

FinMatch: Multimodal Embedding and Vector Search Pipeline

Mohamed Amine Rachid Nour Zouari Karim Dami Eya Azouzi

January 26, 2026

Contents

1 Overview	2
2 Pipeline Description	2
3 Embedding Functions	2
3.1 Audio Embedding	2
3.2 Text Embedding	3
3.3 Image Embedding	3
4 Multimodal Fusion Logic	3
5 Main Script	4
6 CSV Product Embedding Preparation	4
6.1 Reading Product Data	4
6.2 Text Embedding Preparation	5
7 Vector Storage with Qdrant	5
7.1 Connection	5
7.2 Point Construction and Insertion	5
8 Conclusion	5

1 Overview

FinMatch is a **multimodal recommendation system** that combines:

- Text
- Images
- Audio

All inputs are converted into a **single semantic vector** using pretrained models. This vector can be stored in a vector database such as **Qdrant** for similarity search.

Main technologies used:

- CLIP — for text and image embeddings
- Whisper — for audio transcription
- PyTorch — for tensor operations
- Qdrant — for vector storage and retrieval

2 Pipeline Description

The system follows these steps:

1. Transcribe audio into text using Whisper
2. Generate embeddings for text using CLIP
3. Generate embeddings for images using CLIP
4. Normalize all embeddings
5. Fuse embeddings using weighted combination
6. Store vectors in Qdrant

3 Embedding Functions

All embeddings have dimension **512** and are L2-normalized.

3.1 Audio Embedding

Audio is first converted into text, then embedded using CLIP.

```
1 def audioVector(audio_path):
2     result = whisper_model.transcribe(audio_path)
3     text = result["text"]
4
5     text_inputs = processor(
6         text=[text],
7         return_tensors="pt",
8         padding=True
9     )
10
11    text_emb = model.get_text_features(**text_inputs)
12
```

```

13     # Normalize embedding
14     text_emb = text_emb / text_emb.norm(dim=1, keepdim=True)
15
16     return text_emb

```

Explanation:

- Whisper converts speech to text
- CLIP embeds the resulting text
- Normalization ensures cosine similarity works correctly

3.2 Text Embedding

```

1 def txtVector(text):
2     text_inputs = processor(
3         text=[text],
4         return_tensors="pt",
5         padding=True
6     )
7
8     text_emb = model.get_text_features(**text_inputs)
9     text_emb = text_emb / text_emb.norm(dim=1, keepdim=True)
10
11    return text_emb

```

This function directly embeds textual input into a semantic vector.

3.3 Image Embedding

```

1 def imVector(image):
2     image_inputs = processor(
3         images=image,
4         return_tensors="pt"
5     )
6
7     image_emb = model.get_image_features(**image_inputs)
8     image_emb = image_emb / image_emb.norm(dim=1, keepdim=True)
9
10    return image_emb

```

Images are encoded using CLIP's image encoder.

4 Multimodal Fusion Logic

Each modality is assigned a weight:

$$x_{\text{image}} = 0.4, \quad x_{\text{text}} = 0.4, \quad x_{\text{audio}} = 0.2$$

If a modality is missing, its weight is redistributed to the remaining ones.

```

1 q = v_image * x1 + v_text * x2 + v_audio * x3
2 q = q / q.norm(dim=1, keepdim=True)

```

Result: A single normalized vector representing all available modalities.

5 Main Script

The main script:

- Loads available inputs
- Computes embeddings
- Applies fusion logic

```
1 if __name__ == "__main__":
2     x1, x2, x3 = 0.4, 0.4, 0.2
3     v0 = torch.zeros(1, 512)
4     v1=v0.clone()
5     v2=v0.clone()
6     v3=v0.clone()
7
8     if torch.all(v1==v0):
9         x2=x2+x1/2
10        x3=x3+x1/2
11        x1=0
12
13    if torch.all(v2==v0):
14        x1=x1+x2/2
15        x3=x3+x2/2
16        x2=0
17
18    if torch.all(v3==v0):
19        x1=x1+x3/2
20        x3=0
21
22    v_image = v_text = v_audio = v0.clone()
23
24    if txt_path.exists():
25        v_text = txtVector(text)
26
27    if img_path.exists():
28        v_image = imVector(image)
29
30    if audio_path.exists():
31        v_audio = audioVector(audio_path)
32
33    q = v_image*x1 + v_text*x2 + v_audio*x3
34    q = q / q.norm(dim=1, keepdim=True)
35
36    print(q)
```

6 CSV Product Embedding Preparation

Products are read from a CSV file and embedded using text only.

6.1 Reading Product Data

```
1 def read_products(path):
2     with open(path, newline="", encoding="utf-8") as f:
3         reader = csv.reader(f)
```

```

4     next(reader)
5     return list(reader)

```

6.2 Text Embedding Preparation

```

1 def prepareVector(product):
2     text = product[2] + " " + product[5]
3     v = txtVector(text)
4     return v.squeeze(0)

```

7 Vector Storage with Qdrant

7.1 Connection

```

1 client = QdrantClient(
2     url="https://<QDRANT_URL>",
3     api_key=""
4 )

```

7.2 Point Construction and Insertion

```

1 points.append(
2     PointStruct(
3         id=int(p[0]),
4         vector=v,
5         payload={
6             "title": p[2],
7             "price": p[3],
8             "ram": p[8],
9             "os": p[16]
10        }
11    )
12 )
13
14 client.upsert(
15     collection_name="finmatch",
16     points=points
17 )

```

8 Conclusion

This system:

- Combines multiple modalities into one vector
- Enables semantic similarity search
- Scales efficiently using Qdrant

It forms a solid foundation for intelligent product recommendation systems.