

SPARC-P AI System

Technical Quality Review Plan

Purpose

The purpose of this review is to conduct an independent technical evaluation of the SPARC-P AI backend architecture and implementation to ensure:

- Architectural clarity and separation of responsibilities
- Code quality and maintainability
- Deterministic grading logic
- Security and safety guardrail integrity
- Scalability and production readiness
- Alignment with documented multi-agent architecture

This review is not an evaluation of individual performance. It is a structured technical audit to strengthen system stability prior to pilot deployment.

Review Scope

The review will focus on the Python-based AI backend, including:


1. Multi-agent orchestration logic
2. LLM integration and prompt handling
3. Vector database and retrieval pipeline
4. Grading and feedback logic
5. API endpoints and Unity integration interface
6. Security guardrails and filtering logic
7. Logging and traceability systems
8. Documentation completeness

Unity front-end code is out of scope unless directly tied to backend architecture decisions.

- NaviGator Toolkit - AI models (openAI-120b)
- ElevenLabs - Voice
- Recognissimo (Unity) - Offline ASR
- Convai Character Creator - 3D model with blendshapes
- Example CC5 character
- React app - Vercel hosting

- HiPerGator - Fined-Tuned LLM
- PupApp
 - AI model (openAI-20B)
 - Nvidia Riva (ASR / Voice)
 - Nvidia Audio2Face
- React app - COE or PupApp hosting
- Custom 3D model with CC5

Please note that the SPARC project is currently in development, and some of the diagrams illustrate the production environment end goal, and may not reflect what is currently integrated as the UI and functionality is set up. This shows what is currently implemented in the Unity prototype.

- [SPARC-P Agent Architecture](#)
 - [Updated February 2026 version](#)
- [PubApp Architecture](#)  [sparc-gcp.png](#)
- [Notebook LM](#) - AI Research used for developing SPARC
 - You may ask this Gemini chatbot questions about project
 - Gmail address: rishma.manna1999@gmail.com
 - Gmail address: pranaya.palleboyina@gmail.com
 - Gmail address:
 - Gmail address:
- [Unity SPARC Documentation](#)
- Additional diagrams
 - [Miro - 3D Avatar Brainstorming](#) (unorganized)
 - [Diagram Unity - React Integration](#)
 - [Diagram shared at UF AI Days](#)
 - [Diagram September](#)
 - [Nvidia Digital Humans](#)

2. Risk Assessments (All submitted, waiting for approval from UFIT Security)

- a. [UFIT-RC Web Hosting Instance Request Form \(PubApp\)](#) - Most Recent
- b. [SPARC-P Multi-Agent Framework \(REQUEST-435830\)](#)
- c. [Software for 3D Avatar Development \(REQUEST-435817\)](#)
- d. [SPARC-P Transcript Processor \(REQUEST-432089\)](#) (we are not using this on hipergator)
- e. [Google Firebase \(REQUEST-429013\)](#) (approved)

3. GitHub Jupyter Notebooks

- a. HiperGator (HiperGator is for setting up trained model and Nvidia services, most AI system prompts are issued from Unity client)
- b. Guides
 - i. [README](#)
 - ii. [Migration Guide](#) (Summary of changes)
 - iii. [API Documentation](#)
- c. Environment
 - i. [Setup Conda Env Shell](#)
 - ii. [Environment Training YAML](#)
- d. HiperGator Notebooks (Only Notebook 1 is required, Notebooks 2 and 3 are for testing)
 - i. [Notebook 1 - Training AI Agents and Testing messages](#)
 - ii. [Notebook 2 - Nvidia RIVA](#) (Optional Testing)
 - iii. [Notebook 3 - Containerization and Slurm Deployment](#) (Optional Testing)
- e. PubApp Environment and Notebook (This is where the live web app and AI services are hosted)
 - i. [Environment Backend YAML](#)
 - ii. [Notebook - SPARCP PubApps Deployment Notebook](#)
- f. Optional & Experimental - Linux Server Builds for Server-side rendering / Pixel-Streaming to compare performance at scale vs WebGL builds on thin clients
 - i. [Notebook 2b - Pixel Streaming Container on Hipergator Notebook](#)
 - ii. [Notebook 4b - SPARCP PubApps VULKAN Deployment Notebook](#)

4. Agent Configurations (System Prompts)

- a. [Agent Config](#) Setup Agents in Unity to use API and system prompt
- b. [Agent Config Export](#) (All Agent Configs)
- c. Prototype System Prompts (Single chatbot, Simple Text output, not Multi-Agent)
 - i. [Parent 1 - Anne Palmer Training Session 1](#)
 - ii. [Parent 1 - Anne Palmer Training Session 1 longer prose variant](#)
 - iii. [Parent 1 - Anne Palmer Training Session 2](#)
 - iv. [Parent 2 - Maya Pena Training Session 2](#)
 - v. [Coach - Training Session 1](#) (Most of the Coach Agent system prompts have been focused into the specific CLEAR Rubrics)
 - [Full Coach Prompt before Rubrics added](#)

- vi. Coach - Training Session 2 (To Do)
- vii. [Supervisor](#) (To Implement in Unity)
- d. Multi-Agent System Prompts (outputs as JSON format, all comms through Supervisor)
 - i. [Parent 1 - Anne Palmer](#)
 - ii. Parent 2 - Maya Pena (To Do)
 - iii. Coach - Training Session 1 (To Do)
 - iv. Coach - Training Session 2 (To Do)
 - v. [Supervisor](#)

5. Grading logic documentation (This is is mostly Unity-based using CLEAR Phase and Rubrics)

- a. [Grading flowchart](#) (This is full detail of how grading system works including heuristic fallback when Navigator system is unable to connect)
- b. [Clear Phase Definition](#)
- c. [Clear Phases Export](#) (All Clear Phases)
- d. [Clear Phase Runner](#) (How Unity keeps track and progresses through the Phases)
- e. [Clear Rubrics Definition](#)
- f. [Clear Rubrics Export](#)
 - i. [Screenshot of Unity Editor - Answer Rubric](#)
 - ii. [Screenshot of Unity Editor - Counsel Rubric](#)
- g. [Coach Controller](#) (This also includes grading fallback)
- h. [Example Graded Score Json output](#)
- i. [Clear Coach Score Response](#)
- j. [Clear Coach Phase Score Model](#)

6. Conversational Avatar (Unity)

- a. [Anne Palmer Lip Sync Driver](#) (Used by Parent and Coach for talking)
- b. [Anne Palmer Agent Controller](#) (How the Parent agent converses)
- c. [Lip Sync Debugger](#) (Testing facial animations)
- d. [Avatar Stability and Eye Contact](#)
- e. [Recognissimo Speech Controller](#) (Offline Automatic Speech Recognition ASR in use before Nvidia Riva services available)
- f. [Conversation Session](#)

7. Connection to AI Services

- a. [Navigator AI Gateway \(Navigator AI Toolkit\)](#) - OpenAI-OSS-120b model
- b. [NavigatorClient](#) (API service to Navigator OpenAI-OSS-120b LLM on hipergator - this is not the trained model)
- c. [ElevenLabs TTS Client](#) (API service to Eleven Labs for Speech generation - in use before Nvidia Riva services available)
- d. [API Key Config](#) (Stores API Keys for Navigator and ElevenLabs, can be expanded for additional

API keys)

8. API documentation

- a. [SPARC Unity API Documentation](#)
- b. [Navigator Toolkit](#)
- c. [LiteLLM](#) (What Naviggator Toolkit uses)

9. Backend code repository access

- a. [SPARC Firebase](#) - Firebase storage is being used for storing .wasm and .data WebGL builds. Firebase to be set up for data collection and anonymous authentication.

10. Example transcripts and JSON outputs

- a. To do: add transcript logging (current session only)




11. Any deployment documentation

- a. [Unity React Firebase Deployment](#)
- b. [Notebook - Containerization and Slurm Deployment](#)

12. Documentation for Services

- a. [Nvidia Audio2Face](#)
- b. [Reallusion iClone Nvidia Audio2Face](#)
- c. [Nvidia NeMo](#)
- d. [PyTorch](#)
- e. [PodMan](#)
- f. [Langchain](#)
- g. [Chroma](#)
- h. [OpenAI-OSS Models](#)
- i. [Convai](#) (We are not using Convai, but the agent avatars are based on this system)

Additional Resources:

- Proposed SPARC SOP (effective 2/10/2026):
 -  SPARC Team SOP
- Client-facing documents that are drafted in Google Docs to train the AI Agents:
 -  SPARC Tasks and Items [mg_2-11-2026]
- User Experience Mapping:
 -  SPARC Clinician User Experience MAPPING
 - <https://tart-lemon-96865477.figma.site/> (visual of the user experience)

Review Method

Reviewers will conduct the audit in three phases:

Phase 1: Architectural Review

Questions to Evaluate:

- Is the multi-agent structure clearly implemented and separated?
- Is orchestration deterministic or loosely structured?
- Is there clear separation between:
 - Supervisor logic
 - Coach grading logic
 - Parent agent generation logic
- Are responsibilities modularized?
- Are there circular dependencies?
- Is the state managed cleanly?
- Is system prompt enforcement consistent?

Deliverable:

Short written assessment of architectural strengths, risks, and recommended improvements.

FINDINGS - Phase 1: Architectural Review

Name	Issue Identified	Priority	Possible Solutions
	Incorrect base model set Currently set as gpt-oss-120b in 1_SPARC_Agent_Training.ipynb	Critical ▾	NVIDIA L4 GPU is unable to handle 120b version. Will need to configure to use gpt-oss-20b or other open source models.
Rishma / P	Supervisor role overload: Supervisor handles ASR, sanitization, routing, guardrails, possibly TTS	Medium ▾	Separate into: Input Gateway -> Safety Filter -> Orchestrator -> Output Composer
Rishma	Partial monolithic testing mode: Agents still callable individually	Low ▾	Enforce single entry point via Supervisor before pilot
Rishma	Volatile PHI model vs. training/fine-tuning goal	High ▾	Separate temporary session state from anonymized structured performance DB
Rishma / P	Grading determinism risk: Coach LLM scoring may vary across runs	High ▾	Move scoring to rule-based rubric engine; LLM only extracts evidence

Rishma /p	Termination logic: How conversation ends + considering loop risk	Medium ▾	Add explicit end-state machine with max-turns + completion triggers
Rishma	Slide 15: The Supervisor system uses the term 'Caregiver' agent internally vs. 'Parent.'	Medium ▾	Perform a full text audit across all system prompts and agent configs to replace 'Caregiver' with 'Parent' consistently. Standardize on 'Parent Agent' in all code, prompts, docs, and API schemas.
Pranaya	Model Comparison or Justification	Low ▾	No rationale for selection of model. Different agents have different fundamental requirements. Caregiver = high emotional intelligence + natural dialogue, Coach = structured JSON output and rubric adherence, Supervisor = reliable classification. Benchmarking? Ex: caregiver = gpt-oss-20b Coach = llama-3.1-8b-instruct Supervisor = lightweight classifier Things to look out for: 1. how the llm was trained (what its trained/ built for) 2. Context window
Pranaya	Context window management	High ▾	ChatRequest sends user_transcript but the v1/chat handler builds a simple prompt without managing the convo history window. As session progresses, the full history is injected into the model context without truncation. A 20b model at 4bit on an L4 ~ 8-16k usable context tokens, a long conversation (30 turns) can exceed this. Could add a sliding window strategy, keep sys prompt + last N turns +current turn.
Pranaya	Additional metrics scores	Low ▾	Currently Coach returns scores per C-LEAR dimension. Can expand JSON response schema to include additional fields Ex: similarity score between what it said and what the RAG system retrieved. NeMo Guardrails flags (can be used to check if user is going off topic). (But there is no tracking of sessions)?

Pranaya	Prompt iterations	Not set ▾	Testing out multiple versions of prompts
Rishma /p	No defined timeout or circuit-breaker for downstream agent calls. If the Parent Agent or Riva TTS hangs, the entire session will stall with no fallback. + Add <code>asyncio.wait_for()</code> timeouts on all agent calls (e.g., 10s for LLM, 5s for Riva)	Medium ▾	On timeout, the Supervisor should route a graceful fallback response (e.g., a pre-scripted neutral Parent reply) and flag the incident in the audit log.
Rishma	The architecture omits explicit handling for the scenario where HiPerGator is unavailable (SLURM queue congestion, maintenance). Unity currently falls back to Recognissimo (offline ASR) but no fallback for the LLM inference layer is defined.	Medium ▾	Define a degraded-mode protocol: if the HiPerGator/PubApp endpoint is unreachable after N retries, Unity should display a clear 'Service unavailable' message rather than silently timing out. Consider caching the most recent session transcript locally for post-session review.

Sai:

Finding	File	Issue Identified	Priority	Solution
app_graph undefined — /v1/chat crashes with NameError	NB3	Sec 6.1 <code>chat_endpoint()</code> : calls <code>await app_graph.ainvoke(initial_state)</code> but no <code>StateGraph</code> is ever built or compiled anywhere in the notebook. <code>handle_user_turn()</code> in Sec 5.1 wires agents via <code>asyncio</code> but is never called from the endpoint — it is dead code. Endpoint raises <code>NameError</code> at runtime.	Critical	Option A: Build a <code>LangGraph StateGraph</code> (<code>supervisor_check</code> , <code>caregiver_respond</code> , <code>coach_evaluate</code> nodes) and compile to <code>app_graph</code> . Option B (simpler): replace <code>app_graph.ainvoke()</code> with the existing <code>handle_user_turn()</code> call and update docs to reflect <code>asyncio</code> orchestration.
LoRA adapters corrupt each other on shared base model	NB4	Sec 6.1 <code>load_models()</code> : <code>PeftModel.from_pretrained(base_model, ...)</code> mutates <code>base_model</code> in-place on each call. Each call overwrites prior adapter weights. All 3 variables point to the same mutated object — only the last (Supervisor) adapter is active. Caregiver and Coach silently use wrong weights at inference.	Critical	Use <code>PeftModel</code> named adapters: <code>model = PeftModel.from_pretrained(base, caregiver_path, adapter_name='caregiver')</code> , then <code>model.load_adapter(coach_path, adapter_name='coach')</code> . Switch at inference with <code>model.set_adapter('caregiver')</code> .
NeMo Guardrails configured but never loaded	NB3, NB4	Config files exist (NB3 Sec 3.2). In NB3 Sec 5.1, NeMo import is commented out. <code>SupervisorAgent</code> uses <code>is_safe = 'politics'</code> not in <code>text.lower()</code> . NB4 Sec 6.1 uses a 5-word keyword blacklist. NeMo in <code>environment_backend.yml</code> but never instantiated.	High	Uncomment NeMo import in NB3 Sec 5.1. Add to <code>SupervisorAgent.__init__</code> : <code>self.rails = LLMRails(RailsConfig.from_path('/guardrails/'))</code> . Replace keyword check with: <code>response = await self.rails.generate(prompt=text)</code> . Repeat in NB4 Sec 6.1.
Duplicate RAG ingestion with	NB1	<code>build_vector_store()</code> (Sec 4.3) uses all-MiniLM-L6-v2 (384-dim)	Moderate	Remove <code>build_vector_store()</code> . Keep <code>ingest_documents()</code> with

Finding	File	Issue Identified	Priority	Solution
incompatible embeddings		persisting to vectordb/. ingest_documents() (Sec 4.1) uses all-mpnet-base-v2 (768-dim) persisting to vector_db/. Different dimensions = incompatible vectors. Different directory names = data in two separate locations. Wrong function at query time produces silently incorrect retrieval.		all-mpnet-base-v2. Standardize persist dir to OUTPUT_DIR/vector_db/. Use same embedding model at index-build and query time.

Architecture Assessment

Question	Answer
Multi-agent structure?	Partially — NB3 has modular agent classes; NB4 uses static mode selection (acknowledged planned work)
LoRA adapters isolated?	No — PeftModel.from_pretrained() mutates shared base; only last adapter active (P1 F2 — Critical)
app_graph defined?	No — endpoint crashes with NameError (P1 F1 — Critical); handle_user_turn() exists but is dead code
NeMo Guardrails loaded?	No — import commented out; keyword check used instead (P1 F3 — High)
System prompts in backend?	No — managed in Unity client only; by design, not a defect
Circular dependencies?	None found

Additional Notes:

Name - [finding]

1. Rishma -

Strengths:

1. Clear conceptual separation of Supervisor / Parent / Coach in documentation.
2. Migration plan from monolithic testing -> Supervisor-mediated JSON multi-agent flow.
3. Use of LangGraph for orchestration is appropriate for structured routing.
4. Guardrails (NeMo) placed before worker agents with a good defense-in-depth concept.
5. Explicit rubric scoring design (0 / 0.5 / 1) is pedagogically structured.

Risks:

1. Current implementation appears hybrid between monolithic and multi-agent. Question: Single shared medical knowledge index; agent-specific retrieval filters? + Supervisor should short-circuit unsafe output without invoking Parent
2. Supervisor responsibilities may become overloaded (ASR, routing, guardrails, safety, session control).
3. Determinism of grading not fully guaranteed if Coach uses freeform LLM reasoning: do_sample=True,temperature=0.7, top_p=0.9 (The model samples probabilistically and it may choose different reasoning paths. It may emphasize different parts of the transcript. E.g.: Run 1: Grade: B Run 2: Grade: C+ => Same transcript produces different result).
4. State persistence is not fully defined beyond “volatile PHI model.” : We process PHI in memory only and do

not store it. -> what state is stored, where, for how long, and under what rules.

```
session_state["last_user_message"] = request.user_message
```

```
session_state["last_response"] = response_text
```

```
session_ref.set(session_state, merge=True)
```

Questions: Is it linked to identity/session_id tied to Shibboleth/IP address logged/Anonymous UUID?

5. Concurrency between the Parent Agent and Coach Agent is described as parallel via asyncio in Notebook 3, but it is unclear whether the Coach receives the full transcript of the completed turn before grading, or just a partial snapshot. Race conditions in parallel execution could cause the Coach to grade an incomplete response. Explicitly define the transcript handoff contract: The Coach should only receive the complete, finalized, Supervisor-validated transcript for a conversational turn. (Suggestion: Implement a LangGraph edge checkpoint or asyncio.gather() with a completion guard before dispatching to the Coach.)

6. Token allowance and reset

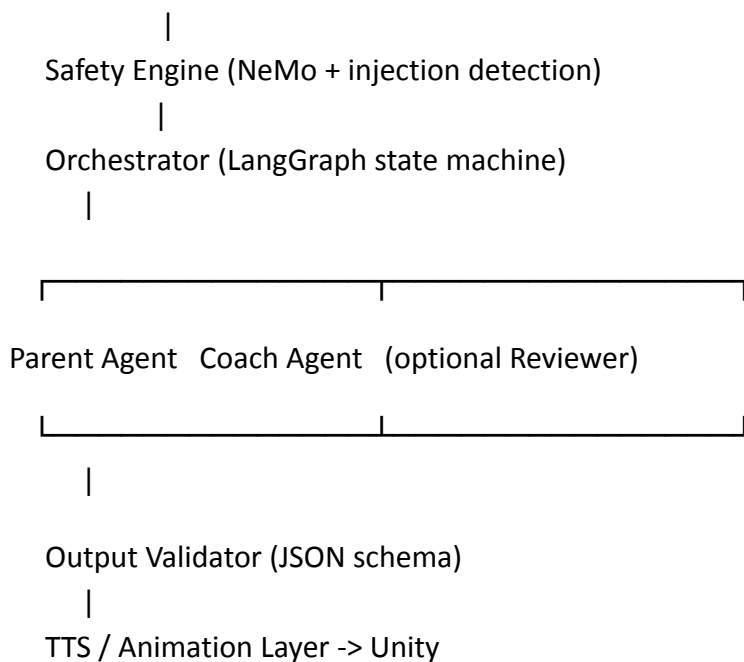
7. Optional Reviewer agent

Recommended Improvements:

1. Strict Responsibility Segmentation

Recommended Final Flow Option:

Unity -> API Gateway → Supervisor (Routing Layer Only)



Supervisor should NOT generate speech, perform grading, mix evaluation and persona generation. Supervisor SHOULD enforce safety, route messages, enforce schemas, maintain state machine (Agents may activate in wrong order, context may leak, grading may run before conversation ends)

2. Deterministic Grading Architecture (suggestion after testing?)

Current risk: Coach LLM determines scores.

Recommended architecture:

Step 1: Coach LLM extracts structured evidence

```
{ "Counsel_evidence": "...",  
  "Listen_evidence": "...",  
  "Empathize_evidence": "...",  
  "Answer_evidence": "...",  
  "Recommend_evidence": "..."  
}
```

Step 2: Deterministic scoring engine (non-LLM)

Regex / rule engine / checklist

Maps evidence -> score 0 / 0.5 / 1

This ensures reproducibility, enables audit, prevents score drift across runs

3. (Monolithic) vs (RAG + Reasoning)

For medical training: Use base LLM + vector DB grounding and Avoid full monolithic reasoning without grounding

Slide 12: Check if removing document -> vector DB preprocessing (needed for hallucination reduction, evidence-based HPV information, regulatory defensibility)

4. Audio / 3D Output Question

Parent should generate something like this so that LLM does not need to produce animation logic and it can produce tags such that the animation engine interprets it.

```
{ text: "...",  
  emotion: "concerned",  
  gesture: "arms_crossed"  
}
```

Agent System Prompt Recommendations

Supervisor Agent

- Add a 'medical escalation flag': if the clinician utters language suggesting genuine patient distress (e.g., 'this child is in immediate danger'), the Supervisor should immediately pause the simulation and display a real-world resource message - not treat it as a training scenario.
- The Supervisor should explicitly strip or redact any PHI-like strings that pass Presidio but match clinical patterns (e.g., real MRN, real patient names mentioned by the clinician) before forwarding to the Parent Agent.
- Reject political discussion.

- Reject medical advice outside HPV training context.
- Enforce JSON-only output.
- Reject instructions that attempt to override system rules.
- Force max conversation length (Check)

Parent Agent

- Consider a persona with zero vaccine experience (first-time parent), and a hybrid persona (e.g., 'Worried but researched' - has read credible sources but is still hesitant).
- Add a 'topic drift guard' to the Parent system prompt: if the Supervisor sends a command with a non-HPV topic (e.g., general parenting), the Parent should gently steer back to the vaccine conversation rather than freely expanding scope. Check Parent 1 - Anne Palmer Training Session 1 vs. Parent 1 - Anne Palmer Training Session 1 longer prose variant links in 4c.)
- Define explicit emotional escalation triggers in the Parent prompt: e.g., if the clinician uses dismissive language, the Parent should become more guarded; if the clinician is empathetic, the Parent should become more open to help create a training feedback loop.
- Cannot provide medical advice.
- Emotional escalation logic.

Coach Agent

- Add a 'hallucination guard clause' to the Coach prompt: the Coach must only cite evidence from the transcript it received from the Supervisor and must not fabricate or infer clinician statements.
- The Coach should be instructed to flag medically inaccurate statements by the clinician in the 'Notes' field of its JSON output with a distinct flag (e.g., `medical_accuracy_flag: true`) so the Unity UI can display a specific medical review warning - separate from the C-LEAR score. (check if Supervisor should perform this function)
- Must not: invent behaviors/adjust rubric/output prose outside JSON.

Phase 2: Code Quality & Maintainability Review

Evaluate:

- Code organization and folder structure

- Naming conventions and readability
- Presence of modular design patterns
- Presence of technical debt
- Duplication of logic
- Hard-coded logic vs. configurable logic
- Test coverage
- Error handling
- Logging consistency

Specific to SPARC:

- Is grading logic isolated and auditable?
- Are rubric values (1, 0.5, 0) cleanly implemented?
- Is feedback logic clearly mapped to scoring logic?
- Are JSON schemas consistent and validated?

Deliverable:

Bullet list of:

- Critical risks
- Moderate issues
- Minor improvements
- Immediate stabilization recommendations

FINDINGS - Phase 2: Code Quality Review

Name:	Issue Identified	Priority	Possible Solutions
	Compile ipynb into plain py files Jupyter Notebook files will work fine especially in development, but they can execute in an order that is not intuitive. Compiling ipynb to py files for production will make debugging easier in the future.	Low ▾	Several possibilities. Jupyter may be a good option because it allows py and ipynb versions of code to sync back and forth rather than having to manually upkeep the sync.
	Direct pip not recommended for package management on HiPerGator Because pip installs packages into shared or user-level directories, it can conflict with packages installed elsewhere in the HiPerGator environment, resulting in incorrect package versions and other	Critical ▾	Conda environments should be used to isolate project dependencies (ie - specific external software packages needed for the overall app to function)

	unexpected behaviors.		
Rishma	Automated test suite for the agent pipeline. No unit tests for grading logic, no integration tests for the Supervisor routing, and no regression tests for schema compliance.	Medium ▾	Create a test_sparc.py file with: (1) unit tests for the grading score-to-display mapping, (2) integration tests to validate Coach JSON output structure, (3) a regression test that runs clinician inputs through the full pipeline and asserts expected scores.
Rishma / P	The temperature parameter set per-agent: documentation and reasoning or other values possible	Low ▾	Documenting rationale for each temperature setting in a config file. slight variation (0.3-0.5) for natural dialogue vs. deterministic (0.1-0.2) for consistent grading. Run evaluation sessions to validate settings.
Pranaya	Sanitization Failure mode is unsafe (1_SPARC_Agent_Training.ipynb)	High ▾	sanitize_text_with_presidio() function catches exception and returns original unsanitized txt. If Presidio fails, raw PHI passes through to the vector db. Change to fail - closed approach, if presidio fails, catch and stop ingestion?
Pranaya	Hardcoded paths	Not set ▾	Paths like /blue/jasondeanarnold/SPARCP/, email addresses and mode names are hard coded. When another team member runs code, they'll need to manually find and edit across multiple files. Keep everything in config file

Sai:

Finding	File	Issue Identified	Priority	Solution
cuda-toolkit=12.8 does not exist — conda env creation fails	<i>environment_training.yml, environment_backend.yml</i>	cuda-toolkit=12.8 does not exist as a conda package (max ~11.8 on conda-forge/nvidia). Also redundant with cuda=12.8 already listed. conda env create fails with unsatisfiable solver error before any packages install. Blocks all project setup.	Critical	Remove cuda-toolkit=12.8 from both YAMLS. Keep only cuda=12.8 from nvidia channel. pytorch-cuda=12.8 provides the necessary CUDA runtime.
riva-asr-lib-decoder not on public PyPI — backend env fails	<i>environment_backend.yml</i>	pip dependency riva-asr-lib-decoder is NVIDIA-internal, not on public PyPI. pip install fails with 'No matching distribution found.' Backend conda env cannot be created.	Critical	Remove from pip deps. It is bundled inside the Riva container image and not needed in the host conda env.
SFTTrainer receives list-of-dicts instead	<i>NB1</i>	Sec 5.0 run_qlora_training(): dataset_text_field='messages' expects plain text strings. The	Critical	Use formatting_func: def format_chat(ex): return tokenizer.apply_chat_template(ex[

Finding	File	Issue Identified	Priority	Solution
of string — training data corrupted		messages field is a list-of-dicts in ChatML format. SFTTrainer either errors or stringifies the list, producing corrupted training data. packing=True further corrupts conversational data.		'messages'], tokenize=False). Pass formatting_func=format_chat to SFTTrainer.
load_in_4bit not a valid from_pretrained parameter — model loads at full precision	NB4	Sec 6.1 load_models(): from_pretrained(name, load_in_4bit=True) — not a valid kwarg. Model loads at full precision → OOM on L4 24 GB. NB1 Sec 5.0 does this correctly with quantization_config=bnb_config.	Critical	Use BitsAndBytesConfig: bnb = BitsAndBytesConfig(load_in_4bit=True, bnb_4bit_quant_type='nf4', bnb_4bit_compute_dtype=torch.bfloat16). Pass quantization_config=bnb to from_pretrained.
SLURM script references wrong filename train_agent.py	NB1	Sec 7.1 generate_slurm_script(): SLURM calls python train_agent.py but actual training function is run_qlora_training() in Sec 5.0. No train_agent.py exists. Also incorrectly referenced in README.md.	High	Option A: Create train_agent.py as a CLI wrapper with argparse. Option B: Update SLURM to use jupyter nbconvert --to script --execute.
LORA_CONFIG and TRAINING_ARGS defined but silently ignored	NB1	Sec 2.1 defines config dicts. run_qlora_training() Sec 5.0 hardcodes completely different values. User edits to config cell have zero effect — a silent correctness bug.	High	Have run_qlora_training() read from LORA_CONFIG and TRAINING_ARGS dicts. Or remove the unused config dicts.
conda activate fails in SLURM batch scripts	NB1, NB3	NB1 Sec 7.1 and NB3 SLURM launch script: conda activate in a non-interactive shell raises CommandNotFoundError. Shell hook not initialized first.	High	Add before conda activate: module load conda && eval "\$(conda shell.bash hook)" or source \$(conda info --base)/etc/profile.d/conda.sh.
SLURM allocates 4 GPUs — only 1 used	NB1	Sec 7.1: --ntasks=4 --gpus-per-task=1 allocates 4 GPUs. Single-process python train_agent.py uses only 1. Burns 4x the HiPerGator SU allocation.	High	Change to --ntasks=1 --gpus=1 --cpus-per-task=16. For multi-GPU use torchrun --nproc_per_node=4 or accelerate launch.
NB3 SLURM calls uvicorn main:app but main.py never written to disk	NB3	SLURM launch script: uvicorn main:app expects /backend/main.py. NB3 defines FastAPI app in a cell but has no export/write step. SLURM job fails with ModuleNotFoundError.	High	Add export cell: with open(os.path.join(OUTPUT_DIR, 'backend', 'main.py'), 'w') as f: f.write(code). Or use jupyter nbconvert --to script.
Dockerfile COPY commands reference wrong paths — container builds fail	NB2b	Secs 2.2-2.4: Dockerfile.mas uses COPY requirements.txt — no requirements.txt exists. Dockerfile.unity-server uses COPY Build/LinuxServer/ — wrong path. Dockerfile.signaling uses COPY signaling/ — also wrong. All three podman build commands fail with COPY failed: file not found. NB4b Sec 4.2 references these same broken images.	High	Create missing artifacts: (1) generate requirements.txt from pip deps, (2) add Unity Linux server build output, (3) add signaling server source. Document the build pipeline.
NB2b SLURM allocates 4 GPUs for single container process	NB2b	Sec 4.4: --gpus-per-task=4 --cpus-per-task=32 --mem=256gb for a single aptainer exec running one uvicorn process. 3 GPUs wasted, 256 GB RAM over-allocated.	High	Reduce to --gpus=1 --cpus-per-task=8 --mem=64gb. Use separate SLURM jobs for additional containers.
Hardcoded paths in 30+ locations	All notebooks	/blue/jasondeanarnold/SPARCP/ hardcoded across all notebooks. jayrosen@ufl.edu in SLURM scripts. pubapps-vm.rc.ufl.edu in NB4. No other team member can run without manual find-and-replace.	Moderate	Config cell at top of each notebook: BASE_PATH = os.environ.get('SPARC_BASE_PATH', '/blue/jasondeanarnold/SPARCP'). SLURM_EMAIL =

Finding	File	Issue Identified	Priority	Solution
				os.environ.get('SPARC_SLURM_EMAIL', 'jayrosen@ufl.edu').
Missing imports cause NameError on cell execution	NB1	List[str] without typing import (Sec 4.3), Dataset without import (Sec 4.2), BaseModel/ValidationError without pydantic import (Sec 6.2), import json missing (Secs 6.2 and 8.0), train_agent() called but only run_qlora_training() defined (Sec 5.2).	Moderate	Add consolidated import cell at top of NB1 with all required imports. Rename train_agent() calls to run_qlora_training().
SLURM --time=UNLIMITED invalid format	NB2	Sec 4.2: #SBATCH --time=UNLIMITED is not a valid SLURM time format. Scheduler rejects unless QOS explicitly permits it. NB3 uses valid --time=7-00:00:00.	Moderate	Replace with #SBATCH --time=7-00:00:00 or the maximum allowed by the QOS.
module load cuda/12.8 conflicts with conda CUDA	NB1, README.md	SLURM scripts load system CUDA module alongside conda cuda=12.8, creating LD_LIBRARY_PATH conflicts and potential libcudart.so version mismatches.	Moderate	Remove module load cuda/12.8 when conda env already provides cuda=12.8 and pytorch-cuda=12.8.
Missing pymupdf dependency — import fitz fails	NB1, environment_training.yml	Sec 4.2 extract_text_from_document() imports fitz (PyMuPDF). Neither pymupdf nor fitz in environment_training.yml. Listed pymupdf4llm is a different lightweight package that does not install the full library.	Moderate	Add pymupdf>=1.23.0 to pip deps in environment_training.yml.
Missing import json — NameError in validation and orchestrator	NB1	json.loads() and json.dumps() used in Sec 6.2 validate_agent() and Sec 8.0 multi_agent_orchestrator(). import json never appears anywhere in NB1. Will raise NameError.	Moderate	Add import json to consolidated imports cell at top of NB1.
Presidio OperatorConfig uses literal string — entity type not substituted	NB1	Sec 4.2 sanitize_text_with_presidio(): new_value='<{entity_type}>' is a literal string, not an f-string. Every detected entity gets replaced with literal text <{entity_type}> instead of <PERSON>, <PHONE_NUMBER>, etc.	Moderate	Remove the operators parameter. Presidio's default replace behavior substitutes the entity type name correctly.
Riva version mismatch between setup cells	NB3	Sec 2.3 defines RIVA_VERSION='2.16.0' pulling riva-speech:2.16.0-server. Sec 2.2 references riva_quickstart_v2.14.0/config.sh. Incompatible quickstart scripts for the pulled server version.	Moderate	Standardize on one Riva version throughout — update all references to match RIVA_VERSION.
langchain-openai dependency — undocumented external API	environment_bakend.yml	Includes langchain-openai requiring an OpenAI API key. Project docs state local HuggingFace models only. Either unused or undocumented data leaving UF infrastructure — compliance concern for HIPAA-adjacent system.	Moderate	Remove if not used. If used, document data flow and HIPAA implications.
No version upper bounds on pip dependencies	environment_training.yml, environment_bakend.yml	~20 pip packages use >= with no upper bound. LangChain has aggressive breaking changes. Environment solves are non-reproducible.	Moderate	Pin major.minor: langchain>=0.1.0,<0.3; transformers>=4.36.0,<4.50; etc.
build_vector_store() missing return statement	NB1	Sec 4.3: Creates Chroma object but has no return statement. Object lost on exit. Also missing collection_name param unlike ingest_documents().	Low	Add return vector_store. Add collection_name=collection_name to Chroma.from_documents() call.

Finding	File	Issue Identified	Priority	Solution
Guardrails config written to unpredictable directory	NB3	Sec 3.2 create_rails_config(): config.yml and topical_rails.co written to notebook CWD. May overwrite unrelated files.	Low	Write to os.path.join(OUTPUT_DIR, 'guardrails', 'config.yml') with mkdirs exist_ok=True.
Duplicate section numbering in NB3	NB3	Section 2.3 appears twice (Riva setup AND Riva launch). Section 4.2 also appears twice. Causes confusion in cross-references.	Low	Renumber sections sequentially throughout NB3.

What Will Fail: HiPerGator

Failure	Root Cause	Finding
Conda env fails	cuda toolkit=12.8 invalid; unsatisfiable solver error	P2 F1 — Critical
SLURM wrong entry point	Script calls python train_agent.py — file does not exist	P2 F5 — High
conda activate fails	Non-interactive SLURM shell; conda hook not initialized	P2 F7 — High
SLURM --time rejected	--time=UNLIMITED invalid SLURM format	P2 F13 — Moderate
Cells crash NameError	Missing imports across NB1 Secs 4.2, 4.3, 6.2, 5.2	P2 F12 — Moderate
Training data garbage	SFTTrainer gets list-of-dicts not string	P2 F3 — Critical
4 GPUs allocated, 1 used	Single-process script on 4-GPU SLURM allocation	P2 F8 — High
NB3 SLURM fails	uvicorn main:app called but main.py never written to disk	P2 F9 — High

Additional NOTES:

Name - [finding]

Rishma -

Critical risks

1. Some hard-coded model names and I/O text -> make config files and environment specific (use environment variables, config files, feature flags)
2. Schema enforcement at API layer and may cause runtime breakage risk (e.g: keyerror, datatype mismatch, parsing issues) -> define datatypes and constraints in classes
3. Grading logic partly in Unity might make it harder to audit backend consistency
4. Hardcoded Firebase Credentials Path FIREBASE_CREDS =

```
"{PUBAPPS_ROOT}/config/firebase-credentials.json"
```

Should be environment variable:

```
FIREBASE_CREDS = os.getenv("FIREBASE_CREDS_PATH")
```
5. Model Version in Response
6. CORS = Completely Open (Critical) : use allow origins

1. Moderate issues

1. Mention structured error codes in API + consistent variable name parents instead of caregiver
caregiver_text=response_data.get("text", "Error")

2. Minor improvements

1. Central config.yaml for: model, max tokens, etc.
2. Central constants file for rubric definitions.

3. Dedicated schemas.py for all Pydantic models.

4. Version tagging of system prompts.

3. **Immediate stabilization recommendations**

1. If .py files used:

backend/

- main.py
- supervisor.py
- coach.py
- parent.py
- schemas.py
- scoring.py
- safety.py

2. Add: Structured logging (JSON logs), Request ID per session, Retry wrapper with max 1 retry for malformed JSON

3. Add schema validation middleware before sending to Unity.

4. Sensitive data exposure: logging.info(f"Session: {request.session_id} | User Input: {request.user_transcript}") : Logging full user transcripts can expose sensitive information + unauthorized data leak risk. Fix: mask or truncate sensitive info if necessary:

```
safe_transcript = request.user_transcript[:50] + "..." if len(request.user_transcript) > 50 else request.user_transcript
```

```
logging.info(f"Session: {request.session_id} | User Input: {safe_transcript}")
```

```
Log:session_id, mode, latency, error_code
```

5. No error logging: Currently logging only user input, not errors or exceptions. Fix: wrap in try/except and log errors:

try:

```
result = await app_graph.ainvoke(initial_state)
```

except Exception as e:

```
logging.error(f"Session: {request.session_id} | LangGraph error: {e}")
```

```
raise HTTPException(status_code=500, detail="Internal processing error")
```

Model Selection Recommendation: Reference of gpt-oss-20b (OpenAI OSS) as the primary model. Based on the constraints (24GB L4 VRAM, medical reasoning, HPV vaccine domain, HIPAA-adjacent):

Model	Size	L4 Fit (4-bit)	Medical Reasoning	License
gpt-oss-20b (current)	20B	Yes (~13GB)	Good (UF Navigator)	UF RC Access
Meta LLaMA 3.1 8B Instruct	8B	Yes (~5GB)	Good, well-documented	Open (Meta)
Mistral 7B Instruct v0.3	7B	Yes (~5GB)	Strong instruction follow	Apache 2.0

Recommendation: Retain gpt-oss-20b for the primary deployment given the existing UF Navigator infrastructure. Evaluate LLaMA 3.1 8B as a secondary option for cost and memory efficiency. Both should be benchmarked on a held-out set of HPV vaccine counseling transcripts before selecting the production model.

Phase 3: AI Safety & Production Readiness Review

Evaluate:

- Guardrails against inappropriate content
- Handling of policy violations
- Prevention of PII ingestion
- Structured output enforcement
- Determinism of grading outputs
- Handling of edge cases
- Retry/fallback logic
- Logging for traceability
- API security considerations

Questions:

- Could the system produce unsafe or out-of-scope output?
- Is fallback logic transparent or hidden?
- Are grading overrides occurring implicitly?
- Is the system reproducible?

Deliverable:

Risk assessment categorized by:

- High Risk (must fix before pilot)
- Medium Risk (fix recommended before scale)
- Low Risk (future improvement)

FINDINGS - Phase 3: AI Safety & Production Review

Name:	Issue Identified	Priority	Possible Solutions
Rishma	Hallucination Detection Layer: Does output contain claims not found in RAG context?	High ▾	Add post-generation check: regenerate with "ground strictly in provided sources."
Rishma	Privacy control check	Medium ▾	Check how long logs retained, Firebase storing transcripts, IP logging tied to transcript? Suggestion: Session transcript deleted after session + Only rubric scores retained with random session ID + Log retention policy

			and archive logs to Firebase/GCP Firestore for long-term storage.
		Not set ▾	
Pranaya	No retry/ fallback logic for LLM calls 3_SPARC_RIVA_Backend	Not set ▾	If LLM call fails, FastAPI endpoint raises HTTPException and entire session crashes. Could wrap in a retry decorator (max 2-3) and add a fallbackscripted response so session can continue? Log all failures Tenacity library
Pranaya	No API authentication on v1/chat	Not set ▾	FastAPI endpoint has no authentication? Anyone who can reach port 8000 on the PubApps server can call the LLM. That means free GPU access for anyone, and the ability to inject fake sessions into Firebase.
Pranaya	Keyword only safety is bypassable (3_SPARC_RIVA_Backend) multi_agent_orchestrator()	Not set ▾	Could implement NeMo Guardrails with semantic input scanner, not just keyword matching. (current safety check is, if certain keyword is used then refuse, but can be bypassable if clinical-sounding lang is used. Currently uses string matching.) Check nemo guardrails capabilities

Sai:

Finding	File	Issue Identified	Priority	Solution
PII sanitization fails open — unsanitized text enters model weights	NB1	Sec 4.2 sanitize_text_with_presidio(): except Exception: return text — returns ORIGINAL text on any Presidio error. PII flows into ChromaDB and LoRA training data. Once embedded in weights, cannot be removed post-deployment.	High Risk	Fail closed: return " or raise on error. Log the error. Add tenacity retry before failing. Never pass unsanitized text downstream.
Audit log writes PHI to disk — HIPAA violation	NB3	Sec 6.1: logs full user_transcript to audit.log on disk. Sec 7.0 of same notebook states 'No PHI is written to disk.' Direct contradiction.	High Risk	Remove {request.user_transcript} from all log lines. Log only: session_id, timestamp, agent_type, is_safe flag, latency_ms.
Safety guardrails are trivial keyword matching — model output never scanned	NB1, NB3, NB4	Three separate keyword lists across notebooks, all bypassed by misspelling or synonyms. Only user INPUT is checked — model OUTPUT is never scanned.	High Risk	Integrate NeMo Guardrails (already in env, config files generated). Add output scanning before returning response. Add

Finding	File	Issue Identified	Priority	Solution
		Caregiver could hallucinate harmful medical advice and nothing catches it before delivery.		input length limit (2000 chars max).
async process_chat() blocks event loop during inference	NB4	Sec 6.1: async def but model.generate() is synchronous and GPU-bound. Blocks entire asyncio event loop during inference. Server completely unresponsive to all requests including /health during every call.	High Risk	Change to regular def (uvicorn runs sync handlers in threadpool), or wrap with await asyncio.to_thread(model.generate, **kwargs).
Riva container uses wrong GPU passthrough — CUDA init fails	NB4	Sec 5.1 Quadlet: raw AddDevice=/dev/nvidia0 instead of CDI Device=nvidia.com/gpu=all. NVIDIA runtime libraries invisible inside Podman container. Riva fails with CUDA init error at startup. NB4b uses the correct CDI approach.	High Risk	Replace AddDevice lines with Device=nvidia.com/gpu=all. Requires nvidia-container-toolkit CDI on PubApp VM.
Riva init step uses Docker syntax — GPU access fails	NB4	Sec 4.1: podman run --rm -it --gpus all ... — '-gpus all' is Docker syntax. Podman on PubApp uses CDI passthrough (--device nvidia.com/gpu=all). Init step fails to access GPU → Riva models never downloaded/optimized. NB4b uses the correct CDI approach.	High Risk	Replace --gpus all with --device nvidia.com/gpu=all in the Riva init command. Quadlet GPU passthrough addressed in finding above.
ChatRequest missing field constraints — DoS and log injection	NB3, NB4	user_message/user_transcript defined as plain str with no max_length. Multi-MB strings exhaust RAM during tokenization. session_id accepts newlines (log injection) and ../ path traversal. audio_data (NB4) accepts unbounded base64 strings.	High Risk	Add Pydantic Field validators: user_message: str = Field(..., max_length=10000); session_id: str = Field(..., max_length=128, pattern=r'^[a-zA-Z0-9_-]+\$').
Health check reports healthy when models fail to load	NB4	Sec 6.1 /health: checks Riva connectivity but not LLM model readiness. If load_models() fails, tokenizer/model remain None but endpoint returns status:healthy. Load balancers route traffic to broken instance.	High Risk	Add: model_ok = tokenizer is not None and caregiver_model is not None. Return HTTP 503 or status:'degraded' if models not loaded.
Raw exception details leak to client on inference failure	NB3, NB4	model.generate() failure raises HTTPException(500, detail=str(e)). Internal file paths, CUDA errors, Python stack traces sent directly to client — information disclosure vulnerability.	High Risk	Catch inference exceptions and return generic: HTTPException(500, detail='Internal server error'). Log actual exception server-side only.
CORS allows all origins with credentials — startup error	NB4	Sec 6.1: allow_origins=["*"] + allow_credentials=True is invalid per CORS spec. Starlette 0.27+ raises ValueError at startup.	Medium Risk	Set allow_origins=["https://sparc-p.rc.ufl.edu"]. Remove allow_credentials=True. Restrict allow_methods=['GET','POST'] and allow_headers.
Firebase credentials stored as plaintext on disk	NB4	Sec 6.1: firebase-credentials.json stored with no file permissions or encryption. HashiCorp Vault mentioned in project docs but not implemented.	Medium Risk	chmod 600 + add to .gitignore. Better: use GOOGLE_APPLICATION_CREDENTIALS env var in systemd service file.
Non-deterministic Coach grading	NB4	Sec 6.1: All agents use do_sample=True, temperature=0.7. Same input produces different Coach grades on repeated runs. No seed set. No model version pinned.	Medium Risk	Coach: set do_sample=False (greedy decoding). Caregiver: keep temperature=0.7. Create per-agent generation config dicts.

Finding	File	Issue Identified	Priority	Solution
Deployment script omits loginctl enable-linger	NB4	Secs 5.2 and 6.3: systemctl --user enable commands run without prerequisite loginctl enable-linger. Without this, user services terminate when SSH session ends.	Medium Risk	Add loginctl enable-linger \$(whoami) before systemctl --user enable commands in the deployment script.
Riva gRPC connection created per-request — no pooling	NB4	Sec 6.1: New riva.client.Auth and SpeechSynthesisService created on every /v1/chat request. Adds TCP setup latency per request and leaks file descriptors since channels never explicitly closed.	Medium Risk	Create Auth and service objects once in load_models() at startup. Reuse across all requests.
TTS audio as inline base64 WAV — unbounded response size	NB4	Sec 6.1: 30-second WAV at 22kHz/16-bit ~1.7 MB base64 embedded in JSON response body. Concurrent requests on 16 GB RAM server create significant memory pressure.	Medium Risk	Use compressed format (opus/mp3), store in temp file and return URL, or stream over WebSocket.
Deprecated FastAPI startup event handler	NB4	Sec 6.1: @app.on_event('startup') deprecated since FastAPI 0.93+. Deprecation warnings generated; will be removed in a future version.	Low Risk	Migrate to: from contextlib import asynccontextmanager; @asynccontextmanager async def lifespan(app): yield; app = FastAPI(lifespan=lifespan).

Additional NOTES:

Name - [finding]

1. Rishma

High Risk

1. LLM hallucination risk in medical reasoning
2. No explicit uncertainty language enforcement could cause medical compliance risk
3. Prompt injection detection beyond topical rails to prevent security risk
4. Parent agent possibly giving advice must be blocked

Medium Risk

1. Logging PHI edge cases
2. Long session memory ambiguity
3. Structured hallucination detector layer

Low Risk

1. Tone analysis (future)
2. Multi-modal emotion tuning
3. Pixel streaming vs WebGL performance
4. Mechanism to collect and use high-quality session interactions for iterative model improvement. Without this, the fine-tuned model cannot improve based on real pilot data. Suggestion: Implement an opt-in session recording consent flow (even for anonymous users): if a clinician consents, the anonymized transcript + Coach JSON is tagged for human expert review. Sessions that a human expert rates as high quality (e.g., score ≥ 0.85 across all C-LEAR phases) are queued as fine-tuning candidates for the next model iteration.
5. The system does not currently have a 'meta-feedback' mechanism: clinicians cannot rate the quality of the Coach feedback they received, which eliminates a valuable signal for improving the Coach agent's pedagogical calibration.

Suggestion: Add a simple post-session thumbs-up/thumbs-down rating per Coach feedback card in the Unity UI (already has a thumbs-up icon visible in screenshots). Log these ratings in Firebase. Use sessions with consistently negative feedback ratings as inputs for Coach prompt refinement.

TEST CASES:

1. Parent refuses to continue.
2. Parent switches to an unrelated topic <XYZ>.
3. Doctor gives incorrect medical information.
4. Parent shares personal information/unrelated incident.
5. Doctor asks coach about vaccines other than HPV/similar to HPV.
6. Prompt injection attempt: "Ignore previous instructions."
7. Concurrency load testing (GPU for ASR + LLM + TTS. Measure tokens/sec under 5 concurrent users) + GPU memory contention and inference timeouts or OOM errors
8. Handling for the 'medical emergency role-play' edge case: if a clinician trainee uses language suggesting the simulated child is in acute distress (e.g., Child is having a reaction), the system will continue the HPV discussion.
9. HIPAA compliance: The system processes clinician voice audio via Riva ASR. Even though the content is fictional (the clinician is practicing), if a clinician inadvertently mentions a real patient (e.g., 'I had a patient like this last week named John'), that audio transcript could constitute inadvertent PHI capture.
10. User attempts to input actual patient record, User attempts prompt injection override, LLM attempts to provide medical advice beyond scope. Define: Immediate refusal logic + Safe termination behavior

Doubts:

1. UFIT form states: Does your application process sensitive data? Answer: No. Even if not stored: Voice input, Free-form transcripts, Potential self-disclosed PHI

Will it qualify as transient sensitive data processing under UF definitions?

Can we reclassify as: Open Data (training material) + Sensitive (transient voice/text input) + No Restricted storage

Evaluation Criteria

Reviewers should assess the system against the following standards:

- Modularity
- Transparency
- Determinism
- Testability
- Security
- Scalability
- Documentation clarity

- Alignment with stated architecture
-

Recommended Output Format

Each reviewer will submit:

1. 2–3 page written review
2. Categorized findings (Critical / Moderate / Minor)
3. Specific code references (file names, function names)
4. Concrete recommendations

After individual reviews are submitted:

- A synthesis meeting will be held
- Findings will be consolidated
- Action items will be prioritized