




Software Development Project

Multilingual: Text to Speech
Presentation 6



Albert Millert
Shalini Priya
Wenjun Sun
Soklay Heng

Primary goal

Develop a web application that uses Grad-TTS Model for Text to Speech conversion. Languages supported by the TTS converter app are:

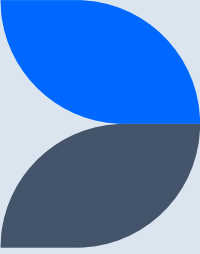
- English
- French



In Progress



Completed



Planning

1

Grid5000
Account

Testbed for
experiment

2

Preprocessing

Generating
phoneme from input
text

3

RNN LSTM
Classifier

For Identifying
language

4

Training

To train a model

5

Deployment

Combine all
modules and deploy
working model on
WebApp

Last time:

1. Language classifier + prediction
2. Server + (ready) modules integration
3. Server + frontend integration
4. English TTS

Completed

- Included EnglishTTS model in the environment
- Evaluation
- Server + (ready) modules integration
- Server + frontend integration
- French model Training

Backend

- vocalization function provided to the server

```
1 def say(sent: str):
2     fout_path = os.path.abspath('./Grad-TTS/out/sample.wav')
3
4     with torch.no_grad():
5         x = torch.LongTensor(intersperse(text_to_sequence(sent, dictionary=cmu), len(symbols)))[None]
6         x_lengths = torch.LongTensor([x.shape[-1]])
7         _, y_dec, _ = generator.forward(x, x_lengths, n_timesteps=TIMESTEPS, temperature=1.5, stoc=False, spk=SPEAKER, length_scale=0.91)
8
9         audio = (vocoder.forward(y_dec).cpu().squeeze().clamp(-1, 1).numpy() * 32768).astype(np.int16)
10
11         write(fout_path, 22050, audio)
12
13     return fout_path
```

Backend

- vocalization function provided to the server
- server handles the output and serves audio file as a response

```
1 from flask import Flask, request, jsonify, send_file
2
3 @app.route("/synthesize", methods=["POST"])
4 @cross_origin(origin="*", headers=["Content-Type", "Authorization"])
5 def synthesize():
6     if request.method == "POST":
7         content = request.json
8         sentence = content["sentence"]
9
10        print(f"Processing: {sentence} in progress...")
11        lang = classify(sentence)
12        print(f"Language: {lang} detected", end="\n\n")
13        print(f"Let's synthesize")
14
15        if lang == "EN":
16            out_path = say(sentence)
17            playsound(out_path)
18        else:
19            default_response = "Sorry, you have to wait for the french model; only english one available"
20            out_path = say(default_response)
21            playsound(out_path)
22
23        return send_file(out_path, mimetype="audio/wav", as_attachment=True, attachment_filename="sample.wav")
24    else:
25        return jsonify({"speech": "Can't touch this"})
```

Backend

- vocalization function provided to the server
- server handles the output and serves audio file as a response
- Server dockerized
 - files copied, environment set up, and server spinned up at the end

```
1 FROM continuumio/miniconda3
2
3 RUN apt-get -y update && apt-get install -y libzbar-dev
4
5 WORKDIR /src
6
7 # copy conda, environment configuration files
8 COPY ./env/ .
9
10 SHELL ["/bin/bash", "--login", "-c"]
11
12 # recreate conda environment
13 RUN conda env create -f environment.yml \
14     && conda init bash \
15     && conda activate tts-env
16
17 # copy language model
18 COPY ./classifierModel .
19
20 # copy project source files
21 COPY ./src/ .
22
23 RUN conda activate tts-env \
24     && cd Grad-TTS/model/monotonic_align \
25     && python setup.py build_ext --inplace \
26     && cd ../../..
27
28 EXPOSE 5000
29
30 CMD conda activate tts-env \
31     && ./server.py
32
```


Backend Demo

Client (English)

HearOut

Generate Text to Speech

Hi , nice to meet you

Let's synthesize

127.0.0.1 - - [13/Jan/2022 08:52:20] "POST /synthesize HTTP/1.1" 200 -

127.0.0.1 - - [13/Jan/2022 08:52:28] "OPTIONS /synthesize HTTP/1.1" 200 -

Processing: Hi , nice to meet you in progress...

Language: EN detected

Client (French)

HearOut

Generate Text to Speech

Bonjour

Let's synthesize

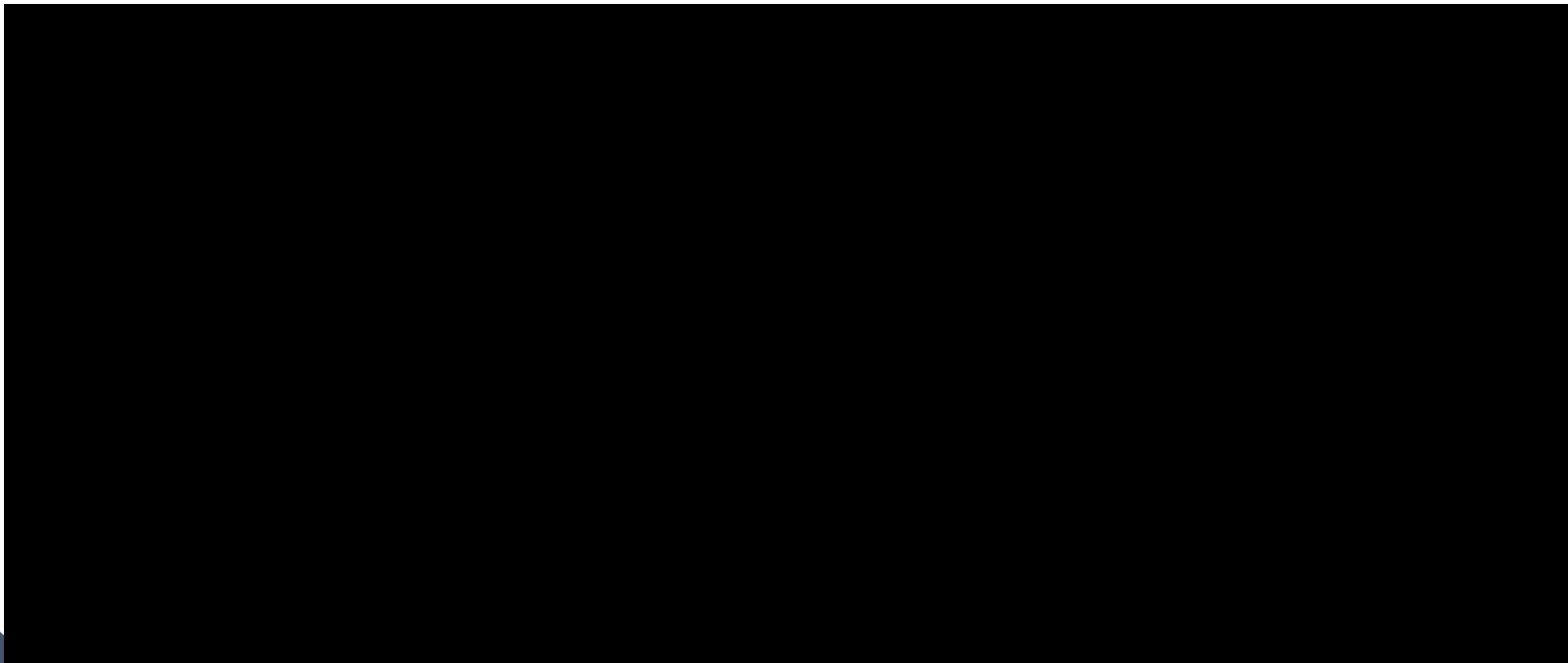
127.0.0.1 - - [13/Jan/2022 08:54:04] "POST /synthesize HTTP/1.1" 200 -

127.0.0.1 - - [13/Jan/2022 08:54:32] "OPTIONS /synthesize HTTP/1.1" 200 -

Processing: Bonjour in progress...

Language: FR detected

Client (Demo)



Evaluation



1. iFLYTEK speech evaluation platform
2. China's largest intelligent voice technology provider
3. Used by multiple schools in China as a scoring platform for spoken English exams

Evaluation Dimension

1. accuracy_score

2. fluency_score

3. Integrity_score

4. Standard_score

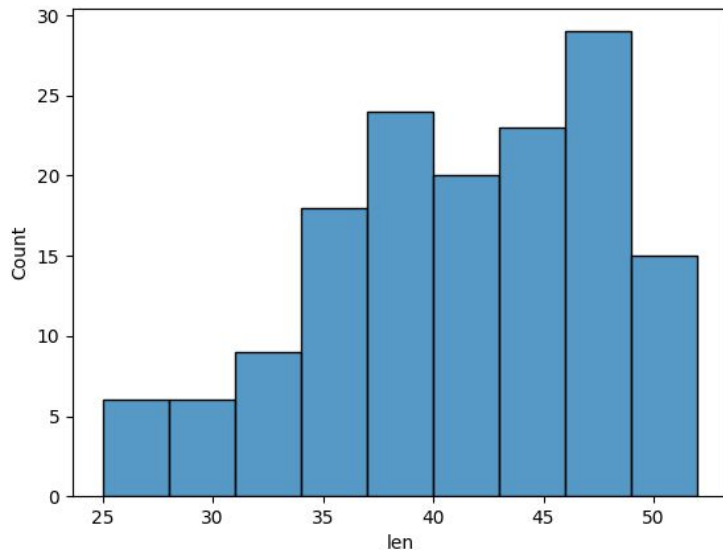
5. Total_score:

$$(0.5 * a_score + 0.3 * f_score + 0.2 * s_score) * i_score / 100$$

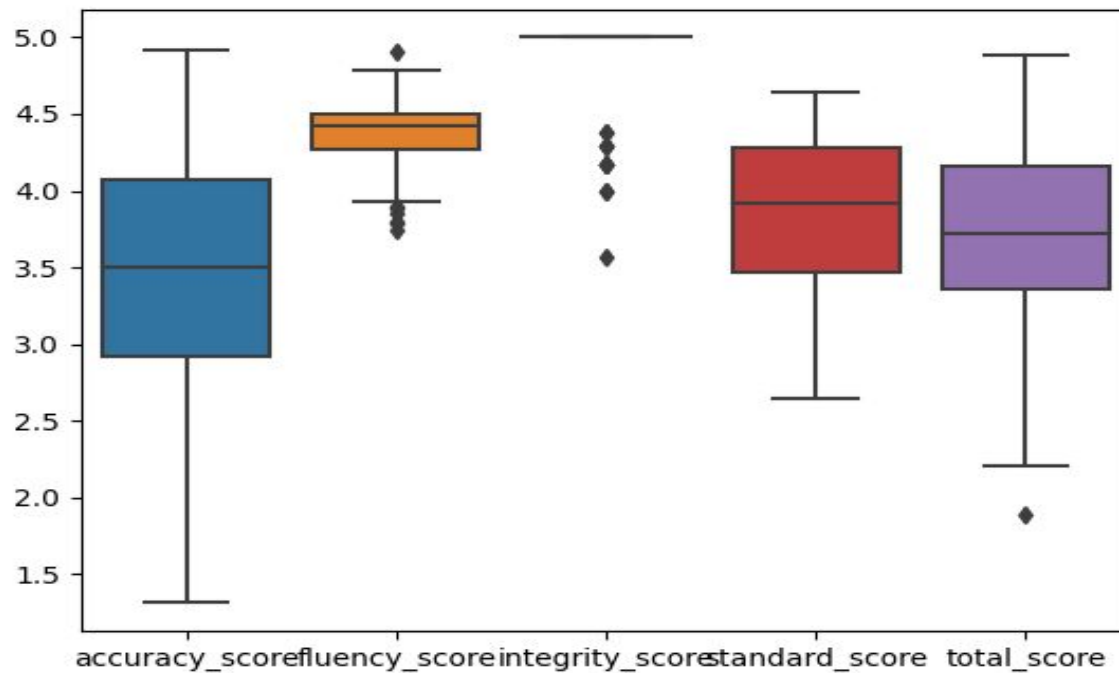
Evaluation Data

1. Corpus: BBC Headlines

2. Amount: 150



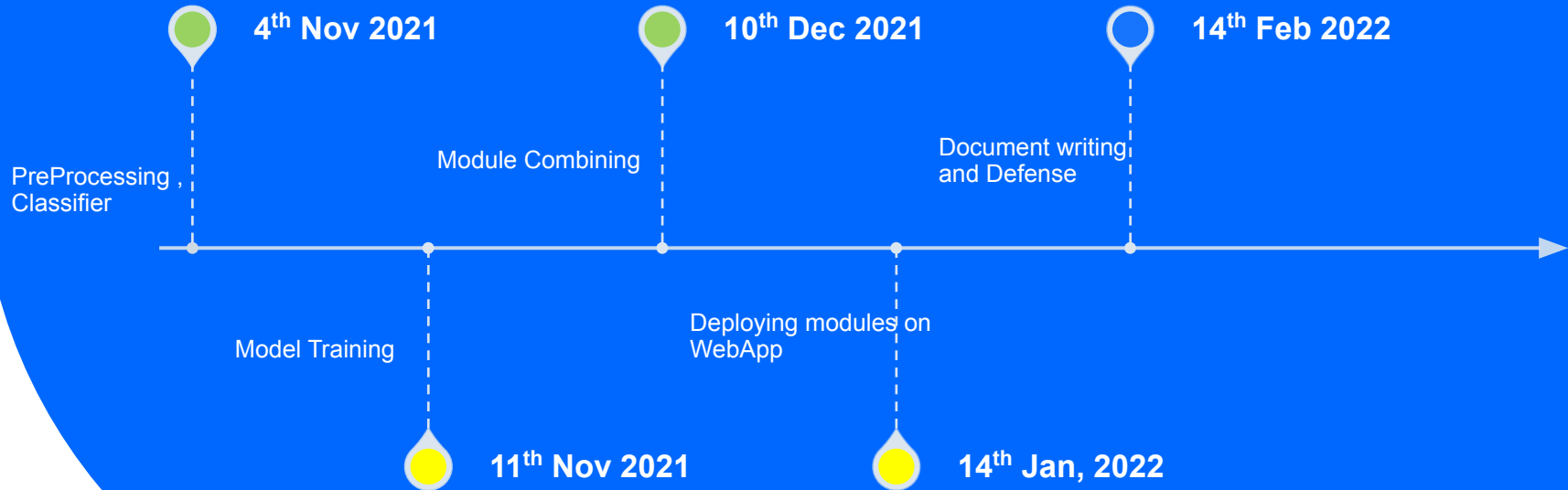
Evaluation Result



Next Time

- Include French model in the environment
- Evaluation
- Server + (ready) modules integration
- Server + frontend integration
- Defense

Timeline





Thank you