

# Projet INF8225 – Modèle text to speech

Hugo Petrilli 2306643  
Antoine Leblanc 2310186

## 1 Introduction

Le développement récent des transformers a révolutionné l'intelligence artificielle dans de nombreux domaines, aussi bien dans les NLP, que la computer vision ou les modèles génératifs. Nous cherchons à savoir si le développement de ses transformers a aussi permis de révolutionner les modèles Text To Speech.

Dans ce projet, nous présenterons tout d'abord l'évolution des modèles text to speech à travers le temps. Nous détaillons les étapes pour transformer un texte en audio et nous expliquerons ensuite en détail le fonctionnement des modèles text to speech les plus récents. Nous implémenterons un modèle de text to speech classique à l'aide de Speech T5 et un modèle de multi speaker et nous observerons les spectrogrammes et les audios de sortie.

## 2 L'histoire du Text to Speech

Le Text To Speech est une tâche qui consiste à transformer un texte en un audio. Pour ce faire, de nombreuses méthodes ont été développées au cours du temps.

### 2.1 Le Vocoder

Le premier pas vers la synthèse vocale électronique est le développement du Voder et du Vocoder dans les années 1930. En 1939, l'ingénieur acoustique et électronique Homer Dudley crée le premier synthétiseur vocal électronique. Cette machine, appelé le Voder est manipulée à l'aide d'un clavier et de pédales qui permettent de moduler les effets sonores et de sortir un son de parole. Homer Dudley développe également le vocoder, qui est un système de codage et de décodage de la parole. Le signal vocal reçu par la machine est divisé en bandes de fréquences, et les caractéristiques de la voix sont extraites à l'aide de filtres. L'utilisation des filtres inverses permet de transformer le signal extrait en un audio. Le vocoder est donc le premier mécanisme permettant de transformer un audio en signal et inversement.

### 2.2 La synthèse à formants

Dans les années 1970, une autre technique pour le text to speech est développée. Il s'agit de la synthèse à formants. C'est une méthode permettant de simuler les caractéristiques

des cordes vocales humaines de manière très technique et de prévoir les sons qui en sortent. Cela permet d'obtenir une voix synthétique.

### 2.3 La synthèse par concaténation

Dans les années 1980, la synthèse par concaténation est développée. Celle-ci repose sur l'enregistrement de petites unités de parole, les phonèmes, qui sont enregistrés. Pour former des phrases entières, il suffit alors de concaténer tous les phonèmes des phrases. Cela améliore la qualité de l'audio.

### 2.4 La synthèse statistique paramétrique

Dans les années 2000 apparaît le modèle de synthèse statistique paramétrique. Il s'agit d'un modèle reposant sur des modèles statistiques, notamment les modèles de Markov pour générer des ondes sonores à partir d'entrées textuelles. Il s'agit donc de modèles text to speech.

### 2.5 L'ajout de réseaux de neurones

Dans les années 2010, des modèles basés sur les réseaux de neurones se développent. Des modèles comme Tacotron chez Google ou WaveNet, chez DeepMind se développent. Ils utilisent des réseaux de neurones profonds avec des CNN et des LSTM. Cela permet d'avoir des voix plus naturelles et fluides et expressives.

### 2.6 Les transformers

Dans les dernières années, avec la découverte des transformers, des modèles comme Transformer TTS ou Tacotron 2 se développent. Ils utilisent des transformers généralement de type Encoder-Decoder et permettent d'augmenter la qualité des audios. C'est ces types de modèles que nous cherchons à étudier dans notre projet.

Les modèles les plus récents sont également basé sur les transformer mais peuvent permettre plus de chose. Ils peuvent synthétiser la voix de quelqu'un à partir d'audios de cette personne, ils peuvent traiter des tâches de multi speaker ou peuvent détecter la langue d'un texte pour fournir un audio avec le bon accent.

### 77 3 Les modèles TTS avec transformer

78 Le premier modèle qui se rapproche le plus du modèle Trans-  
79 former est Transformer TTS proposé par une équipe de Mi-  
80 crosoft en 2019. L'objectif était d'améliorer le modèle Taco-  
81 tron proposé par Google, qui était un modèle Seq2Seq basé  
82 sur des LSTM et de l'attention, en intégrant l'auto-attention  
83 multi-tête. L'objectif de ce modèle, comme tous les autres  
84 modèles TTS, est de prendre une phrase en entrée et de ren-  
85 voyer un spectrogramme. Ce spectrogramme est ensuite con-  
86 verti en audio par un vocodeur.

#### 87 3.1 Le spectrogramme de Mel

88 Le spectrogramme de Mel est un outil très important dans le  
89 TTS puisque c'est lui qui fait le lien entre sortie du modèle et  
90 entrée du vocodeur afin de générer l'audio. Tout d'abord, re-  
91 gardons ce qu'est un spectrogramme.  
92 Un audio est échantillonné (avec un taux d'échantillonnage)  
93 pour pouvoir le numériser et le rendre compréhensible à l'uti-  
94 lisateur, plus le taux est élevé, meilleure la numérisation de  
95 l'audio sera.

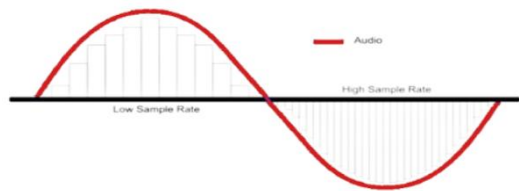


Figure 1 : Audio échantillonné

107 A partir de cet échantillonnage on peut construire le spectre  
108 du signal avec les fréquences calculées par une transformée  
109 de Fourier. Le spectrogramme va ainsi être seulement la  
110 concaténation des différents spectres calculés sur des petites  
111 durées de temps.

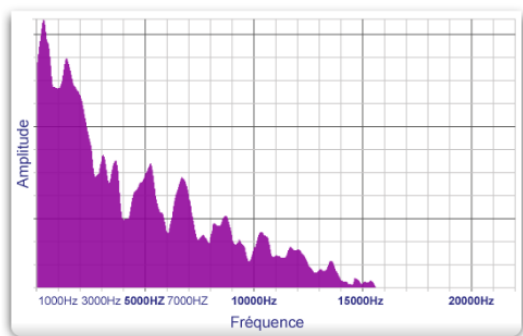


Figure 2 : Spectre du signal

128 Enfin, la technologie TTS utilise un autre type de  
129 spectrogramme, qui est le spectrogramme de Mel. Ce  
130 spectrogramme a été créé car l'oreille humaine détecte mieux  
131 les sons de basse fréquence que les hautes. Pour construire ce  
132 spectrogramme, on vient seulement filtrer chaque spectre et

133 projeter les fréquences sur le l'échelle de Mel (en Hz) qui est  
134 plus adapté pour l'oreille humaine.

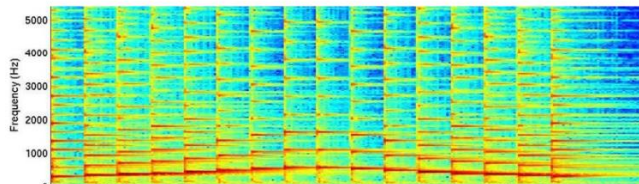


Figure 3 : Spectrogramme de Mel

143 Ainsi, on comprend que les sorties de notre Transformer vont  
144 être les différentes valeurs des fréquences pour un instant  
145 donné (un instant est égale à une très courte période de  
146 temps), et cela pour le nombre d'instant souhaité. On peut  
147 aussi noter que l'on a ici un hyperparamètre lié à notre  
148 modèle qui est lié à la qualité de l'audio : le taux  
149 d'échantillonnage. Un modèle est entraîné sur un taux  
150 d'échantillonnage unique et ne pourra être utilisé pour du fine  
151 tuning ou de l'inférence avec un taux d'échantillonnage  
152 différent.

154 Après cela, le spectrogramme en sortie passe dans un vocoder  
155 qui transforme l'image en audio. Le vocodeur utilisé est  
156 WaveNet qui est réseau type CNN.

#### 157 3.2 L'architecture de Transformer TTS

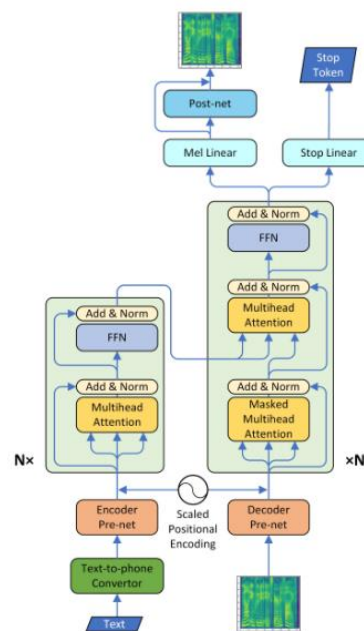


Figure 4 : Architecture de Transformer TTS

185 Le transformeur TTS reçoit un texte, et renvoie le spectro-  
186 gramme de Mel de l'audio correspondant. Il est formé autour  
187 d'un transformeur de type Encoder-Decoder. L'architecture  
188 de l'encodeur et du décodeur est similaire à celle des enco-  
189 deurs et décodeurs dans les modèles NLP. On observe un

190 plongement du texte et un plongement positionnel avant  
191 d'envoyer les données dans l'encodeur. Une fois dans l'en-  
192 codeur, l'attention multi tête est calculée puis dans le déco-  
193 der, il y a calcul de l'attention multi-tête masquée et de la  
194 cross attention. En sortie du décodeur, il y a une couche li-  
195 néaire.

196  
197 Il y a cependant des différences notables entre les modèles de  
198 NLP et le transformer TTS. Le texte est converti en phonème  
199 au début du modèle. Les phonèmes sont la manière de pro-  
200 noncer les mots. Cela permet d'être plus précis sur le texte à  
201 convertir en audio.

202  
203 De plus, le décodeur prend et renvoie des images de la forme  
204 du spectrogramme de Mel. Cela signifie que les blocs Deco-  
205 der Pre Net et Decoder Post Net sont des blocs de pré traite-  
206 ment et post traitement de l'image. De manière plus précise,  
207 Decoder Pre Net prend en entrée le spectrogramme de Mel de  
208 l'input donné par le décodeur et est un réseau de neurones de  
209 3 réseaux connectés, avec une activation ReLu. Le bloc Mel  
210 linear et Decoder Post Net est composé de deux modules. Le  
211 premier est un layer linear avec comme entrée l'output du dé-  
212 codeur, qui prédit le spectrogramme de Mel. Le deuxième  
213 bloc est composé de 5 layers de convolution 1D pour affiner  
214 la prédiction du spectrogramme de Mel.

215  
216 Pour résumer, le modèle Transformer TTS prend en entrée  
217 un audio qu'il convertit en phonème et qu'il envoie dans le  
218 transformer de type encodeur/décodeur et qui est centré sur le  
219 calcul de l'attention multi-tête. Le modèle prédit alors le  
220 spectrogramme de Mel du texte en entrée. Après cela, il faut  
221 utiliser un vocoder pur transformer le spectrogramme de Mel  
222 en audio.

## 223 4 Les modèles implémentées

224 Pour ce projet, nous avons implémentées différents modèles.

### 225 4.1 Speech T5

226 Le modèle Speech T5 est un modèle développé par Microsoft  
227 et inspiré du modèle T5 (Text To Text Transfer Transformer).  
228 Le modèle T5 est un modèle de traitement de texte reposant  
229 sur un transformer de type encodeur-décodeur. Il repose éga-  
230 lement sur le principe d'attention. Speech T5 est donc un type  
231 de modèle text to speech utilisant les transformers. Il peut gé-  
232 rer tout type de tâches lié à la parole, aussi bien le text to  
233 speech que la reconnaissance vocale ou la traduction vocale.  
234 Dans ce projet, nous nous intéresserons à l'utilisation text to  
235 speech de Speech T5.

236  
237  
238  
239  
240  
241  
242  
243  
244  
245

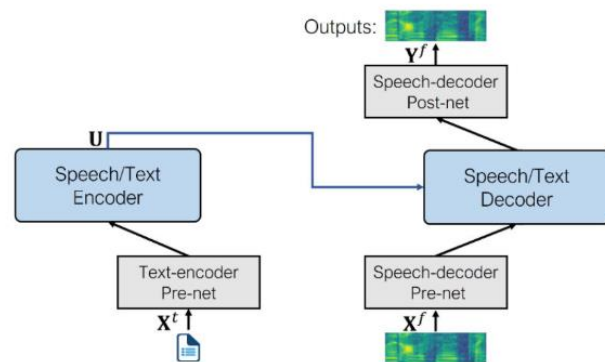


Figure 5 : Architecture du modèle TTS de Speech T5

La version text to speech de Speech T5 est similaire, pour ne pas dire quasi identique à celle du transformer TTS. C'est ce modèle que nous allons implémenter et tester dans la suite du projet.

### 266 4.2 Speech T5 multi speaker

267 Un aspect intéressant des modèles Text-To-Speech est de  
268 pouvoir générer des audios mais pour plusieurs voix, accents  
269 ou langues différents. Le Transformer TTS est un premier  
270 modèle simple qui ne peut générer qu'une seule voix, accent  
271 ou langue à la fois. Pour pouvoir générer des variantes il faut  
272 récupérer le modèle pré-entraîné et le réentraîner sur la va-  
273 riant de notre choix.

274 Concernant le modèle SpeechT5 il est possible avec un  
275 seul modèle de générer différentes voix, accents ou même  
276 langue. La technologie fonctionne en récupérant un plonge-  
277 ment de la voix de la personne qui parle et de concaténer ce  
278 plongement à la sortie et de passer ce nouveau vecteur dans  
279 une couche linéaire pour générer un nouveau spectrogramme  
280 de Mel.

281

### 282 4.3 Entraînement de notre voix

283 Nous nous sommes appuyés sur le tutoriel offert par Hugging  
284 Face sur SpeechT5 qui propose l'entraînement du modèle sur  
285 une base de données d'audio néerlandais. Nous allons plutôt  
286 essayer de faire du voice cloning avec la voix d'A. Leblanc,  
287 auteur de ce rapport. Pour constituer la base de données, An-  
288 toine a effectué 116 enregistrements entre 5 et 8 secondes de  
289 sa voix à partir de phrases générées par ChatGPT. Il s'est en-  
290 registé à l'aide de son téléphone personnel et chez lui avec  
291 le bruit environnant qu'il pourrait y avoir.

292 Notre objectif est de faire du voice cloning avec du zero  
293 shot, few shots et un entraînement un peu plus conséquent.  
294 Le dataset a été divisé en deux : un petit contenant 16 données  
295 et un plus grand contenant 102 données.

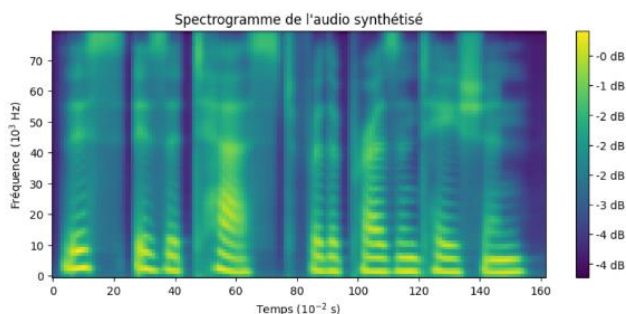
## 296 5 Résultats de nos modèles

297 Nos tests effectués sont disponibles au lien github suivant :  
298 <https://github.com/AntoineLeblancFr/INF8225/tree/main>

## 5.1 Le modèle text to speech classique (modèle pré-entraîné)

Pour créer notre modèle de text to speech, nous utilisons des modèles pré-entraînés. Nous utilisons la librairie transformers, avec des modules comme SpeechT5Processor, SpeechT5ForTextToSpeech ou SpeechT5HifiGan, ce qui permet d'avoir un modèle ou un vocoder. Nous utilisons également des datasets importés.

Après, nous choisissons le texte à transformer en audio. Le modèle nous renvoie un spectrogramme.



En utilisant le vocoder, nous transformons alors ce spectrogramme en audio, que nous pouvons enregistrer.

## 5.2 Le modèle text to speech multi-speaker

Pour le multi-speaker nous réutilisons le même modèle, seulement il faut générer pour chaque exemple un plongement permettant de représenter la voix du speaker pour cet exemple. Le modèle se charge de faire la concaténation seule. C'est la même chose pour l'inférence, il suffit de transmettre un plongement de la voix que l'on désire.

## 5.3 Entraînement de notre voix

Les hyperparamètres choisis pour l'entraînement de notre voix sont disponibles dans notre notebook présent à notre lien github.

On peut voir, tout d'abord, que pour le zero shot la voix est compréhensible avec un accent anglais. Cela s'explique par le fait que le modèle a été pré-entraîné sur un dataset d'audio en anglais. Ainsi, le zero-shot fonctionne si on est s'exprime dans la bonne langue.

Pour le few shot, l'audio ressort presque inaudible. Cela peut s'expliquer par le fait que l'on entraîne le modèle sur un petit dataset français et que l'on a trop peu d'exemples pour pouvoir obtenir un résultat convenable.

Enfin, pour le fine tuning sur le plus gros dataset, l'audio reste de mauvaise qualité mais on discerne tout de même les sonorités les plus importantes et on reconnaît un peu la voix d'A. Leblanc. Tous les audios sont disponibles sur le notebook.

## 6 Conclusion

En conclusion, nous pouvons dire que le voice cloning n'a pas été réussi dans notre cas car premièrement le modèle ne

connaissait pas le français et deuxièmement car on ne disposait pas d'assez d'exemples d'audio. En effet, avec le gros dataset on arrive à un peu plus de 10 minutes d'audio alors que pour faire du bon voice cloning il faut au minimum entre 20 et 25 minutes d'audio, d'après la littérature. Une idée qui permettrait de rendre notre voice cloning meilleur même avec si peu d'exemples, serait d'entraîner le modèle sur une grande base de données d'audio français et de venir ajuster à la fin les paramètres du modèle avec nos quelques exemples.

On comprend maintenant pourquoi il est facile de copier la voix de politicien comme on a pu le voir sur les réseaux sociaux, puisque l'on dispose d'un grand nombre de données de leur voix via leur discours avec aucune musique de fond et aucun bruit. A contrario, il est plus compliqué d'effectuer du voice cloning sur quelqu'un dont on a pas assez d'audio, même si aujourd'hui des techniques sont développées en utilisant la technologie auto-encodeur pour faire des représentations plus efficaces des voix et ainsi n'avoir besoin que de peu d'exemples pour faire du voice cloning.

Cette pratique soulève aussi des questions car on voit le danger et les arnaques que l'on peut faire grâce à cette technologie. Cependant elle peut aussi permettre à des personnes qui ont perdu leur voix de la retrouver avec un ancien audio d'eux.

## Références

- [Li et al., 2019] Naihan Li, Shujie Liu, Yanqing Liu, Sheng Zhao, Ming Liu and Ming Zhou. *Neural speech synthesis with transformer network* Guizhou, China, 2019.
- [Ao et al., 2022] Junyi Ao, Rui Wang, Long Zhou, Chengyi Wang, Shuo Ren, Yu Wu, Shujie Liu, Tom Ko, Qing Li, Yu Zhang, Zhihua Wei, Yao Qian, Jinyu Li and Furu Wei. *SpeechT5: Unified-Modal Encoder-Decoder Pre-Training for Spoken Language Processing* Hong Kong, 2022.