
Cross Lingual Speaker Adaptation for TTS Applications

— Software Project-2021/2022 —
Rasul, Anna, Claesia, Sharmila

Outline

- ❖ Project Idea
- ❖ Datasets
- ❖ Subtasks
- ❖ Preprocessing Task
- ❖ Web App
- ❖ Timeline
- ❖ References

Project Idea

Grad-TTS modified for multilingual TTS

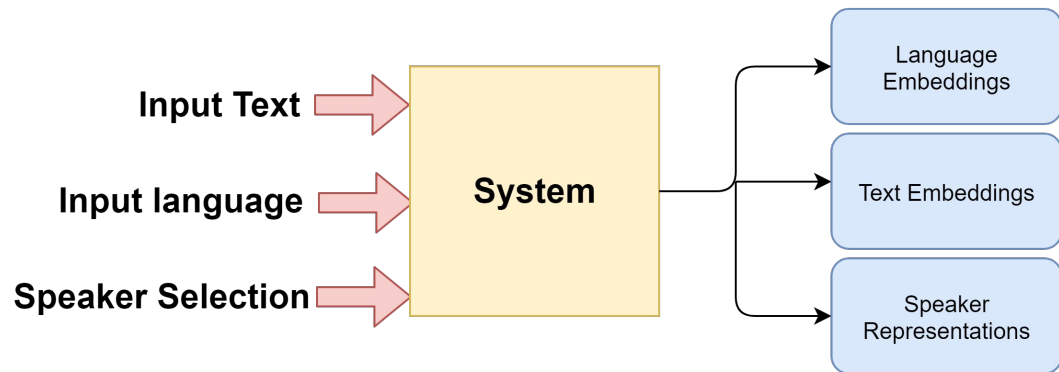
Input Text 

Input language 

Speaker Selection 

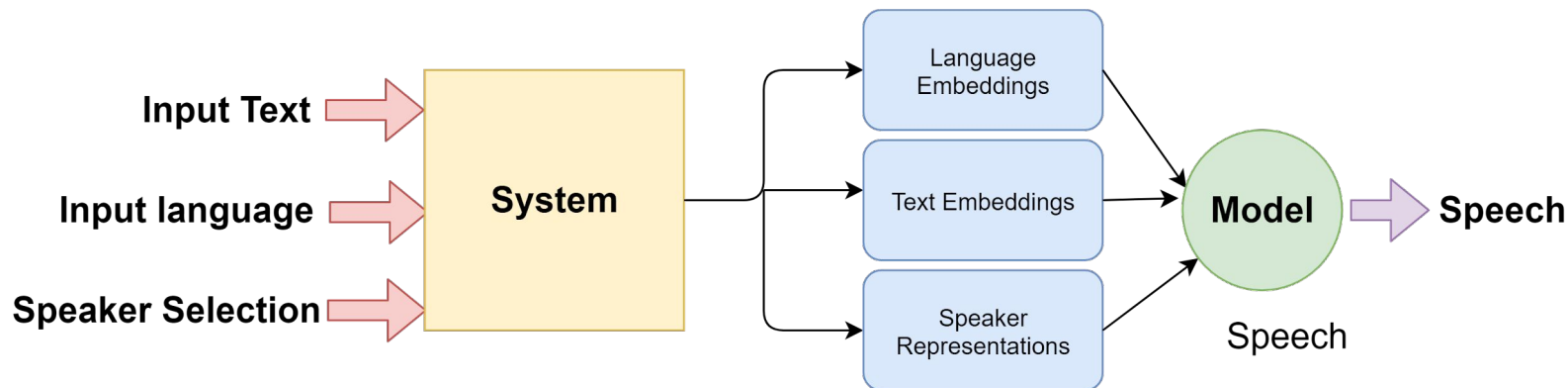
Project Idea

Grad-TTS modified for multilingual TTS



Project Idea

Grad-TTS modified for multilingual TTS



Project Idea

Grad-TTS modified for multilingual TTS

Text_english, english, english embedding , French speaker \Rightarrow TTS \Rightarrow English speech but in French speakers voice

Text_french, french, french embedding , English Speaker \Rightarrow TTS \Rightarrow French Speech but in English Speakers voice

Datasets

- English: LJS Speech dataset (<https://keithito.com/LJ-Speech-Dataset/>)
 - 13,100 short audio clips
 - single speaker
 - passages from non-fiction books
 - length from 1 to 10 seconds
 - total length of ~24 hours
- French: Siwis Speech dataset (<https://infoscience.epfl.ch/record/225946?ln=en>)
 - multiple styles and emphasis
 - about ten hours of speech

Subtasks

Tasks we are working on this week:

- Preprocessing of French corpus (phoneme representation)
- Web Application (simple layout)
- Explore Docker
- Set working environment in Grid5000

Preprocessing Task

From French texts to the sequence of phonemes.

- For each line from SIWIS dataset we will have a sequence of phonemes.
- Encode each phoneme by a unique ID.
- Data for the training:

```
path/to/wav/file.wav | [12, 52, 07, 14, 34, 20, 46, 42] |  
speaker_id | lang_id
```

Web App

Input Elements:

Text box

Dropdown to select speaker,
language

Generate button

Result: Audio

Multilingual Multispeaker TTS

Bonjour!

Language Options

French

English

Language of Text

Speaker

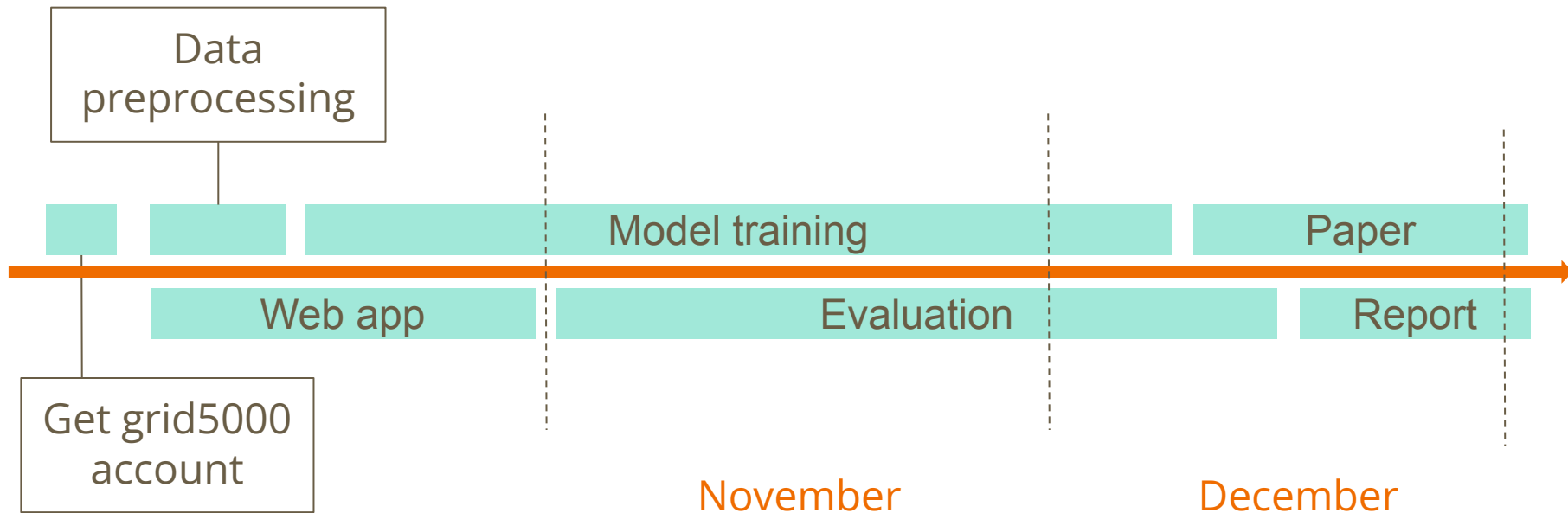
Please provide the text you wish to hear.

synthesize

{{output}}

Claésia Costa, Rasul Dent, Anna Kriukova, Sharmila Upadhyaya Last updated: -

Timeline



References

- Grad-TTS: A Diffusion Probabilistic Model for Text-to-Speech (<https://arxiv.org/pdf/2105.06337.pdf>)
- Glow-TTS: A Generative Flow for Text-to-Speech via Monotonic Alignment Search (<https://arxiv.org/pdf/2005.11129.pdf>)

Thank you ! Any Questions?