

Local LLM - Student Facing Rubrics		
How Student will be Graded		
<p>Your grade will be based on both your code implementation and your reflection/reporting. Some parts will be automatically graded (your code must run correctly), while others will be manually graded (clarity of your report, depth of reflection).</p>		
Section	What's Expected	Weight
<b>Environment Setup &amp; Dependencies</b>	Set up Python environment correctly, install all required libraries, place models properly, and ensure the model loads.	15%
<b>Core Functions Implementation</b>	Implement all required functions (chunk_text, embeddings, FAISS, search, run_llm). Functions should return correct results and demonstrate your understanding.	30%
<b>PDF Context Integration</b>	Successfully extract text from the PDF, generate embeddings, and use them to improve query responses. Show differences between queries with and without PDF context.	20%
<b>Reflection &amp; Reporting</b>	Submit a short (1-2 pages) report that: <ol style="list-style-type: none"> <li>1. Summarizes your setup steps</li> <li>2. Describes errors and how you solved them</li> <li>3. Compares results with and without context</li> <li>4. Reflects on how LLMs can be applied in healthcare.</li> </ol>	20%
<b>Submission Quality</b>	Your Jupyter Notebook runs end-to-end without errors, includes clear outputs, and your report/screenshots are well-organized and concise.	15%
<b>Notes for Students</b>		
Autograding: Many functions will be tested automatically (e.g., chunking, embeddings, FAISS index). If your notebook doesn't run, you'll lose those points.		
Manual review: Reports, reflections, and the clarity of your responses will be graded by the professor/TA.		
Report requirement: 1-2 pages, concise and structured. Include examples of queries (with and without context).		
Evidence: Include screenshots or copied outputs from your notebook showing that your code works.		
Pro tip: Make sure your functions run as expected with simple test cases before moving to the full PDF workflow.		

Local LLM - Professor Facing Rubrics						
How to Grade						
Grade only what's demonstrably shown: code that runs end-to-end, clear outputs/screenshots, and a concise report. Combine autograder results (does it execute and produce required artifacts?) with manual judgment (clarity, correctness, depth of analysis). Apply the four levels consistently—Exceeds / Meets / Approaches / Missing mapped to 100 / 75 / 50 / 25% of each criterion's weight—and favor reproducibility, no missing files, brittle paths, or hidden dependencies. If evidence is missing or unverifiable, assign the lower level.						
Section	Criterion	Weight	Excellent (100%) - Exceeds Expectations	Proficient (75%) - Meets Expectations	Basic (50%) - Partially Meets Expectations	Inadequate (25%) - Does Not Meet
<b>Environment Setup &amp; Dependencies (15%)</b>	Python env & libraries	5 pts	All dependencies installed; imports run without error.	All dependencies installed; imports run without error.	Minor missing library/import, easily fixed.	Environment fails to run.
	Model download & placement	5 pts	Both models correctly placed and accessible in ./llms/local_llm.	Models downloaded but misplaced/misnamed.	Only one or incomplete model.	No usable models found.
	Model loading	5 pts	Model loads without error; run_llm("Hello") produces valid coherent response.	Model loads without error; run_llm("Hello") produces warnings but produces output.	Loads with minor warnings but produces output.	Fails to load or crashes.
<b>Core Functions Implementation (30%)</b>	chunk_text()	5 pts	Correct chunking & overlap; outputs usable segments.	Minor overlap/size errors but mostly works.	Inconsistent splitting or incorrect overlap.	Missing or non-functional.
	llama_embed_text()	5 pts	Returns correct embeddings as np.ndarray with proper shape.	Returns embeddings but with inconsistent shape/type.	Embeddings returned but with evident errors.	Missing or fails completely.
	Embedding loop / FAISS index	5 pts	All chunks embedded; FAISS index created and populated.	Most chunks embedded; index partially populated.	Few chunks embedded or index corrupted.	No embeddings or FAISS integration.
	search_similar_chunks()	5 pts	Retrieves top-k relevant chunks with clear similarity logic.	Retrieves chunks but sometimes irrelevant.	Results mostly irrelevant to query.	Function missing or fails.
	run_llm()	5 pts	Generates coherent, relevant responses to test prompts.	Responses mostly coherent; minor weaknesses.	Incomplete or inconsistent responses.	No meaningful response.
<b>PDF Context Integration (20%)</b>	run_llm_with_pdf_knowledge()	5 pts	Effectively integrates retrieved context in responses.	Context partially integrated; some relevance.	Weak integration; context barely affects answer.	No integration attempted.
	PDF extraction	5 pts	>500 chars extracted accurately from PDF.	Text extracted but slightly incomplete.	Only partial/limited extraction.	Extraction fails or missing.
	PDF embeddings	5 pts	Embeddings generated for all text chunks.	Most embeddings generated successfully.	Partial embeddings only.	No embeddings created.
<b>Reflection &amp; Reporting (20%)</b>	Query without context	5 pts	Shows LLM's baseline knowledge limits clearly.	Some limits visible but inconsistent.	Results unclear about LLM's limitations.	Misleadingly confident or missing comparison.
	Query with context	5 pts	Responses significantly improved with PDF context.	Partial improvement visible.	Limited or unclear improvement.	No observable improvement.
	Setup summary	5 pts	Clear, complete, step-by-step documentation.	Mostly clear; minor gaps.	Limited detail/clarity.	Missing or very unclear.
<b>Submission Quality (15%)</b>	Difficulties/errors	5 pts	Thorough documentation of errors & solutions.	Documents errors but not all solutions.	Mentions errors with no resolution.	No discussion of errors.
	Query comparison	5 pts	Strong analysis of context vs. no-context differences.	Adequate comparison; some insights.	Weak or superficial comparison.	No comparison provided.
	Insights & applications	5 pts	Description on LLM limitations & healthcare use cases.	Some relevant insights.	Surface-level or generic insights.	No meaningful reflection.
	Notebook	5 pts	Runs end-to-end with clean outputs; reproducible.	Minor execution/formatting issues.	Multiple errors, requires fixing to run.	Notebook incomplete/unusable.
	Screenshots/outputs	5 pts	Clear, labeled screenshots/outputs included.	Some outputs missing/unclear.	Outputs incomplete or poorly presented.	No outputs/screenshots.
	Report formatting	5 pts	Well-structured, concise, within 1-2 pages.	Clear but slightly long/short.	Formatting issues or unclear writing.	Report missing.

Performance levels map to:
100% = full 5 pts
75% = 3.75 pts
50% = 2.5 pts
25% = 1.25 pts