

Impute to Improve

AI-Powered Attribute Imputation for Retail Catalog





Faculty Advisor: Yu Ma

DSG Team: Sunanda Parthasarathy, Kyle Bennison, Jie Luo

Problem Statement and Objective

Dick's Sporting Goods (DSG) has an online catalog of 2.4M+ SKUs, but key product attributes are only ~40% complete. This limits search and filter performance, product discoverability, and business analytics.

Our goal is to leverage a Retrieval-Augmented Generation (RAG)-based LLM pipeline to:





Improve recommendation quality and user engagement

Data







Product Catalog

Structured data including product ID, category, and known attribute values

Product Description

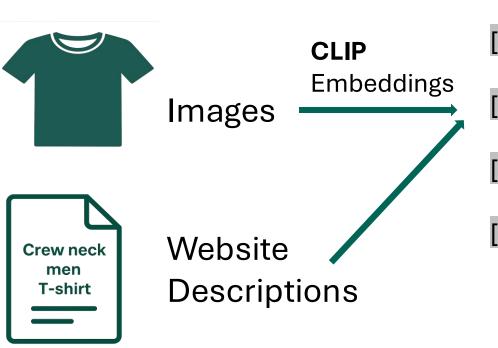
Short text containing cues like fit, material, brand, and use-case

Product Image

Visual input used to infer style-related details such as sleeve length, neckline, and cut

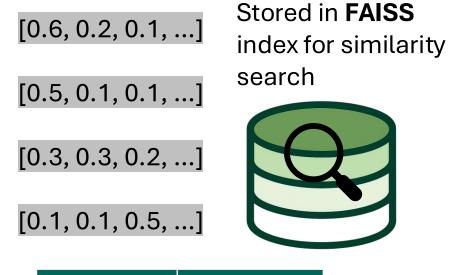
Methodology

1.Gather Inputs



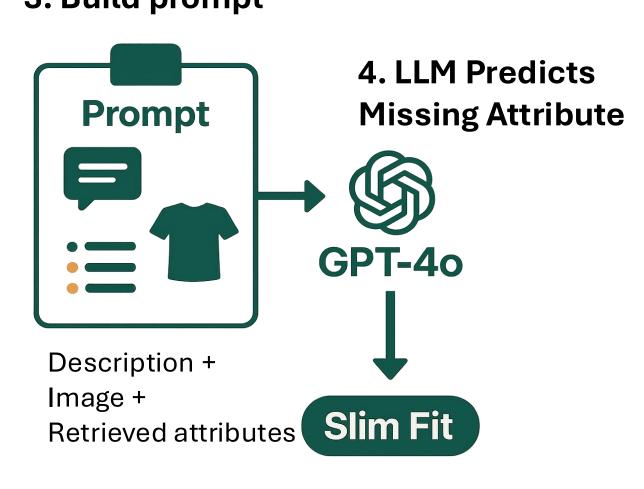


2. Retrieve neighbors

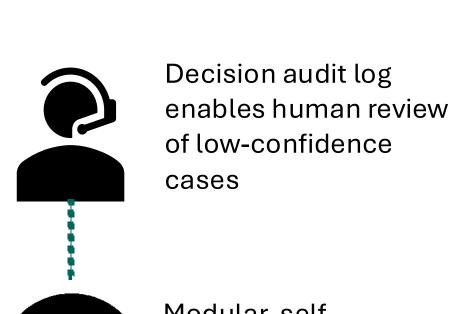


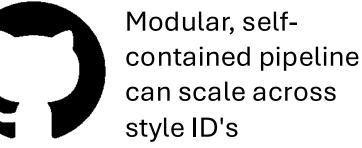
Potential	Count
Value	(# styles)
Slim	3
Regular	2
Relaxed	0

3. Build prompt



5. Ensure write-back and QA pipeline





Retrieval Pipeline

We retrieve similar styles using text and image embeddings to provide relevant context for LLM prediction. This improves prediction accuracy and reduces hallucination.



Embed product text and image using CLIP (Contrastive Language-Image Pretraining) into a shared vector space.



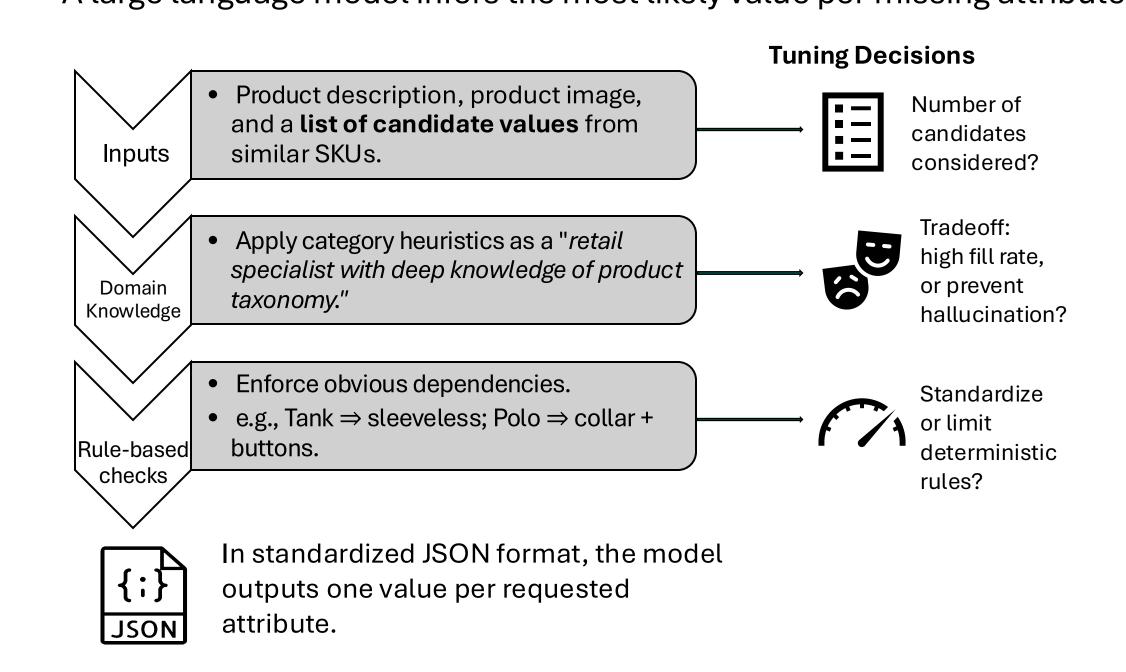
Use FAISS (Facebook Al Similarity Search) to index combined embeddings for fast, scalable similarity retrieval.



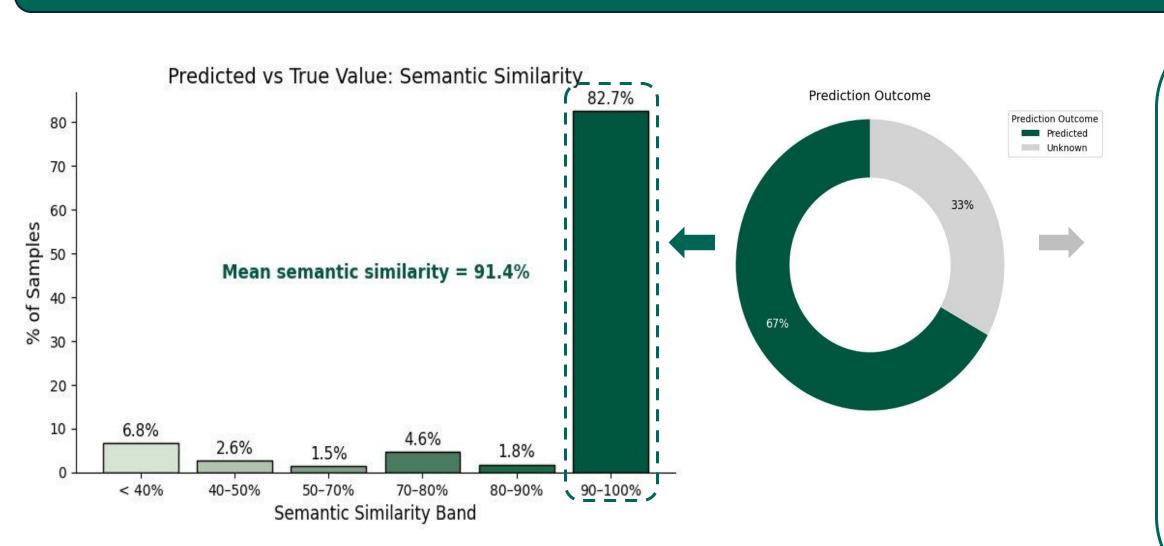
For each product with missing attributes, retrieve the top-k most similar styles with known attribute values.

LLM Inference and Guardrails

A large language model infers the most likely value per missing attribute.



Result and Impact



Why the Model Predicted "Unknown"



Ambiguous product info



Missing context from similar items

ا المنا

No clear visual/ textual cues



Attribute too subtle to detect

Estimated annual incremental revenue

\$3.0M



S derivation: 12-14M daily impressions \div 48 results x 365 \approx 91-106M. Assumption-based; to be validated via A/B. Figures illustrative.