

PROBLE

Problem

- Marketplaces usually support text search, not image search.
- Clothing is a high-value category where users want to search from a photo.
- Retrieval must distinguish not only object type, but also print, texture, and color.
- Catalogs are large and contain many categories and subcategories.

Goal

- Add image-based search for marketplace products.
- Return the most visually similar items, especially for clothing.
- Build a compact on-device demo for the course project.

Key challenge

- CLIP / SigLIP capture semantic similarity, but may miss fine-grained pattern differences.
- A real marketplace also needs category-aware retrieval and efficient search.



Visual Search for Marketplace

On-device visual retrieval for e-commerce clothing search

RESULT

Output

- Top-10 visually similar products.
- Returned as product cards.
- Designed as an Android demo

Why on-device?

- Demonstrates a compact mobile computer vision pipeline.
- Useful for prototyping and showing end-to-end functionality.
- Highlights latency and resource constraints.

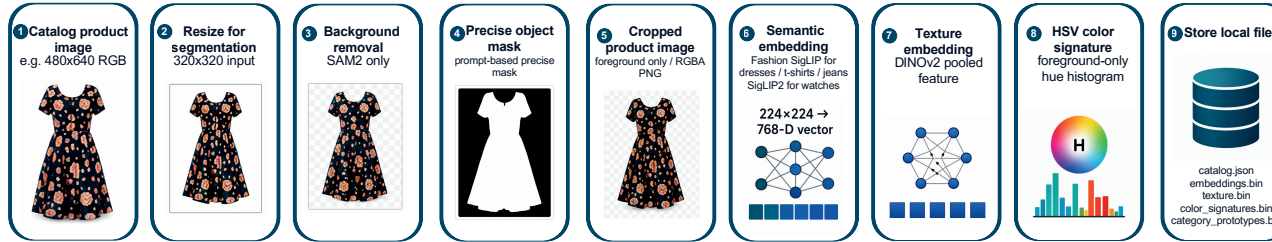
Production note

- Real marketplace deployment would likely use server-side retrieval.
- Large catalogs require scalable indexing, server-side vector retrieval, and stronger infrastructure.
- The on-device version is a research demo, not the final production architecture.



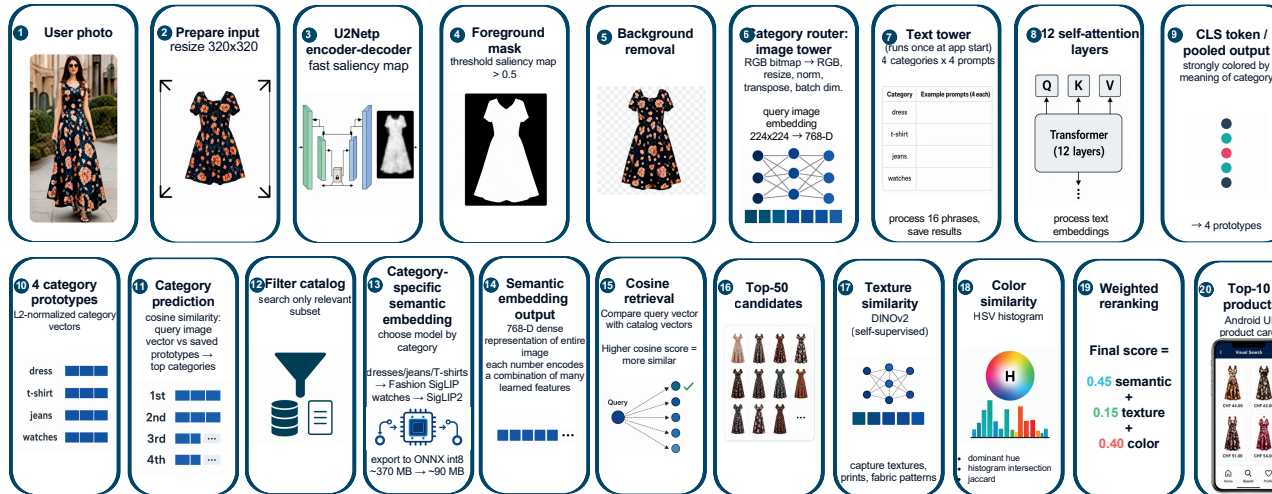
Student: Iuliia Medvednikova
 Course: Computer Vision
 Project: Visual Search for Marketplace
 Supervisor: Prof. Francis Engelmann
 Teaching Assistant: Nicolai Hermann
 USI Master's in Informatics
 May 2026
 Contact: iuliia.medvednikova@usi.ch

OFFLINE CATALOG PREPROCESSING



SAM 2 is used offline for higher-quality catalog preprocessing.

ON-DEVICE RUNTIME PIPELINE



Runtime uses precomputed local files from Section A (embeddings.bin, texture.bin, color_signatures.bin, category_prototypes.bin, catalog.json).

Why U2Netp on-device?

- small, fast, prompt-free
- good trade-off for mobile inference
- approx. 4-5 MB class model footprint

Why SAM 2 offline?

- higher-quality masks
- heavier model, better for preprocessing than mobile runtime

Category router details

- text tower runs once during app start
- processes 4x4 phrases
- stores 4 category prototypes
- image tower runs for each user photo
- category via cosine similarity to prototypes

Why DINOv2 for texture?

- self-supervised learning
- good for textures, prints, fabric patterns
- helps distinguish floral vs polka-dot vs solid

Why reranking?

- semantic embeddings find similar items in general
- texture and color recover details left by generic embeddings