# Multimodal Video-Only RAG - Detailed Development Plan



## Project Scope

Goal: Build a fully-local, video-only Retrieval-Augmented Generation (RAG) pipeline that answers natural-language queries and returns (a) a textual answer citing timestamps and (b) clipped video segments corresponding to those timestamps.

All code in, managed in virtualenv.

## Directory Layout (Monorepo)

```
repo_root/
  - data/
                         # raw + processed
      - raw_videos/
       frames/
                       # extracted .jpg
      - clips/
                       # returned snippets
      - phase1_audio/
       - phase2 visual/
       phase3_db/
       - phase4_retriever/
       phase5_generation/
       - phase6_clipper/
                         # .env, model paths, settings.yaml
    conf/
    README.md
```

Each sub-package has its own `\*\* entry point\*\* so phases can be run as standalone modules`.

## Phase 1 – Audio Processing & Embedding

## 1-A | Audio Extraction & Transcription

Item	Detail	
Tool	OpenAl Whisper (CPU, medium model) via openai-whisper	
Script	<pre>src/phase1_audio/extract_transcribe.py</pre>	
Algorithm	<pre>1.ffmpeg -i video.mp4 -vn -acodec pcm_s16le temp.wav</pre>	

2. Run Whisper with `to obtain JSON with per-word timings. | | \*\*Output\*\* | data/transcripts/{video\_id}.json- list of words withstart,end,word`.|

### 1-B | 10-Second Segmentation & Normalisation

- segment\_transcript.py converts word list → 10 s buckets.
- Normalise text: lower-case, strip punctuation, collapse whitespace. **Silent buckets** (≤ 2 words) retain empty string but keep times.

```
{
  "video_id": "demo.mp4",
  "segments": [
     {"start": 0.0,      "end": 10.0,      "text": "introduction of the device ..."},
     {"start": 10.0,      "end": 20.0,      "text": ""},
     ...
  ]
}
```

### 1-C | Embedding

- embed\_text.py loads CLIP ViT-B/32 text encoder (open\_clip).
- Batched forward pass (≤ 32 segments/batch, CPU).
- Persist embeddings → phase3\_db via gRPC interface (see Phase 3 API).

#### **Deliverables Phase 1**

Deliverable	Acceptance Criteria
extract_trans cribe.py	JSON transcript matches Whisper CLI output; unit test verifies ≥ 95 % coverage for timing fields.
<pre>segment_trans cript.py</pre>	Exactly [duration/10] segments with correct start/end; silent buckets kept.
embed_text.py	Embeds N segments in $\leq$ ( N / 8 ) s on dev laptop (benchmark).

## Phase 2 – Visual Frame Extraction & Embedding

## 2-A | Frame Sampling (10 s Grid)

- ffmpeg -ss {t} -i video.mp4 -frames:v 1 frames/{id}\_{t}.jpg for t = 0,10,20,...
- Metadata: start = t, end = t + 10.

### 2-B | Frame Embedding

- Use same CLIP model image encoder.
- Images resized to 224×224, center crop.
- Batch size = 64.
- Store embedding + metadata.

#### **Deliverables Phase 2**

Deliverable	Criterion
sample_frames.py	All frames exist; naming matches regex videoid_\d+.jpg.
embed_frames.py	Embedding file or direct DB push validated by checksum.

## Phase 3 – Vector Store Service (ChromaDB)

### 3-A | Containerised DB

Docker service in docker-compose.yml – Chroma with
 CHROMA\_DB\_IMPL=duckdb+parquet, volume-mounted under data/chroma/.

### 3-B | Schema & API

One **collection** video\_segments with schema:

Expose a thin gRPC layer (src/phase3\_db/server.py) so Phase 1 & 2 import DbClient without raw Chroma dependency.

## 3-C | Batch Ingestion

- DbClient.add\_batch(vectors, metadatas) with batch\_size=20.
- Unit test: ingest 100 dummy vectors, assert count==100.

**Deliverables Phase 3**: Docker compose, DbClient, ingestion tests.

#### Phase 4 – Retrieval Service

### 4-A | Query Embedding

• embed\_query.py → CLIP text encoder.

### 4-B | Search Endpoint

Retriever.search(query: str, k:int=10) -> List[Document]

- Executes single cosine-similarity search in Chroma.
- Returns rank-ordered list; each Document holds text (for audio) or "<IMAGE\_FRAME>" placeholder plus metadata.

#### **Deliverables Phase 4**

• src/phase4\_retriever/retriever.py with pydantic models and coverage tests.

### Phase 5 – LLM Generation Micro-service

### 5-A | Prompt Template (LangChain)

### 5-B | API Wrapper

 FastAPI endpoint /ask → returns JSON: { "answer": str, "sources": List[metadata] } (sources are the same objects given to the LLM).

**Deliverables Phase 5**: prompt files, qa\_service.py, e2e test with canned video and expected timestamp in answer.

## Phase 6 - Clip Builder

## 6-A | Timestamp Parser

- Regexes for HH:MM:SS, MM:SS, SS.S → seconds.
- Merge intervals with ≤ 2 s gap.

### 6-B | Clip Extraction

 Use moviepy.VideoFileClip(video).subclip(start,end).write\_videofile(outfile , codec="libx264"). • Runs asynchronously (asyncio) so multiple clips export in parallel.

### 6-C | Clip Registry

• SQLite table clips(id, video\_id, start, end, path, created\_at) for re-use.

Deliverables Phase 6: clipper.py, integration test - feed answer JSON, verify clips exist and duration matches.



# Milestone Timeline & Team Allocation

W ee

k	Phase / Component	Owner	Key Review
1	Phase 1-A/B	Backend 1	Transcript JSON contract
2	Phase 1-C + Phase 2-A	Backend 1 / ML 1	Embedding fidelity test
3	Phase 2-B + Phase 3	ML 1 / DevOps 1	DB ingestion benchmark
4	Phase 4	Backend 2	Retrieval precision@10
5	Phase 5	ML 2	Answer quality & token costs
6	Phase 6	Backend 2	Full e2e demo



## Readiness Checklist per Phase

## Future Enhancements (Post-MVP)

1. **UI**: React front-end with embedded <video> + timeline markers.

© 2025 KR x OP Dev Team