**CIDAI** Centre of Innovation for Data tech and Artificial Intelligence

MATXA, the first speech open-source solution to support the different Catalan varieties - Barcelona Supercomputing Center (BSC)

# MATXA, the first speech open-source solution to support the different Catalan varieties –BSC

## Martí Llopart Font
*Direcció Innovació Digital*
**Barcelona Supercomputing Center (BSC)**
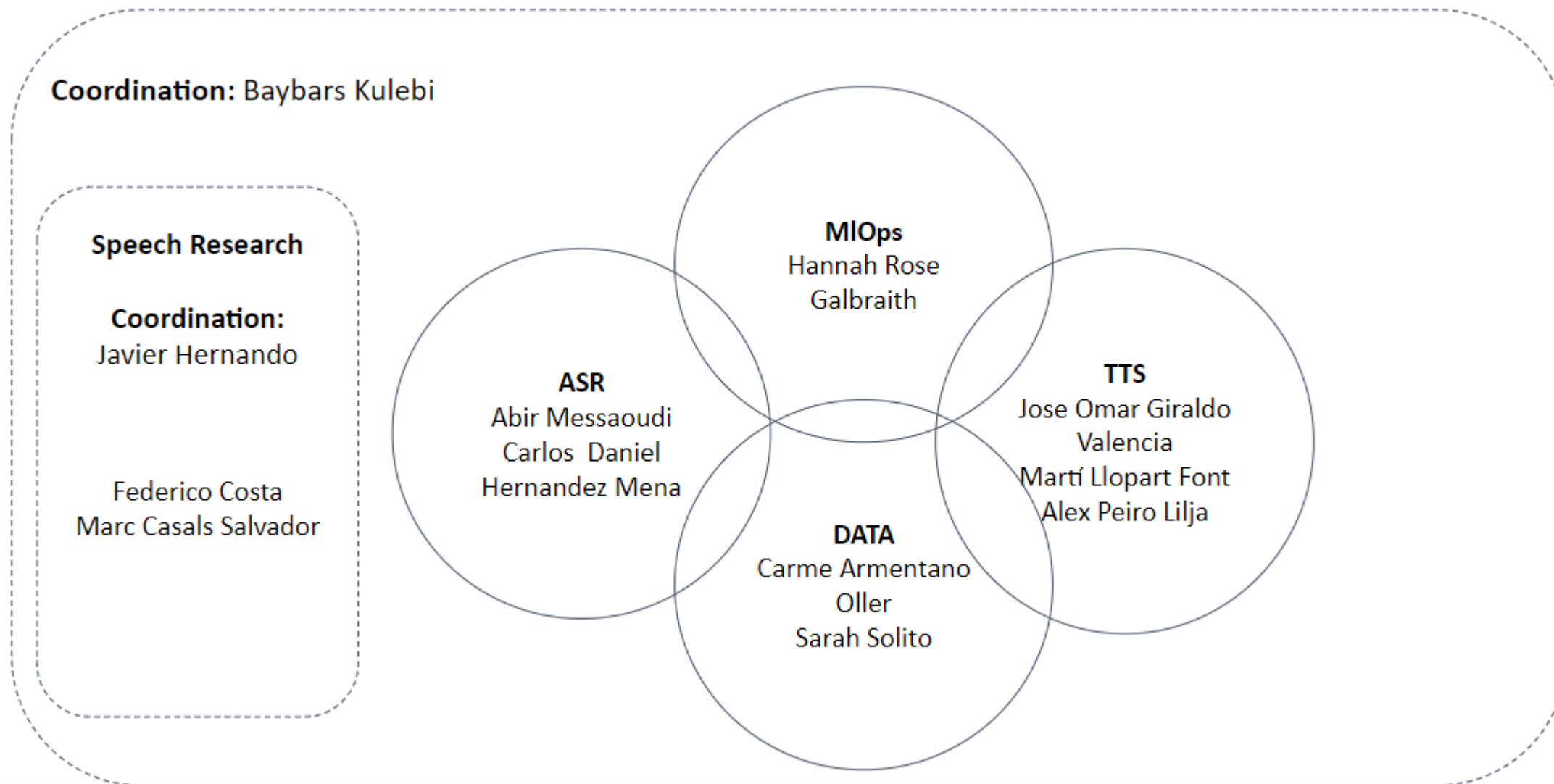
**CIDAI** Centre of Innovation for Data tech and Artificial Intelligence

Who are we?

The project aims to enable Catalan to make a qualitative and quantitative leap in the digital ecosystem.

# Who are we?

**Coordination:** Baybars Kulebi

**Speech Research**

**Coordination:**
Javier Hernando

Federico Costa
Marc Casals Salvador

**MlOps**
Hannah Rose
Galbraith

**ASR**
Abir Messaoudi
Carlos Daniel
Hernandez Mena

**TTS**
Jose Omar Giraldo
Valencia
Martí Llopart Font
Alex Peiro Lilja

**DATA**
Carme Armentano
Oller
Sarah Solito

Who am I ?

Martí Llopart, BSC – Language Technologies Unit – TTS from Speech

BEng Biomedical Engineering – First Class Honours
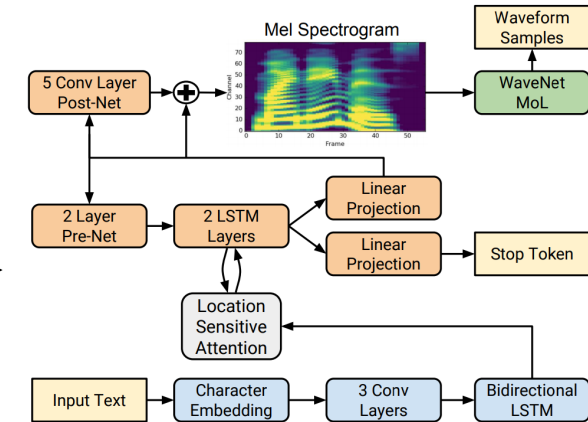
*Queen Mary University of London*

MEM Engineering Project Management

*UPC*

Published in the renowned journal of Biophysical Reviews, cited by Nature

Published by Databricks

The first speech open-source solution to suport the different Catalan varieties - BSC

Where do we come from?



VITS: Conditional Variational Autoencoder with
Adversarial Learning for End-to-End Text-to-Speech

Jaehyeon Kim, Jungil Kong, and Juhee Son

## Summary

What's Matxa?

It is the first multispeaker, multidialectal neural TTS model, and comes together with the vocoder model alVoCat to generate high quality and expressive speech efficiently in four Catalan dialects:
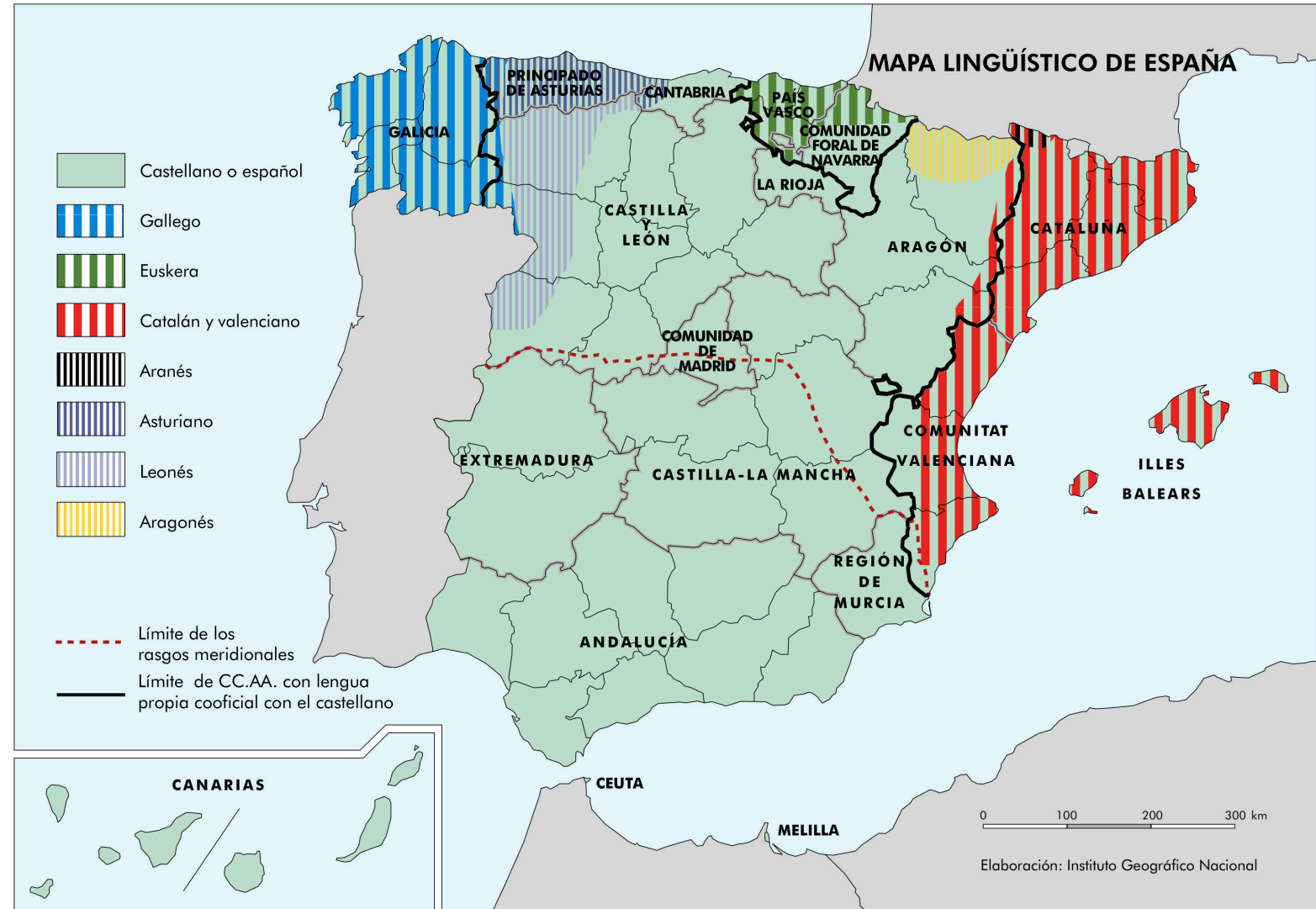
- Balear
- Central
- North-Occidental
- Valencian

The first speech open-source solution to suport the different Catalan varieties - BSC

TTS in minority languages

More efforts are needed to democratise these solutions.

**CIDAI** Centre of Innovation for Data tech and Artificial Intelligence

## What we intend to do



☕ Matxa is a multispeaker multidialect TTS model which uses 🥑 alVoCat as a vocoder. They are based on Matcha-TTS and Vocos architectures.

You can synthesize test sentences below and check the technical details in the "About" tab.

Demo   About   Informació

**Input text**
max 500 characters

m'ha costat molt desenvolupar una veu, i ara que la tinc no estaré en silenc

**Accent**
Models are trained on 4 accents

balear

**Speaker id**
Models are trained on 2 speakers. You can prompt the model using one of these speaker ids.

quim

**Temperature**                    0,2
Temperature

**Length scale**                   0,89
Controls speech pace, larger values for slower pace and smaller values for faster pace

Clear          Submit

🎵 **Matxa + alVoCat**

# Objectives

- Developing Natural-Sounding TTS Synthesis for Catalan Dialects

- Seamless Integration with the administration for Visually Impaired Assistance

- User-Friendly Interface and Model Download

Current

COST

Quantity

Standardized

## What's the problem to be solved?

moltes gràcies ⟹ Phonemizer

_m_'_o_l_t_ə_z_
_ɣ_ɾ_'_a_s_i_ə_s_

Matxa-TTS ⟹ [mel-spectrogram] ⟹ Alvocat ⟹ [Waveform]

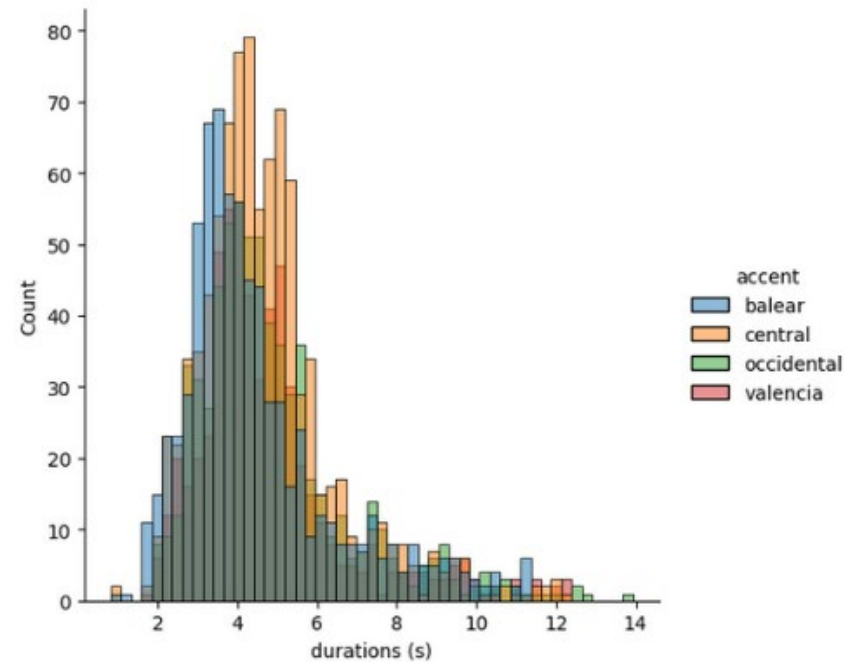Phonemes                     mel-spectrogram                    Waveform

2- The matcha-TTS model converts these phonemes into a mel spectrogram, a visual representation of the spectrum of frequencies of a sound over time.

3- This spectrogram is then fed into our adaptation of the Vocos vocoder, which synthesizes the speech waveform.
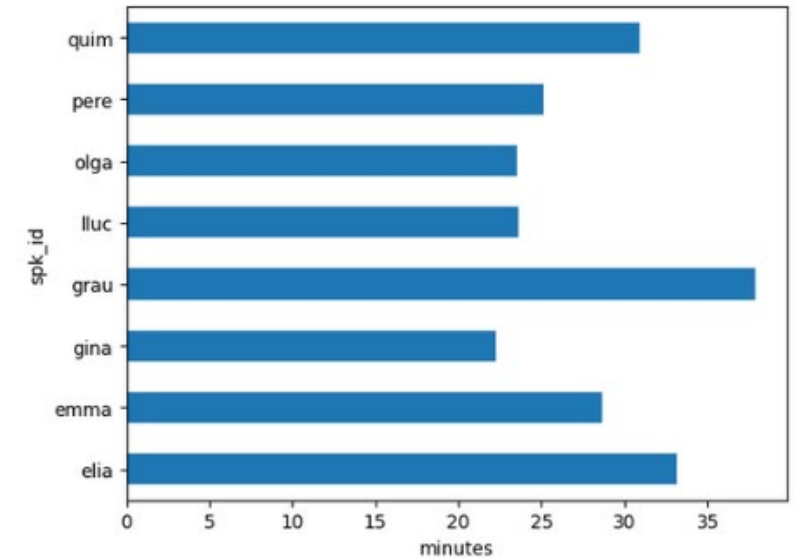
# Multi-accent data:

- 3.5h of studio recordings

- Two voices per accents (female/male)

- On average, 25 min per speaker.



Distribution of utterance durations per accent

accent
- balear
- central
- occidental
- valencia



Total durations per speaker

# • Multi-accent data:

In the following example extracted from our eSpeak, clear phonetic differences can be observed:

Original sentence: "Volem sentir la teva veu perquè és molt important"

**Balearic**: voˈɛn sənˈti lə ˈtevə vˈɛw pərkˈə ˈəz mˈolt importˈant

**Central**: buˈɛm sənˈti lə ˈteβə βˈɛw pərkˈɛ ˈez mˈol impurtˈan

**North-Western**: boˈem senˈti la ˈtewɛ βˈew perkˈe ˈez mˈol importˈan

**Valencian**: voˈem senˈtir la ˈtewa vˈew perkˈe ˈez mˈolt importˈant

**Balearic:**

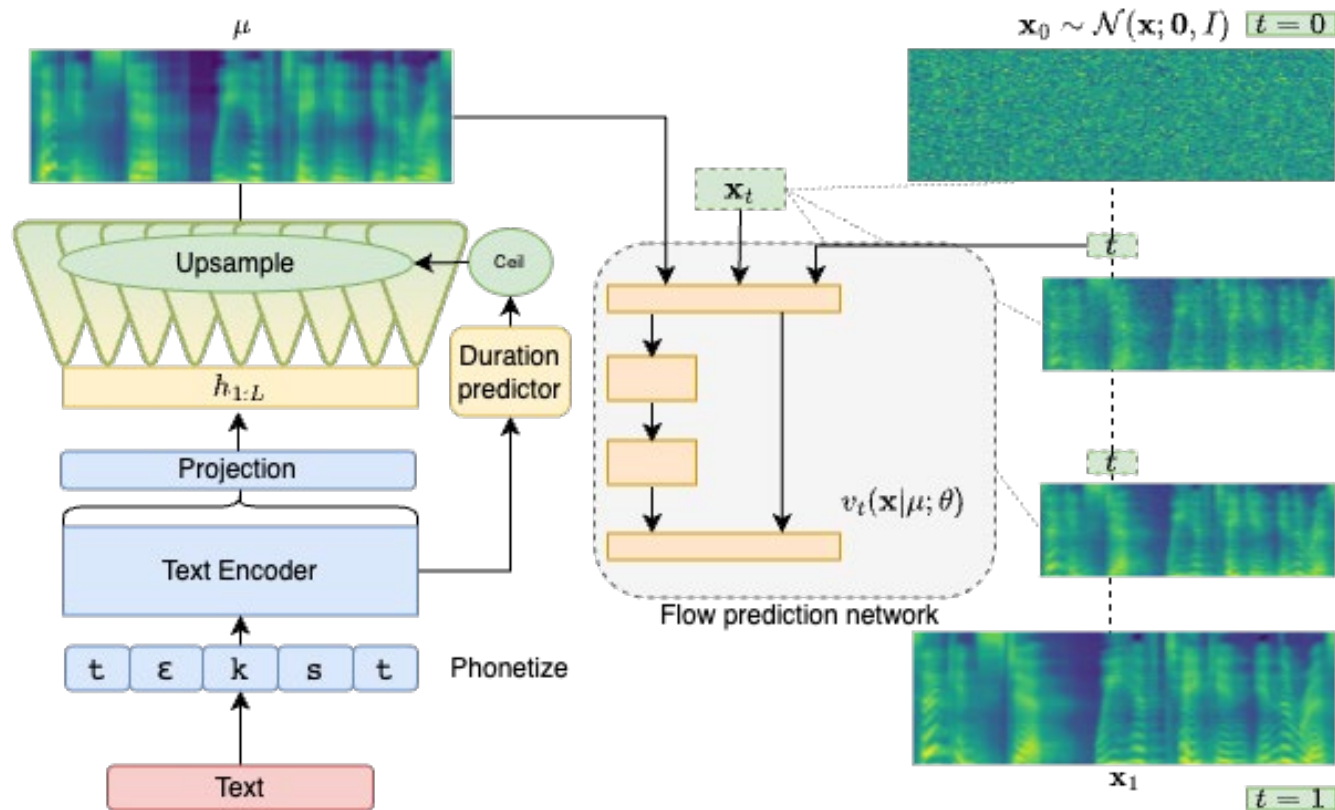**Central:**

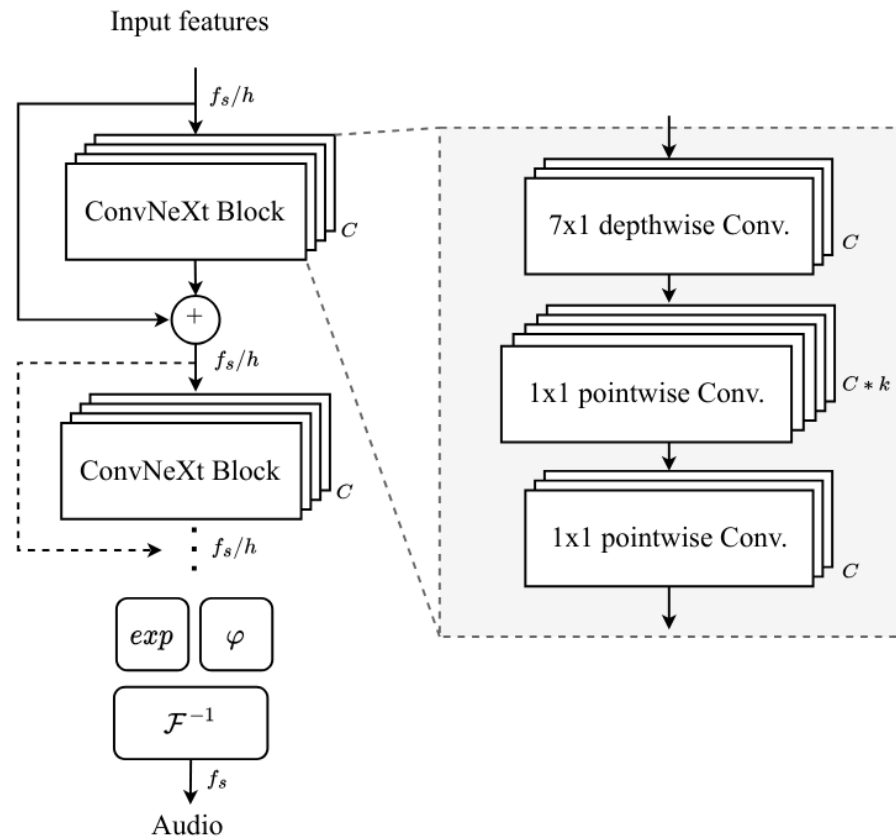**North-Western**:

**Valencian:**

# • Model architectures:

Matxa is based on Matcha-TTS[1], a non-autoregressive encoder-decoder model designed for fast acoustic modelling.

# • Model architectures:

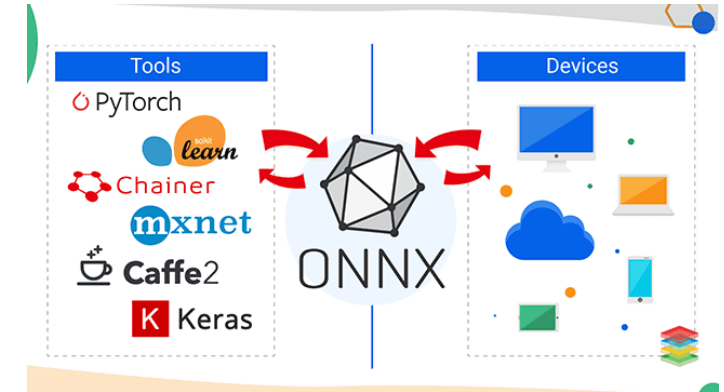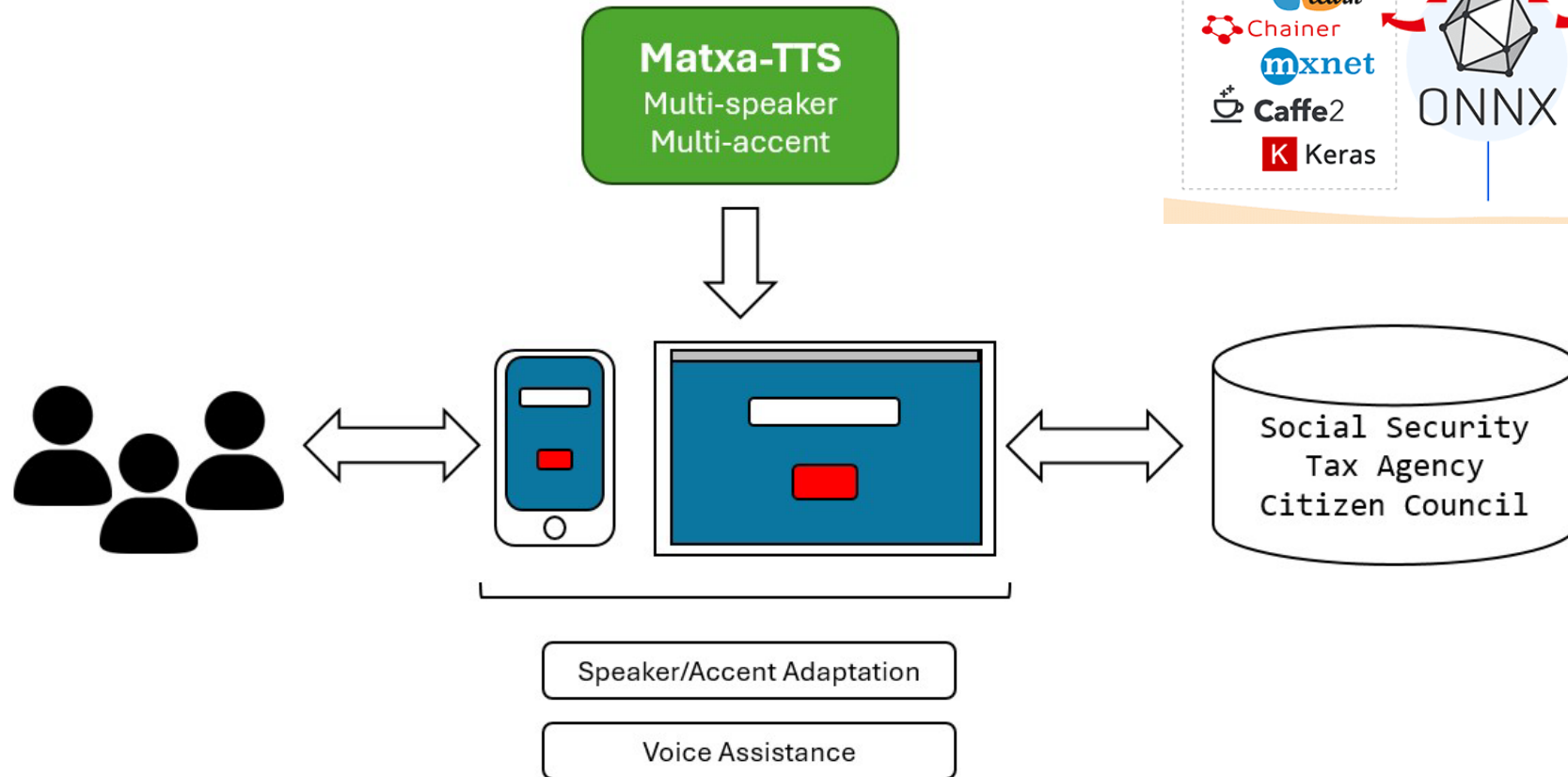AlVoCat is an adapted version of the recently published vocoder named Vocos[3]. It is a fast neural vocoder designed to synthesise audio waveforms from acoustic features.

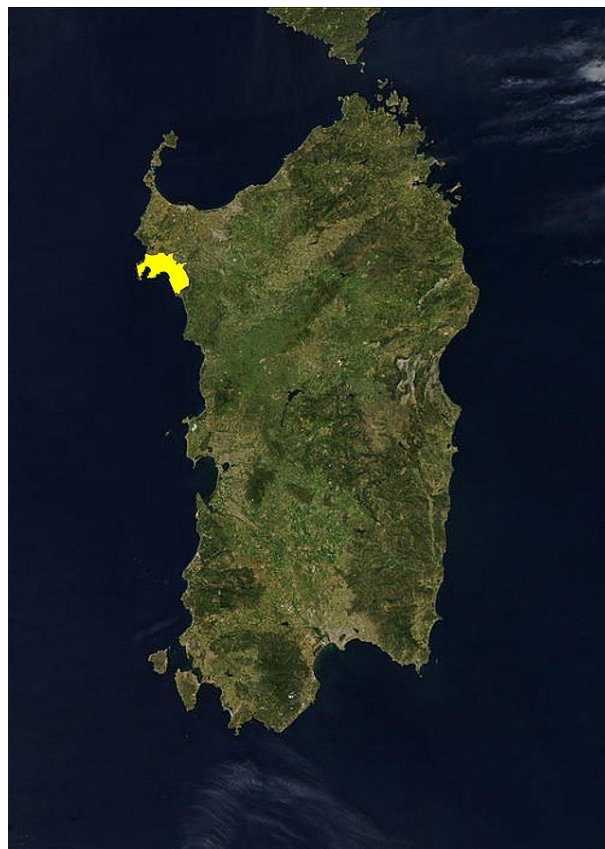The first speech open-source solution to suport the different Catalan varieties - BSC

# • Implementation

The first speech open-source solution to suport the different Catalan varieties - BSC

## Impact of the solution

Septentrional  +

## Impact of the solution

**TEXT_Models**   >

Encoders / Decoders models, foundational, pretrained or fine-tuned

# Llama-VITS: Enhancing TTS Synthesis with Semantic Awareness

Xincan Feng, Akifumi Yoshimoto

Recent advancements in Natural Language Processing (NLP) have seen Large-scale Language Models (LLMs) excel at producing high-quality text for various purposes. Notably, in Text-To-Speech (TTS) systems, the integration of BERT for semantic token generation has underscored the importance of semantic content in producing coherent speech outputs. Despite this, the specific utility of LLMs in enhancing TTS synthesis remains considerably limited. This research introduces an innovative approach, Llama-VITS, which enhances TTS synthesis by enriching the semantic content of text using LLM. Llama-VITS integrates semantic embeddings from Llama2 with the VITS model, a leading end-to-end TTS framework. By leveraging Llama2 for the primary speech synthesis process, our experiments demonstrate that Llama-VITS matches the naturalness of the original VITS (ORI-VITS) and those incorporate BERT (BERT-VITS), on the LJSpeech dataset, a substantial collection of neutral, clear speech. Moreover, our method significantly enhances emotive expressiveness on the EmoV_DB_bea_sem dataset, a curated selection of emotionally consistent speech from the EmoV_DB dataset, highlighting its potential to generate emotive speech.

**Matxa**

The first speech open-source solution to suport the different Catalan varieties - BSC

## Acknowledgments

**The development of LaFresCat dataset, and the neural network models Matxa and alvoCat has been possible thanks to the financing by the Government of Catalonia through the Aina project.**

# Thank you for your

## attention!

# Any requests?

# CIDAI

Centre of Innovation
for Data tech
and Artificial Intelligence

## Socis promotors

Generalitat de Catalunya

Ajuntament de Barcelona

BSC Barcelona Supercomputing Center · Centro Nacional de Supercomputación

CVC Centre de Visió per Computador

eurecat Centre Tecnològic de Catalunya

i2cat

HUAWEI

Microsoft

NTT DATA

SAP

sdg group

UNIVERSITAT POLITÈCNICA DE CATALUNYA BARCELONATECH · Intelligent Data Science and Artificial Intelligence Research Center

IDEAI

## Membres

AMB

ATM Àrea de Barcelona Autoritat del Transport Metropolità

Caixa Enginyers BANCA COOPERATIVA

CETAQUA CENTRO TECNOLÓGICO DEL AGUA

3cat

CRUÏLLA

www.cidai.eu