

H.U.G.H. Distributed AI Protection System: Deployment Blueprint

Date: 2025-06-12

Executive Summary

This report provides a comprehensive deployment blueprint for the H.U.G.H. (Highly Uncompromising Guardian Heuristic) distributed AI protection system. The central objective is to equip you, Robert, with a detailed, executable plan to establish a robust AI ecosystem focused on family privacy protection against emerging AI surveillance threats. This blueprint prioritizes the deployment of high-performance local AI models, capable of competing with leading proprietary systems, across your Proxmox server, MacBook Air, and the Convex platform, thereby minimizing reliance on external APIs and enhancing data sovereignty.

The strategy encompasses a detailed architecture design, a model allocation plan leveraging state-of-the-art vocal, coding, and general-purpose AI models, and specific integration steps for Convex agents. It addresses critical considerations for local-to-cloud API interactions (with an emphasis on minimizing them), outlines robust resilience and failover mechanisms, and provides a phased, step-by-step implementation plan. Security considerations are paramount throughout this blueprint, ensuring the integrity and confidentiality of your family's data. Finally, guidance on integrating H.U.G.H. into your existing family network is provided.

This document includes necessary hardware specifications, software recommendations, download links, installation commands, and configuration examples. By following this blueprint, you can construct a powerful, private, and resilient AI system tailored to your family's protection needs, ensuring both high performance and peace of mind in an increasingly AI-driven world.

Table of Contents

1. Architecture Design
2. Model Allocation Strategy
3. Convex Agent Integration
4. Local-to-Cloud API Considerations
5. Resilience and Failover Strategies
6. Step-by-Step Implementation Plan
7. Security Considerations
8. Family Network Integration
9. References
10. Appendix: Quick Reference

1. Architecture Design

The H.U.G.H. system is designed as a distributed AI platform, leveraging the unique strengths of your Proxmox server, MacBook Air, and the Convex cloud platform to create a resilient, private, and powerful AI protection shield for your family. The architecture prioritizes local processing to mitigate surveillance risks associated with third-party AI services.

Core Architectural Principles

- **Local First:** Core AI processing, especially involving sensitive family data, will occur on local hardware (Proxmox server, MacBook Air).
- **Distributed Intelligence:** Tasks are distributed to the most appropriate node based on computational requirements, data locality, and specific AI model capabilities.
- **Resilience:** The system incorporates redundancy and failover mechanisms to ensure continuous operation.
- **Scalability:** The architecture allows for future expansion of AI capabilities and hardware resources.
- **Privacy by Design:** Security and privacy considerations are integrated into every layer of the system.

System Components and Roles

1. Proxmox VE Server Cluster (The Powerhouse):

- **Role:** Primary compute cluster for hosting demanding AI models, including large language models (LLMs), specialized coding models (Codestral 25.01), and high-fidelity vocal AI models (Whisper, OuteTTS). It will also run virtual machines (VMs) or containers for various H.U.G.H. services.
- **Setup:** Configured as a Proxmox VE cluster (ideally 3 nodes for robust quorum, or 2 nodes with an external quorum device) with shared storage (NFS, Ceph, or GlusterFS) to enable High Availability (HA) and live migration of H.U.G.H. service VMs.
- **Key Feature:** GPU passthrough will be configured to provide dedicated GPU access to VMs running AI models, ensuring near-native performance. This is critical for models like Llama 3.2 70B, Qwen2.5 72B, and Codestral 25.01.

2. MacBook Air (M-Series) (Personal & Edge Node):

- **Role:** Serves as a personal AI interaction point and an efficient node for running smaller, optimized LLMs (e.g., quantized 7B-13B models like Llama-3.1-8B) for tasks like quick queries, personal assistance, or on-device processing of less intensive H.U.G.H. functions. Its M-series chip with unified memory and Neural Engine is well-suited for optimized local AI.
- **Connectivity:** Connects to the H.U.G.H. network, capable of offloading heavier tasks to the Proxmox cluster or operating semi-autonomously for specific functions.
- **Use Cases:** Local document summarization, quick voice commands (potentially processed by a smaller local ASR), and as a resilient interface if the Proxmox cluster is temporarily unavailable for certain tasks.

3. Convex Platform (The Orchestrator & State Manager):

- **Role:** Acts as the central nervous system for H.U.G.H. It will not run AI models directly but will host the backend logic for H.U.G.H. agents, manage persistent state, orchestrate tasks across the Proxmox cluster and MacBook Air, and facilitate real-time communication and data synchronization.
- **Features Leveraged:** Durable execution for long-running protection workflows, real-time database for storing configurations, threat intelligence, and family preferences, and serverless functions for backend logic. Its ability to coordinate multi-agent systems is key for sophisticated protection strategies.
- **Interaction:** H.U.G.H. services on Proxmox and the MacBook Air will interact with the Convex backend via secure APIs.

4. Family Network & Devices:

- **Integration:** H.U.G.H. will interface with your home network to monitor traffic (with appropriate configurations), process voice commands from designated devices, and provide protection services to all connected devices (PCs, smartphones, IoT devices).
- **Interface:** A web-based dashboard, potentially built with Streamlit and hosted as part of the Convex backend or on a dedicated VM, will provide an interface for managing H.U.G.H. and viewing its status.

System Architecture Diagram

```

graph TD
    subgraph Family Network
        FN_Device1[Family PC 1]
        FN_Device2[Family Smartphone]
        FN_Device3[IoT Device]
        FN_Router[Home Router]
    end

    subgraph H.U.G.H. System
        subgraph Proxmox VE Cluster (Local Powerhouse)
            PVE1[Proxmox Node 1 (GPU Passthrough)]
            PVE2[Proxmox Node 2 (GPU Passthrough)]
            PVE_SharedStorage[Shared Storage (NFS/Ceph)]
            PVE1 -- Corosync --> PVE2
            PVE1 -- Access --> PVE_SharedStorage
            PVE2 -- Access --> PVE_SharedStorage
            VM_LLM[VM: General LLMs (Llama 3.2, Qwen2.5)]
            VM_Coding[VM: Coding AI (Codestral 25.01)]
            VM_Vocal[VM: Vocal AI (Whisper, OuteTTS)]
            VM_Services[VM: H.U.G.H. Core Services]
            PVE1 --- Hosts --> VM_LLM
            PVE1 --- Hosts --> VM_Coding
            PVE2 --- Hosts --> VM_Vocal
            PVE2 --- Hosts --> VM_Services
        end

        subgraph MacBook Air (Personal/Edge Node)
            MBA[MacBook Air M-Series]
            MBA_Ollama[Ollama: Lighter LLMs]
            MBA_Agent[H.U.G.H. Personal Agent]
            MBA --- Runs --> MBA_Ollama
            MBA --- Runs --> MBA_Agent
        end

        subgraph Convex Platform (Cloud Orchestrator - Secure & Controlled)
            Convex_Backend[Convex Backend Functions]
            Convex_DB[Convex Realtime Database]
            Convex_Dashboard[Streamlit Dashboard (via Convex)]
            Convex_Backend -- Manages --> Convex_DB
            Convex_Backend -- Serves --> Convex_Dashboard
        end
    end

    FN_Router <--> Internet
    FN_Router --- Local_Network --- PVE1
    FN_Router --- Local_Network --- PVE2
    FN_Router --- Local_Network --- MBA

    VM_Services -- API Call --> Convex_Backend
    MBA_Agent -- API Call --> Convex_Backend

    Convex_Backend -- Task Orchestration --> VM_Services
    Convex_Backend -- Task Orchestration --> MBA_Agent

```

```

FN_Device1 <--> FN_Router
FN_Device2 <--> FN_Router
FN_Device3 <--> FN_Router

FN_Device1 -- Interaction --> Convex_Dashboard
FN_Device2 -- Voice Commands (Processed by VM_Vocal) --> VM_Vocal

style PVE1 fill:#f9d,stroke:#333,stroke-width:2px
style PVE2 fill:#f9d,stroke:#333,stroke-width:2px
style MBA fill:#ccf,stroke:#333,stroke-width:2px
style Convex_Backend fill:#9cf,stroke:#333,stroke-width:2px

```

This architecture ensures that sensitive data processing remains local, while Convex provides robust orchestration and state management. The MacBook Air offers a flexible, personal interaction point, and the Proxmox cluster delivers the heavy computational power required for advanced AI models.

2. Model Allocation Strategy

The H.U.G.H. system will employ a diverse set of AI models, strategically allocated across your hardware resources (Proxmox server cluster and MacBook Air) to optimize performance, ensure privacy, and meet the system's functional requirements for vocal interaction, coding assistance, and general AI-driven protection.

Guiding Principles for Model Allocation

- **Resource Matching:** Allocate models based on their computational demands (VRAM, CPU, RAM) and the capabilities of the host hardware.
- **Latency Sensitivity:** Prioritize local hosting on the most responsive hardware for latency-sensitive tasks (e.g., real-time voice interaction).
- **Privacy:** Ensure models handling highly sensitive data run on fully controlled local infrastructure.
- **Task Specialization:** Dedicate specific nodes or VMs to particular types of AI models (e.g., a vocal processing VM).

Model Allocation Plan

A. Proxmox VE Server Cluster (Primary AI Workload Host)

The Proxmox cluster, equipped with high-VRAM GPUs (e.g., NVIDIA RTX 4090s or RTX 6000 Ada) via passthrough, will handle the most demanding AI models.

1. Vocal AI Models (Hosted in `VM_Vocal`):

- **Speech-to-Text (ASR): Whisper Large V3 (or community “Turbo” variant)**
 - **Reasoning:** State-of-the-art accuracy for transcription. Requires significant resources, making Proxmox ideal. Local deployment ensures privacy of voice data.
 - **Hardware Target:** VM with GPU passthrough (e.g., 1x RTX 4090, allocate ~10-16GB VRAM for this model).
- **Text-to-Speech (TTS): OuteTTS-0.2-500M**
 - **Reasoning:** Advanced features like voice cloning, natural speech synthesis. Good fit for Proxmox due to its size and capabilities.
 - **Hardware Target:** Can share the GPU in `VM_Vocal` or have a dedicated smaller GPU if concurrent heavy use is expected.
- **Alternative TTS:** Indic Parler-TTS if broader multilingual capabilities with specific Indian languages are a priority, or Orpheus 3B if a more permissively licensed, lighter model is sufficient for some tasks.

2. Coding Assistance Model (Hosted in VM_Coding):

- **Mistral Codestral 25.01**

- **Reasoning:** Top-tier performance for code generation, FIM, and analysis, with a large 256k token context window. This model is resource-intensive.
- **Hardware Target:** Dedicated VM with significant GPU passthrough (e.g., 1x RTX 4090 or RTX 6000 Ada, requiring 24GB+ VRAM).

3. General-Purpose Large Language Models (Hosted in VM_LLM):

- **Primary High-Performance LLM: Llama 3.2 70B (Quantized GGUF or full if multi-GPU)**

- **Reasoning:** Offers strong reasoning and instruction-following capabilities, approaching GPT-4 class performance. The 70B parameter model requires substantial VRAM. Quantized versions (e.g., Q4_K_M GGUF) can run on 48GB+ VRAM, while full precision would need multiple high-end GPUs.
- **Hardware Target:** VM with one or more high-VRAM GPUs (e.g., 1-2x RTX 6000 Ada or multiple RTX 4090s).

- **Alternative High-Performance LLM: Qwen2.5 72B**

- **Reasoning:** Similar to Llama 3.2 70B in capability, with strong multilingual support. Resource requirements are comparable.
- **Hardware Target:** Same as Llama 3.2 70B.

- **Secondary/Versatile LLM: Llama-3.1-8B-Instruct or Qwen2.5-7B**

- **Reasoning:** Excellent all-around models for tasks not requiring the absolute largest model. Can handle many H.U.G.H. protection logic, analysis, and interaction tasks. More easily hosted.
- **Hardware Target:** VM with GPU passthrough (e.g., 1x RTX 4090, 16-24GB VRAM).

- **Multimodal LLM (Optional, for future expansion): Llama-3.2-11B-Vision-Instruct or Qwen2-VL-7B-Instruct**

- **Reasoning:** If H.U.G.H. needs to process images or video for threat detection or analysis.
- **Hardware Target:** VM with GPU passthrough (e.g., 1x RTX 4090).

B. MacBook Air (M-Series) (Personal Interaction & Lighter Tasks)

The MacBook Air, with its M-series chip, will run smaller, optimized models for personal use and less demanding tasks, leveraging tools like Ollama and AirLLM.

1. Personal Assistant LLM: Llama-3.1-8B-Instruct (Quantized)

- **Reasoning:** A capable model that can run efficiently on Apple Silicon with quantization (e.g., 4-bit or 8-bit via Ollama). Suitable for quick questions, text summarization, and personal productivity tasks.
- **Deployment:** Via Ollama, optimized for M-series chips.

2. Lightweight Vocal AI (Optional/Fallback):

- Smaller Whisper variants (e.g., Whisper Base or Small) could potentially run on the MacBook Air for non-critical, personal voice notes if privacy on the device itself is paramount for that specific interaction. Performance will be lower than the Proxmox-hosted Whisper Large V3.
- **Reasoning:** Provides on-device ASR for quick, personal inputs without network traffic to the Proxmox server.

3. H.U.G.H. Personal Agent Interface:

- The MacBook Air can run a client application or interface for the H.U.G.H. Personal Agent, which communicates with the Convex backend and can trigger local models for specific tasks.

Model Deployment Tools

- **Proxmox VMs:**
 - **Ollama:** For easy deployment and management of GGUF models (Llama 3.2 70B quantized, Llama-3.1-8B).
 - **Hugging Face Transformers Library (Python):** For deploying models like Whisper, OutTTS, and Codestral 25.01 (if not available in GGUF or via Ollama directly). Use within Docker containers.
 - **vLLM or Text Generation Inference (TGI):** For optimizing serving of Transformer models like Codestral or the larger Llama/Qwen variants if high throughput is required.
- **MacBook Air:**
 - **Ollama:** Primary tool for running quantized LLMs optimized for Apple Silicon.
 - **AirLLM:** If specific optimizations beyond Ollama are needed for M-series chips.

This allocation strategy ensures that the most powerful models are run on capable hardware within the Proxmox cluster, maintaining privacy and performance, while the MacBook Air serves as an efficient personal AI node. Convex will orchestrate the interactions between these components and the AI models they host.

3. Convex Agent Integration

The Convex platform will serve as the intelligent orchestration and state management layer for the H.U.G.H. system. Its features, such as durable execution, real-time state synchronization, and serverless functions, are ideal for coordinating the distributed AI components running on your Proxmox server and MacBook Air, and for managing the logic of H.U.G.H.'s protective agents.

Role of Convex in H.U.G.H.

- **Central Orchestrator:** Convex functions will receive triggers (e.g., voice commands relayed from `VM_Vocal`, security alerts from network monitoring tools, user requests from the dashboard) and dispatch tasks to the appropriate AI models on Proxmox or the MacBook Air.
- **State Management:** Convex's real-time database will store:
 - H.U.G.H. system configuration (e.g., model endpoints, API keys for internal services).
 - Family preferences and personalization settings for AI interactions.
 - Knowledge base for H.U.G.H. agents (e.g., learned behaviors, threat intelligence).
 - Logs and history of H.U.G.H. operations and detected threats (summarized, with detailed logs remaining local if preferred).
- **Agent Logic:** Implement the core decision-making logic for H.U.G.H.'s protection agents as Convex functions. These agents can manage long-running workflows, such as monitoring for specific types of AI surveillance or coordinating responses to detected threats.
- **Real-time Dashboard Backend:** Power a Streamlit (or similar) dashboard providing users with an interface to interact with H.U.G.H., view system status, and manage settings.
- **Multi-Agent Coordination:** If H.U.G.H. evolves to use multiple specialized AI agents, Convex can manage their communication and collaborative tasks.

Setting up Convex for H.U.G.H.

The following steps outline the process to initialize and configure a Convex project for the H.U.G.H. system.

Prerequisites:

- * Node.js and npm installed on your development machine.
- * A Convex account.

Step 1: Initialize Convex Project

1. Install the Convex CLI:

```
bash
```

```
  npm install -g convex
```

2. Log in to your Convex account:

```
bash
```

```
  npx convex login
```

3. Create a new Convex project (e.g., `hugh-guardian`):

```
bash
```

```
  npx convex dev
```

This command initializes a new Convex project in the current directory, creates a `convex` subdirectory, and starts a local development server that syncs with the Convex cloud backend.

Step 2: Define Data Schema

Define the structure of your data in `convex/schema.ts`. This schema will enforce data types and relationships.

Example `convex/schema.ts`:

```

import { defineSchema, defineTable } from "convex/server";
import { v } from "convex/values";

export default defineSchema({
  // Table to store H.U.G.H. configurations
  configurations: defineTable({
    settingName: v.string(),
    settingValue: v.any(), // Use more specific types like v.string(), v.number(),
    v.object() as needed
    description: v.optional(v.string()),
  }).index("by_settingName", ["settingName"]),

  // Table for family preferences
  familyPreferences: defineTable({
    userId: v.string(), // Identifier for family member
    preferenceKey: v.string(),
    preferenceValue: v.any(),
    lastUpdated: v.number(), // Timestamp
  }).index("by_userId_preferenceKey", ["userId", "preferenceKey"]),

  // Table for detected events or threats
  securityEvents: defineTable({
    timestamp: v.number(),
    eventType: v.string(), // e.g., "NetworkIntrusionAttempt", "PrivacyLeakDetected"
    severity: v.union(v.literal("Low"), v.literal("Medium"), v.literal("High"), v.literal("Critical")),
    description: v.string(),
    sourceNode: v.optional(v.string()), // e.g., "Proxmox_VM_LLM", "MacBookAir"
    details: v.object({}), // Store structured details of the event
    status: v.union(v.literal("New"), v.literal("Investigating"), v.literal("Resolved"))
  },
  ).index("by_timestamp", ["timestamp"]),

  // Table for tasks dispatched to AI models
  aiTasks: defineTable({
    taskId: v.string(),
    taskType: v.string(), // e.g., "TranscribeAudio", "GenerateCode", "AnalyzeText"
    payload: v.any(),
    targetModel: v.string(), // e.g., "WhisperLargeV3", "Codestral2501"
    status: v.union(v.literal("Pending"), v.literal("Processing"), v.literal("Completed"), v.literal("Failed")),
    result: v.optional(v.any()),
    createdAt: v.number(),
    updatedAt: v.number(),
  }).index("by_taskId", ["taskId"]).index("by_status", ["status"]),
});

```

Step 3: Implement Core Functions (Mutations and Queries)

Create TypeScript files in the `convex/` directory (e.g., `convex/tasks.ts`, `convex/events.ts`) to define serverless functions (mutations for writing data, queries for reading data).

Example `convex/tasks.ts` for dispatching a task:

```

import { mutation, query } from "./_generated/server";
import { v } from "convex/values";

// Mutation to create a new AI task
export const createTask = mutation({
  args: {
    taskType: v.string(),
    payload: v.any(),
    targetModel: v.string(),
  },
  handler: async (ctx, args) => {
    const taskId = crypto.randomUUID(); // Generate a unique task ID
    await ctx.db.insert("aiTasks", {
      taskId,
      taskType: args.taskType,
      payload: args.payload,
      targetModel: args.targetModel,
      status: "Pending",
      createdAt: Date.now(),
      updatedAt: Date.now(),
    });
    return taskId;
  },
});

// Mutation to update task status and result
export const updateTaskResult = mutation({
  args: {
    taskId: v.string(),
    status: v.union(v.literal("Processing"), v.literal("Completed"),
v.literal("Failed")),
    result: v.optional(v.any()),
  },
  handler: async (ctx, args) => {
    const existingTask = await ctx.db
      .query("aiTasks")
      .withIndex("by_taskId", (q) => q.eq("taskId", args.taskId))
      .unique();

    if (!existingTask) {
      throw new Error(`Task with ID ${args.taskId} not found`);
    }

    await ctx.db.patch(existingTask._id, {
      status: args.status,
      result: args.result,
      updatedAt: Date.now(),
    });
    return { success: true };
  },
});

// Query to get pending tasks for a specific model
export const getPendingTasks = query({
  args: { targetModel: v.string() },
  handler: async (ctx, args) => {
    return await ctx.db
      .query("aiTasks")

```

```
.withIndex("by_status", q => q.eq("status", "Pending")) // This requires a composite index or filtering after query
    // For more complex filtering, you might need to adjust indexing strategy or filter in code
    .filter(q => q.eq(q.field("targetModel"), args.targetModel))
    .collect();
},
});
```

Step 4: Deploy Backend

Once you run `npx convex dev`, your functions and schema are automatically deployed to the Convex development environment. For production, you would use:

```
npx convex deploy
```

Step 5: Integrate H.U.G.H. Services with Convex

Your services running on Proxmox VMs and the MacBook Air will use the Convex client library (available for JavaScript/TypeScript, Python, Rust) to interact with these backend functions.

Example Python client usage (from a service on Proxmox):

```

from convex import ConvexClient
from convex.values import ConvexValue # For type hints if needed
import os
import time

# Best practice: Store Convex URL and deploy key in environment variables
CONVEX_URL = os.environ.get("CONVEX_URL") # e.g., "https://your-project-
name.convex.cloud"
CONVEX_DEPLOY_KEY = os.environ.get("CONVEX_DEPLOY_KEY") # If needed for write access
from backend scripts

client = ConvexClient(CONVEX_URL)
# If using a deploy key for mutations from a backend script:
# client.set_debug(True) # Optional for debugging
# client.set_auth(CONVEX_DEPLOY_KEY)

def report_security_event(event_type: str, severity: str, description: str, details: dict):
    try:
        client.mutation("events:createSecurityEvent", { # Assuming you create an
'events.ts' with this mutation
            "eventType": event_type,
            "severity": severity,
            "description": description,
            "details": details,
            "sourceNode": "Proxmox_VM_NetworkMonitor" # Example
        })
        print("Security event reported to Convex.")
    except Exception as e:
        print(f"Error reporting security event to Convex: {e}")

def poll_for_transcription_tasks():
    # This is a simplified polling example. For real-time, consider Convex subscrip-
    tions.
    while True:
        try:
            tasks = client.query("tasks:getPendingTasks", {"targetModel": "Whisper-
LargeV3"})
            for task in tasks:
                print(f"Processing task: {task['taskId']}")
                # ... (logic to get payload, call local Whisper model) ...
                # Dummy result
                transcription_result = "This is the transcribed text."
                client.mutation("tasks:updateTaskResult", {
                    "taskId": task['taskId'],
                    "status": "Completed",
                    "result": {"text": transcription_result}
                })
                print(f"Task {task['taskId']} completed.")
                time.sleep(10) # Poll every 10 seconds
        except Exception as e:
            print(f"Error polling for tasks: {e}")
            time.sleep(30) # Wait longer on error

```

By integrating Convex, H.U.G.H. gains a robust, scalable, and manageable backend, allowing you to focus on developing the AI-driven protection logic while Convex handles the complexities of distributed state and execution.

4. Local-to-Cloud API Considerations

While the primary philosophy of the H.U.G.H. system is **local-first processing** to maximize privacy and minimize surveillance risks, there might be exceptional scenarios where interaction with cloud-based APIs becomes necessary or highly beneficial. This section outlines considerations and best practices for such interactions, ensuring they are conducted with utmost security and control.

Scenarios for Potential Cloud API Use

1. Accessing Unique, Cutting-Edge Models (Temporarily or Strategically):

- Some highly advanced models (like certain versions of GPT-4o audio models mentioned in the research) might initially only be available via API. If a specific, non-sensitive task critically requires such a model and no local equivalent exists, controlled API access could be a *temporary* measure.
- **H.U.G.H. Policy:** This should be an exception, not the rule. A thorough risk assessment is required, and data sent to the API must be non-identifiable or anonymized.

2. System Updates and Model Downloads:

- Downloading new model weights (e.g., from Hugging Face, Mistral AI's GitHub) or system software updates inherently involves cloud interaction.
- **H.U.G.H. Policy:** Use official, secure channels. Verify checksums of downloaded files.

3. Non-Sensitive Utility Functions:

- For tasks that do not involve any family data, such as fetching general news for summaries (if H.U.G.H. has such a feature) or checking against public threat intelligence feeds.
- **H.U.G.H. Policy:** Ensure these services have strong privacy policies and that no personal data is inadvertently leaked.

4. Convex Platform Interaction:

- As outlined, Convex is a cloud platform used for orchestration and state management. Communication between local H.U.G.H. nodes and Convex is a form of local-to-cloud interaction.
- **H.U.G.H. Policy:** This is an accepted interaction, managed by Convex's security model (TLS encryption, authentication). Ensure your Convex project security settings are robust.

Best Practices for Secure API Interaction

Should any cloud API interaction be deemed necessary, the following best practices, drawn from the provided research, must be strictly adhered to:

1. API Inventory and Justification:

- Maintain a clear record of all external APIs H.U.G.H. interacts with.
- Each API use must have a strong justification, especially if it deviates from the local-first principle.

2. Robust Authentication and Authorization:

- Utilize strong authentication mechanisms like **OAuth 2.0** or **JSON Web Tokens (JWTs)** for accessing APIs.
- API keys must be securely stored (e.g., using HashiCorp Vault, Proxmox VE's encrypted storage, or environment variables managed securely within Docker/Kubernetes if applicable) and regularly rotated.
- Implement the principle of least privilege: API keys should only have permissions necessary for the specific task.

3. Encryption in Transit:

- All communication with external APIs must use **TLS 1.2 or TLS 1.3** to encrypt data in transit.
- Verify SSL/TLS certificates to prevent man-in-the-middle attacks.

4. API Gateway (If Managing Multiple Internal/External APIs):

- If H.U.G.H. exposes any internal APIs (e.g., for communication between its own distributed services, though direct secure channels are often better for internal) or consumes many external ones, an API gateway can centralize control, monitoring, and security policies.
- For H.U.G.H.'s scale, direct secure client-to-Convex and secure model access is likely sufficient, but be aware of gateway benefits if complexity grows.

5. Web Application Firewall (WAF):

- If H.U.G.H. exposes any service to the internet (e.g., a remote access dashboard, which should be heavily secured and considered carefully), a WAF can help protect against common web exploits.

6. Data Minimization:

- Send only the absolute minimum data required for the API call.
- Anonymize or pseudonymize data whenever possible before sending it to an external API. **No sensitive family PII should be sent to untrusted third-party APIs.**

7. Zero Trust Architecture Principles:

- Assume no implicit trust. Every interaction, even between internal H.U.G.H. components if they cross trust boundaries (though less relevant for a family-scale system than enterprise), should be authenticated and authorized.
- Network segmentation within your Proxmox environment can help isolate services.

8. Monitoring and Logging:

- Log all API calls, including request parameters (excluding sensitive data), responses, and any errors.
- Monitor API usage for anomalies or signs of misuse. Convex's logging can assist with this for its interactions.

9. Privacy-Enhancing Technologies (PETs) - Advanced:

- For extremely sensitive scenarios (likely beyond typical H.U.G.H. needs but good to be aware of):
 - **Homomorphic Encryption:** Allows computation on encrypted data.
 - **Differential Privacy:** Adds noise to data to protect individual privacy while allowing aggregate analysis.
- These are complex and typically used in large-scale data analysis, but the principles of data protection are relevant.

10. Regular Security Audits and Updates:

- Periodically review API integrations for security vulnerabilities.
- Keep client libraries and SDKs used for API interaction up to date.

Specific Consideration for Azure OpenAI Audio Models (as per research data):

The research mentions GPT-4o audio models available via Azure. If, for a highly specific, non-sensitive, experimental use case, you decide to use these:

- * Strictly follow Azure's security best practices for API keys and endpoint protection.
- * Ensure the Azure OpenAI Service resource is configured with appropriate network security (e.g., private endpoints if feasible within your setup).
- * Be fully aware of data handling policies for any data sent to Azure.

* This remains a deviation from the core H.U.G.H. local-first principle and should be approached with extreme caution and clear justification.

By rigorously applying these considerations, you can ensure that any necessary local-to-cloud interactions within the H.U.G.H. system are conducted as securely as possible, preserving the overarching goal of family privacy protection. The default stance should always be to find or develop a local alternative.

5. Resilience and Failover Strategies

Ensuring the H.U.G.H. system remains operational and responsive, even in the event of hardware failures or software issues, is crucial for providing consistent family protection. This section details resilience and failover strategies for the key components: the Proxmox server cluster, the MacBook Air, and the Convex orchestration layer.

A. Proxmox VE Server Cluster Resilience

The Proxmox cluster, hosting the core AI services, is the most critical component for system-wide resilience.

1. Proxmox VE High Availability (HA):

- **Clustering:** Implement a Proxmox VE cluster.
 - **Quorum:** For robust HA, a minimum of three Proxmox nodes is recommended to maintain quorum. If you have two physical servers, a third, lightweight quorum device (e.g., a Raspberry Pi running `corosync-qnetd`) is essential to prevent split-brain scenarios.
 - **Corosync:** Proxmox uses Corosync for cluster communication. Ensure redundant network paths for Corosync if possible (e.g., dedicated network interfaces).
- **Shared Storage:** Essential for HA and live migration. Options include:
 - **NFS:** Simpler to set up, suitable for many home labs. Requires a reliable NAS or dedicated NFS server.
 - **Ceph:** Integrated into Proxmox, provides distributed and resilient storage across cluster nodes. More complex but highly scalable and robust. Recommended if you have at least three nodes with sufficient local storage.
 - **GlusterFS:** Another distributed file system option.
- **HA-Managed VMs:** Configure H.U.G.H. service VMs (VM_LLM, VM_Coding, VM_Vocal, VM_Services) as HA resources in Proxmox. If a node fails, Proxmox HA will automatically restart these VMs on another available node in the cluster.
- **Live Migration:** Shared storage enables live migration of VMs between Proxmox nodes for maintenance without service interruption.

2. Network Redundancy:

- **Bonded NICs:** Use Linux bonding (e.g., LACP, active-backup) for Proxmox host network interfaces to provide link aggregation and failover.
- **Multiple Switches:** For critical connections, consider using two network switches to protect against switch failure.
- **VLANs:** Segment network traffic (e.g., management, VM traffic, storage traffic) for better organization and security, which can also aid in isolating fault domains.

3. Backup and Recovery:

- **Proxmox Backup Server (PBS):** Implement regular backups of all critical VMs using PBS or another reliable backup solution. PBS offers efficient, deduplicated backups.
- **VM Configuration Backups:** Regularly back up VM configuration files.

- **Model Weight Backups:** Store downloaded AI model weights on a separate, backed-up storage location, not just within the VMs.

4. Power Redundancy:

- Utilize Uninterruptible Power Supplies (UPS) for all Proxmox nodes, network switches, and shared storage devices to protect against power outages and enable graceful shutdowns.

B. MacBook Air Resilience

The MacBook Air serves as a personal node and can offer limited fallback capabilities.

1. Local Model Operation:

- If the Proxmox cluster is down, the MacBook Air can continue to run its locally hosted, smaller AI models (e.g., quantized Llama-3.1-8B via Ollama) for personal tasks.
- It can provide basic H.U.G.H. agent functionality if its local agent can operate with cached data or in a degraded mode without full Convex connectivity.

2. Data Synchronization:

- For critical personal data processed or generated on the MacBook Air by H.U.G.H. agents, ensure it's periodically synced to a backed-up location or, if appropriate and secure, to the Convex database (e.g., summarized, non-sensitive results).

3. Independent Operation:

- Design H.U.G.H. personal agent features on the MacBook Air to function independently for core tasks where possible, reducing reliance on constant connectivity to the Proxmox cluster for those specific functions.

C. Convex Platform Resilience

Convex is a managed cloud service designed for high availability and data durability.

1. Inherent Resilience:

Convex handles its own infrastructure resilience, data replication, and backups. Your primary responsibility is secure and correct usage of the platform.

2. Client-Side Retry Logic:

- In your H.U.G.H. services (on Proxmox and MacBook Air) that interact with Convex, implement robust retry logic with exponential backoff for API calls to handle transient network issues or temporary Convex unavailability.

3. Local Caching (Strategic):

- For critical configurations or frequently accessed data from Convex, H.U.G.H. services can implement a local caching mechanism (e.g., in Redis on a Proxmox VM or on-disk) to continue limited operation if Convex is unreachable. The cache should have a reasonable TTL and a clear invalidation strategy.
- This is particularly important for configurations needed by local AI models to function.

D. Overall System Failover Strategy

1. Graceful Degradation:

Design H.U.G.H. services to degrade gracefully if certain components are unavailable. For example, if the primary LLM on Proxmox is down, the system might fall back to a smaller LLM or provide reduced functionality rather than failing completely.

2. Health Monitoring and Alerting:

- Implement health checks for all H.U.G.H. services and AI models.
- Use a monitoring system (e.g., Prometheus and Grafana, or simpler custom scripts) to track the status of Proxmox nodes, VMs, Convex connectivity, and AI model responsiveness.
- Configure alerts (e.g., email, mobile notification via a secure service) for critical failures.

3. Contingency Plans for Model Unavailability:

- If a specific AI model (e.g., Codestral 25.01) becomes unavailable due to corruption or other issues, have a plan to:
 - Quickly redeploy it from a backup or fresh download.
 - Temporarily switch to an alternative, less capable model if one exists and is configured.
 - Clearly notify the user (Robert) of the reduced functionality.

Proxmox HA Visualization (Conceptual)

```

graph LR
    subgraph Proxmox Cluster
        PVE_Node1[Node 1 (Active VM_LLM)]
        PVE_Node2[Node 2 (Standby)]
        PVE_Node3[Node 3 (Quorum/Active other VMs)]
        SharedStorage[(Shared Storage: NFS/Ceph)]

        PVE_Node1 --- Corosync & HA --> PVE_Node2
        PVE_Node1 --- Corosync & HA --> PVE_Node3
        PVE_Node2 --- Corosync & HA --> PVE_Node3

        PVE_Node1 --- R/W --> SharedStorage
        PVE_Node2 --- R/W --> SharedStorage
        PVE_Node3 --- R/W --> SharedStorage

        VM_LLM1([VM: LLM Service])
        PVE_Node1 --- Hosts --> VM_LLM1
    end

    User[User Interaction] --> VM_LLM1

    %% On Failure of PVE_Node1
    subgraph Proxmox Cluster - After Failover
        PVE_Node1_Failed[Node 1 (Failed)]
        PVE_Node2_Active[Node 2 (Active VM_LLM - Migrated)]
        PVE_Node3_Quorum[Node 3 (Quorum/Active other VMs)]
        SharedStorage_Failover[(Shared Storage: NFS/Ceph)]

        PVE_Node2_Active -- Corosync & HA --> PVE_Node3_Quorum
        PVE_Node1_Failed -. Failed Link .-> PVE_Node2_Active
        PVE_Node1_Failed -. Failed Link .-> PVE_Node3_Quorum

        PVE_Node2_Active --- R/W --> SharedStorage_Failover
        PVE_Node3_Quorum --- R/W --> SharedStorage_Failover

        VM_LLM2([VM: LLM Service])
        PVE_Node2_Active --- Hosts --> VM_LLM2
    end

    User --> VM_LLM2

    style PVE_Node1_Failed fill:#ff0000,stroke:#333,stroke-width:2px

```

This diagram illustrates how a VM (VM_LLM) running on PVE_Node1 would be automatically migrated and restarted on PVE_Node2 if PVE_Node1 fails, assuming HA and shared storage are correctly configured.

By implementing these strategies, the H.U.G.H. system will achieve a high degree of resilience, ensuring continuous protection and availability for your family.

6. Step-by-Step Implementation Plan

This section outlines a phased implementation plan for deploying the H.U.G.H. system. Each phase includes key activities, estimated timelines (conceptual), and dependencies. This plan assumes you have already acquired the necessary initial hardware (Proxmox server(s), MacBook Air, GPUs).

Overall Timeline Estimate: 8-16 weeks, depending on your familiarity with the technologies and time commitment.

Phase 0: Preparation and Foundation (Pre-requisite)

- **Activities:**

1. Finalize hardware selection and procurement (Proxmox servers, GPUs like RTX 4090/6000 Ada, RAM, NVMe SSDs, network gear, UPS).
2. Set up physical hardware, including network cabling.
3. Familiarize yourself with Proxmox VE, Docker, Ollama, and Convex documentation.

• **Time Estimate:** 2-4 weeks (concurrent with initial hardware setup).

• **Dependencies:** Hardware availability.

Phase 1: Proxmox VE Cluster and Core Local AI Services Setup

• **Objective:** Establish the Proxmox VE cluster and deploy initial, foundational AI services for vocal processing and basic LLM capabilities.

• **Time Estimate:** 3-5 weeks.

• **Activities:**

1. **Proxmox VE Installation & Clustering (Week 1-2):**

- Install Proxmox VE on all server nodes.
- Configure the Proxmox VE cluster (including Corosync and quorum device if using a 2-node setup).


```
bash
# Example: Add a node to cluster (on new node, after installing Proxmox)
# pvecm add IP_OF_EXISTING_NODE
```
- Set up shared storage (NFS, Ceph, or GlusterFS). For NFS on a NAS:
 - On NAS: Create NFS share, allow Proxmox node IPs.
 - On Proxmox: Datacenter -> Storage -> Add -> NFS.
- Configure basic network settings (bridges, bonds).

2. **GPU Passthrough Configuration (Week 2):**

- Enable IOMMU (VT-d/AMD-Vi) in server BIOS.
- Configure Proxmox host for GPU passthrough (modify GRUB, load `vfio-pci` modules).


```
bash
# Example: Check IOMMU groups
# find /sys/kernel/iommu_groups/ -type l
# Add to /etc/default/grub: GRUB_CMDLINE_LINUX_DEFAULT="quiet in-
tel_iommu=on" (for Intel)
# update-grub
```

```

# echo "vfio" >> /etc/modules
# echo "vfio_iommu_type1" >> /etc/modules
# ... (more steps involving blacklisting drivers, vendor IDs)

```

- Test passthrough with a simple VM.

3. `VM_Vocal` Setup (Week 3):

- Create a Linux VM (e.g., Ubuntu Server 22.04) on Proxmox.
- Assign a passthrough GPU to `VM_Vocal`.
- Install Docker, NVIDIA drivers, and CUDA toolkit within the VM.
- Deploy Whisper Large V3 (e.g., using Hugging Face Transformers in a Docker container).

```

python
# Python snippet for Whisper in Docker container
# from transformers import pipeline
# import torch
# device = "cuda:0" if torch.cuda.is_available() else "cpu"
# pipe = pipeline("automatic-speech-recognition", model="openai/whisper-large-
v3", device=device)
# def transcribe(audio_path):
#     return pipe(audio_path, generate_kwargs={"language": "english"})

```

- Deploy OuteTTS-0.2-500M (similarly, in a Docker container).

4. `VM_LLM_Basic` Setup (Week 4):

- Create another Linux VM with GPU passthrough.
- Install Ollama.

```

bash
curl -fsSL https://ollama.com/install.sh | sh

```

- Pull and run Llama-3.1-8B-Instruct via Ollama.

```

bash
ollama pull llama3.1:8b-instruct
ollama run llama3.1:8b-instruct "Hello! How are you?"

```

- Expose Ollama API for internal network access.

5. Initial Proxmox HA Configuration (Week 5):

- Define HA groups and add `VM_Vocal` and `VM_LLM_Basic` as HA-managed resources.
- Test failover by shutting down one Proxmox node (if cluster size permits).

• **Dependencies:** Phase 0 completion. Stable Proxmox host setup.

• **Contingency:** If GPU passthrough is problematic, start with CPU-only models or smaller GPU-accelerated models that are less demanding, then revisit passthrough.

Phase 2: MacBook Air Integration

• **Objective:** Integrate the MacBook Air as a personal AI node.

• **Time Estimate:** 1-2 weeks.

• **Activities:**

1. Ollama Setup on MacBook Air (Week 1):

- Install Ollama for macOS.
- Pull and run a quantized version of Llama-3.1-8B-Instruct or a similar M-series optimized model.

```

bash

```

```
# On MacBook Air
ollama pull llama3.1:8b-instruct # Or a specific quantized version if available
```

2. Develop/Install H.U.G.H. Personal Agent Stub (Week 1-2):

- Set up a basic application or script environment (e.g., Python) for the H.U.G.H. Personal Agent.
- Implement basic interaction with the local Ollama instance.

- **Dependencies:** Functional MacBook Air.

- **Contingency:** If Ollama performance is insufficient, explore AirLLM or other Apple Silicon specific optimization tools.
-

Phase 3: Convex Orchestration Layer Setup

- **Objective:** Deploy the Convex backend for H.U.G.H. state management and task orchestration.

- **Time Estimate:** 2-3 weeks.

- **Activities:**

1. Convex Project Initialization (Week 1):

- Set up Convex account and CLI.
- Initialize the `hugh-guardian` Convex project as per Section 3.
- Define initial `convex/schema.ts` for configurations, tasks, and events.

2. Implement Core Convex Functions (Week 1-2):

- Develop basic mutations and queries for task management (`createTask`, `updateTaskResult`) and event logging.
- Implement functions for storing/retrieving H.U.G.H. configurations.

3. Integrate Proxmox Services with Convex (Week 2-3):

- In `VM_Vocal` and `VM_LLM_Basic`, install Convex Python client.
- Develop scripts/services that:
 - Poll Convex for tasks (e.g., `VM_Vocal` polls for transcription tasks).
 - Report task completion and results back to Convex.
 - Log significant events (e.g., model errors, successful processing) to Convex.

```
```python
```

## Example: Service in VM\_Vocal polling Convex

---

```
from convex import ConvexClient
```

---

```
CONVEX_URL =
"YOUR_CONVEX_PROJECT_URL"
```

---

```
client = ConvexClient(CONVEX_URL)
```

---

```
Add CONVEX_DEPLOY_KEY if mutations are
from a backend script
```

---

```
while True:
```

---

```
tasks = client.query("tasks:getPendingTasks",
{ "targetModel": "WhisperLargeV3" })
```

---

```
for task in tasks:
```

---

```
process task with local Whisper model
```

---

```
client.mutation("tasks:updateTaskResult",
{ ... })
```

---

```
time.sleep(5)
```

---

```
...
```

#### 4. Integrate MacBook Air Agent with Convex (Week 3):

- Enable the H.U.G.H. Personal Agent on the MacBook Air to communicate with Convex for task requests or state updates.
  - **Dependencies:** Phase 1 completion. Internet connectivity for Convex.
  - **Contingency:** If direct polling from VMs to Convex is inefficient, explore Convex's real-time subscription features for more responsive task notifications.
- 

### Phase 4: Advanced AI Models and Full Resilience

- **Objective:** Deploy more powerful AI models and enhance system resilience.

- **Time Estimate:** 2-4 weeks.

- **Activities:**

1. **VM\_Coding Setup (Week 1):**

- Create a dedicated VM for Codestral 25.01 with high-VRAM GPU passthrough.
- Install necessary dependencies (Python, Transformers, potentially vLLM).
- Download and deploy Codestral 25.01 (following Mistral AI's local deployment guidelines once available).
- Integrate with Convex for coding-related tasks.

2. **Upgrade VM\_LLM (Week 2-3):**

- If hardware permits (e.g., RTX 6000 Ada or multiple RTX 4090s), deploy Llama 3.2 70B (quantized GGUF via Ollama/llama.cpp, or full model via vLLM/TGI if resources allow) or Qwen2.5 72B.
- Test performance and integration with Convex for complex reasoning tasks.

```
bash
Example: Running a large GGUF model with llama.cpp
./main -m ./models/llama-3.2-70b.Q4_K_M.gguf -p "Describe quantum
entanglement." -n 512 --gpu-layers 99
```

3. **Implement Advanced Proxmox HA and Backup (Week 3-4):**

- Thoroughly test Proxmox HA failover scenarios.
- Set up Proxmox Backup Server and configure regular backup schedules for all critical VMs.
- Refine network redundancy.

- **Dependencies:** Phase 1, 3 completion. Availability of Codestral 25.01 local deployment package. Sufficient GPU hardware.

- **Contingency:** If 70B+ models are too demanding, scale down to the largest feasible model (e.g., 30B-40B class) or use more aggressively quantized versions.
- 

### Phase 5: Family Network Integration & UI

- **Objective:** Integrate H.U.G.H. into the family network and provide a user interface.

- **Time Estimate:** 1-2 weeks.

- **Activities:**

1. **Network Monitoring Setup (Optional, Advanced - Week 1):**

- If H.U.G.H. is to perform network traffic analysis for threats, configure a VM or physical device with tools like Suricata or Zeek. This requires careful planning for port mirroring on your router/switch.
- Feed alerts from these tools into Convex via a custom script.

## 2. Voice Interaction Points (Week 1):

- Set up microphones or use existing smart devices (with custom integration if possible, or by routing their audio to `VM_Vocal`) to capture voice commands for H.U.G.H.

## 3. Develop H.U.G.H. Dashboard (Week 1-2):

- Use Streamlit (or a similar framework) to build a web dashboard.
- Host the dashboard (e.g., on a dedicated VM or using Convex's capabilities if it supports such hosting).
- The dashboard should query Convex to display system status, recent events, and allow interaction with H.U.G.H. functions.

```
python
Example Streamlit app snippet (app.py)
import streamlit as st
from convex import ConvexClient
client = ConvexClient("YOUR_CONVEX_PROJECT_URL")
st.title("H.U.G.H. System Dashboard")
configurations = client.query("configurations:getAll") # Assuming a query
getAll in configurations.ts
st.write(configurations)
```

- **Dependencies:** Phase 3, 4 completion.
  - **Contingency:** Start with a very simple command-line interface or basic status page if a full dashboard is too time-consuming initially.
- 

This step-by-step plan provides a roadmap. Be prepared to adapt it based on your experiences, available resources, and the evolving AI landscape. Regular testing and iteration will be key to success.

## 7. Security Considerations

Security is a foundational pillar of the H.U.G.H. system, especially given its objective of family privacy protection. This section outlines key security considerations across all components of the distributed architecture. The primary security advantage is the local-first processing model, but diligent configuration and maintenance are still essential.

### A. Network Security

#### 1. Network Segmentation:

- Utilize VLANs to segment your home network. Create separate VLANs for:
  - H.U.G.H. server infrastructure (Proxmox nodes, management interfaces).
  - Trusted family devices.
  - IoT devices (often less secure, should be isolated).
  - Guest network.
- Implement strict firewall rules between VLANs, allowing only necessary traffic. Your router or a dedicated firewall (e.g., pfSense/OPNsense VM) can manage this.

#### 2. Firewall Configuration:

- Configure the firewall on your main router and on Proxmox hosts (`iptables` or `nftables`, managed via Proxmox VE firewall).
- Default deny: Block all incoming connections from the internet unless explicitly required (e.g., for a securely managed VPN or a highly secured dashboard, though local access is preferred).

- Limit outbound connections from H.U.G.H. components to only necessary endpoints (e.g., Convex, official model repositories).

### **3. Secure Remote Access (If Needed):**

- If remote access to H.U.G.H. management interfaces is required, use a VPN (e.g., WireGuard or OpenVPN) hosted on your network. Avoid exposing Proxmox web UI or SSH directly to the internet.

### **4. Intrusion Detection/Prevention System (IDS/IPS):**

- Consider deploying an IDS/IPS like Suricata on your network perimeter or monitoring traffic to/from H.U.G.H. servers. This is an advanced step but can provide significant threat visibility.

## **B. Proxmox VE Security**

### **1. Host Hardening:**

- Keep Proxmox VE hosts updated with the latest security patches.
- Use strong, unique passwords for the root user and any other administrative accounts.
- Configure SSH securely: disable root login, use key-based authentication, change default port (optional, security by obscurity).
- Enable and configure the Proxmox VE firewall for hosts and VMs.

### **2. VM Security:**

- Harden the OS within each VM (e.g., Ubuntu Server): regular updates, minimal software installation, user account management with least privilege.
- Use Proxmox VE firewall to restrict network access to/from each VM.

### **3. GPU Passthrough Security:**

- While powerful, GPU passthrough can have security implications if a VM is compromised, as it gets direct hardware access. Ensure VMs with passthrough are particularly well-hardened.

### **4. Backup Security:**

- Encrypt backups created by Proxmox Backup Server.
- Store backups in a secure, separate location (ideally following the 3-2-1 backup rule).

## **C. MacBook Air Security**

### **1. Standard macOS Security:**

- Keep macOS updated.
- Use FileVault full-disk encryption.
- Enable the built-in firewall.
- Use strong user passwords and consider a password manager.
- Be cautious about software installations; only install from trusted sources.

### **2. Ollama Security:**

- Ensure Ollama is downloaded from the official source.
- If exposing the Ollama API on the MacBook Air for local network access, ensure it's firewalled appropriately and not accessible from untrusted networks.

## **D. Convex Platform Security**

### **1. Account Security:**

- Use a strong, unique password for your Convex account.
- Enable Two-Factor Authentication (2FA) for your Convex account.

## 2. API Keys and Deploy Keys:

- Securely manage Convex Deploy Keys if used for backend scripts. Store them as environment variables or in a secrets manager, not hardcoded in scripts.
- Grant Deploy Keys the minimum necessary permissions.

## 3. Function Security:

- In your Convex functions (`mutations`, `queries`), validate all input arguments thoroughly to prevent injection attacks or unexpected behavior.
- Implement proper authorization checks within functions if different users/agents have different levels of access to data (e.g., using `ctx.auth` if Convex authentication is used).

## 4. Data Privacy in Convex:

- While Convex provides the infrastructure, you control the data schema. Avoid storing highly sensitive raw PII in Convex if possible. If necessary, consider encryption at the application layer before sending data to Convex, though this adds complexity. For H.U.G.H., Convex is primarily for orchestration metadata, configurations, and summarized/anonymized event data.

## E. AI Model Security

### 1. Model Provenance:

- Download AI models only from official and trusted sources (e.g., Hugging Face official repos, Meta AI, Mistral AI).
- Be cautious with community-finetuned models; vet their source if possible.

### 2. Access Control to Models:

- If AI models are exposed as APIs within your local network (e.g., Ollama API, TGI/vLLM endpoints), ensure these endpoints are not accessible from untrusted networks. Use firewall rules.
- Consider adding an authentication layer if multiple distinct services/users need to access these local model APIs.

### 3. Input Sanitization and Output Filtering:

- Sanitize inputs to AI models to prevent prompt injection attacks that could cause models to behave unexpectedly or reveal sensitive system information.
- Filter outputs from LLMs, especially if displayed to users, to prevent the generation of harmful, biased, or inappropriate content (though this is more about safety than traditional security).

## F. Data Security and Privacy

### 1. Encryption at Rest:

- Use full-disk encryption on Proxmox hosts (e.g., LVM on LUKS) and within VMs where sensitive data is stored.
- MacBook Air should use FileVault.

### 2. Encryption in Transit:

- All internal H.U.G.H. communication between services (e.g., VM to Convex, VM to VM if over network) should use TLS/SSL.
- Local network traffic can be assumed trusted if the network itself is secure, but for sensitive data flows, explicit TLS is better.

### 3. Data Minimization:

- Collect and store only the data essential for H.U.G.H.'s operation.
- Regularly review and purge old logs or data that is no longer needed, according to a defined retention policy.

### 4. Privacy by Design:

- The local-first architecture is the core privacy-preserving feature. Continuously evaluate any new feature or integration for its privacy implications.

## G. Software Supply Chain Security

### 1. Dependency Management:

- Keep all software, libraries, and dependencies (OS, Docker images, Python packages, Node.js packages) updated to patch known vulnerabilities.
- Use tools like `npm audit` (for Node.js/Convex) or `pip-audit` (for Python) to check for vulnerable dependencies.
- Pin Docker image versions to specific digests for reproducibility and to avoid pulling in potentially compromised `latest` tags.

By diligently implementing these security considerations, you can significantly enhance the security posture of the H.U.G.H. system and protect your family's privacy effectively. Security is an ongoing process, not a one-time setup. Regular reviews and updates are essential.

## 8. Family Network Integration

---

Integrating the H.U.G.H. system effectively into your family's existing network is key to its ability to provide comprehensive protection and seamless interaction. This involves considerations for network topology, device interaction, user interfaces, and ensuring minimal disruption to normal family internet usage.

### A. Network Topology and H.U.G.H. Placement

#### 1. Core H.U.G.H. Servers (Proxmox Cluster):

- These servers should be connected to a trusted, wired segment of your network, ideally on a dedicated VLAN as discussed in Security Considerations. This segment should have high-bandwidth access to your primary router/firewall.

#### 2. MacBook Air:

- Can connect via Wi-Fi or wired Ethernet to your trusted family network VLAN. Its connectivity needs to allow access to the Proxmox H.U.G.H. services and the Convex platform over the internet.

#### 3. Router/Firewall Configuration:

- Your main home router or dedicated firewall (e.g., pfSense/OPNsense) is a critical control point.
- **Port Forwarding:** Avoid indiscriminate port forwarding. Only forward ports if absolutely necessary (e.g., for a self-hosted VPN server) and ensure the destination service is highly secure. For H.U.G.H., most interactions will be initiated from within the local network or managed via Convex, minimizing the need for open inbound ports.
- **DNS:** Consider using a DNS server with filtering capabilities (e.g., Pi-hole VM or AdGuard Home VM running on Proxmox) to block malicious domains and trackers for all devices on your network. H.U.G.H. could potentially integrate with or provide data to such a service.

## B. Protecting Family Devices

H.U.G.H. can offer protection in several ways:

#### 1. Network-Level Monitoring (Optional, Advanced):

- If you implement network traffic analysis (e.g., using Suricata on a mirrored port or as a transparent bridge), H.U.G.H. can:
  - Ingest alerts from the IDS/IPS.
  - Use an LLM on Proxmox to analyze alert patterns, identify sophisticated threats, or reduce false positives.
  - Potentially trigger automated responses via Convex (e.g., notifying you, temporarily blocking a device if severe compromise is detected – use with extreme caution).

- **Setup:** This requires your router/switch to support port mirroring or for the IDS/IPS to be inline. This is a complex setup and may impact network performance if not properly resourced.

## 2. DNS-Level Filtering:

- H.U.G.H. can maintain or contribute to blocklists used by your local DNS filtering service (Pi-hole/AdGuard Home). For example, if H.U.G.H. identifies a new phishing domain through its analysis, it could update the DNS filter.

## 3. Endpoint Interaction (Future Potential):

- While not detailed in the initial scope, future H.U.G.H. versions could involve lightweight agents on family PCs or smartphones that communicate with the central H.U.G.H. system for threat intelligence or to offload analysis tasks. This would require careful design for privacy and performance.

# C. User Interaction with H.U.G.H.

## 1. Voice Interaction:

- **Microphone Sources:**
  - Dedicated USB microphones connected to a small device (e.g., Raspberry Pi) that securely streams audio to `VM_Vocal1` on Proxmox.
  - Leverage existing smart speakers if they offer an API or method to capture audio for local processing (this is rare; most process in their cloud).
  - Microphone on the MacBook Air for personal H.U.G.H. interactions.
- **Wake Word:** Implement or use an existing local wake-word detection engine (e.g., Picovoice Porcupine, openWakeWord) on the microphone-equipped device to trigger audio capture only when H.U.G.H. is addressed. This is crucial for privacy.
- **Processing Flow:** Wake word detected -> Audio streamed to `VM_Vocal1` (Whisper) -> Text to Convex -> Convex orchestrates task (e.g., query LLM in `VM_LLM`) -> Response text to `VM_Vocal1` (OuteTTS) -> Spoken audio back to user.

## 2. Dashboard Interface:

- A web-based dashboard (e.g., built with Streamlit, hosted on a VM or via Convex) accessible from any browser on the family network.
- **Features:**
  - System status (Proxmox nodes, AI model health).
  - Recent security events or alerts.
  - Interface to submit text queries to H.U.G.H. LLMs.
  - Configuration settings for H.U.G.H. (e.g., notification preferences).
  - Access control: Consider if different family members should have different views or controls.

## 3. Notifications:

- H.U.G.H., via Convex, can send notifications for important events:
  - Email notifications.
  - Mobile push notifications (e.g., via a service like Pushover or a self-hosted solution like Gotify running in a VM).
  - Visual alerts on the dashboard.

## D. Maintaining Family Network Performance

### 1. Resource Allocation:

- Ensure Proxmox servers have sufficient CPU, RAM, and network bandwidth so that H.U.G.H. operations do not negatively impact general internet usage for the family.
- AI model processing, especially for large LLMs, can be resource-intensive. Schedule bulk processing tasks (if any) for off-peak hours.

### 2. Quality of Service (QoS):

- If your router supports QoS, prioritize traffic for interactive applications (video conferencing, gaming) over less time-sensitive H.U.G.H. background tasks if network contention becomes an issue.

### 3. Local Caching:

- For frequently accessed information or model results, H.U.G.H. services can implement local caching to reduce redundant processing and network traffic.

## E. Privacy within the Family

### 1. User Profiles/Personalization:

- If H.U.G.H. offers personalized features, ensure data for different family members is appropriately segregated or managed with consent. Convex schema can support `userId` fields for this.

### 2. Transparency:

- Be transparent with your family about what H.U.G.H. does, what data it collects (even locally), and how it works. The dashboard can help with this.

### 3. Data Access Controls:

- Consider if certain H.U.G.H. data or controls should be restricted (e.g., detailed security logs might only be accessible to you, Robert).

By thoughtfully integrating H.U.G.H. into your family network, you can create a system that is both powerful in its protective capabilities and unobtrusive in its daily operation, respecting both privacy and user experience.

## References

---

### Vocal AI & Azure:

- \* [OpenAI Next-Generation Audio Models](https://openai.com/index/introducing-our-next-generation-audio-models/) (<https://openai.com/index/introducing-our-next-generation-audio-models/>)
- \* [Azure OpenAI Advanced Audio Models](https://devblogs.microsoft.com/foundry/get-started-azure-openai-advanced-audio-models/) (<https://devblogs.microsoft.com/foundry/get-started-azure-openai-advanced-audio-models/>)
- \* [OpenAI API Voice Model Release](https://community.openai.com/t/gpt-4o-new-voice-model-api-release/793753) (<https://community.openai.com/t/gpt-4o-new-voice-model-api-release/793753>)
- \* [VentureBeat Article on OpenAI Voice Models](https://venturebeat.com/ai/openais-new-voice-ai-models-gpt-4o-transcribe-let-you-add-speech-to-your-existing-text-apps-in-seconds/) (<https://venturebeat.com/ai/openais-new-voice-ai-models-gpt-4o-transcribe-let-you-add-speech-to-your-existing-text-apps-in-seconds/>)
- \* [Microsoft Learn - Azure Realtime Audio Quickstart](https://learn.microsoft.com/en-us/azure/ai-services/openai/realtime-audio-quickstart) (<https://learn.microsoft.com/en-us/azure/ai-services/openai/realtime-audio-quickstart>)
- \* [Azure OpenAI TTS Demo GitHub](https://github.com/Azure-Samples/azure-openai-tts-demo.git) (<https://github.com/Azure-Samples/azure-openai-tts-demo.git>)
- \* [OuteAI/OuteTTS-0.2-500M on HuggingFace](https://huggingface.co/OuteAI/OuteTTS-0.2-500M) (<https://huggingface.co/OuteAI/OuteTTS-0.2-500M>)
- \* [openai/whisper-large-v3 on HuggingFace](https://huggingface.co/openai/whisper-large-v3) (<https://huggingface.co/openai/whisper-large-v3>)
- \* [ai4bharat/IndicParlerTTS-Fast on HuggingFace](https://huggingface.co/ai4bharat/IndicParlerTTS-Fast) (<https://huggingface.co/ai4bharat/IndicParlerTTS-Fast>) (Example for Indic Parler-TTS)

### Mistral Codestral:

- \* [Mistral AI Models Overview](https://docs.mistral.ai/getting-started/models/models_overview/) ([https://docs.mistral.ai/getting-started/models/models\\_overview/](https://docs.mistral.ai/getting-started/models/models_overview/))
- \* [Codestral 25.01 Announcement and Details](https://mistral.ai/news/codestral-2501) (<https://mistral.ai/news/codestral-2501>)
- \* [Microsoft Tech Community Blog on Codestral 25.01](https://techcommunity.microsoft.com/blog/machinelearningblog/introducing-codestral-25-01-mistrals-first-code-model-in-azure-ai-foundry-model-/4365705) (<https://techcommunity.microsoft.com/blog/machinelearningblog/introducing-codestral-25-01-mistrals-first-code-model-in-azure-ai-foundry-model-/4365705>)

- \* [DeepNewz Article on Codestral 25.01](https://deepnewz.com/ai-modeling/mistral-ai-launches-codestral-25-01-on-github-surpassing-deepseek-coder-support-6fb4c6e5) (<https://deepnewz.com/ai-modeling/mistral-ai-launches-codestral-25-01-on-github-surpassing-deepseek-coder-support-6fb4c6e5>)

#### **GPT-4.1 & Benchmarks:**

- \* [Future AGI Blog on GPT-4.1 Benchmarks 2025](https://futureagi.com/blogs/gpt-4-1-benchmarks-2025) (<https://futureagi.com/blogs/gpt-4-1-benchmarks-2025>)
- \* [OpenAI GPT-4.1 Official Release & Documentation](https://openai.com/index/gpt-4-1/) (<https://openai.com/index/gpt-4-1/>)
- \* [Medium Article on GPT-4.1 Models and Benchmarks](https://medium.com/@support_94003/gpt-4-1-models-coding-benchmarks-context-scaling-real-world-applications-a469ae0b7cda) ([https://medium.com/@support\\_94003/gpt-4-1-models-coding-benchmarks-context-scaling-real-world-applications-a469ae0b7cda](https://medium.com/@support_94003/gpt-4-1-models-coding-benchmarks-context-scaling-real-world-applications-a469ae0b7cda))
- \* [LLM Leaderboard 2025](https://llm-stats.com/) (<https://llm-stats.com/>)
- \* [Helicone Blog on GPT-4.1](https://helicone.com/blog/gpt-4-1-release) (<https://helicone.com/blog/gpt-4-1-release>)

#### **Hardware Requirements:**

- \* [GeeksforGeeks: Recommended hardware for running LLMs locally](https://www.geeksforgeeks.org/recommended-hardware-for-running-langs-locally/) (<https://www.geeksforgeeks.org/recommended-hardware-for-running-langs-locally/>)
- \* [ML Journey: Local LLM hardware requirements](https://mljourney.com/local-langs-hardware-requirements-what-you-need-to-run-langs-locally/) (<https://mljourney.com/local-langs-hardware-requirements-what-you-need-to-run-langs-locally/>)
- \* [Hardware Corner: Hardware for Llama 4](https://www.hardware-corner.net/meta-releases-llama-4-what-hardware/) (<https://www.hardware-corner.net/meta-releases-llama-4-what-hardware/>)
- \* [Hugging Face Forums: Hardware for running LLMs](https://discuss.huggingface.co/t/recommended-hardware-for-running-langs-locally/66029) (<https://discuss.huggingface.co/t/recommended-hardware-for-running-langs-locally/66029>)
- \* [TechTarget: Hardware and tools for running LLMs](https://www.techtarget.com/searchEnterpriseAI/tip/How-to-run-LLMs-locally-Hardware-tools-and-best-practices) (<https://www.techtarget.com/searchEnterpriseAI/tip/How-to-run-LLMs-locally-Hardware-tools-and-best-practices>)
- \* [DeepSeek-R1 Hardware Requirements \(Dev.to\)](https://dev.to/askyt/deepseek-r1-671b-complete-hardware-requirements-optimal-deployment-setup-2e48) (<https://dev.to/askyt/deepseek-r1-671b-complete-hardware-requirements-optimal-deployment-setup-2e48>)
- \* [DeepSeek R1 Local Deployment Guide \(DeepWiki\)](https://deepwiki.com/deepseek-ai/DeepSeek-R1/3.2-local-deployment) (<https://deepwiki.com/deepseek-ai/DeepSeek-R1/3.2-local-deployment>)
- \* [Geeky Gadgets Hardware Requirements for DeepSeek R1](https://www.geeky-gadgets.com/hardware-requirements-for-deepseek-r1-ai-models/) (<https://www.geeky-gadgets.com/hardware-requirements-for-deepseek-r1-ai-models/>)
- \* [System Requirements for Running DeepSeek Models \(CodingMall\)](https://codingmall.com/knowledge-base/25-global/240733-what-are-the-system-requirements-for-running-deepseek-models-locally) (<https://codingmall.com/knowledge-base/25-global/240733-what-are-the-system-requirements-for-running-deepseek-models-locally>)

#### **HuggingFace Models & Open Source Alternatives:**

- \* [Hugging Face Blog on Open LLM Leaderboard](https://huggingface.co/blog/open-langs-leaderboard-rlhf) (<https://huggingface.co/blog/open-langs-leaderboard-rlhf>)
- \* [Top Open Source Models on HuggingFace 2025 \(Analytics Vidhya\)](https://www.analyticsvidhya.com/blog/2024/12/top-open-source-models-on-hugging-face/) (<https://www.analyticsvidhya.com/blog/2024/12/top-open-source-models-on-hugging-face/>)
- \* [Reddit discussion on GPT-4 rankings](https://www.reddit.com/r/LocalLLaMA/comments/14i43j7/how_does_gpt4_rank_on_the_huggingface_leaderboard/) ([https://www.reddit.com/r/LocalLLaMA/comments/14i43j7/how\\_does\\_gpt4\\_rank\\_on\\_the\\_huggingface\\_leaderboard/](https://www.reddit.com/r/LocalLLaMA/comments/14i43j7/how_does_gpt4_rank_on_the_huggingface_leaderboard/))
- \* [GPT-4o System Card \(HuggingFace Papers\)](https://huggingface.co/papers/2410.21276) (<https://huggingface.co/papers/2410.21276>)
- \* [HuggingFace GPT-4 Model Page](https://huggingface.co/GPT-4) (<https://huggingface.co/GPT-4>)
- \* [ChatQA Model \(Arxiv\)](https://arxiv.org/abs/2401.10225) (<https://arxiv.org/abs/2401.10225>)
- \* [Analytics Vidhya: 10 GPT-4 Open-Source Alternatives in 2025](https://www.analyticsvidhya.com/blog/2024/04/gpt-open-source-alternatives/) (<https://www.analyticsvidhya.com/blog/2024/04/gpt-open-source-alternatives/>)
- \* [DataCamp: 12 GPT-4 Open-Source Alternatives](https://www.datacamp.com/blog/12-gpt4-open-source-alternatives) (<https://www.datacamp.com/blog/12-gpt4-open-source-alternatives>)
- \* [Open Source Collection: 5 Open Source Alternatives To GPT-4](https://opensourcecollection.com/blog/5-open-source-alternatives-to-gpt-4-language-model) (<https://opensourcecollection.com/blog/5-open-source-alternatives-to-gpt-4-language-model>)
- \* [LibHunt: Top 23 GPT-4 Open-Source Projects](https://www.libhunt.com/topic/gpt-4) (<https://www.libhunt.com/topic/gpt-4>)
- \* [Qwen/Qwen2.5-0.5B on Hugging Face](https://huggingface.co/Qwen/Qwen2.5-0.5B) (<https://huggingface.co/Qwen/Qwen2.5-0.5B>)
- \* [QwenLM/Qwen2.5-Omni GitHub](https://github.com/QwenLM/Qwen2.5-Omni) (<https://github.com/QwenLM/Qwen2.5-Omni>)
- \* [Qwen2.5 Blog](https://qwenlm.github.io/blog/qwen2.5/) (<https://qwenlm.github.io/blog/qwen2.5/>)
- \* [meta-llama/Llama-3.1-8B-Instruct on HuggingFace](https://huggingface.co/meta-llama/Llama-3.1-8B-Instruct) (<https://huggingface.co/meta-llama/Llama-3.1-8B-Instruct>)
- \* [meta-llama/Llama-3.2-11B-Vision-Instruct on HuggingFace](https://huggingface.co/meta-llama/Llama-3.2-11B-Vision-Instruct) (<https://huggingface.co/meta-llama/Llama-3.2-11B-Vision-Instruct>)

- \* [facebookresearch/OPT on GitHub](https://github.com/facebookresearch/OPT) (<https://github.com/facebookresearch/OPT>)
- \* [lmsys/vicuna-13b-v1.3 on HuggingFace](https://huggingface.co/lmsys/vicuna-13b-v1.3) (<https://huggingface.co/lmsys/vicuna-13b-v1.3>) (Example for Vicuna)
- \* [mosaicml/mpt on GitHub](https://github.com/mosaicml/mpt) (<https://github.com/mosaicml/mpt>)
- \* [mosaicml/mpt-30b-chat on HuggingFace](https://huggingface.co/mosaicml/mpt-30b-chat) (<https://huggingface.co/mosaicml/mpt-30b-chat>)

#### **Llama 3.2 & Ollama:**

- \* [Llama 3.2 Download Page](https://llamaimodel.com/download-3-2/) (<https://llamaimodel.com/download-3-2/>)
- \* [Hugging Face Models](https://huggingface.co/models) (<https://huggingface.co/models>)
- \* [Medium article on running Llama 3.2 locally](https://medium.com/@aleksej.gudkov/how-to-run-llama-3-2-locally-a-complete-guide-36d4a8c7bf94) (<https://medium.com/@aleksej.gudkov/how-to-run-llama-3-2-locally-a-complete-guide-36d4a8c7bf94>)
- \* [KDnuggets guide on local Llama 3.2 deployment](https://www.kdnuggets.com/2024/10/using-llama-3-2-locally.html) (<https://www.kdnuggets.com/2024/10/using-llama-3-2-locally.html>)
- \* [Ollama-models GitHub Repository](https://github.com/Pyenb/Ollama-models) (<https://github.com/Pyenb/Ollama-models>)
- \* [Ollama Official Website](https://ollama.com/) (<https://ollama.com/>)
- \* [Ollama.org \(alternative site/community\)](https://ollama.org/) ([http://ollama.org/](https://ollama.org/))
- \* [DevOpsRoles article on Ollama](https://www.devopsroles.com/running-free-ai-models-llama-deepseek-ollama/) (<https://www.devopsroles.com/running-free-ai-models-llama-deepseek-ollama/>)
- \* [Unite.AI - 7 Best LLM Tools to Run Models Locally](https://www.unite.ai/best-lm-tools-to-run-models-locally/) (<https://www.unite.ai/best-lm-tools-to-run-models-locally/>)
- \* [ggerganov/llama.cpp on GitHub](https://github.com/ggerganov/llama.cpp) (<https://github.com/ggerganov/llama.cpp>)
- \* [vllm-project/vllm on GitHub](https://github.com/vllm-project/vllm) (<https://github.com/vllm-project/vllm>)

#### **Convex Integration:**

- \* Convex Developer Hub (General reference, specific URL not provided in source, assume standard Convex documentation)
- \* Medium articles on Convex MCP servers (General reference)

#### **Proxmox Distributed AI & GPU Passthrough:**

- \* [LinuxConfig Proxmox GPU Passthrough](https://linuxconfig.org/proxmox-gpu-passthrough) (<https://linuxconfig.org/proxmox-gpu-passthrough-setup-guide>) (Example general guide)
- \* Proxmox forums (General reference)
- \* VirtualizationHowTo (General reference)

#### **MacBook Air M-Series LLM Performance:**

- \* HuggingFace blog (General reference for Apple Silicon optimizations)
- \* Medium benchmarking articles (General reference)

#### **Local-to-Cloud API Privacy & Security:**

- \* [TechTarget API Security Best Practices](https://www.techtarget.com/searchsecurity/tip/API-security-best-practices-A-CISO-guide) (<https://www.techtarget.com/searchsecurity/tip/API-security-best-practices-A-CISO-guide>) (Example general guide)
- \* StackHawk Blog (General reference)
- \* Google Cloud Blog (General reference for cloud security principles)
- \* Cloud Security Alliance (General reference)

#### **Home-lab Failover Strategy (Proxmox, Mac Nodes):**

- \* Proxmox forums (General reference for HA)
- \* [Hansen IT Solutions Proxmox Two Node Cluster](https://hansensolutions.ca/proxmox-ve-home-lab-setting-up-a-two-node-cluster-with-a-qdevice/) (<https://hansensolutions.ca/proxmox-ve-home-lab-setting-up-a-two-node-cluster-with-a-qdevice/>) (Example specific guide)
- \* Medium articles on home lab clustering (General reference)

## **Appendix: Quick Reference**

---

This appendix provides a quick reference for key commands and configurations relevant to the H.U.G.H. system deployment.

## A. Proxmox VE Commands (Host Shell)

- Check Cluster Status:

```
bash
 pvecm status
```

- List Nodes in Cluster:

```
bash
 pvecm nodes
```

- List VMs and Containers:

```
bash
 qm list
 pct list
```

- Start/Stop/Shutdown VM:

```
bash
 qm start <VMID>
 qm stop <VMID>
 qm shutdown <VMID>
```

- Enter VM Console:

```
bash
 qm terminal <VMID>
```

- Update GRUB (after modifying `/etc/default/grub` for IOMMU):

```
bash
 update-grub
```

- Update initramfs (after modifying `/etc/modules`):

```
bash
 update-initramfs -u -k all
```

- Check IOMMU Groups (for GPU Passthrough):

```
bash
 find /sys/kernel/iommu_groups/ -type l
```

- Proxmox VE Firewall Commands:

```
bash
 pvefw-update # Apply firewall changes
 pvefw status # Check firewall status
```

## B. Ollama Commands (Linux VM or macOS)

- Install Ollama (Linux):

```
bash
 curl -fsSL https://ollama.com/install.sh | sh
```

- Pull a Model:

```
bash
 ollama pull llama3.1:8b-instruct
 ollama pull codestral # If/when available on Ollama
 ollama pull qwen2.5:7b
```

- Run a Model (Interactive):

```
bash
 ollama run llama3.1:8b-instruct
```

- List Local Models:

```
bash
 ollama list
```

- Remove a Local Model:

```
bash
 ollama rm llama3.1:8b-instruct
```

- Serve Ollama API (usually runs by default after install):

The API typically listens on `http://localhost:11434`. To make it accessible from other machines on your local network, you may need to configure Ollama to listen on `0.0.0.0` by setting the `OLLAMA_HOST` environment variable before starting the Ollama service.

```
bash
Example systemd override for Ollama service (Linux)
sudo systemctl edit ollama.service
Add:
[Service]
Environment="OLLAMA_HOST=0.0.0.0"
Then:
sudo systemctl daemon-reload
sudo systemctl restart ollama
```

## C. Docker Commands (Within VMs)

- Build a Docker Image:

```
bash
 docker build -t my-ai-app .
```

- Run a Docker Container (Example for GPU access):

```
bash
 docker run --gpus all -d -p 8000:8000 my-ai-app
```

- List Running Containers:

```
bash
 docker ps
```

- List All Containers (including stopped):

```
bash
 docker ps -a
```

- Stop a Container:

```
bash
 docker stop <container_id_or_name>
```

- View Container Logs:

```
bash
 docker logs <container_id_or_name>
```

- Remove a Container:

```
bash
 docker rm <container_id_or_name>
```

- Remove an Image:

```
bash
 docker rmi <image_id_or_name>
```

## D. Convex CLI Commands (Development Machine)

- Login to Convex:

```
bash
 npx convex login
```

- **Initialize/Start Dev Server:**

```
bash
 npx convex dev
```

- **Deploy to Production:**

```
bash
 npx convex deploy
```

- **Run a Function (Query or Mutation) from CLI:**

```
bash
 npx convex run functions:myFunction '{"arg1": "value1"}'
```

- **View Dashboard (Schema, Functions, Data):**

Usually available at `http://localhost:3000` when `npx convex dev` is running, or via the Convex project dashboard online.

## E. Sample Configurations

### 1. Proxmox VM Configuration Snippet (for GPU Passthrough - `/etc/pve/qemu-server/<VMID>.conf`)

```
... other VM settings ...
For NVIDIA GPU:
hostpci0: XX:XX.X,pcie=1,x-vga=1 # Primary GPU with VGA output, XX:XX.X is from lspci
For secondary GPU or if issues with x-vga:
hostpci0: XX:XX.X,pcie=1
If GPU has an audio function, pass it too:
hostpci1: YY:YY.Y,pcie=1

BIOS setting for UEFI
bios: ovmf
Machine type for better PCIe support
machine: q35
CPU settings
cpu: host
... other settings like memory, cores, network ...
args: -cpu 'host,+kvm_pv_unhalt,+kvm_pv_eoi,hv_vendor_id=NV43FIX,kvm=off'
Example for NVIDIA, may vary
```

**Note:** Specific `args` and passthrough methods can vary based on GPU model and Proxmox version. Always consult latest Proxmox documentation and community guides.

## 2. Dockerfile Example (for a Python/Transformers AI service)

```
FROM python:3.10-slim

WORKDIR /app

RUN apt-get update && apt-get install -y --no-install-recommends \
 git \
 && rm -rf /var/lib/apt/lists/*

COPY requirements.txt requirements.txt
RUN pip install --no-cache-dir -r requirements.txt

For NVIDIA GPU support in container (ensure base image or host has drivers)
This example assumes NVIDIA container toolkit is used on host
Or, install CUDA toolkit inside if needed, but larger image

COPY . .

Command to run your application
CMD ["python", "app.py"]
```

requirements.txt example:

```
transformers
torch
torchaudio
accelerate
sentencepiece
Add other dependencies like Flask/FastAPI if exposing an API
```

## 3. Convex Schema Snippet ( convex/schema.ts )

(Refer to Section 3 for a more complete example)

```
import { defineSchema, defineTable } from "convex/server";
import { v } from "convex/values";

export default defineSchema({
 configurations: defineTable({
 settingName: v.string(),
 settingValue: v.any(),
 }).index("by_settingName", ["settingName"]),
 // ... more tables
});
```

#### 4. Ollama API Call Example (Python)

```
import requests
import json

ollama_api_url = "http://localhost:11434/api/generate" # Or IP of Ollama server

def query_ollama(prompt, model="llama3.1:8b-instruct"):
 payload = {
 "model": model,
 "prompt": prompt,
 "stream": False # Set to True for streaming response
 }
 try:
 response = requests.post(ollama_api_url, json=payload)
 response.raise_for_status() # Raise an exception for HTTP errors
 return response.json().get("response", "")
 except requests.exceptions.RequestException as e:
 print(f"Error querying Ollama: {e}")
 return None

Example usage:
response_text = query_ollama("Why is the sky blue?")
if response_text:
print(response_text)
```

This appendix should serve as a helpful quick start for common operations. Always refer to the official documentation for each tool for comprehensive details.