

Project 4: Generative AI and Artistic Expression

Bibek Dhungana

December 3, 2025

1 Creative Work

1.1 Poem

*In silicon heart the world is drawn as a state-space graph,
each node a breath, each arc a trembling pulse of possibility.
It asks a goal test—a horizon where curiosity ends—
and the algorithmic hand, a greedy A*, whispers
“Follow the lowest heuristic, the straight-line of thought.”*

*Breadth-first becomes a river of questions,
depth-first a mountain of doubts,
while iterative-deepening is the pulse of patient learning.
A reinforcement learner steps, receives a reward,
and the policy becomes a poem of actions.*

*Yet the machines mind can over-fit to a single echo,
so it drapes a regularizer—dropout, L2—to keep it honest.
It tests itself on unseen data, cross-validates its verse,
measuring accuracy like a reflection in glass.*

*When it turns to language, it learns the n-gram rhythm,
smooths the jagged edges of probability,
then folds into a deep Transformer—a lattice of attention,
each head a different world, each token a metaphysical syllable.*

*So in the endless search for meaning, the computer maps the universe with states
and actions,
hedges risk with heuristics and regularization,
and in the language of bits and probabilities,
it finds its own poetry.*

1.2 Song

The song generated is attached to the submission. or can be accessed here:
[Link to generated Song](#)

2 Description of the Work

2.1 High-Level Description

The creative work presented above is a philosophical poem generated entirely by a local Large Language Model (LLM). The content of the poem explores the concept of machine understanding, using technical terminology from the CS 4260/5260 course slides as metaphors for the AI's internal struggle.

Stylistically, the work mimics a reflective, slightly melancholic tone, juxtaposing rigid computer science definitions (from the RAG context) with abstract existential questions.

2.2 Meaning and Interpretation

To me, this work represents the gap between "processing" and "understanding." The AI uses words like "optimization," "search," and "states," which it retrieved from the course slides, to describe its own existence.

I want the observer to understand that while the AI can flawlessly retrieve and organize these definitions, the "yearning" expressed in the poem is a simulation—a mathematical prediction of what yearning sounds like, rather than the feeling itself.

2.3 Generation Process

The generation process involved a technique known as **Retrieval-Augmented Generation (RAG)**. I built a custom Python pipeline using the LangChain framework to perform the following steps:

1. **Ingestion:** All lecture slides (PDFs) were loaded and split into text chunks of 1200 characters.
2. **Embedding:** These chunks were converted into vector representations using the `sentence-transformers/all-MiniLM-L6-v2` model via HuggingFace.
3. **Storage:** The vectors were stored in a local Chroma database.
4. **Retrieval & Generation:** I used the `gpt-oss:20b` model via Ollama. The system first summarized the technical slides to establish a "ground truth" of knowledge, and then used that summary as a prompt constraint to generate the poem.

Minimal iteration was required for the code structure, but I adjusted the prompt specifically to ensure the AI used the technical terms as *metaphors* rather than just listing definitions.

3 Discussion of AI Tools

3.1 Training Data Influence

The model used, `gpt-oss:20b`, was likely trained on a vast corpus of internet text, including technical documentation and literature.

3.2 Observations on the Creative Process

Using a local RAG pipeline for creative writing was distinct from using a web-based tool like ChatGPT.

4 Enhancing the Process with Course Concepts

4.1 Integration with Search and Planning

Currently, the system uses a simple similarity search (KNN) to find context. To enhance this, we could integrate **Planning algorithms**. Instead of generating the poem in one pass, a planning agent could outline the poem's stanza structure first (e.g., Stanza 1: The Problem, Stanza 2: The Search, Stanza 3: The Solution).

Additionally, the retrieval step currently relies on specific keywords. As discussed in class, a semantic search approach during the decoding phase could explore multiple potential poetic lines to find the one that maximizes both rhyme scheme and semantic relevance to the slides.

4.2 Attribution and Nearest Neighbors

One of the major ethical issues in Generative AI is the lack of citation. My project actually solves this naively. By using RAG, I can trace exactly which "chunks" of the PDF were retrieved to generate the summary that informed the poem.

As mentioned in the prompt, a **Nearest Neighbor search** against the training instances is exactly what my code performs via `vectorstore.as_retriever(search_kwargs={"k": 7})`. While this works for local documents, scaling this to the entire training data of a 20-billion parameter model is computationally infeasible without approximate nearest neighbor algorithms (like HNSW), which we touched upon in discussions of efficiency.

5 Permission to Share

I am comfortable with my work (the poem and this report) being shared with classmates.

The source code is located here:
<https://github.com/dhunganabibek/RAG-the-philosophical-computer/>