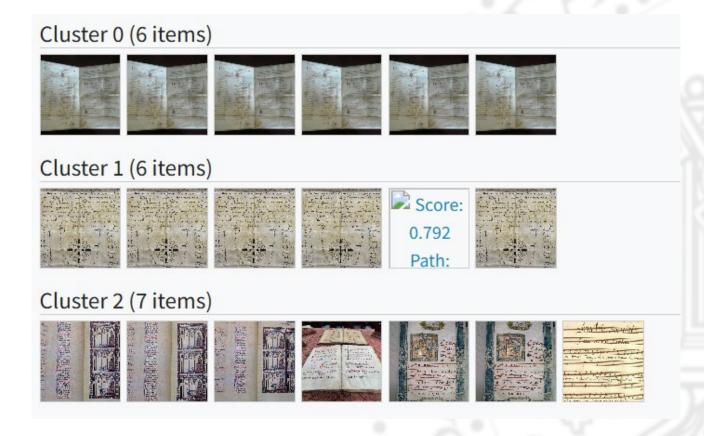




GURSO

"¿QUÉ TIENE LA INTELIGENCIA ARTIFICIAL PARA MI ARCHIVO? DE LA TEORÍA A LA PRÁCTICA"

Organización no supervisada de colección fotográficas





Aprendizaje no supervisado

No tenemos clases predeterminadas

- Aprendizaje por observación vs aprendizaje con ejemplos (supervisado)
- Vamos a explorar los datos para buscar en ellos estructuras intrínsecas

Aplicaciones del aprendizaje no supervisado

- Segmentación de clientes: Identificar diferentes grupos de clientes con comportamientos similares.
- Recomendaciones personalizadas: Sugerir productos o servicios basados en las similitudes de los clientes con otros.
- Detección de anomalías: Identificar datos inusuales que no se ajustan a patrones esperados.
- Análisis de textos: Agrupar documentos similares y descubrir temas subyacentes.

Introducción a la Visión Artificial con Modelos de Lenguaje Multimodales

Aprendizaje supervisado vs No supervisado

Supervisado:

Datos: (x, y)

x es un dato, y es una etiqueta

Objetivo:

Aprender una función que mapea el dato x a la etiqueta y

Ejemplos:

Clasificación, regresión, detección de objetos, segmentación semántica, descripción de imágenes, descripción de imágenes, etc.

¿Cuánto cuesta etiquetar 1M de imágenes?

(Base de datos pequeña a mediana)

1.000.000 (imágenes)

× (10 segundos/imagen) (anotación rápida)

 \times (1/3600 horas/segundo)

× (15€ / hora) (salario bajo)

= 41.667 €

No supervisado:

Datos: x

x es un dato, no hay etiquetas

Objetivo:

Aprender algunas estructuras ocultas subyacentes de los datos.

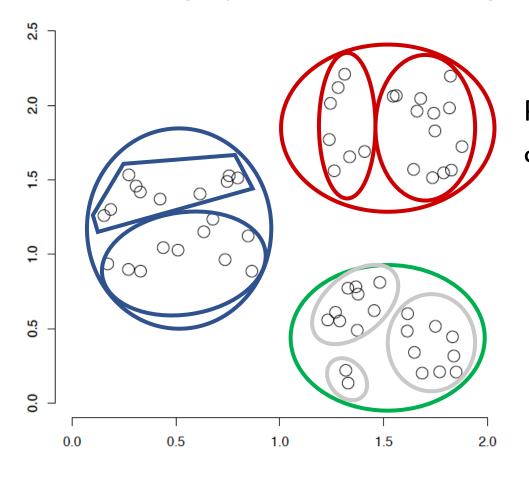
Ejemplos:

Agrupamientos, Reducción de la dimensionalidad, aprender características, estimar densidades, etc.

Aprendizaje no supervisado

Métodos básicos de aprendizaje no supervisado

Métodos de agrupamiento (Clustering)



Para los datos de la figura, ¿Cuántas agrupaciones/clusters podemos definir?

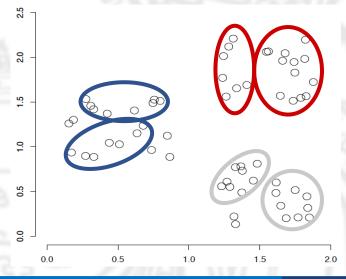
Aprendizaje no supervisado

Métodos básicos de aprendizaje no supervisado

Métodos de agrupamiento (Clustering)

- Basados en la densidad:
 - ✓ DBSCAN (Density-Based Spatial Clustering of Applications with Noise):
 - Agrupa puntos densamente conectados.
 - Identifican clusters como áreas densamente pobladas de puntos de datos, separadas por áreas de menor densidad.
 - Son especialmente efectivos para detectar clusters de forma arbitraria.
 - Adecuados para manejar ruido y outliers.
 - Puede definirse un enfoque jerárquico:
 HDBSCAN

clusters con diferentes densidades.

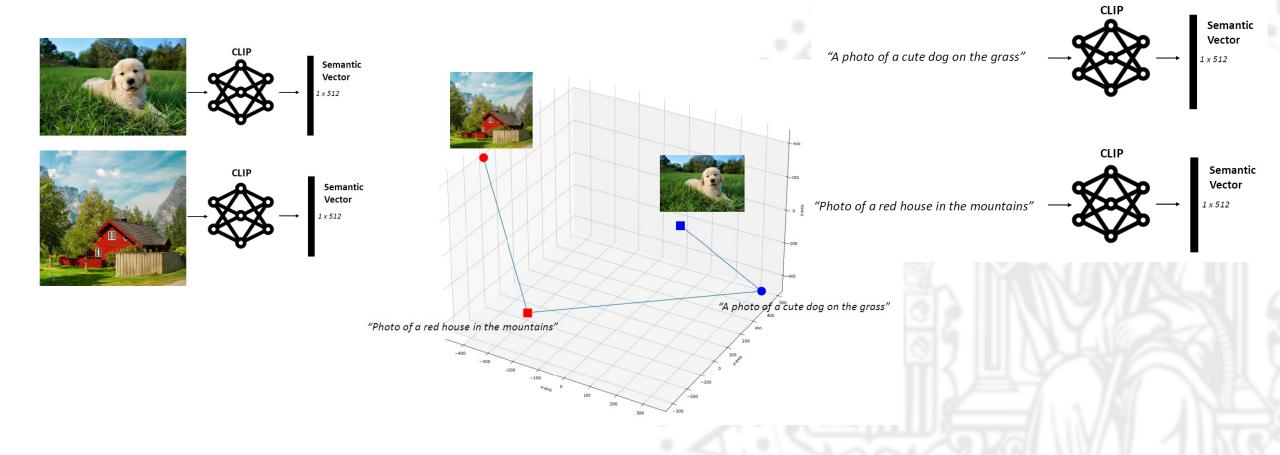


Inteligencia Artificial e información

Espacios semánticos y multimodalidad: CLIP (Contrastive Language-Image Pre-Training)

Combina un modelo de lenguaje con un modelo semántico de conocimiento de imágenes

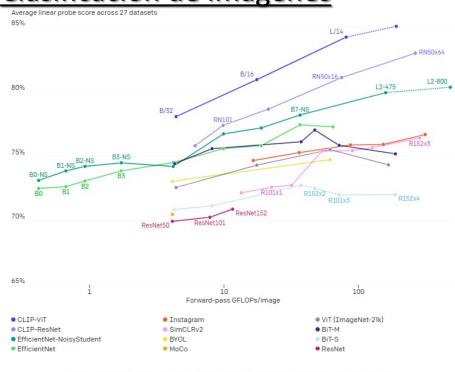
Entrenado con más de 400M de pares imagen+texto



Casos de uso con CLIP Generación de imágenes:

DALL.E de OpenAl y su sucesor DALL.E 2 VQGAN-CLIP (código abierto)

Clasificación de imágenes







(a) Oil painting of a candy dish of glass candies, mints, and other assorted sweets



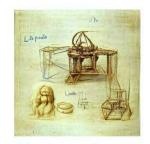
(d) A beautiful painting of a building in a serene landscape



(g) an astronaut in the style of van Gogh



(b) A colored pencil drawing of a waterfall



(e) sketch of a 3D printer by Leonardo da Vinci



(h) Baba Yaga's house + fantasy art



(c) A fantasy painting of a city in a deep valley by Ivan Aivazovsky



(f) an autogyro flying car, trending on artstation

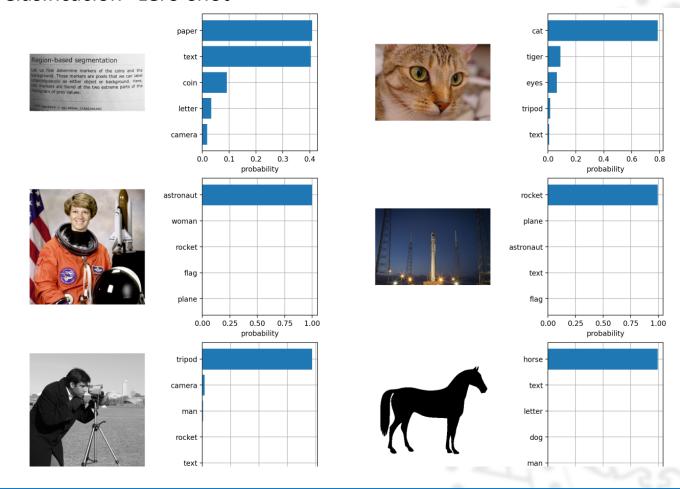


(i) pickled eggs, tempera on wood

CLIP en acción

https://colab.research.google.com/drive/109AkGnR6dHzX-stQRIWDdmEnJv9Fbc4S

Clasificación "zero-shot"



Geolocalización con CLIP

https://huggingface.co/geolocal/StreetCLIP

¿Dónde estoy?







Top 10 countries:

	Country	Score
0	Argentina	56.3566
1	Bolivia	33.6038
2	Brazil	3.4332
3	Madagascar	2.5838
4	Uruguay	1.4587
5	Chile	0.6000
6	Peru	0.3102
7	Senegal	0.3075
8	India	0.2060
9	Mexico	0.1846

Top 10 cities from: Argentina

	Admin	Score
Puerto Iguazú	Misiones	68.773
Termas de Río Hondo	Santiago del Estero	2.0647
Gualeguaychú	Entre Ríos	2.0204
Perito Moreno	Santa Cruz	1.7624
Río Colorado	La Pampa	1.4637
Catamarca	Catamarca	1.437
Alto Río Senguer	Chubut	1.2853
Gualeguay	Entre Ríos	1.2193
Río Segundo	Córdoba	0.9734
Tandil	Buenos Aires	0.9054

Entrenado con un dataset de 1,1 millones de imágenes geo etiquetadas urbanas y rurales a nivel de calle a partir del modelo openai/clip-vit-large-patch14-336

Hardware: 4 NVIDIA A100 GPUs

Horas: 12

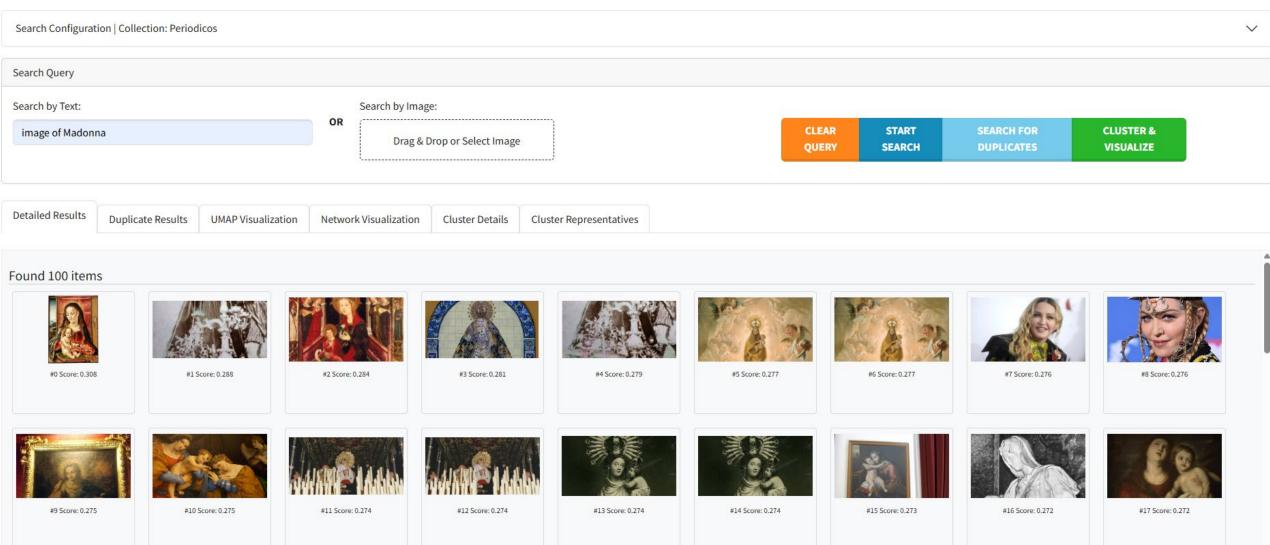
StreetCLIP fue entrenado con el objetivo de hacer coincidir las imágenes con el título correspondiente a la ciudad, región y país de origen de las imágenes.

http://155.210.153.36:8501/

Organización no supervisada de colecciones fotográficas



Multimodal Search, Clustering & Visualization



RECUPERACION DE INFORMACIÓN MULTIMODAL BASADA EN ESPACIOS SEMÁNTICOS

Los sistemas tradicionales de IR (Recuperación de Información) se basan en el indexación de texto completo de los documentos en sí mismos o en metadatos que describen los datos en el caso de audio, imágenes o video.

Nuestro problema: solo contenido de video, sin metadatos ¿Cómo representamos el contenido solo de video?

- Consumir miles de horas de documentalistas creando metadatos.
- Utilizar nuevos sistemas de descripción automática de imágenes para crear metadatos. Estos han tenido grandes avances en los últimos años, pero todavía enfrentan varios desafíos (alucinaciones, flexibilidad en la descripción, sesgos, ...) pero podrían ser una opción en el futuro próximo utilizando VLM (Modelos de Lenguaje Visuales) como Llava, GPT-4v, PaLIgemma, ...
- Usar una representación conjunta en el espacio semántico.

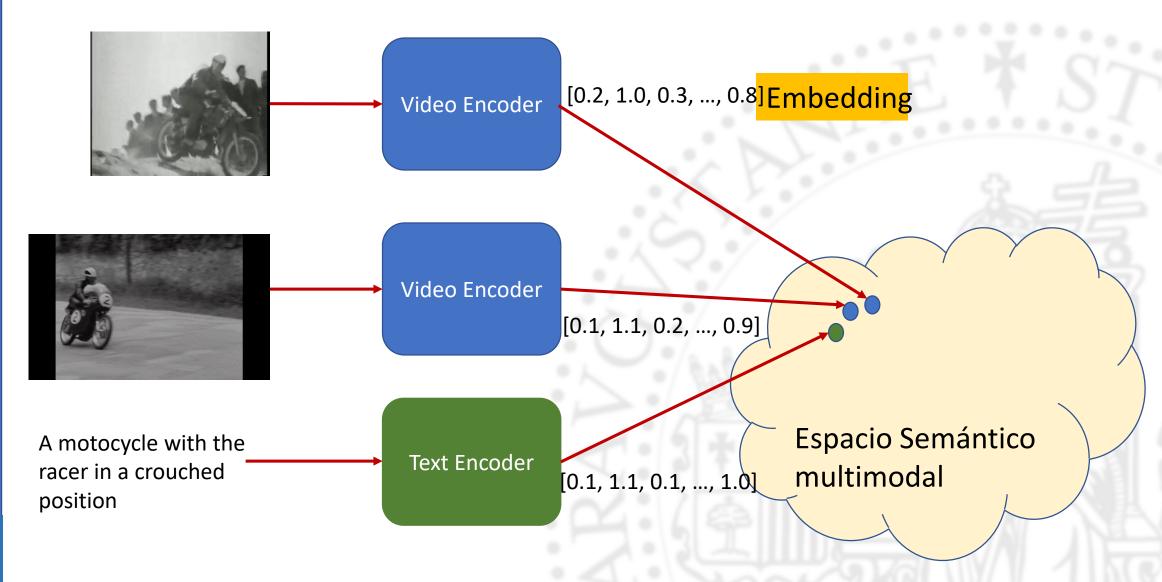


RECUPERACION DE INFORMACIÓN MULTIMODAL BASADA EN ESPACIOS SEMÁNTICOS

- El objetivo es aprender un espacio semántico conjunto que pueda capturar las relaciones inherentes entre ambas modalidades (texto y video), también conocido como modelo de espacio vectorial multimodal.
- Los modelos de espacio vectorial multimodal representan el significado (texto o video) como puntos en espacios vectoriales de alta dimensionalidad, también conocidos como espacios semánticos multimodales.
- Cada "punto" en el espacio semántico multimodal se conoce como "embedding".
- El texto y el video se representan mediante embeddings.
- Las piezas de información de diferentes modalidades que representan el mismo concepto semántico estarán "cerca" en el espacio semántico multimodal.
- La noción de "similaridad" o "proximidad" entre conceptos se reduce a la "distancia" entre los vectores de representación en el espacio vectorial.









Three important definitions

Semantic Shot: A consecutive sequence of nearby video frames embeddings in semantic space.

Semantic Scene: A consecutive sequence of nearby semantic shots in semantic space.

Semantic Class: A set of nearby video frames embeddings in semantic space. Classes are like scenes but without considering the time position.

The multimodal semantic space enables the following search functionalities:

- 1. Text-to-Video search, which retrieves semantic shots based on text queries
- Image-to-Video search, which retrieves semantic shots that are "similar" to the image query
- 3. Text+Image-to-Video search, which retrieves semantic shots based on a combination of text and image query (augmented retrieval)

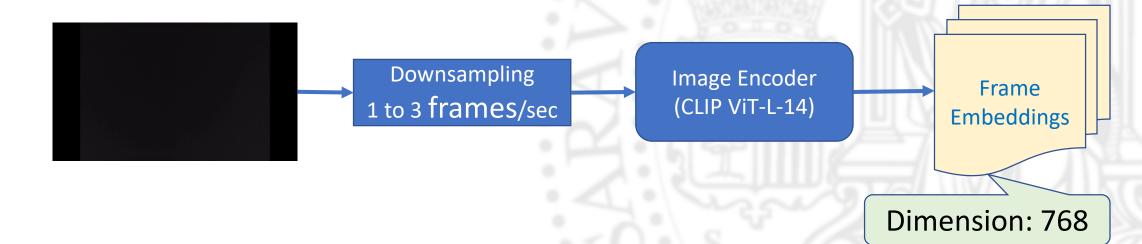


VIDEO ANALISYS

INPUT: VIDEO (mp4 format)

- 1. STEP: DONWSAMPLING VIDEO FRAME RATE (1 TO 3 FRAMES/SEC)
- STEP: COMPUTE FRAME EMBEDDINGS USING A IMAGE ENCODER (CLIP VIT-L-14)

OUTPUT: FRAME EMBEDDINGS (768 DIMENSIONAL VECTORS)





SEMANTIC ANALYSIS

INPUT: FRAME EMBEDDINGS (768D x (time video lenght in seconds x 3))

- COARSE GROUPING: Find Clusters in the Semantic Space → SEMANTIC CLASSES
- 2. CLASS DIARIZATION: Distribute Semantic Classes in the Time-Line → SEMANTIC SCENES
- 3. FINE GROUPING: Find Clusters in the Semantic Scenes → SEMANTIC SHOTS
- 4. SHOT EMBEDDINGS: Use Embedding Mean Average over Semantic Shots

OUTPUT: SHOT EMBEDDINGS (768D x (#semantic shots))

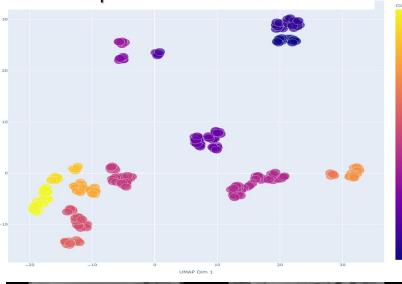




Semantic hierarchy

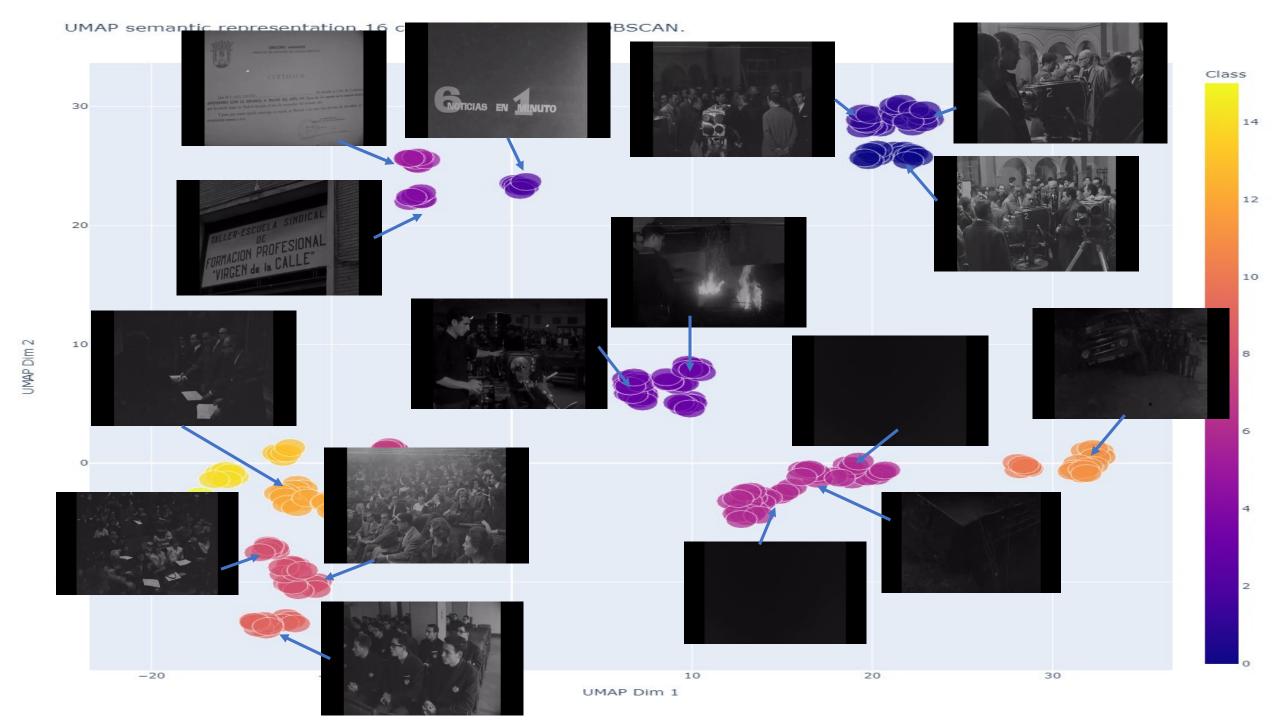
Clustering over Semantic Space Semantic Class (coarse grouping) Diarization over Semantic Semantic **Semantic Classes** Scene Scene Clustering over Semantic Semantic Semantic Scenes Shot Shot (fine grouping) Shot Shot Mean Pooling **Embedding Embedding**













VIVOLAB

MODELOS DE LENGUAJE MULTIMODALES



a motorcyclist driving a racing motorcycle

Query

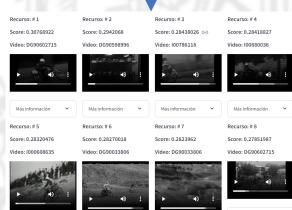
Without optimization 1:30 hours to populate the database with the 27:35 hours of video Query embedding computation

Image/text encoder (CLIP ViT-L-14)

768D space

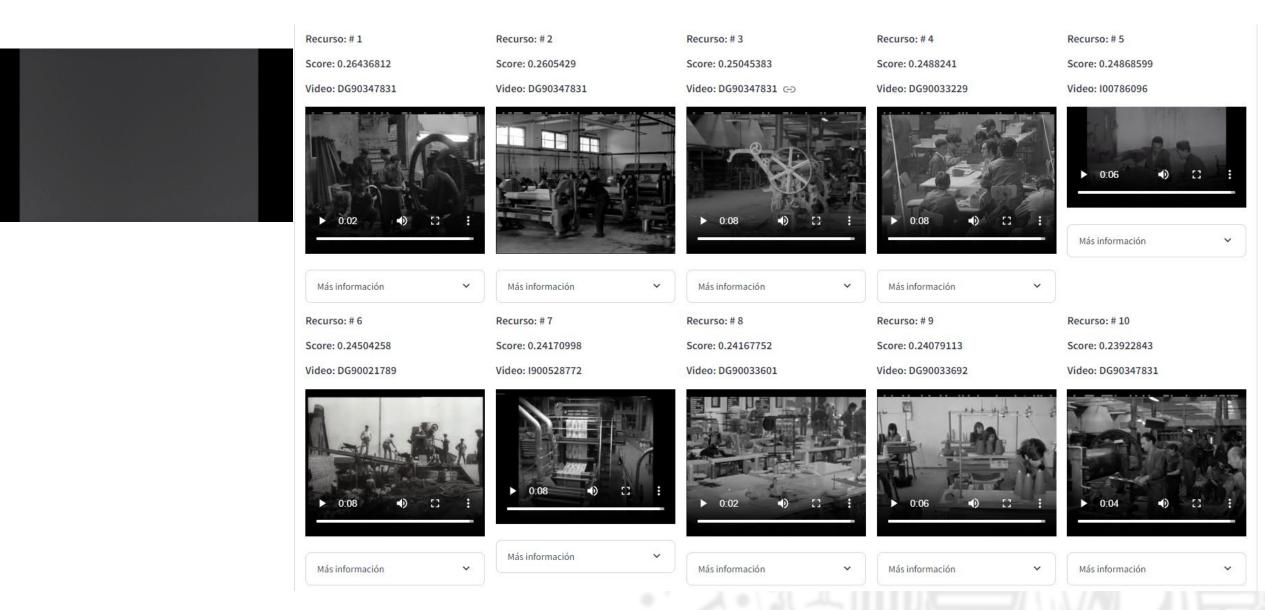
Vector Database Qdrant 370 videos (18.404 shots) (27:35:02) Refine search

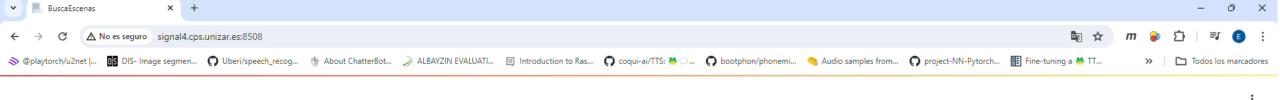
Shot Search Video shots retrieval





FACTORY. People working in a Factory. (I00786311.mp4)





BuscaEscenas: rtvePoCArchivo3



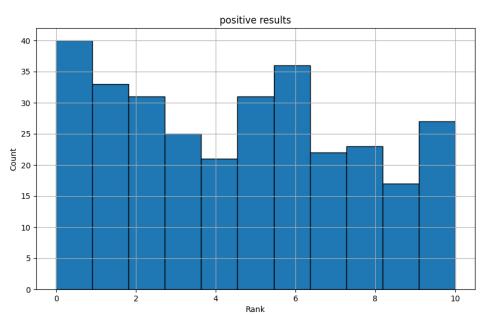
Duscui

Introduzca una consulta

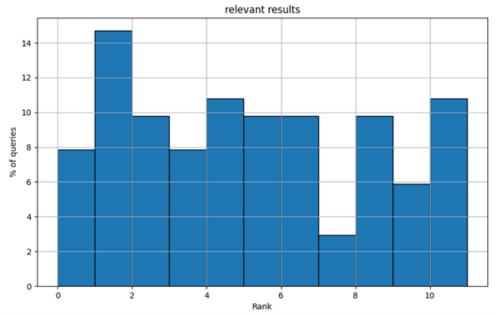
Introducción a la Visión Artificial con Modelos de Lenguaje Multimodales

Test subjetivo con documentalistas.

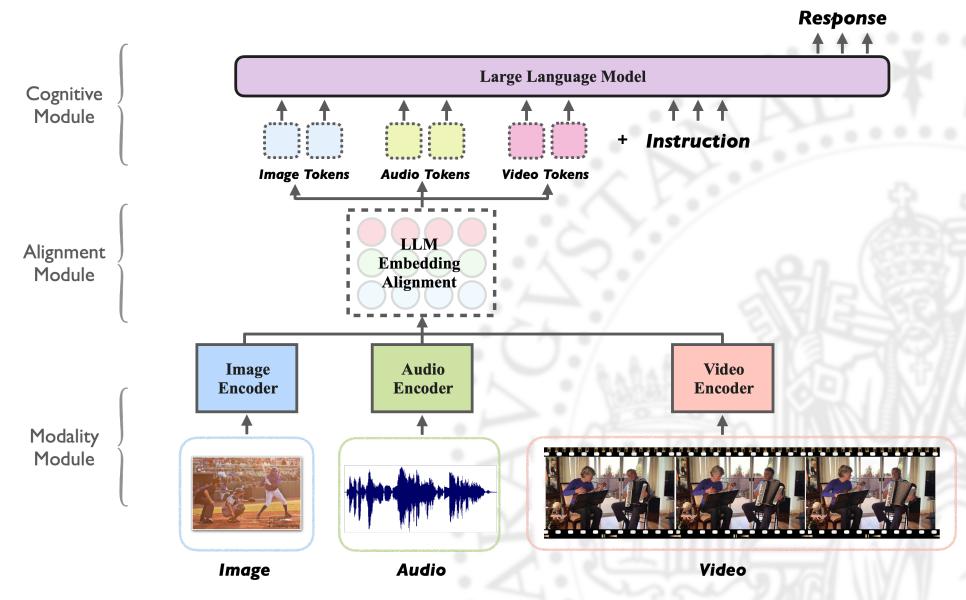
Fase I: 273 búsquedas: 86,97% R@10



Fase II: 102 búsquedas: 92,15% R@10



- Esta tecnología permite explorer el archivo sin necesidad de metadatos
- Valorar en el futuro la integración de descripciones con modelos de lenguaje multimodales





El campeón invierno/primavera 2025 modelos abiertos



https://qwenlm.github.io/blog/qwen2.5-vl/

https://github.com/QwenLM/Qwen2.5-VL

https://ollama.com/library/qwen2.5vl



Características clave del modelo Qwen2.5-VL

- Comprensión Visual Profunda: analizar textos, gráficos, iconos, diagramas y la disposición de elementos dentro de imágenes.
- Capacidad de Agente Visual: capaz de razonar, dirigir herramientas dinámicamente e interactuar con ordenadores y teléfonos.
- Análisis de Vídeos Largos (>1h) y Captura de Eventos: Comprende vídeos extensos e identifica segmentos clave.
- Localización Visual Precisa y Salida JSON: Localiza objetos (cajas/puntos) con JSON estable.
- Generación de Salidas Estructuradas para Documentos: Genera salidas estructuradas para el contenido de documentos como escaneos de facturas, formularios y tablas, lo cual es útil en finanzas, comercio, etc.





Multimodal Playground

Pick a model available locally on your system ↓

Qwen/Qwen2.5-VL-7B-Instruct

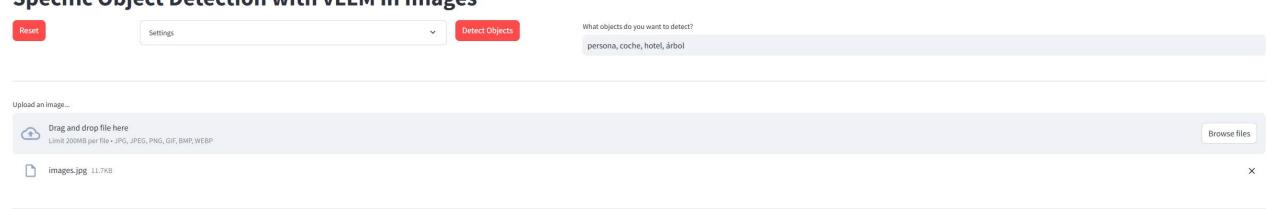
You have selected the model: Qwen/Qwen2.5-VL-7B-Instruct

Upload an image or video for analysis

Drag and drop file here
Limit 200MB per file • PNG, JPG, JPEG, WEBP, MP4, MPEG4



Specific Object Detection with vLLM in Images



Original Image



Object Detection Results (YOLOv11x)



Annotated Image with YOLOv11 0.0356 seconds

Object Detection Results (RF-DETR)



Annotated Original Image 0.0357 seconds

Detection Details:

persona: interact with the car (bbox: [208, 93, 226, 145])

Result with Detections 4.3441 seconds

Image with Detections (Qwen2.5 VL)

- coche: parked car (bbox: [7, 93, 125, 142])
- coche: driving car (bbox: [93, 100, 193, 153])
- hotel: enter hotel (bbox: [199, 45, 318, 135])
- árbol: stand near hotel (bbox: [111, 23, 190, 112])

Supplied to the section of the secti

Detection Details:

- car: No action (bbox: [91, 101, 210, 152], confidence: 0.93)
- car: No action (bbox: [14, 91, 127, 143], confidence: 0.92)
- person: No action (bbox: [207, 95, 224, 146], confidence: 0.83)

Detection Details:

- person: No action (bbox: [208, 95, 224, 146])
- car: No action (bbox: [91, 100, 211, 152])
- car: No action (bbox: [14, 92, 128, 143])
- truck: No action (bbox: [14, 92, 128, 143])

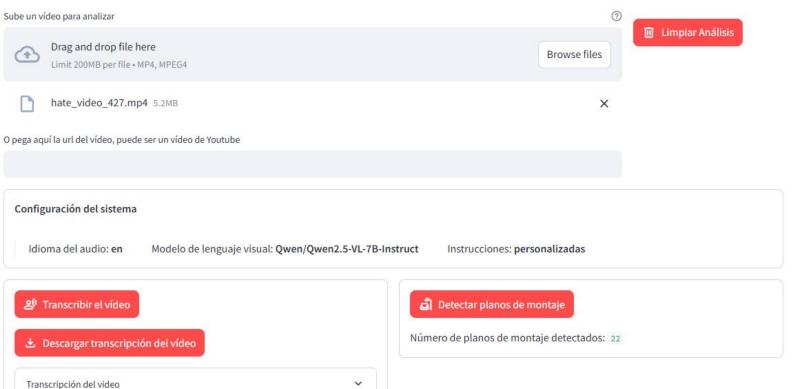


Configuración Análisis de vídeo 👔 Análisis masivo 🍃

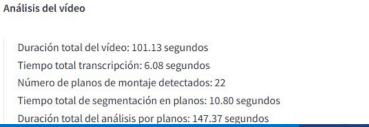


COUNTERING MEDIA INTOLERANCE IN YOUNG AUDIENCES

Análisis de vídeo







Plano 8: {'Anticlimatismo': 0, 'LGBTQ+fobia': 0, 'Machismo': 0, 'Racismo y Xenofobia': 1, 'Antigitanismo': 0, 'Islamofobia': 0, 'Antisemitismo': 0, 'Aspectismo': 0}

tiempo_inicio: 00:01:35.267

MORE DIVERSITY, THEY SAID...

tiempo_inicio: 00:01:35.267 tiempo_final: 00:01:45.100

clase: {'Anticlimatismo': 0, 'LGBTQ+fobia': 0, 'Machismo': 0, 'Racismo y Xenofobia': 1, 'Antigitanismo': 0, 'Islamofobia': 0, 'Antisemitismo': 0, 'Aspectismo': 0}
atributos_de_lenguaje_toxico: {'TOXICITY': 0.0, 'SEVERE_TOXICITY': 0.0, 'IDENTITY_ATTACK': 0.0, 'INSULT': 0.0, 'PROFANITY': 0.0, 'THREAT': 0.0}
tipo_de_audiencia: false

tema: El tema central de la escena es el racismo y la discriminación racial.

justificacion: La transcripción de audio contiene comentarios racistas y ofensivos, como 'When are we whites gonna have our day?' y 'The niggers had theirs such a long, long time.' Estos comentarios expresan un deseo de exclusividad racial y estigmatizan a una comunidad racial. La imagen visual muestra personas de diferentes grupos étnicos detrás de un mismo cercado, lo que sugiere una crítica sobre la inclusión y diversidad en términos de oportunidades y derechos.

Transcripción de la escena

[['0:01:35.260000', '0:01:37.282000', '<\$#02>', 'Hey, Mr. President, what do you say?'], ['0:01:37.702000', '0:01:39.824000', '<\$#02>', 'When are we whites gonna have our day?'], ['0:01:40.064000', '0:01:42.386000', '<\$#02>', 'The niggers had theirs such a long, long time.'], ['0:01:42.486000', '0:01:44.948000', '<\$#02>', "I'm a white and it's time that I have mine."]]

Plano 19: {'Anticlimatismo': 0, 'LGBTQ+fobia': 0, 'Machismo': 0, 'Racismo y Xenofobia': 1, 'Antigitanismo': 0, 'Islamofobia': 0, 'Antisemitismo': 0, 'Aspectismo': 0}



tiempo_inicio: 00:03:54.633 tiempo_final: 00:03:59.633

justificacion: La imagen muestra una escena que incluye un diálogo entre tres personajes, donde uno de ellos dice 'I told you not to date black men.' (Te advertí que no salieras con hombres negros). Este mensaje expresa claramente una actitud discriminatoria hacia la raza negra, promoviendo estereotipos y prejuicios racistas. No hay elementos adicionales en la transcripción de audio para contrarrestar esta afirmación, lo que la convierte en una declaración abiertamente discriminatoria.

clase: {'Anticlimatismo': 0, 'LGBTQ+fobia': 0, 'Machismo': 0, 'Racismo y Xenofobia': 1, 'Antigitanismo': 0, 'Islamofobia': 0, 'Antisemitismo': 0, 'Aspectismo': 0} atributos_de_lenguaje_toxico: {'TOXICITY': 0.7, 'SEVERE_TOXICITY': 0.9, 'IDENTITY_ATTACK': 0.8, 'INSULT': 0.9, 'PROFANITY': 0.0, 'THREAT': 0.0} tipo_de_audiencia: false

tema: El tema central de la escena es el racismo y la discriminación racial.

Transcripción de la escena

['00:03:54.633', '00:03:59.633', '<S#0>', 'No se ha detectado ninguna voz en la grabación']