



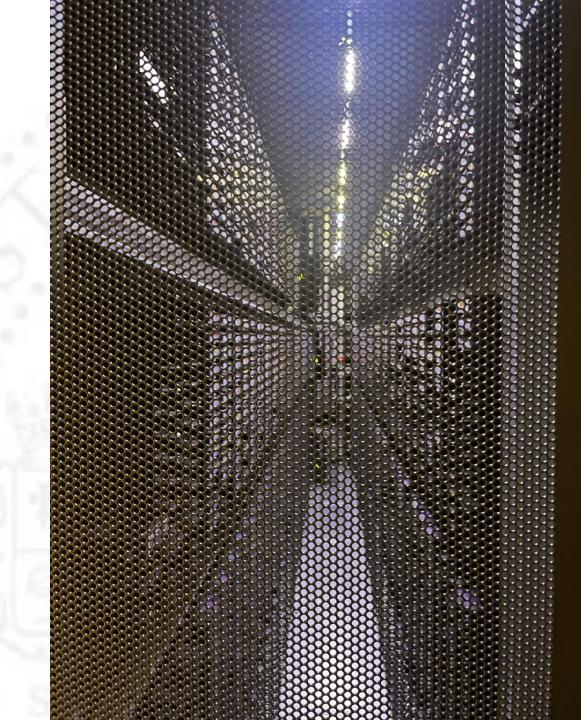
GURSO

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

niversidad ragoza Análisis y extracción de metadatos en contenidos audiovisuales:

Tecnologías del Habla





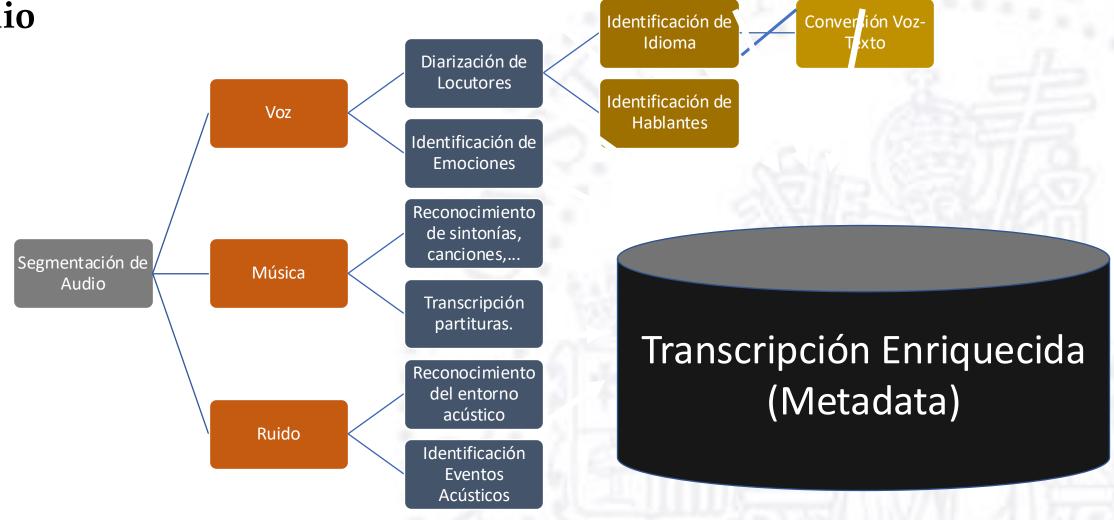


¿Qué información podemos encontrar en un audio?

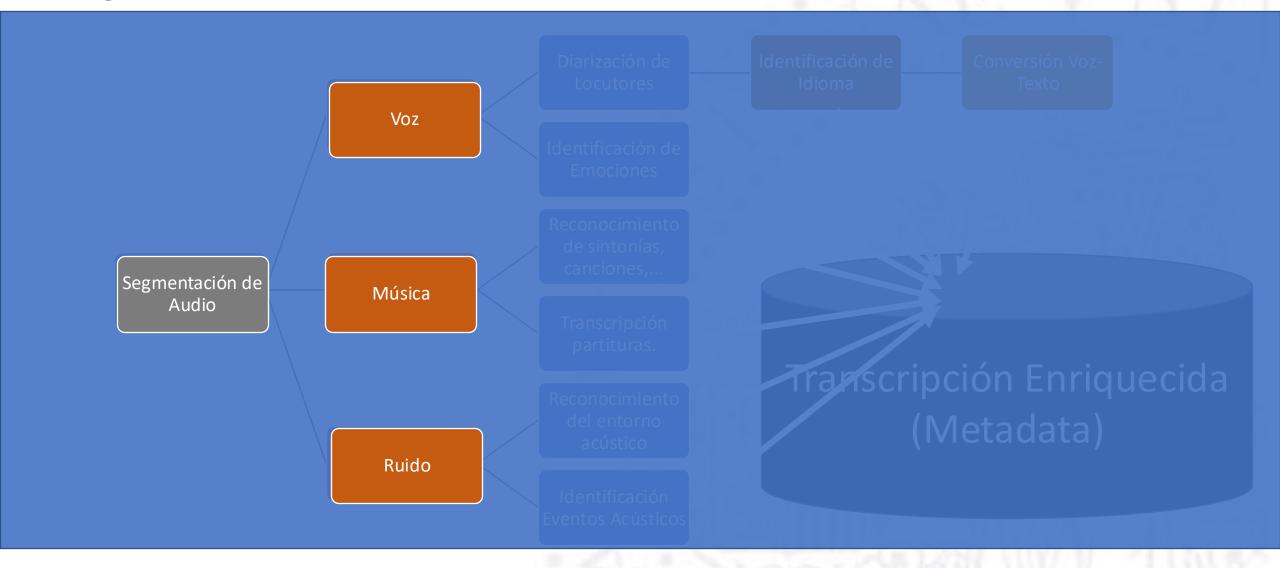


- Hay ruido, música, habla, ...
- Cuántas personas hablan y cuándo habla cada una de ellas
- Cuáles son las identidades de las personas que hablan
- En qué idioma están hablando
- Qué dice cada una de ellas
- · Cuál es el estado emocional de cada una de ellas

Tecnologías Audio



Segmentación de Audio







Segmentación de Audio

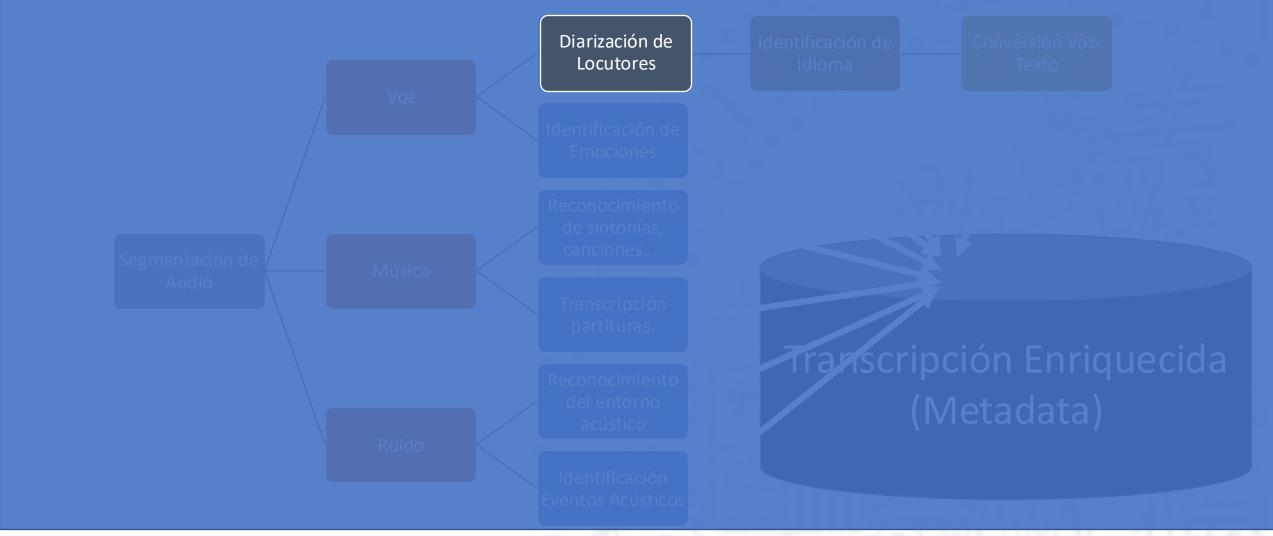
- ¿Qué es?:
 - Dividir el audio de entrada en fragmentos atendiendo al tipo de contenido acústico: Voz / Música / Ruido y combinaciones de estos.







- ¿Para qué sirve?:
 - Da soporte a otras tareas de extracción de información como:
 - Diarización
 - Identificación del hablante
 - Conversión Voz-Texto
 - ...





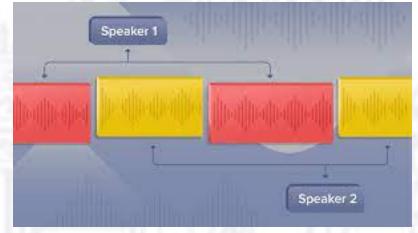


•¿Qué es?:

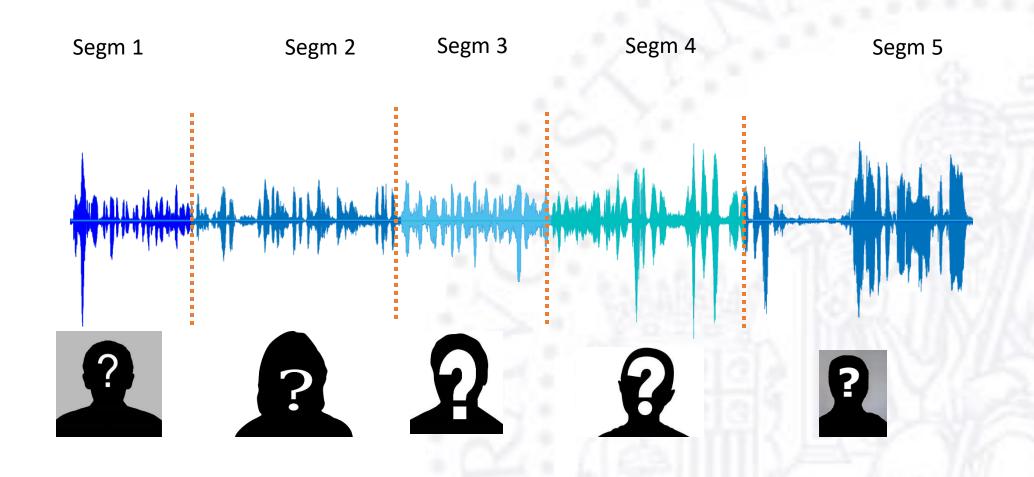
- Dividir en fragmentos atendiendo al interviniente y agrupar dichos fragmentos en función de la identidad del locutor.
- Término usado por la comunidad: Diarización
 - Diarise: (Diarize) to make use of a diary to record past events or those planned for the future.

• ¿Para qué sirve?:

- Tecnología soporte para mejorar prestaciones de:
 - · Reconocimiento automático del habla
 - Reconocimiento del hablante
 - ...



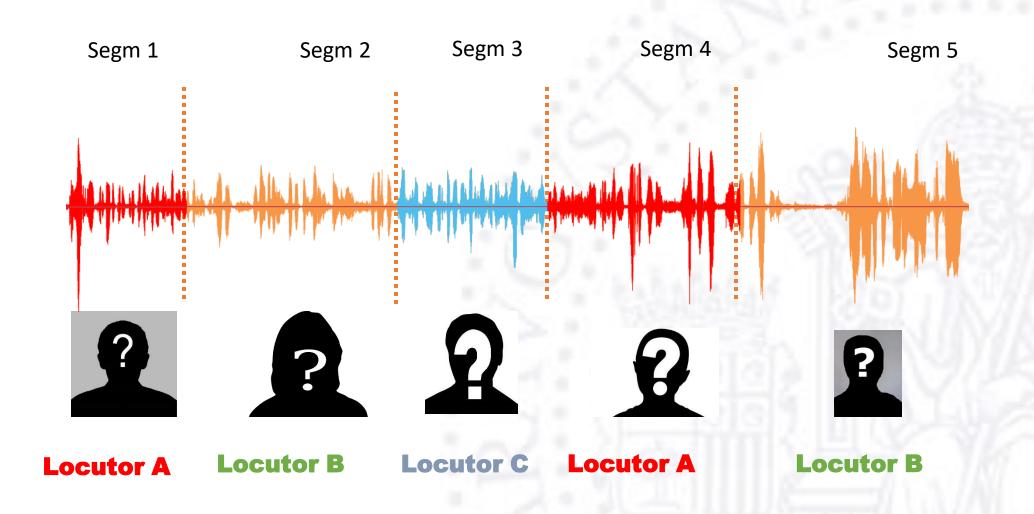






nentación y

Agrupación de Hablantes

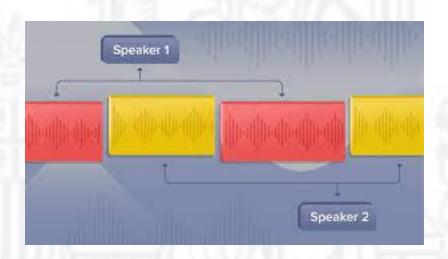




• ¿Para qué sirve?:

- Tecnología soporte para mejorar prestaciones de:
 - Reconocimiento automático del habla
 - Reconocimiento del hablante

• ...

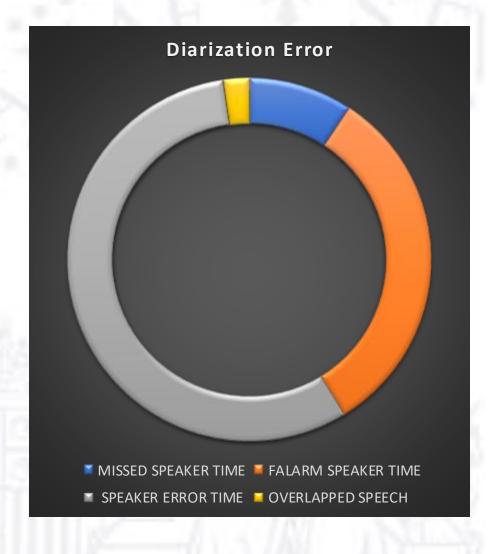


Medida del Error

• Diarization Error Rate:

$$DER = \frac{T_{incorrecto}}{T_{voz}}$$

- Componentes del error:
 - Pérdida:
 - Hay voz pero se ha confundido con silencio.
 - Falsa Alarma:
 - No hay voz pero se ha detectado erróneamente.
 - Error de Locutor:
 - Se ha confundido un locutor con otro.
 - Error por Solape:
 - Dos locutores hablan a la vez, pero solo se ha identificado a uno.



Prestaciones: Albayzín Retos - RTVE







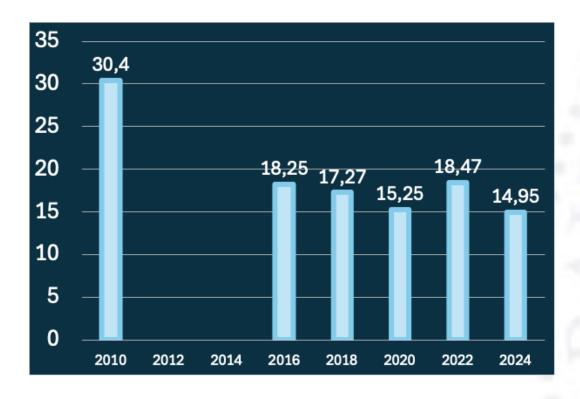


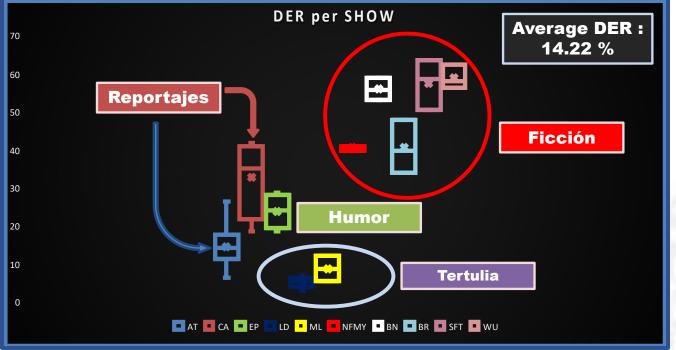




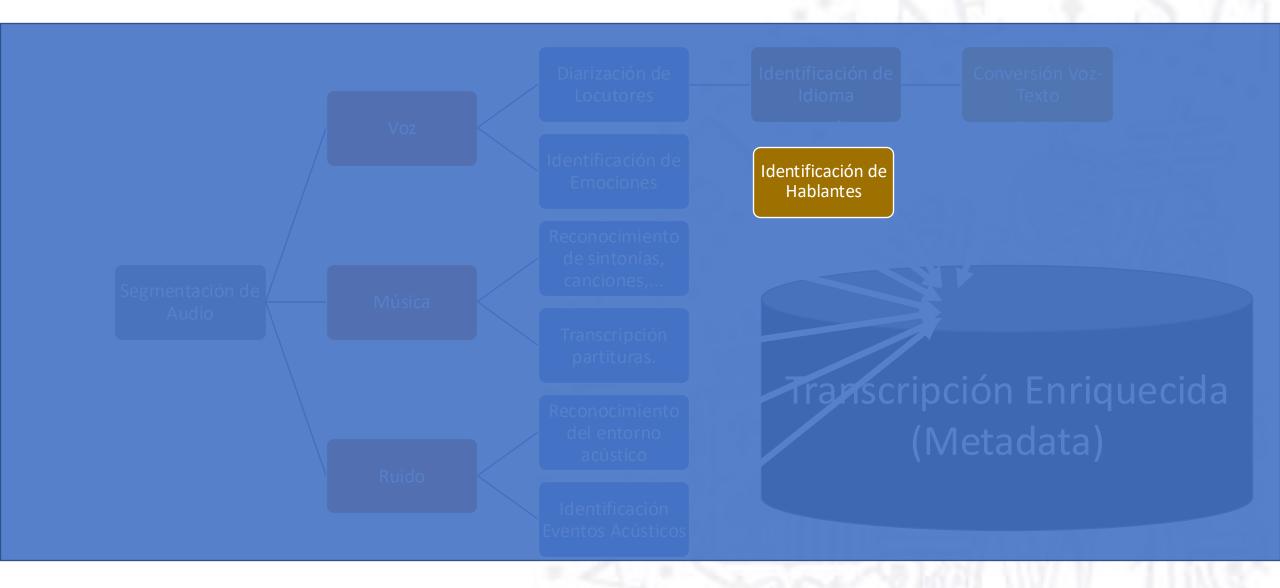








Identificación de Hablantes







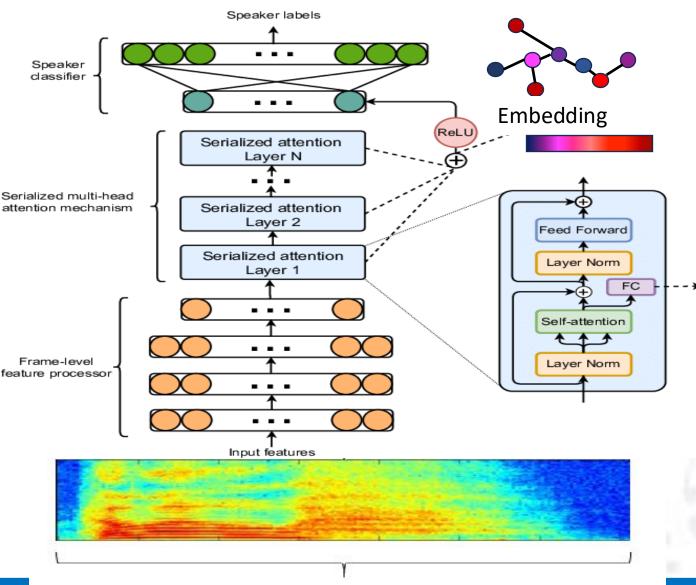
Identificación de Hablantes

• ¿Para qué sirve?:

• Permite asignar identidades concretas a fragmentos de audio de un contenido analizado



Cómo se hace ...



De cada fragmento de voz se extrae una "huella" que se compara con las almacenadas en la fase de registro para saber quién ha pronunciado ese contenido

FASE DE REGISTRO



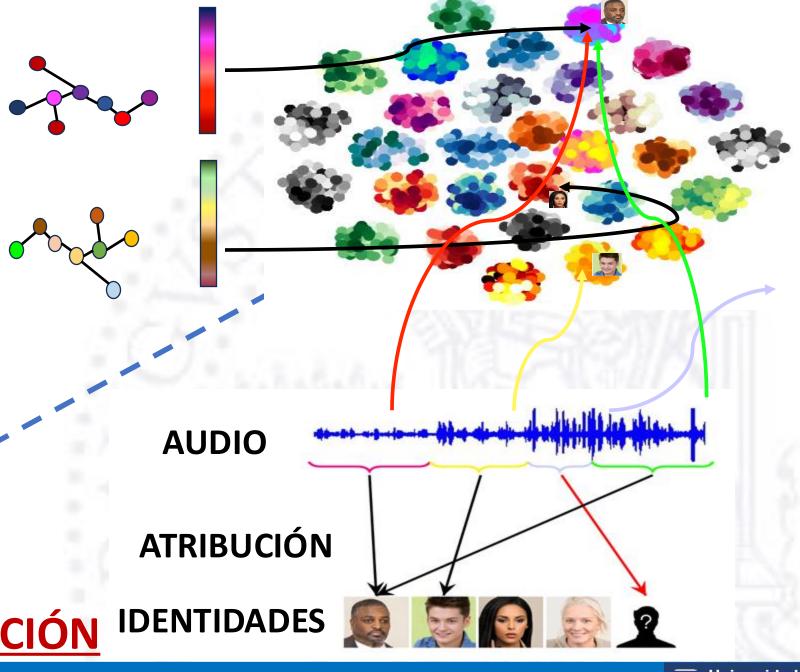
վվիկտա---

Bob

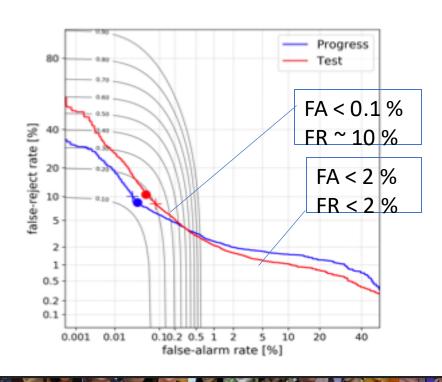
Creación de Modelos



Sally



Precisión ...





The 2019 NIST Speaker Recognition Evaluation CTS Challenge

Seyed Omid Sadjadi¹, Craig Greenberg¹, Elliot Singer^{2,†}, Douglas Reynolds^{2,†}, Lisa Mason³, Jaime Hernandez-Cordero³



utterances

hours

Redes Neuronales con cientos de millones de parámetros consiguen tasas de error por debajo del 2%

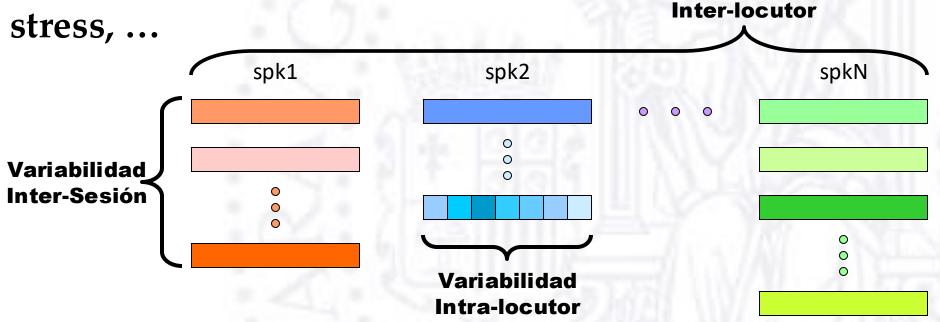
speakers

eligencia Artificial en los Archivos, Jaca 2025



En entornos reales de operación ...

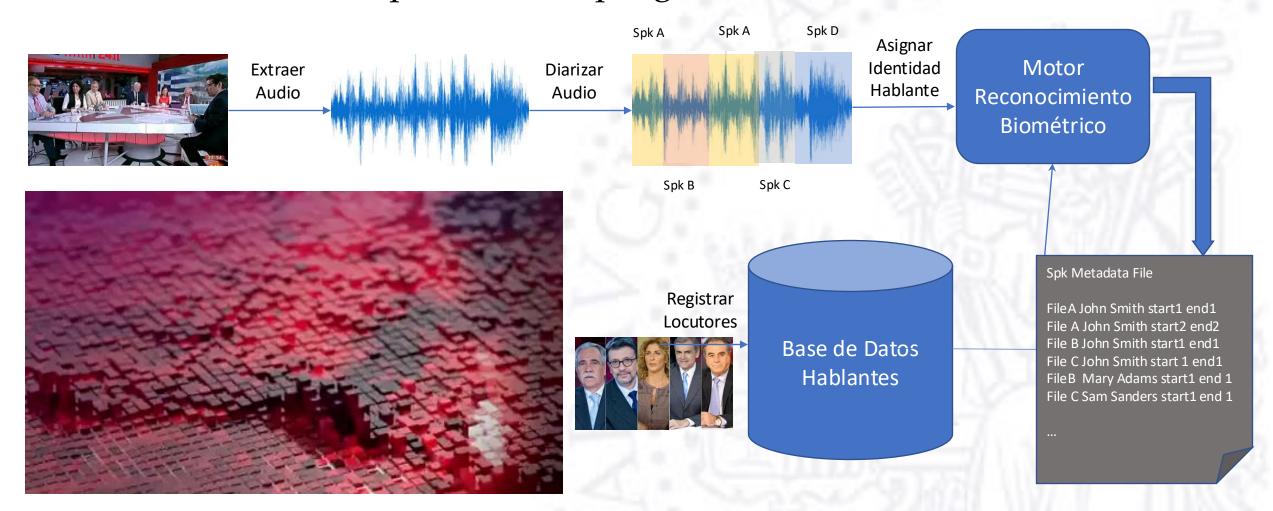
- Alta variabilidad en la voz de los hablantes
- Diversidad de dominios acústicos (Estudio, Calle, Estadio, ...)
- Solape entre hablantes
- Intervenciones de corta duración
- Voz emocional, stress, ...



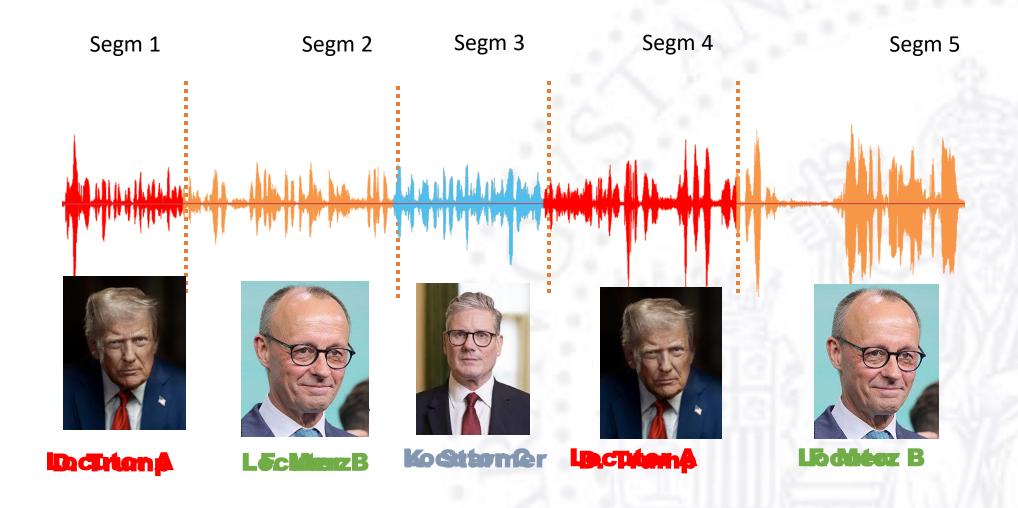
Variabilidad

Identificación de Hablantes

Reconocimiento de personas en programas de TV:



Diarización Junto con Identificación de Hablante:



Prestaciones: Albayzin – Retos RTVE

















	AER
BIOMETRIC VOX	65.09 %
VIVOLAB	72.63 %

	MISS	FA	SPKERR	AER
TEAMIV_p	1,3	229,7	9,5	240,55
TEAMIV_I3	3,7	160,8	20,9	185,42
TEAMIV_I6	8,3	76,5	5,9	90,62
TEAMIV_I9	12,4	15,3	1,2	28,88

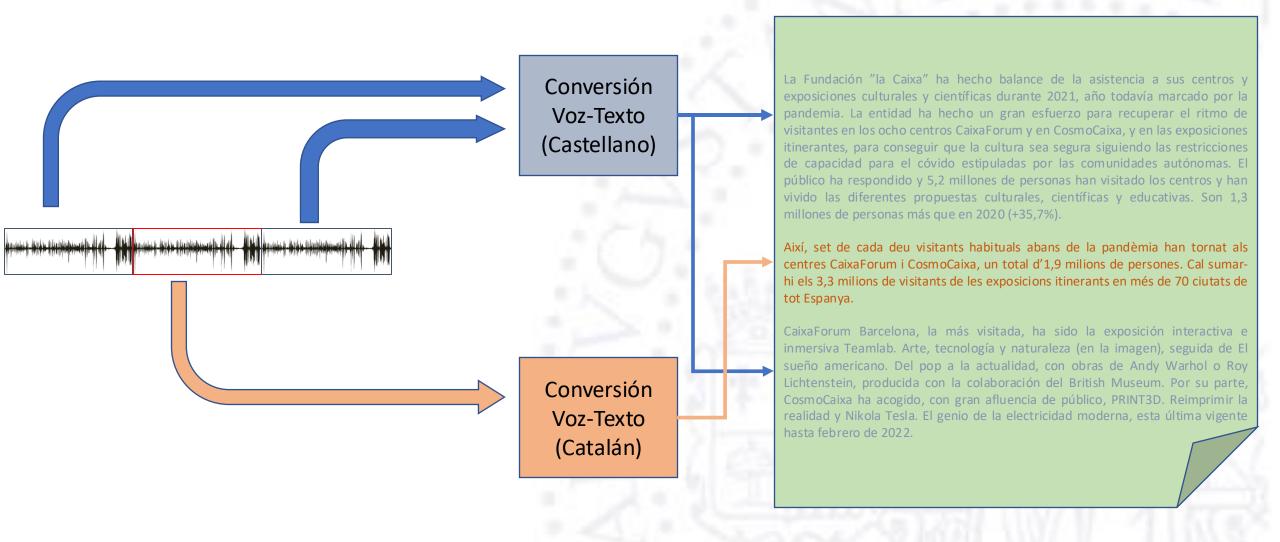
Subset	Closed Condition		Open Condition			
	Direct	Indirect	Hybrid	Direct	Indirect	Hybrid
Dev. subset	13.73	15.27	15.89	41.91	37.45	37.68
Eval. subset	25.11	17.20	16.49	65.31	60.34	31.95

Identificación de Idioma

• ¿Para qué sirve?:

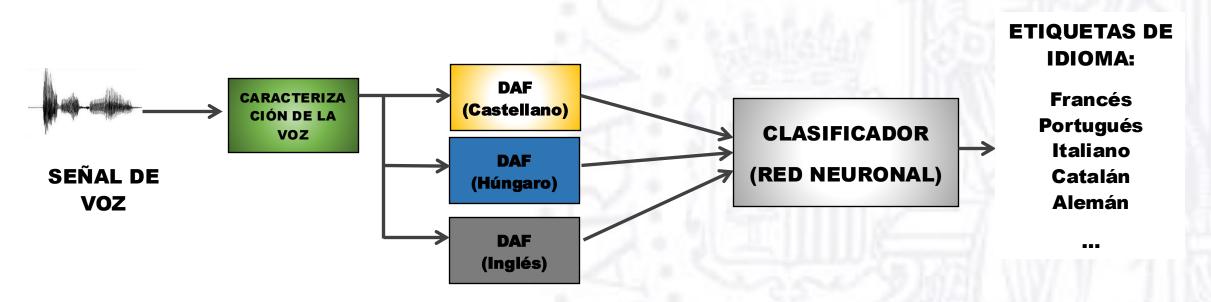
- En entornos multilingües, permite el indexado y la recuperación de documentos:
- Esencial en esos entornos como soporte a:
 - Reconocimiento automático del habla

Identificación de Idioma

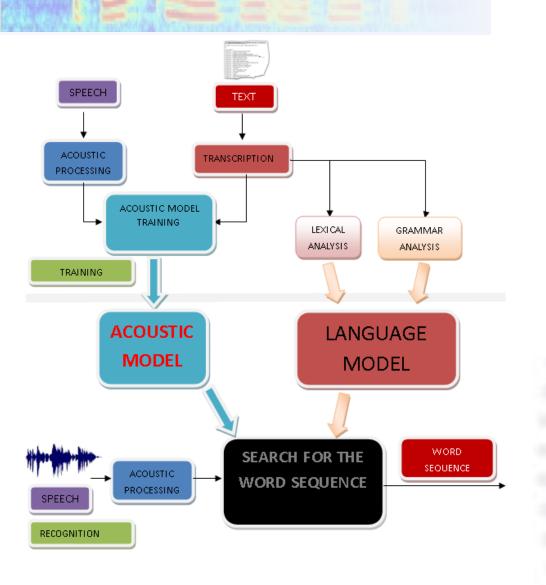


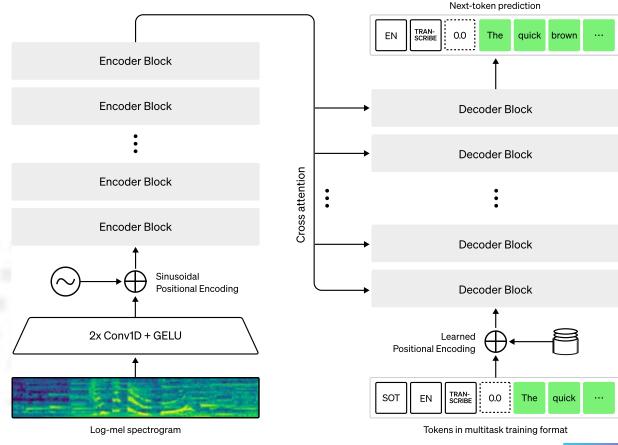
Identificación de Idioma

- Acústicos
 - Tratan de buscar patrones discriminativos directamente sobre la señal de voz
- Fonotácticos (Lingüísticos)
 - Primero procesan la señal de entrada con un (o varios) reconocedor fonético (en varios idiomas) después buscan patrones discriminativos en las secuencias de fonemas de salida



Etapas y Procesos RAH:









Componentes de un Sistema de RAH

MODELO ACÚSTICO

 Describe las características de cada unidad desde el punto de vista de la señal de voz (espectralmente)

MODELO DE LENGUAJE

- Describe las relaciones entre palabras del vocabulario
- Cuantifica la probabilidad de las secuencias de palabras

MODELO LÉXICO

 Describe cómo se forma cada palabra del vocabulario a partir de las diferentes unidades del modelo acústico.

Errores en un Sistema de RAH

Borrados

El locutor dice algo pero el sistema no devuelve nada

Substituciones

 El sistema devuelve a su salida una palabra diferente de la pronunciada por el locutor.

Inserciones

 El locutor no dice nada, pero el sistema devuelve alguna palabra (generalmente debido a artefactos acústicos)

Errores en un Sistema de RAH

Métricas de Precisión y Error:

REF: a las tres y siete minutos de mañana HYP: a las tres diecisiete minutos de la mañana

CORRECTO (C) ERRORES:

Substitutiones (S), Borrados (B), Inserciones (I)

% ACC =
$$\frac{C}{C+S+B+1}$$
 x 100
% WER = $\frac{S+B+1}{C+S+B}$ x 100

Prestaciones: Albayzin – Retos RTVE



Universidad Zaragoza





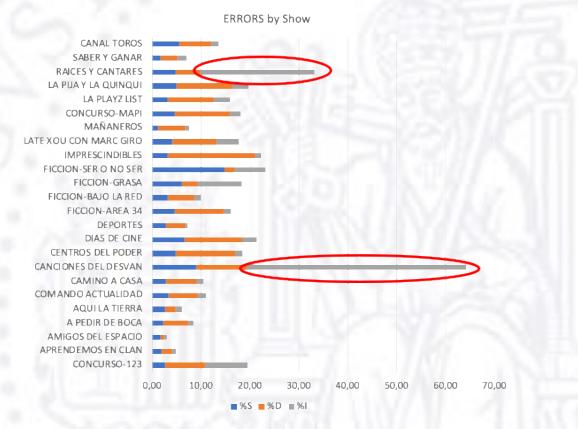












WER Baseline WhisperX large-v3: 2024: 12,10 %



Diarización e Identificación de hablantes Reconocimiento Automático del Habla:





D. Trump: Contenido de la intervención 1 ... <Tcomienzo1> <Tfin1>



F. Merz: Contenido de la intervención 2 ... <Tcomienzo2> <Tfin2>



K. Starmer: Contenido de la intervención 3 ... < Tcomienzo 3 > < Tfin 3 >

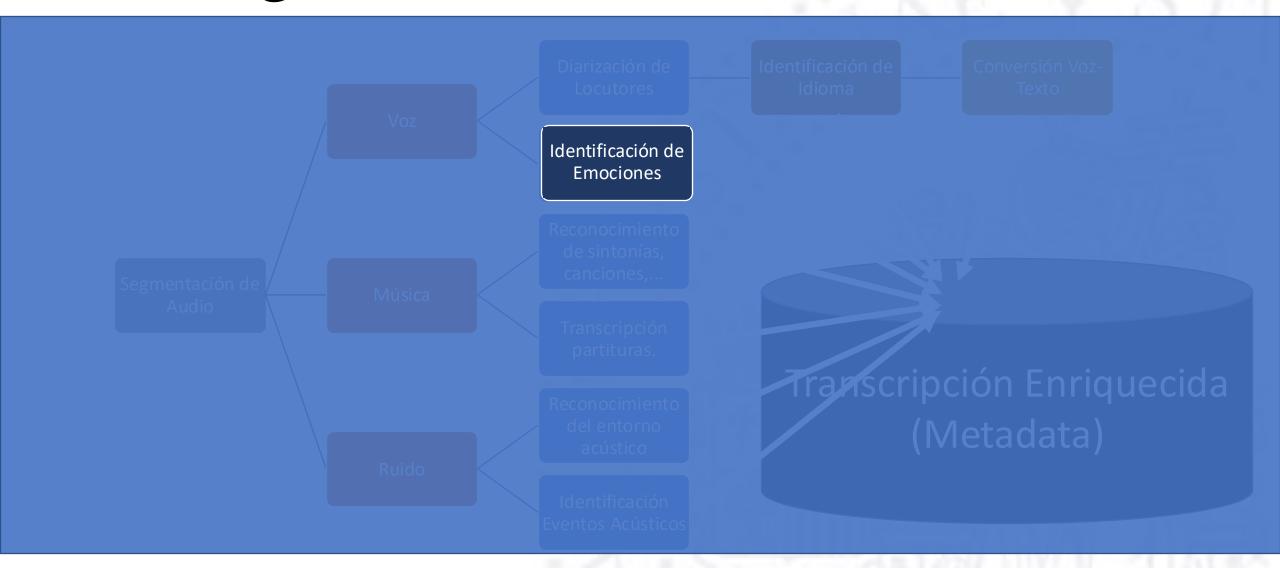


D. Trump: Contenido de la intervención 4 <Tcomienzo4> <Tfin4>



F. Merz: Contenido de la intervención 5 ... <Tcomienzo5> <Tfin5>

Tecnologías







Identificación de Emociones

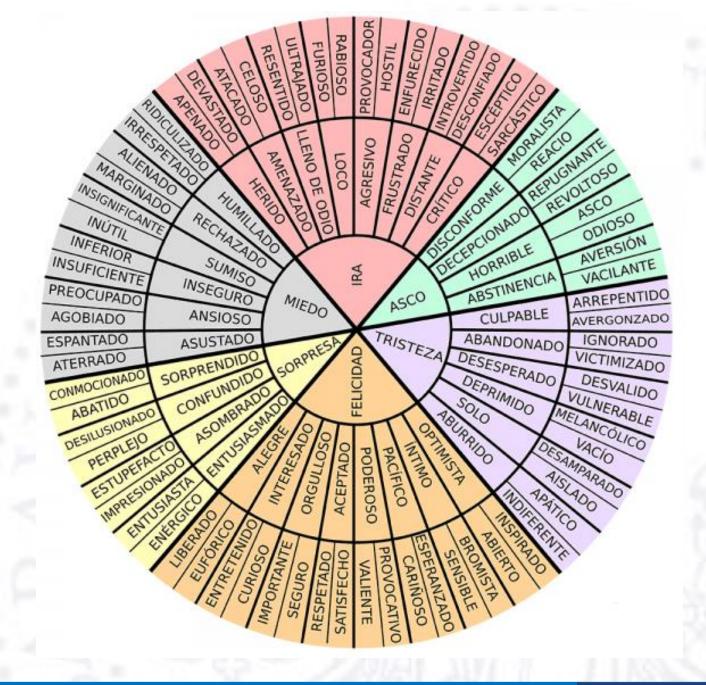
• ¿Para qué sirve?:

• Puede añadir información extra que enriquece el discurso de los protagonistas de un contenido

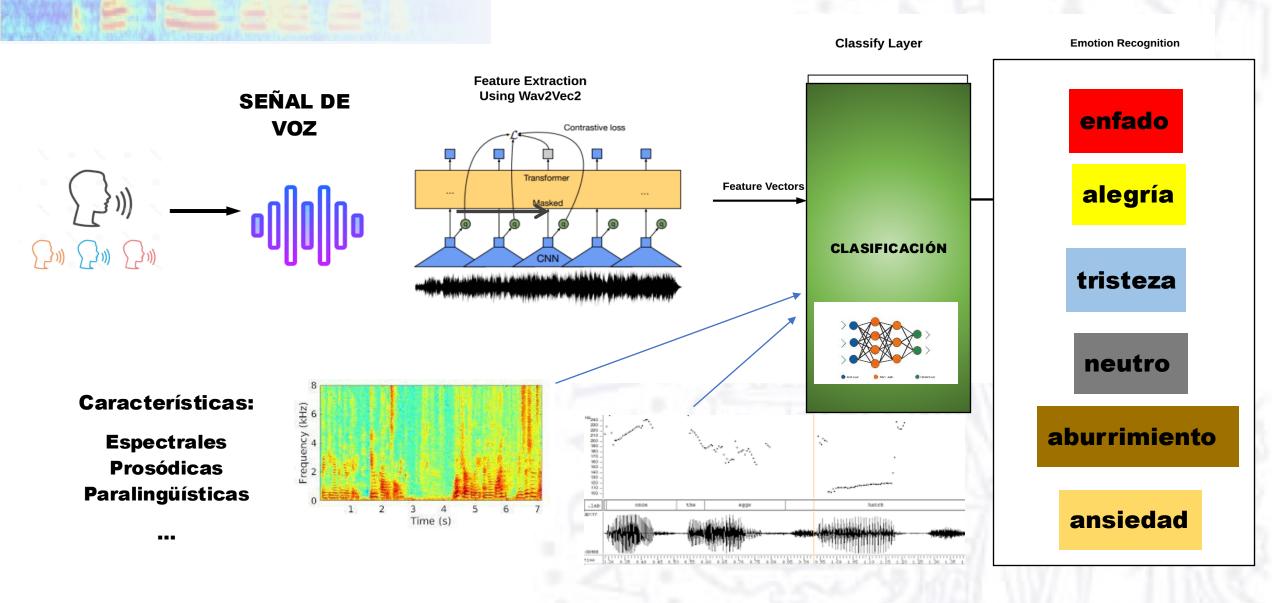


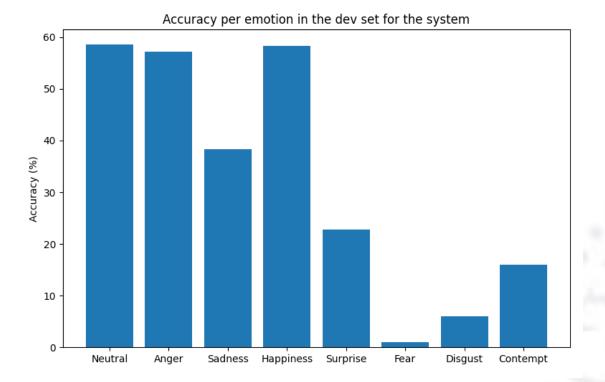
El problema de los datos

- Cantidad limitada de audio
- Emociones simuladas
- Pocos locutores
- Etiquetado subjetivo



Identificación de Emociones





- Diferencias grandes en cuanto a precisión entre emociones
- La diferencia en precisión es consecuente con los estudios realizados en humanos y el consenso del conjunto de datos.

