manual), memory usage grows linearly (or quadratically in older architectures). A standard 70B model might load comfortably on an A6000 (48GB), but if you attempt to utilize its full 32k context window, the KV cache can swell by an additional 10–12 GB, instantly triggering an Out of Memory (OOM) crash.

Architectural Note: For RAG pipelines processing large PDFs, we recommend reserving at least **20-30% of total VRAM** specifically for the KV cache.

3. The "Rule of Thumb" Sizing Guide

For rapid estimation during hardware procurement, use the following tiers based on the **GGUF** (**Q4_K_M**) quantization standard.

Model Class	Example Model	Min. VRAM (Inference Only)	Min. VRAM (RAG / Long Context)	Recommended Hardware
Small (Edge)	Llama-3-8B, Mistral-7B	6 GB	8 GB	NVIDIA RTX 3060 / 4060
Medium (Dev)	Mixtral 8x7B, Qwen-14B	14 GB	24 GB	NVIDIA RTX 3090 / 4090 (24GB)
Large (Pro)	Llama-3-70B, Command R+	40 GB	48 GB	NVIDIA RTX 6000 Ada / A6000 (48GB)
Massive (Cluster)	Grok-1, Llama-3-400B	250+ GB	320+ GB	4x or 8x A100/H100 Cluster

Key Takeaway: For sovereign environments, **VRAM is the bottleneck, not compute**. A slower token generation speed (tokens/sec) is acceptable for a chatbot; an OOM crash is not. When in doubt, prioritize VRAM capacity (e.g., a used 3090 with 24GB) over raw clock speed (e.g., a new 4070 Ti with only 16GB).

4. Network Architecture: The "Zero Trust" Setup

In a connected environment, the "network" is an enabler. In a sovereign environment, it is a constraint. The architecture must be designed to mimic the convenience of the open internet while adhering to strict "No-Egress" policies. We achieve this by building a "Local Internet"—a mirrored ecosystem that lives entirely within the secure boundary.

4.1 The Physical Air-Gap Topology

The foundational requirement is physical and logical isolation. The Sovereign Inference Unit (SIU) does not sit on the corporate WAN; it resides on a dedicated subnet with physically severed uplinks to the outside world.

- **The "Low Side" (Connected):** A sanitization station where initial assets (models, libraries, containers) are downloaded, scanned for CVEs (Common Vulnerabilities and Exposures), and hashed.
- **The "Data Diode" / Transfer Mechanism:** The strictly controlled pathway (often manual "sneaker-net" via encrypted drives or unidirectional network diodes) to move assets to the secure environment.
- **The "High Side" (Disconnected):** The production environment where the SIU operates. This network trusts *nothing* that hasn't passed through the sanitization station.

4.2 Internal Repositories: Building the "Local Internet"

Developers cannot function without dependencies. If a Python script fails because it cannot reach pypi.org, the project stalls. To mitigate this, the architecture must include local mirrors of critical public repositories.

- **Artifact Repository (The "PyPI Mirror"):**
Instead of pip install pandas, developers point their config to an internal instance of Sonatype Nexus or JFrog Artifactory.
 - *Configuration:* pip.ini or pip.conf is modified globally to target <https://nexus.internal/repository/pypi-hosted/>.
 - *Workflow:* A "Librarian" or automated pipeline pulls approved wheels (.whl) from the public internet, scans them, and uploads them to the internal Nexus.
- **Model Registry (The "Hugging Face Mirror"):**
Language models are too large (10GB–100GB) to treat like code. They should not be stored in Git.
 - *Solution:* A dedicated Network Attached Storage (NAS) or S3-compatible object store (e.g., [MinIO](#)) running locally.
 - *Structure:* Organized by Organization/Model-Name/Quantization (e.g., meta-llama/Llama-3-70b/GGUF/).
 - *Access:* Developers mount this drive as read-only volume, allowing instant model loading without duplicating terabytes of data across workstations.

4.3 Containerization Strategy: Docker vs. Apptainer

Shipping code to a disconnected server is risky. "It works on my machine" is not an acceptable standard when you cannot download a missing driver. Containerization is mandatory, but the tool choice matters.

- **Docker (The Development Standard):**
 - *Use Case:* Building images on the "Low Side."

- *The "Save/Load" Workflow:* We utilize `docker save -o image.tar my-app:latest` to create a single file artifact that can be scanned and transferred across the air-gap. Once inside, `docker load -i image.tar` restores it.
 - **Apptainer (formerly Singularity) (The HPC Standard):**
 - *Why use it?* Many secure government labs and High-Performance Computing (HPC) centers forbid Docker because the Docker daemon requires `root` privileges, creating a massive security vulnerability.
 - *The Advantage:* Apptainer runs containers as a standard user (rootless). It encapsulates the entire environment into a single `.sif` file (Singularity Image Format) which is immutable and cryptographically verifiable.
 - *Recommendation:* If deploying to a shared secure cluster, architect for Apptainer. If deploying to a dedicated standalone server, Docker is acceptable.
-

5. Software Stack: The "Sovereign Inference" Layer

Once the hardware is racked and the network is isolated, the final hurdle is the software stack. In a cloud environment, you might just hit an endpoint. In a sovereign environment, you *are* the endpoint. The software stack must be robust, efficient, and capable of operating without any external "phone home" telemetry.

5.1 Operating System & Drivers

Linux is the non-negotiable standard for high-performance AI. Windows (via WSL2) is acceptable for prototyping, but production Sovereign Units should run on **Ubuntu LTS or Red Hat Enterprise Linux (RHEL)**.

- **The Driver Challenge:** You cannot run `apt-get install nvidia-driver`.
- **The Solution:** You must download the **NVIDIA .run installer file** (a self-contained binary) on the low side and transfer it across.
- **Version Pinning:** Strictly pin your CUDA version (e.g., CUDA 12.1) to match your PyTorch/TensorFlow binaries. Mismatched CUDA versions in an offline environment can lead to weeks of debugging "DLL hell" without StackOverflow to help you.

5.2 The Inference Engine: Choosing Your Server

The "Inference Engine" is the software that loads the model weights and exposes an API (usually REST) for your applications to query.

- **Option A: Ollama (The Developer Standard)**
 - *Best For:* Rapid prototyping, single-user workstations, and "chat" interfaces.
 - *Why:* It bundles the model management and server into one binary. It provides an OpenAI-compatible API (`/v1/chat/completions`), meaning you can drop it into existing tools (like VS Code extensions or LangChain scripts) with zero code changes.
 - *Constraint:* Less optimized for high-concurrency (multiple users hitting it at once).

- **Option B: vLLM (The Production Standard)**
 - *Best For:* Multi-user RAG applications, high-throughput batch processing.
 - *Why:* It utilizes **PagedAttention**, a memory management technique that increases throughput by 24x over standard Hugging Face Transformers. If you are serving a department of 50 analysts, use vLLM.
 - *Constraint:* More complex setup; requires precise CUDA version matching.
- **Option C: Llama.cpp (The Compatibility King)**
 - *Best For:* Non-standard hardware (older CPUs, Apple Silicon, AMD cards).
 - *Why:* It runs on almost anything. If your sovereign hardware is an aging Intel Xeon server with no GPU, Llama.cpp is the only viable option.

5.3 Quantization: The GGUF Standard

In the cloud, models run at "BF16" (Brain Float 16). In local environments, this is wasteful. We utilize **Quantization** to compress models with minimal intelligence loss.

- **The Format: GGUF** (GPT-Generated Unified Format) is the file format of choice. It maps model tensors directly to memory (mmap), allowing models to load instantly without parsing.
- **The "Sweet Spot" (Q4_K_M):**
 - We recommend the **4-bit Medium Quantization (Q4_K_M)** for 95% of use cases.
 - *Metric:* A Llama-3-70B model at Q4 retains roughly 99% of the reasoning capability of the uncompressed model but consumes 40GB of VRAM instead of 140GB.
 - *Warning:* Avoid "Q2" or "Q3" quantizations for reasoning tasks; the logic degradation is severe (the model becomes "lobotomized").

6. Conclusion

Building a Sovereign AI capability is not merely a software project; it is a systems engineering challenge. It requires a holistic view that marries **hardware constraints** (VRAM), **network security** (Air-gaps), and **operational discipline** (Local repositories).

While the upfront architectural cost is higher than a simple API subscription, the strategic value is immense: **complete data sovereignty, zero leakage risk, and immunity to external outages.**

By following the architecture laid out in this brief—sizing hardware correctly, establishing local mirrors, and selecting efficient inference engines—organizations can deploy state-of-the-art intelligence into the most sensitive corners of their mission space.

Next in this Series:

- **Whitepaper #02:** *The Disconnected Pipeline: Solving Dependency Management & Containerization.*

About the Author

Dustin J. Ober, PMP, M.Ed. AI Developer & Technical Instructional Systems Designer

Dustin J. Ober is a specialist in the intersection of Artificial Intelligence, Instructional Strategy, and secure systems architecture. With a background spanning over two decades in the United States Air Force and defense contracting (Leidos, Peraton), he focuses on deploying high-impact technical solutions within mission-critical environments.

Unlike traditional developers who focus solely on code, Dustin bridges the gap between **technical capability** and **operational reality**. His expertise lies in architecting "Sovereign AI" systems—designing offline, air-gapped inference pipelines that allow organizations to leverage state-of-the-art intelligence without compromising data security or compliance.

He holds a Master of Education in Instructional Design & Technology and is a certified Project Management Professional (PMP). He actively develops open-source tools for the AI community, focusing on DSPy implementation, neuro-symbolic logic, and verifiable agentic workflows.

Connect:

- **Web:** aiober.com
 - **LinkedIn:** linkedin.com/in/dustinober
 - **Email:** dustinober@me.com
-

Suggested Citation: Ober, D. J. (2025). *Sovereign AI Infrastructure: Architecting Intelligent Systems for Disconnected Environments* (Whitepaper No. 01). AIOber Technical Insights.
<https://aiober.com/research/sovereign-ai-infrastructure>