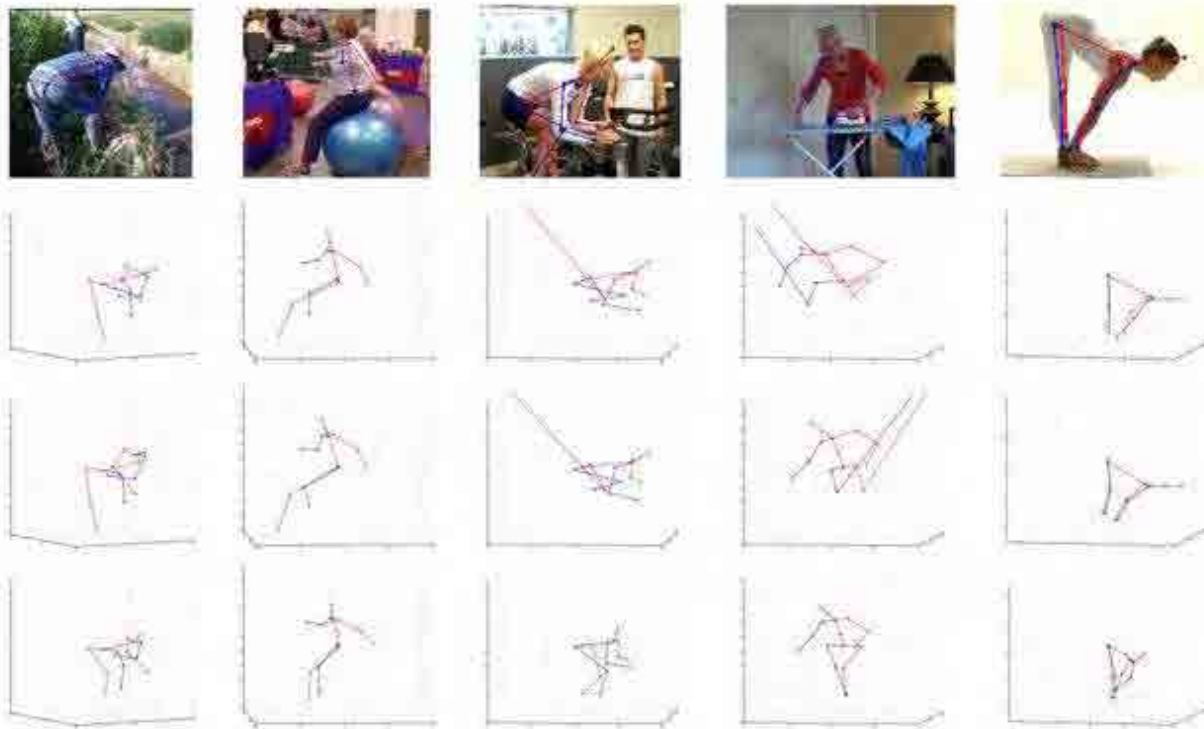# Pose and spatial movement as input for dynamic content search & generation

# Pose and spatial movement as input for dynamic content search & generation

**Main purpose?** To explore the potential of posture correctness analysis and multimodal feedback delivery for different applications (ergonomics, yoga, others).



**Tasks**
1) Task scope definition
2) Dataset search and evaluation.
3) State of the art on building postural models and postural analysis.
4) Setting up an environment for posture classification from images.
5) Model concept proposal for posture analysis, based on angles. Built from reference datasets and literature. Limited scope.
6) First beta prototype set up.

# Pose and spatial movement as input for dynamic content search & generation

**CLIP as classifier**

**DEEPER ANALYSIS OF CLIP AS CLASSIFIER:**

Yoga-82 dataset

- **Low zero-shot performance**

- **Significant improvement after fine tuning**
    - Metrics on **82 classes**:

| | Precision | Recall | F1-score | MCC | Support |
|---|---|---|---|---|---|
| **Weighted avg** | 0,861 | 0,859 | 0,857 | 0.855 | 3826 |

**BUT STILL MARGIN TO IMPROVE !**

**CLIP as classifier**

- **How can we improve the performance ?**
  - Balancing the classes in the dataset (μ= 186, σ=105)
  - Boosting classes that CLIP encounters issues with. (Already identified thanks to the representation based on hierarchical order groups. Confusion matrices)
    - E.g.:

**Makara Adho Mukha Svanasana**          **Chaturanga Dandasana**



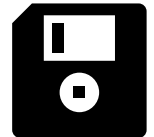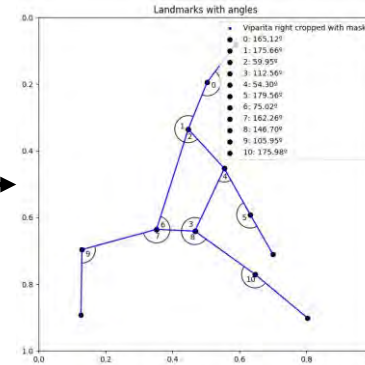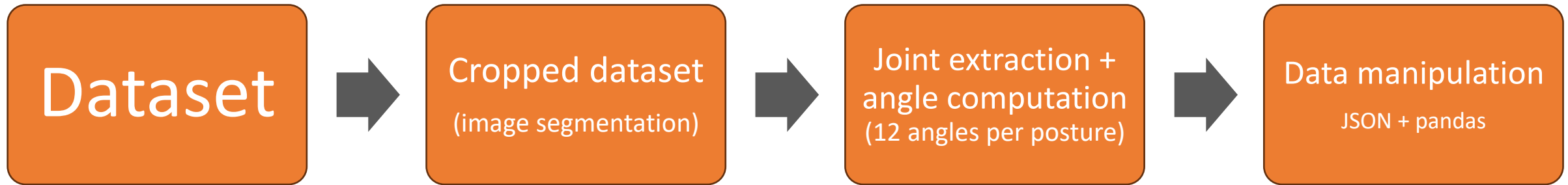F1 score: 0,5 | **36 train imgs**.          F1 score: 0,796 | **165 train imgs**.

- Combining with other types of images (infrared, joints, etc.)
- Trying to fine-tune various models of visual encoders

# Pose and spatial movement as input for dynamic content search & generation

**Mediapipe for pose evaluation**

## Angles extraction pipeline applied to all images of each class.

# Pose and spatial movement as input for dynamic content search & generation

**Posture evaluator model**

## Building the posture evaluator model.

- **Various attempts:**
  - Rules engines
  - KNN
  - XGBoost

  Trained and tested using synthetic data generated from the angle's extraction pipeline.
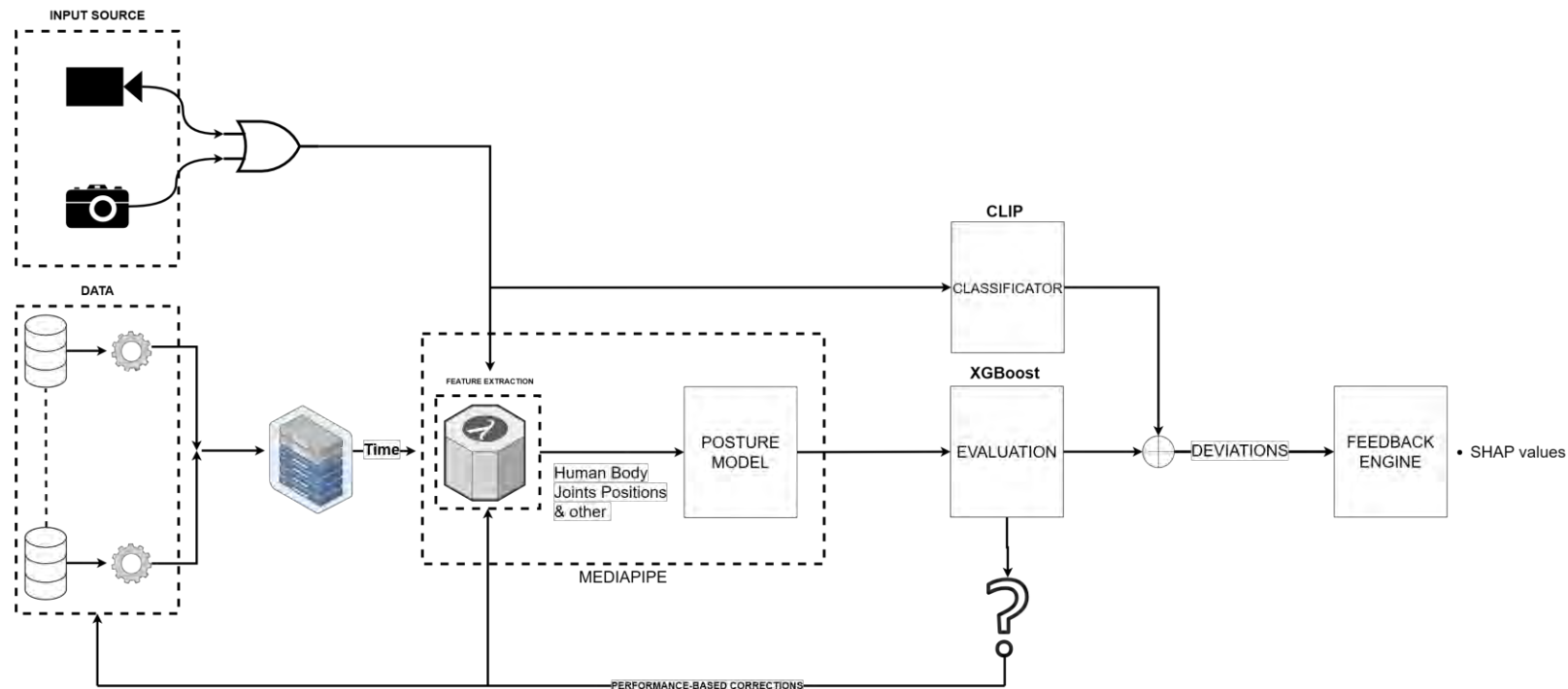
Still necessary to evaluate it on real data.

- **Best results: XGBoost. 1 model by posture.**

| | Precision | Recall | F1-score | Support |
|---|---|---|---|---|
| **Average** | 0,981 | 0,982 | 0,982 | 1000 |

**Results**

- **Results:**
  - <u>Journal article:</u>
    - **Exploring the Use of Contrastive Language-Image Pre-Training for Human Posture Classification: Insights from Yoga Pose Analysis**
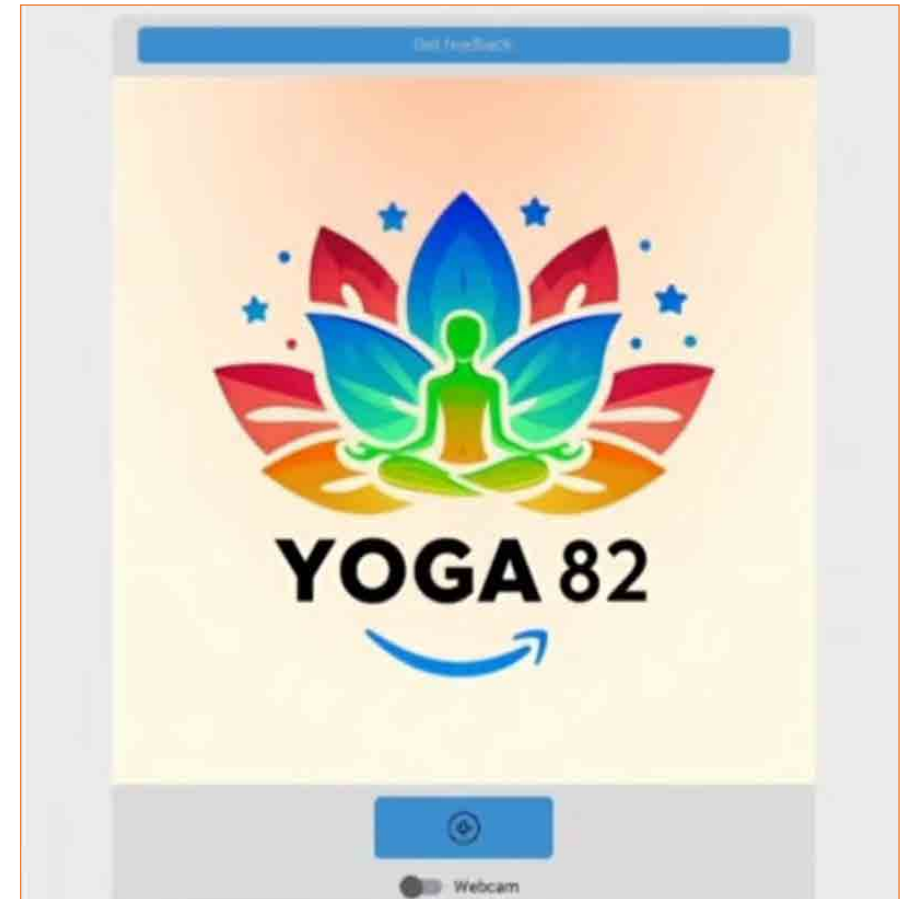  - <u>Yoga-82 app</u>

# Pose and spatial movement as input for dynamic content search & generation
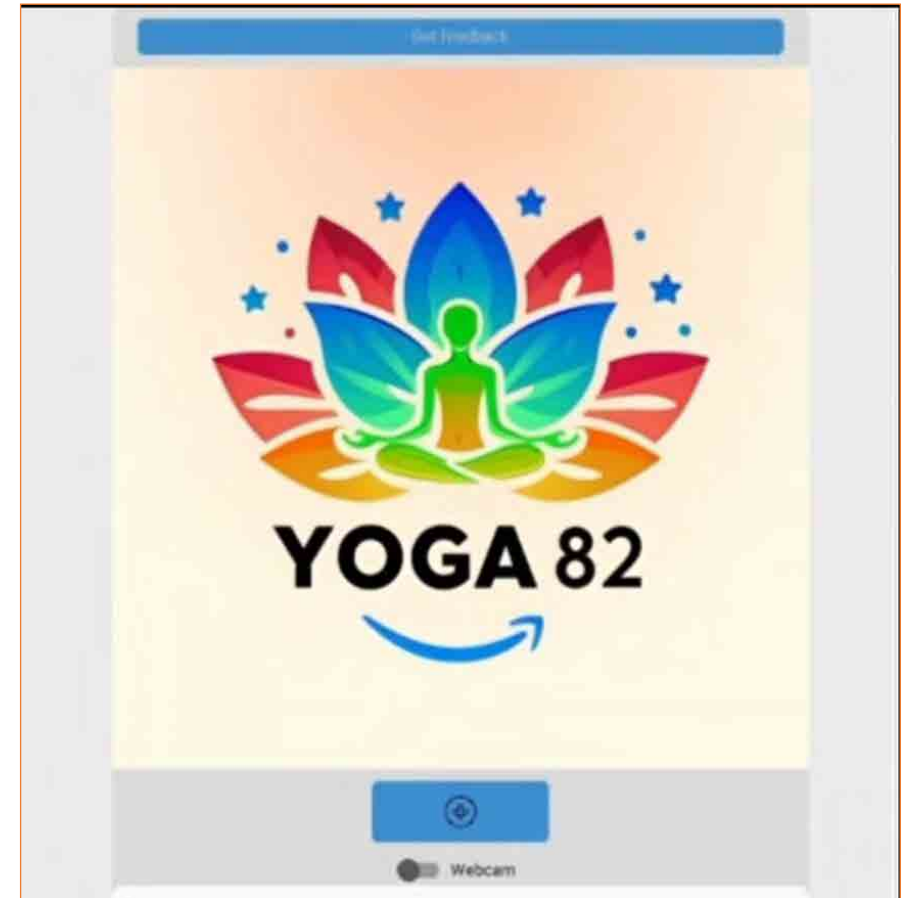
**Results**

## Real time

### Virabhadrasana II posture

# Pose and spatial movement as input for dynamic content search & generation
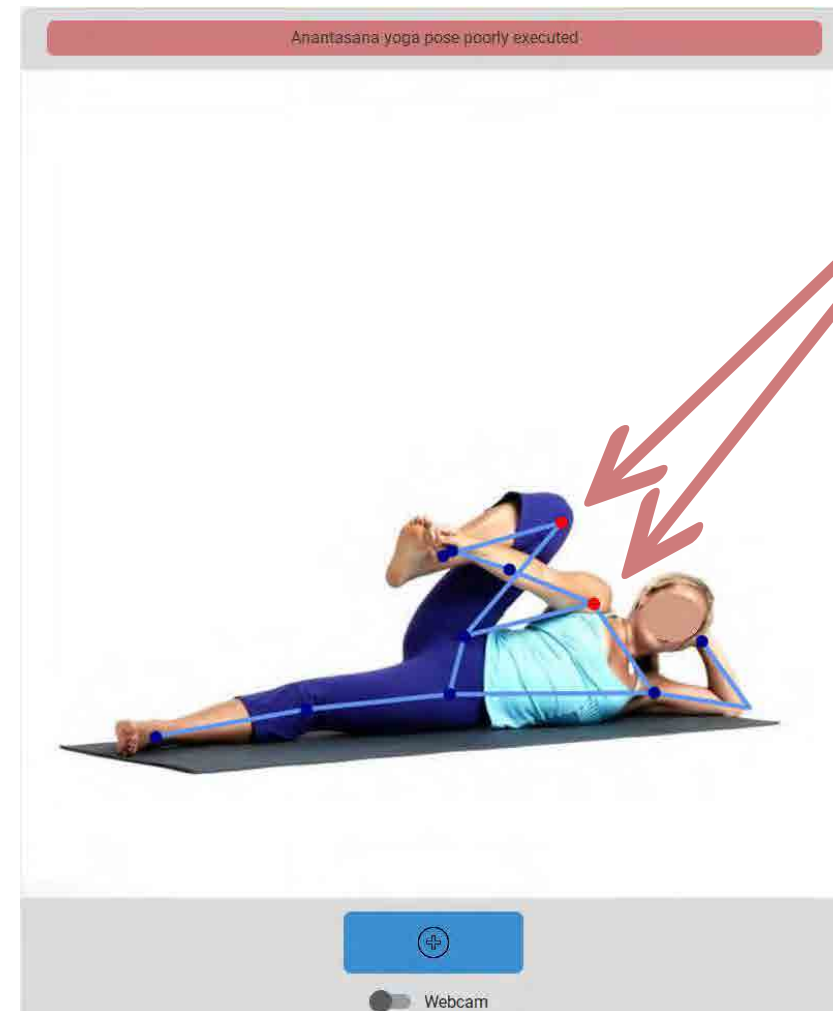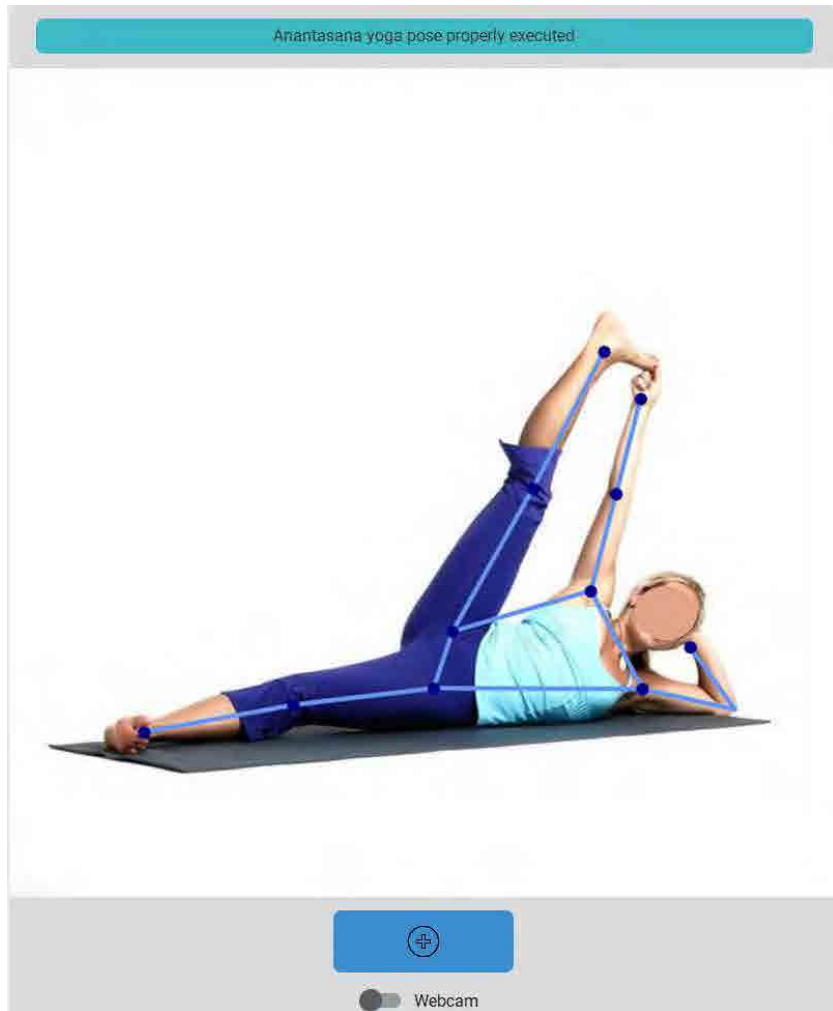
**Results**

## Real time

**Marjaryasana posture**

# Image based feedback



Anantasana yoga pose properly executed

Anantasana yoga pose poorly executed

SHAP VALUES

Webcam

Webcam

# Zero-shot sonorizing of video sequences

# Entangling AI-audio synthesis models and multimodal representations

Signal Processing Applications Group
GAPS – IPTC - UPM

How should this scene sound?

# Should this sound similar?

# The Question
**And a shared goal**

What is **a suitable audio** for a given image or video sequence?

How do we **search** or **create** a matching audio?

How should we **evaluate** if this match is *coherent*?

# The approaches
**Two ways to address the task + a novel way to evaluate**



CLIP
embeddings

...

Image & Text
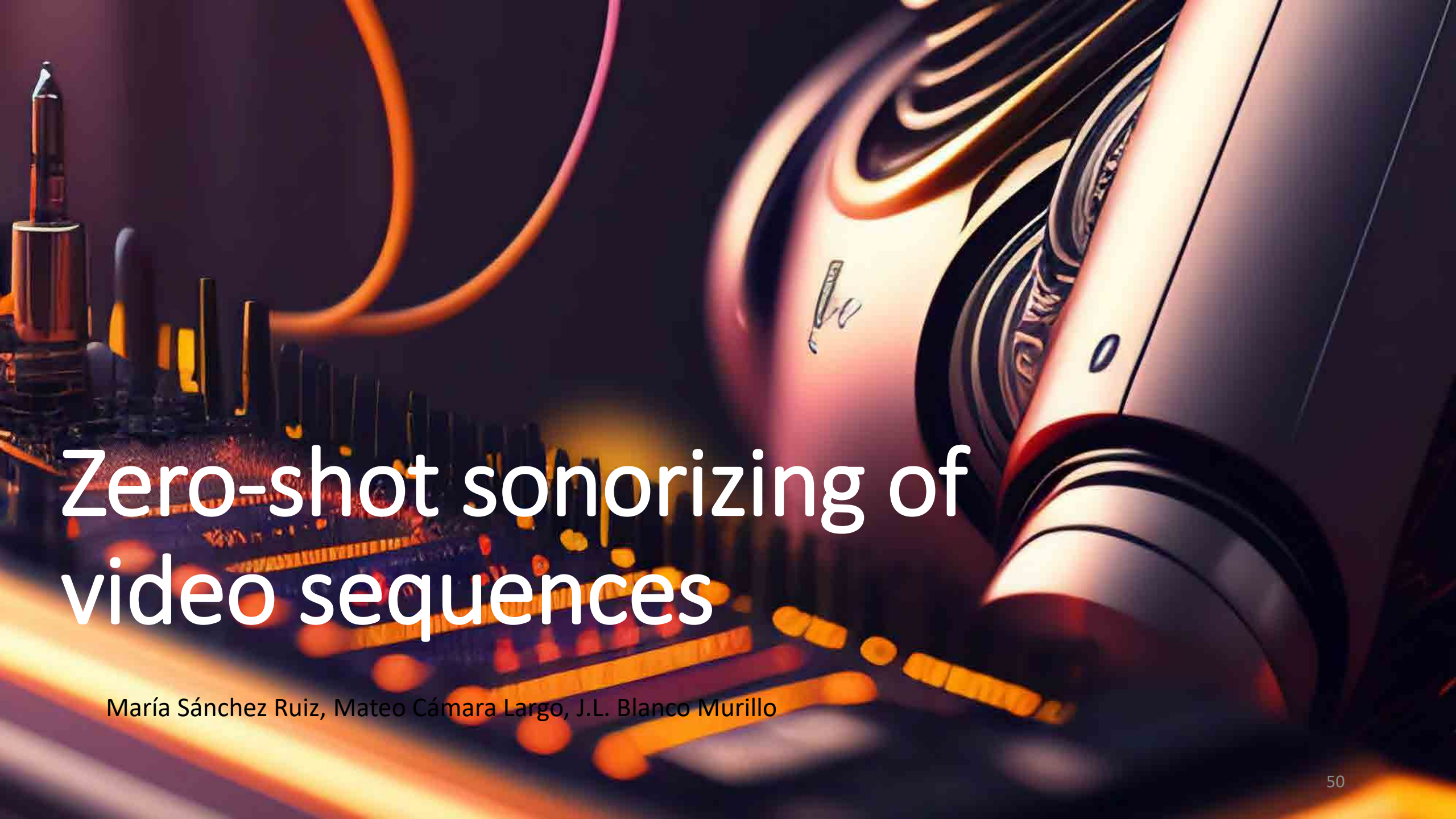encoders

encode

Multimodal representations

on Pretrained Models + Gen AI

**Image** & **Text**
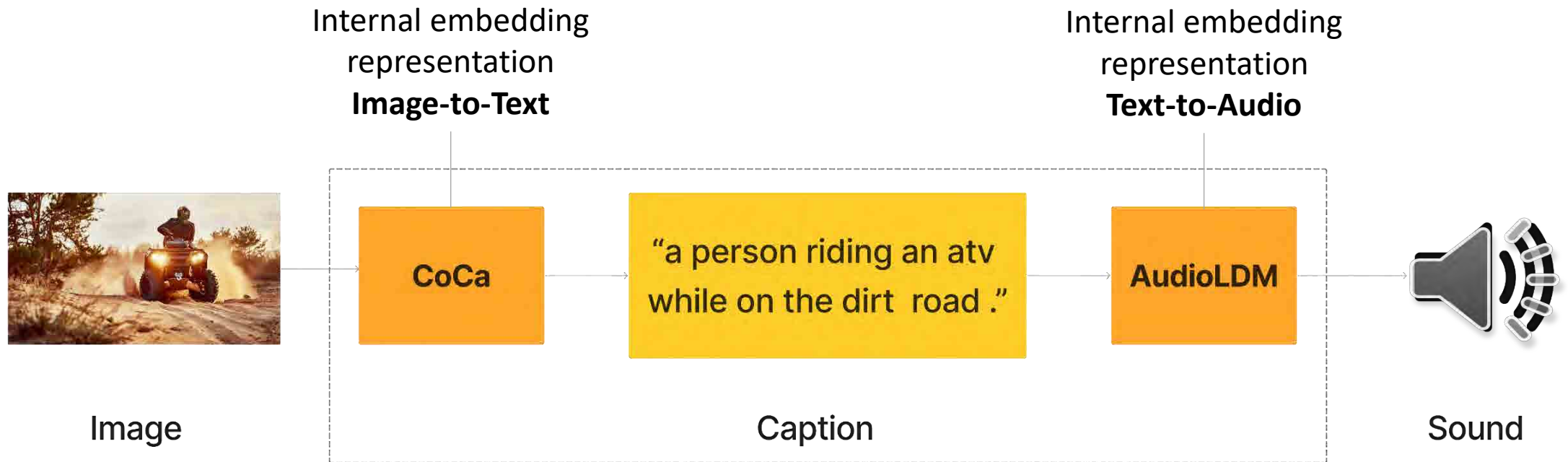
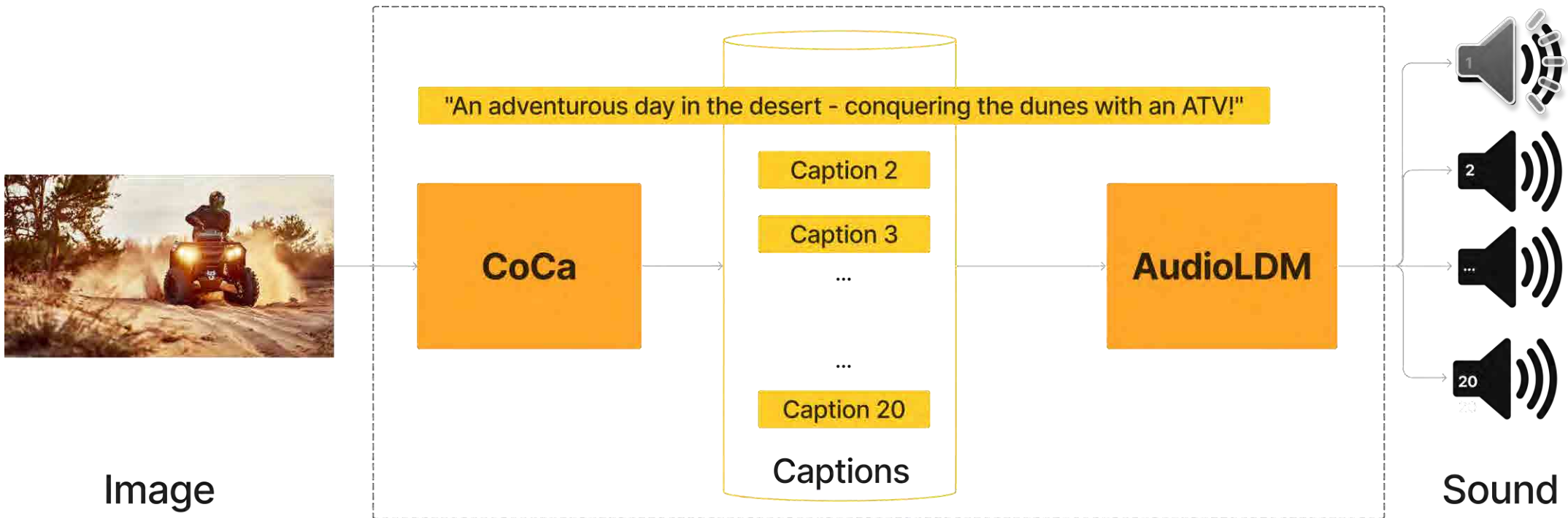**Text** & **Audio**

**Image** & **Audio**

# Zero-shot sonorizing of video sequences

María Sánchez Ruiz, Mateo Cámara Largo, J.L. Blanco Murillo

50

# Text-guided sonirization



Internal embedding representation **Image-to-Text**

Internal embedding representation **Text-to-Audio**

CoCa

"a person riding an atv while on the dirt road ."

AudioLDM

Image

Caption

Sound

# Multi-captioning

# Multi-captioning



Image

Captions

Sound

Caption 1

"Roaring through the sandy wilderness on an ATV, embracing the desert's beauty."

Caption 3

...

...

Caption 20

CoCa

AudioLDM

# Multi-captioning



Image

Captions

"In the heart of the desert, the versatile ATV becomes the ultimate desert explorer."
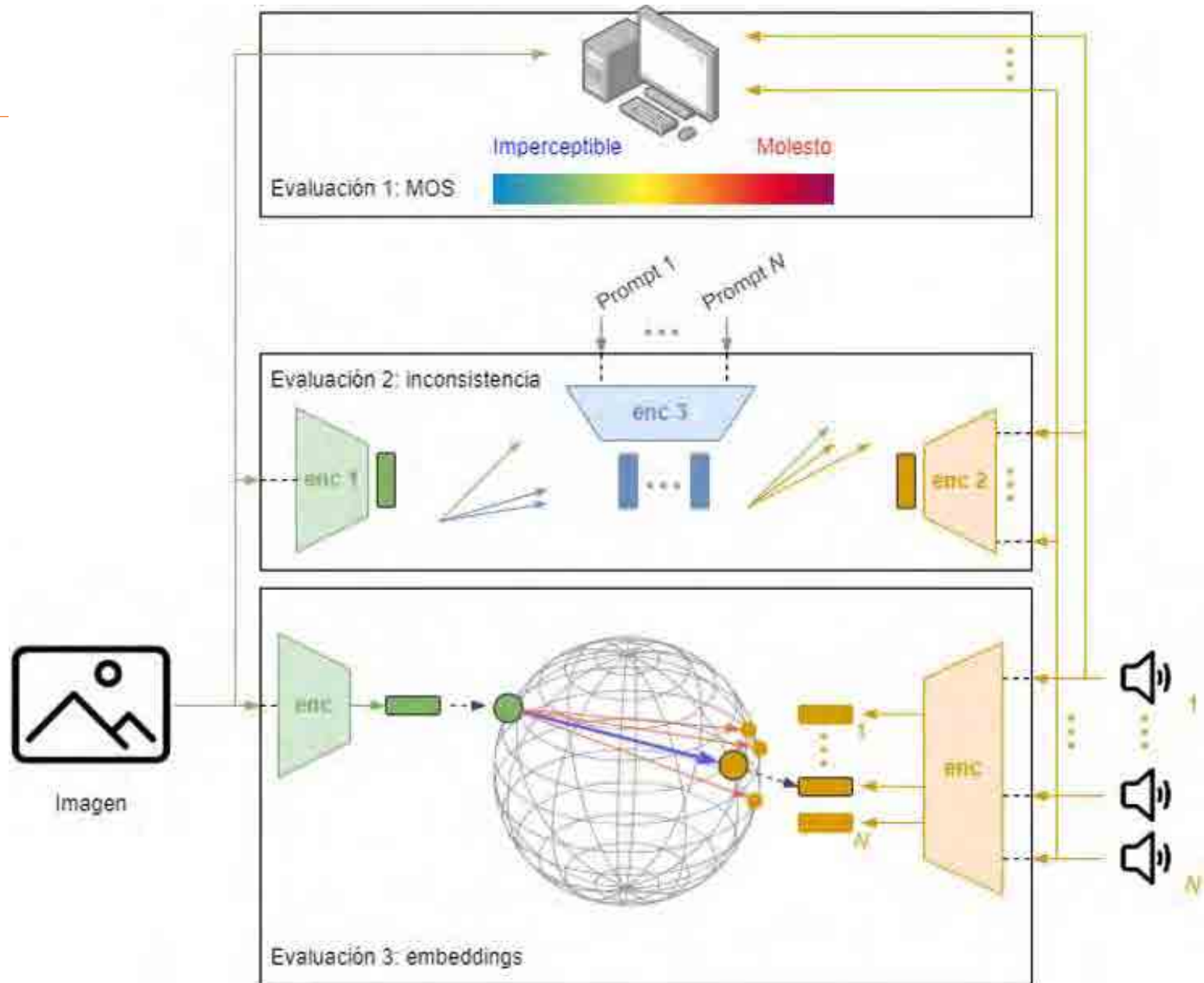
Caption 1
Caption 2
Caption 3
...
...

CoCa

AudioLDM

Sound

54

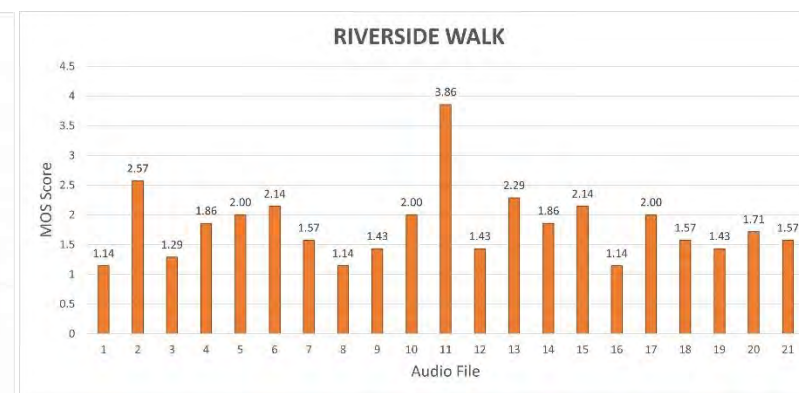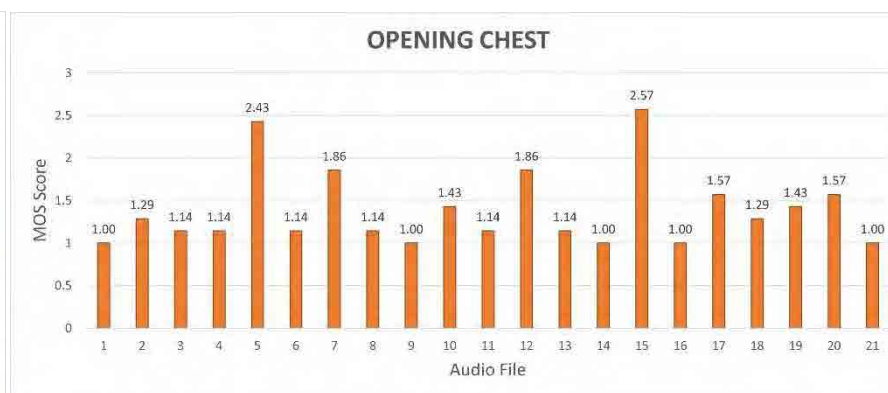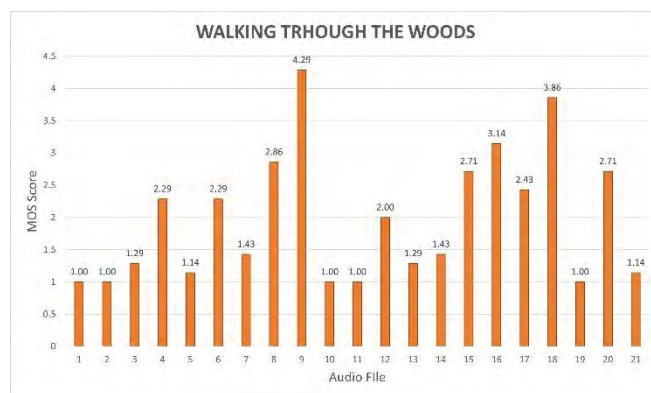# Evaluation

1. **Subjective** user experience.

2. Embeddings **consistency**.

3. Embeddings **projection**.



Evaluación 1: MOS — Imperceptible / Molesto

Evaluación 2: inconsistencia — Prompt 1 ... Prompt N — enc 3, enc 1, enc 2

Evaluación 3: embeddings — Imagen, enc

# Evaluation 1: MOS

# Evaluation 2: Inconsistency
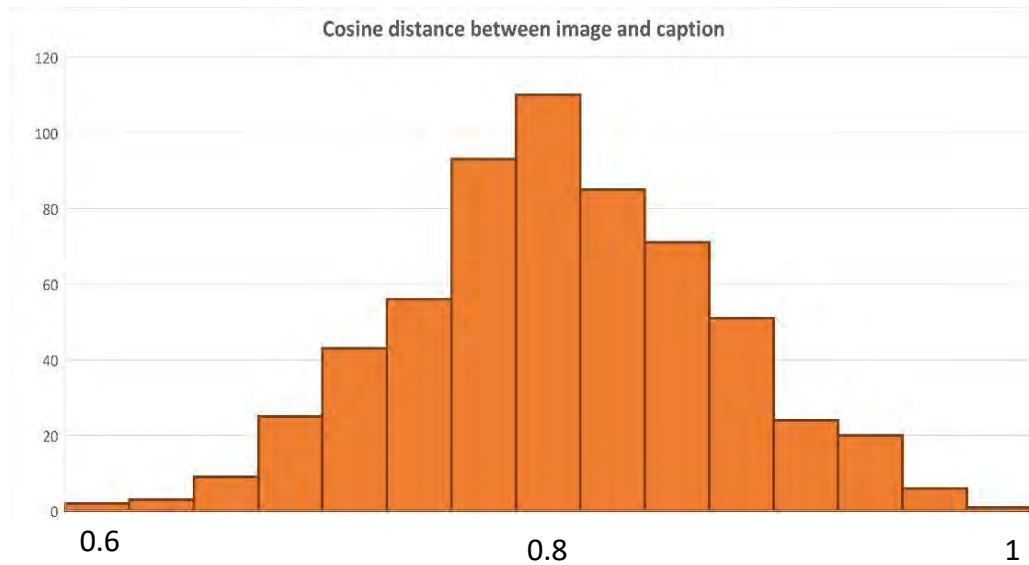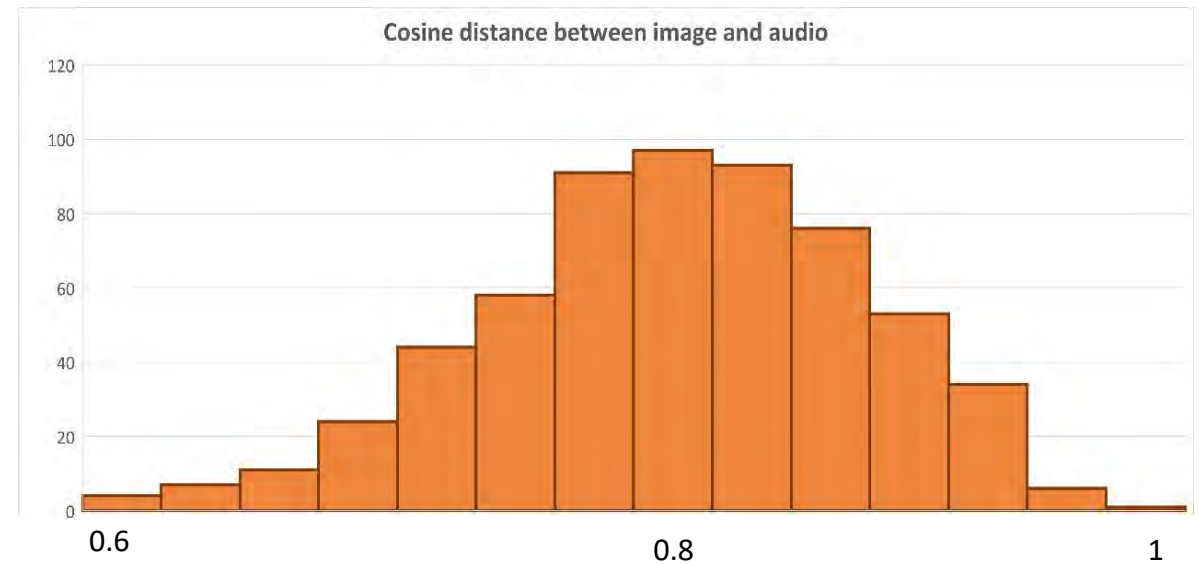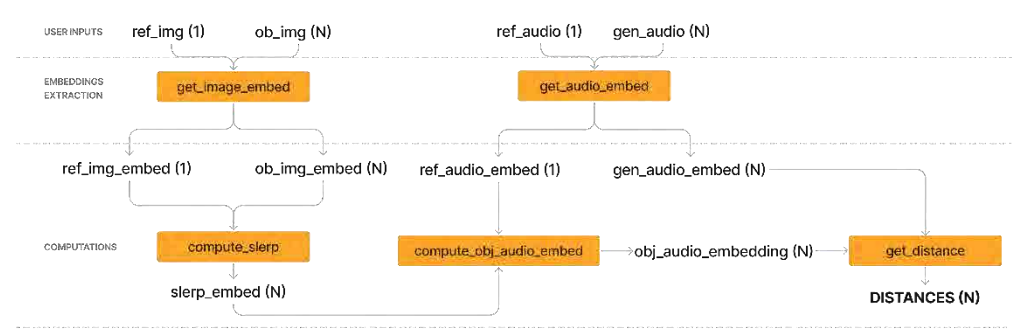
**Image-Caption Cosine Distance**

**Image-Audio Cosine Distance**

# Evaluation 3: Embedding distance

**Embedding Distance Scheme**

**Summary of Validation Test Results**



| Test | Audios | Frames | Mean | STD |
|---|---|---|---|---|
| 1 | ref_audio, gen_audio | ref_img, ob_img | 0 | 0 |
| 2 | ref_audio, gen_audio | ref_img, ob_img | 0 | 0 |
| 3 | ref_audio, gen_audio | ref_img, ob_img | 0 | 0 |
| 4 | ref_audio, gen_audio | ref_img, ob_img | 24.9 | 2.5 |
| 5 | ref_audio, gen_audio | ref_img, ob_img | 24.7 | 2.5 |
| 6 | ref_audio, gen_audio | ref_img, ob_img | 12.5 | 3.2 |
| 7 | ref_audio, gen_audio | ref_img, ob_img | 12.5 | 3.2 |

1. **Valid sonorization approach & evaluation procedure**

2. **Consistency** of the metrics with user subjective assessment.

3. Embeddings **consistency** metric is robust.

Entangling AI-audio synthesis models and multimodal representations

# CoCa-AudioLDM integrated model.



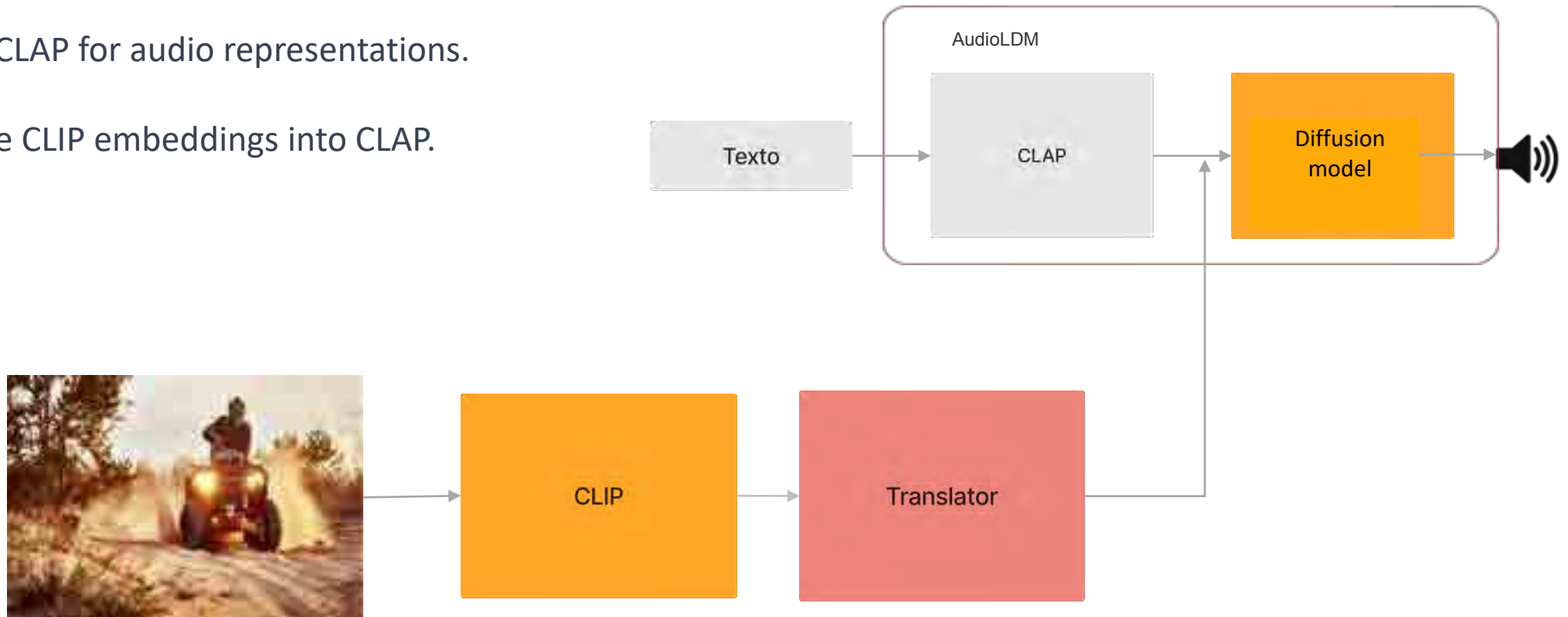An image     CoCa     "a person riding an atv while on the dirt road."     A caption     AudioLDM     A sound
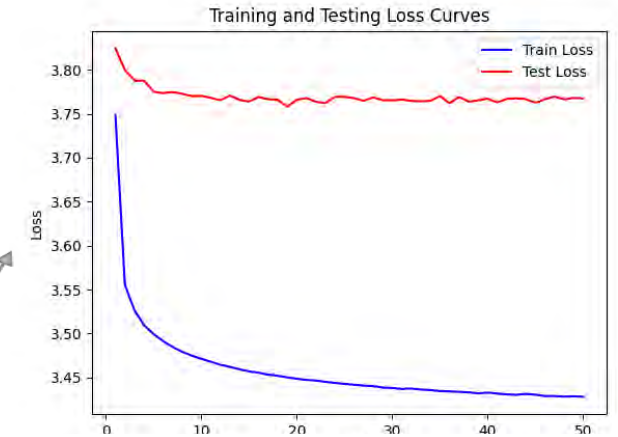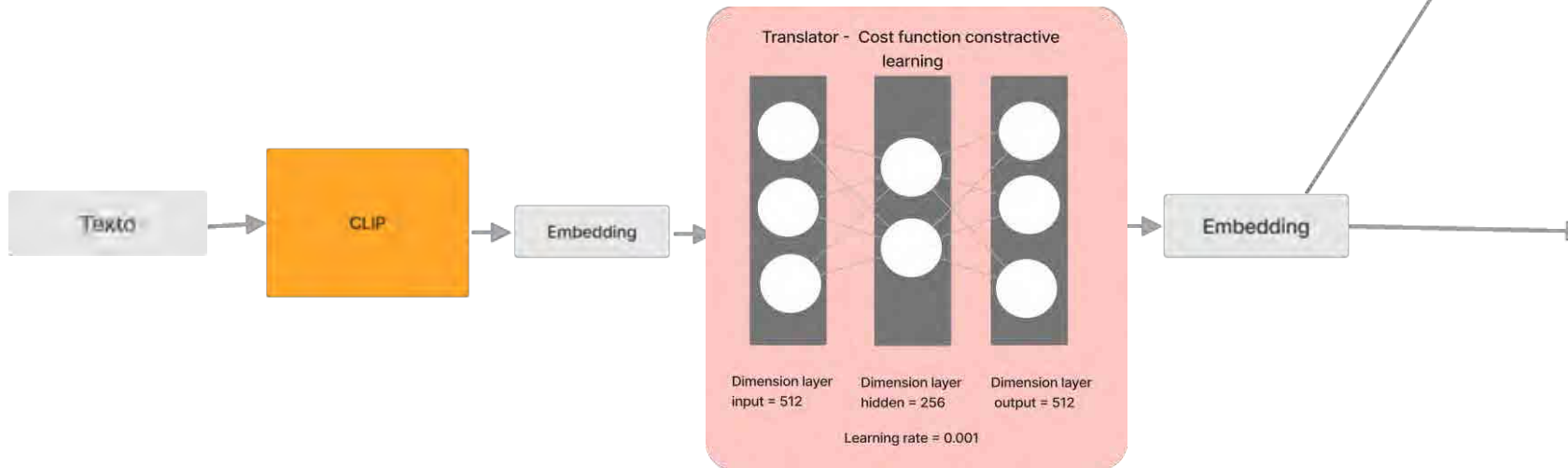
# CLIP-T-AudioLDM

- Eliminate the need for CoCa.

- Utilizes CLAP for audio representations.
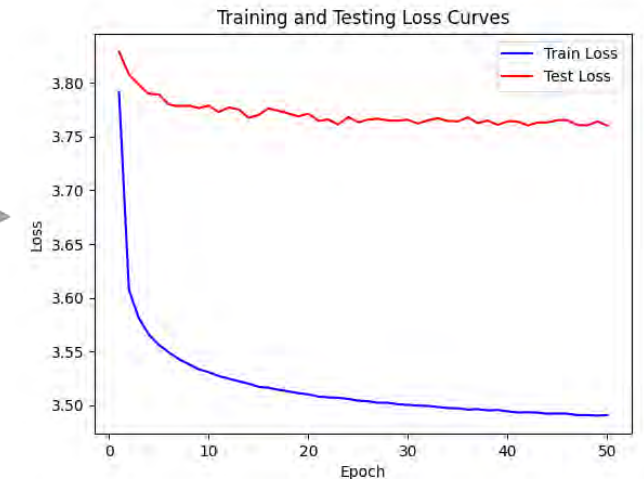
- Translate CLIP embeddings into CLAP.

# CLIP-T-AudioLDM

**18 models implemented:**

- 2 different hidden layer dimensions: 256 and 512.

- 3 different cost functions: MSE, CD and CL.

- 3 training databases: Audiocaps (57K), WIT(10K), Conceptual Captions(10K).



Model trained with all databases.

# Evaluation of the two models

## CoCa-AudioLDM

| Metric | Gigs | Horses | Kids | Piano | Train | Average | Std desviation |
|---|---|---|---|---|---|---|---|
| FAD | 17 | 12,43 | 10,88 | 4,78 | 15,56 | 12,03 | 4,65 |
| DistanceMetric-Pipeline | 23,24 | 24,85 | 24,18 | 24,82 | 28,5 | 25,12 | 2,00 |

FAD and Distance Metric-Pipeline metrics for the CoCa-AudioLDM integrated model.

## CLIP-T-AudioLDM

| Model | DROP | Learning rate | Gigs | Horses | Kids | Piano | Train | Average | Std desviation |
|---|---|---|---|---|---|---|---|---|---|
| CL512 | 0 | 0,001 | 14,12 | 23,85 | **16,25** | **10,5** | 25,8 | **18,11** | **6,51** |
| | 0,5 | | 12,59 | 22,69 | 17,67 | 28,48 | **24,9** | 21,27 | 6,23 |
| | | 0,002 | **11,43** | **22,34** | 17,51 | 30,18 | 27,13 | 21,72 | 7,5 |

FAD metrics for the CLIP-T-AudioLDM integrated model.

| Model | DROP | Learning Rate | Gigs | Horses | Kids | Piano | Train | Average | Std desviation |
|---|---|---|---|---|---|---|---|---|---|
| CL512 | 0 | 0,001 | 14,12 | 23,85 | **16,25** | **10,5** | 25,8 | **18,11** | **6,51** |
| | 0,5 | | 12,59 | 22,69 | 17,67 | 28,48 | **24,9** | 21,27 | 6,23 |
| | | 0,002 | **11,43** | **22,34** | 17,51 | 30,18 | 27,13 | 21,72 | 7,5 |

DistanceMetric-Pipeline metrics for the CLIP-T-AudioLDM integrated model.

# Demo

## Generator of Audio from Images

Welcome to our platform where two powerful models collaborate to generate audio from images. Let's explore the capabilities of these models:

### Contractive Captioner (CoCa)

CoCa is designed to describe the content of an image. It employs an image encoder and a text decoder to obtain unimodal text representations. These representations are then used to create multimodal image and text representations. CoCa captures both global and regional characteristics of images and texts, making it versatile in various tasks such as visual recognition, image caption generation, and more.

localhost:3000
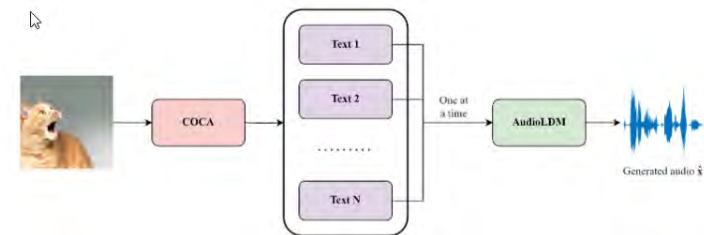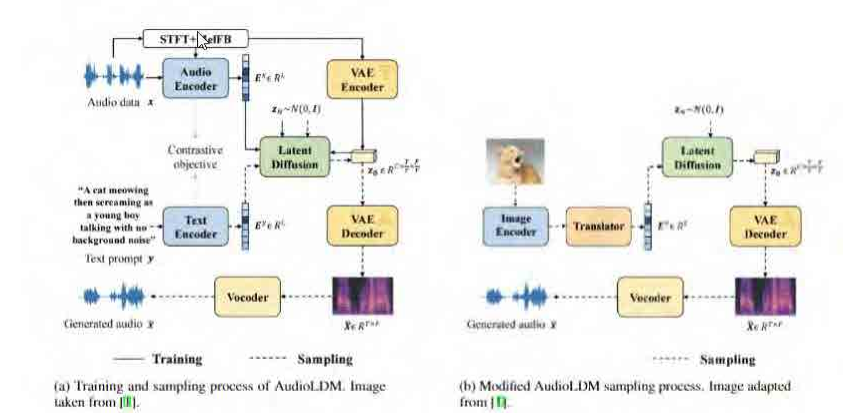
# Generator of audio

Submit

[Insert Image Here]

# Results: submitted & under preparation

International Conference paper:
**Del Visual Al Auditivo: Sonorización De Escenas Guiada Por Imagen**

arxiv working draft:
**Image-conditioned audio generation and evaluation using deep learning models**

# Results: and more than just that…

**María Sánchez Ruiz**
Masters Degree: **DTU+MUIT ETSIT**

12/12

**Laura Fernández Galindo**
Masters Degree: **MUIRST-ETSIT**

10/10

# Thanks

Amazon Team. Particularly to our coworkers:

**Guilia Comini**                    **Adam Gabrys**