

# Program Increment (PI) Plan: LMForge

**Project:** LMForge - Locally Hosted LLM Builder **PI Duration:** Fall Semester 2025 (5 Sprints)

## 1. Executive Summary & Project Evolution

**Vision:** To democratize AI development by providing a seamless, privacy-conscious platform ("LMForge") that enables non-technical users to scrape data, clean it, and fine-tune open-source models for domain-specific tasks without relying on cloud infrastructure.

### Project State: Current Status (Backend & Frontend Focus)

Feature Area	Current State (End of PI)
Backend Architecture	<b>Microservices Architecture:</b> Five distinct backend modules (FastAPI, Django, Flask) that process data independently.
User Interface (UI)	<b>Unified Design System:</b> A specialized cross-functional sub-team integrated a polished, consistent React-based UI across all modules based on a validated Figma prototype.
Data Pipeline	<b>Functional Silos:</b> The system can successfully Scrape (Reddit/Web), Clean (Pandas), and Vectorize (pgvector), though the handoff currently requires manual file transfer.
Infrastructure	<b>Hybrid Containerization:</b> RAG and Database are fully Dockerized; Agentic uses MCP; Fine Tuning & Scraping run in local Python environments due to resource constraints.
Key Metrics	<b>36% Reduction</b> in token usage (Cleaning Team); <b>100% Success rate</b> on Agent tool calls; <b>0% Duplicate headings</b> in processed data.

## 2. Team Roster & Technical Achievements

A deep dive into the specific backend logic, architectural choices, and deliverables defined in the final team documentation.

Team	Key Deliverables & Architecture
<b>UI/UX Design Sub-team</b> <i>(Cross-functional Tiger Team)</i>	<ul style="list-style-type: none"><li>• <b>Deliverable:</b> Created a comprehensive <b>Figma Prototype</b> to address the initial lack of visual polish.</li><li>• <b>Integration:</b> Successfully translated the Figma design into a unified <b>React Frontend</b>, replacing disjointed team-specific interfaces.</li></ul>
<b>Agile Agentic Systems</b>	<ul style="list-style-type: none"><li>• <b>Solution:</b> "LMForge Agent Creator" — A no-code tool to build agents.</li><li>• <b>Architecture:</b> Implemented <b>FastAPI</b> to serve agent logic and adopted <b>MCP (Model Context Protocol)</b> to standardize tool calling.</li><li>• <b>Methodology:</b> Used <b>BPMN</b> to align team goals and resolve initial scope ambiguity.</li><li>• <b>Tech Stack:</b> Python, LangChain, MCP-Use, Figma.</li></ul>
<b>Web Scraping</b>	<ul style="list-style-type: none"><li>• <b>Ingestion Engine:</b> Built a multi-threaded Python scraper using <b>Selenium</b> and <b>Reddit API</b>.</li><li>• <b>Compliance:</b> Focused on "Legal Scraping" by prioritizing official APIs over brute-force methods.</li><li>• <b>Logic:</b> Implemented batch processing for PDFs and URLs.</li></ul>
<b>Data Cleaning</b>	<ul style="list-style-type: none"><li>• <b>Algorithms:</b> Developed "Content-Aware" cleaning scripts that utilize HTML class awareness to strip non-content elements.</li><li>• <b>Optimization:</b> Achieved a <b>63% reduction in lines</b> and <b>36% reduction in characters</b>, directly saving compute costs.</li></ul>

	<ul style="list-style-type: none"> <li>• <b>Logic:</b> Implemented an "Accreditation Block" to preserve author/source metadata while cleaning body text.</li> </ul>
<b>RAG Database</b>	<ul style="list-style-type: none"> <li>• <b>Database:</b> Migrated from MySQL to <b>PostgreSQL with pgvector</b> (IVFFlat indexing).</li> <li>• <b>Models:</b> Standardized on <b>Qwen2.5:0.5b-instruct</b> for chat and <b>all-minilm:33m</b> for embeddings to balance performance with local resource limits.</li> <li>• <b>Architecture:</b> Fully Dockerized the service to ensure reproducibility.</li> </ul>
<b>Fine Tuning</b>	<ul style="list-style-type: none"> <li>• <b>Optimization:</b> Implemented <b>SINQ (Quantization)</b> logic to allow heavy models to run on student laptops.</li> <li>• <b>Implementation:</b> Moved from EconBERTa to a <b>BERT</b> implementation focused on the SQuAD database.</li> <li>• <b>Workflow:</b> Created a Kanban-style automation dashboard to track training jobs.</li> </ul>

### 3. Technical Architecture & Commonality Analysis

*Analysis of high-level components based on final project outcomes.*

#### Commonalities (Strengths)

- **Language:** All five teams standardized on **Python** for backend logic.
- **Frontend:** Convergence on a **Unified React UI** ensured a consistent user experience.
- **Local Focus:** All teams optimized for local execution (e.g., RAG selecting 0.5b parameter models; Fine Tuning using SINQ).

#### Divergences (Friction Points)

- **Databases:** RAG moved to **PostgreSQL**, while Fine Tuning prepares SQuAD data independently.
- **Environment:** RAG is **Docker-based**, while other teams rely on local **requirements.txt**. RAG documentation notes that "Containerization... may limit some future development" if not managed carefully across teams.

- **Data Interfaces:** The Web Scraping team outputs JSON, but the RAG team requires specific chunking formats that are not yet fully aligned.

## 4. Priorities for Future Classes (The Handoff)

*Specific recommendations derived from the "Future Recommendations" section of each team's final report.*

### Priority 1: The "Smart Router" (Web Scraping Recommendation)

- **Objective:** Automate input handling.
- **Deliverable:** Build a "Smart Router" that recognizes pasted content (Text vs. CSV vs. URL vs. Reddit Link) and automatically dispatches it to the correct processor without user selection.

### Priority 2: Advanced RAG Chunking (RAG Recommendation)

- **Objective:** Improve retrieval quality.
- **Deliverable:** Move beyond simple paragraph splits. Implement **Sentence-Window** Chunking and metadata filtering.
- **Action:** RAG team explicitly requests collaboration with Scraping/Cleaning to standardize the JSON input format for these chunks.

### Priority 3: UI Implementation Fidelity (Design Team Recommendation)

- **Objective:** Ensure code matches design.
- **Action:** Future developers must utilize the provided **Figma Layout** as the strict "source of truth." Do not code ad-hoc styles; follow the grid and typography defined in the prototype.

### Priority 4: Accreditation & Metadata (Data Cleaning Recommendation)

- **Objective:** Preserve data provenance.
- **Deliverable:** Ensure the "Accreditation Block" logic (created by the Cleaning team) is preserved and visualized in the final RAG citation output.

### Priority 5: User-Friendly Entry (Fine Tuning Recommendation)

- **Objective:** Lower the barrier to entry.
- **Deliverable:** Implement specific UI graphs below the log terminal to visualize training loss/accuracy in real-time for non-technical users (Business majors).

## 5. Iteration Schedule Summary

- **Sprint 1:** Architecture & Research (Pivots from Amazon/Twitter; RAG explores Vector DBs).
- **Sprint 2:** Prototyping (Figma Design formed; Agentic creates BPMN).
- **Sprint 3:** Core Development (FastAPI, Content-Aware Algorithms, BERT implementation).
- **Sprint 4:** Breakthroughs (MCP Integration, pgvector migration, 36% noise reduction).
- **Sprint 5:** Integration & Polish (React Frontend integration, Dockerization of RAG, Final Documentation).

## 6. Program Board (Dependencies)

Dependent Team	Needs...	From...	Status
All Teams	Unified UI Components & Design System	UI/UX Sub-team	[Resolved]
Data Cleaning	Raw Data Output (JSON/Text)	Web Scraping	[Resolved]
RAG	Defined Data Input Format (Standardized JSON)	Data Cleaning	[Critical Priority]
Agentic Systems	Functional RAG API to query docs	RAG Database	[At Risk]
All Teams	Unified Docker Environment	DevOps/Arch	[Blocked]

## 7. ROAM Board (Risk Management)

Risk Description	Impact	Category	Mitigation Strategy
<b>Environment Inconsistency:</b> requirements.txt differs per machine.	Critical	Owned	<b>Future Action:</b> Mandate Docker for all teams in Sprint 1 of next semester (RAG team has provided the template).
<b>Containerization Limits:</b> RAG team notes full containerization is resource-heavy.	High	Accepted	<b>Mitigation:</b> Use a hybrid approach where only DB/API are containerized, and ML training runs on bare metal if GPU is required.
<b>Knowledge Transfer:</b> Complexity of LangFlow/n8n concepts.	Medium	Mitigated	<b>Action:</b> The "Agent Creator" tool abstracts this complexity, reducing the learning curve for new users.
<b>Ambiguity:</b> Lack of clear task definition (Cleaning Team).	Medium	Resolved	<b>Lesson:</b> Use BPMN and visual workflows (like Agentic team did) to align goals early in Sprint 1.

## 8. Strategic Roadmap Alignment (Year 1)

*Mapping this PI to the 5-Year Vision.*

<b>Strategic Theme</b>	<b>Year 1 Goal (Roadmap)</b>	<b>Status</b>
<b>Ease of Use</b>	Install, bug issues, GUI basics	[Mixed] GUI is <b>Live &amp; Unified</b> , but installation remains complex.
<b>RAG Database</b>	Vector DB Integration	[Done] Successfully migrated MySQL -> Postgres/pgvector.
<b>Data Cleaning</b>	Basic cleaning scripts & processes	[Done] Content-aware scripts operational (36% reduction).
<b>Agentic Systems</b>	Prototype basic agent tool	[Done] "LMForge Agent Creator" with MCP is live.