

# Test du Stage ASR : Optimisation fine basée sur des adaptateurs (Adapter-Based Fine Tuning) pour les langues à faibles ressources

Proposé par :

**NGOM CHRISTINE CARELLE ANGE**

ngom24christine03@gmail.com

+237 651763807 / 697218936

Yaoundé/Douala, Cameroun

26 novembre 2025

## Résumé

Ce rapport présente ma proposition au ASR Fellowship Challenge organisé par Digital Umuganda et IndabaX Cameroun. L'objectif principal de ce défi est d'améliorer les performances d'un modèle de reconnaissance vocale automatique (ASR) pour le Kinyarwanda, une langue nationale du Rwanda sous-représentée en termes de ressources numériques. La particularité de ce challenge réside dans l'utilisation d'adaptateurs : de petits modules entraînables insérés dans un modèle pré-entraîné gelé (ses poids ne changent pas). Cette approche permet d'adapter efficacement le modèle à une nouvelle langue tout en préservant ses capacités générales et en réduisant considérablement les coûts computationnels. Le présent rapport détaille l'architecture utilisée, la méthodologie d'entraînement, les résultats obtenus et fournit des instructions complètes pour reproduire l'expérimentation.

## Introduction

Dans le cadre du ASR Fellowship Challenge organisé par Digital Umuganda en collaboration avec IndabaX Cameroun, ce travail présente une démarche d'amélioration des performances de modèles de reconnaissance vocale automatique (ASR) pour des langues sous-représentées numériquement. Le défi consiste à adapter un modèle de base pré-entraîné à l'aide de petits modules spécialisés appelés **adaptateurs**, tout en respectant les contraintes suivantes : ne pas modifier les poids du modèle de base, n'utiliser que l'ensemble d'entraînement fourni, réaliser l'implémentation en Python avec PyTorch. L'approche adopté permet d'exploiter la puissance d'un modèle existant tout en l'adaptant à la spécificité linguistique de la langue ciblée sans qu'il n'y ait de perte des connaissances préalablement acquises. L'objectif est donc d'intégrer des adaptateurs à un modèle ASR existant, d'entraîner uniquement ces petites couches additionnelles, puis d'évaluer l'amélioration obtenue en termes de précision, de transcription. L'ensemble du processus repose sur une méthodologie organisée en étapes : préparation de l'environnement, téléchargement du dataset, génération des transcriptions du

modèle de base, entraînement des adaptateurs, transcription avec le modèle affiné, puis calcul final du Word Error Rate (WER) afin de mesurer l'impact réel de l'adaptation. Pour rendre ce travail reproductible, pédagogique et compréhensible même pour un débutant, tous les scripts sont documentés afin d'expliquer clairement ce qu'ils font et ce qu'il se passe lorsqu'on les exécute.

## Méthodologie

La démarche méthodologique suivi pour réaliser ce travail est constituer de plusieurs étapes significatives et nécessaire à son bon déroulement. Le pipeline suivant décrit toutes les étapes techniques permettant de reproduire le projet depuis zéro, incluant la création du répertoire, de l'environnement virtuel, du dossier `src`, ainsi que l'exécution de chaque script avec les détails de leur fonctionnement et des fichiers générés.

### 1 Création et organisation du projet

#### 1.1 Création du dossier principal et déplacement à l'intérieur

Créez un dossier que vous nommerez `ASR_test` puis déplacez vous à l'intérieur. En ligne de commande il vous suffira de vous diriger vers le répertoire où vous voulez créer votre projet ( par exemple : `mon_repertoire`) et entrer les commandes suivantes dans votre terminal.

```
1 mkdir ASR_test
2 cd ASR_test
```

Cette étape crée le répertoire contenant tout le projet puis place l'utilisateur à l'intérieur pour exécuter la suite.

#### 1.2 Création de la structure interne du projet

Une fois dans le dossier `ASR_test`, créez les répertoires qui suivent. Dans votre terminal il suffira d'entrer les commandes suivantes :

```
1 mkdir ASR_code
2 mkdir ASR_code/src
3 mkdir ASR_code/src/utils
4 mkdir ASR_code/src/training
5 mkdir ASR_code/src/inference
6 mkdir ASR_code/src/evaluation
7 mkdir ASR_code/src/models
8
9 mkdir weights
10 mkdir weights/adapters
```

Ces dossiers contiendront respectivement :

- le dataset téléchargé,
- les scripts source,
- les poids entraînés,
- les transcriptions générées.

Vous comprendrez bien plus tard.

## 2 Création de l'environnement virtuel

Dirigez vous dans le dossier ASR\_code :

```
1 cd ASR_code
```

Créer l'environement virtuel du projet et l'activer. Il suffit de saisir la commande suivante :

```
1 python3 -m venv venv
2 source venv/bin/activate      # Linux/MacOS
3 venv\Scripts\activate        # Windows
```

Si vous n'avez pas venv vous pouvez installer virtualenv à partir de la commande

```
1 pip install virtualenv      # Windows
```

Ce dossier contiendra les dépendances Python isolées du système. Il ne vous reste plus qu'à installer toutes les dépendances nécessaires de notre projet en exécutant la commande suivante dans votre terminal :

```
1 pip install torch datasets transformers jiwer huggingface_hub
```

## 3 Téléchargement du dataset

Écrivez le code Python suivant dans le fichier nommé data\_utils.py du repertoire ASR\_code/src/utils et l'exécuter :

```
1 from huggingface_hub import snapshot_download
2
3 def download_dataset():
4     snapshot_download(
5         repo_id="DigitalUmuganda/ASR_Fellowship_Challenge_Dataset",
6         repo_type="dataset",
7         local_dir="data"
8     )
9
10 if __name__ == "__main__":
11     download_dataset()
```

Ce script python va :

- créer automatiquement : ASR\_code/data/
- télécharger automatiquement le dataset du challenge,
- le placer dans ASR\_code/data/,
- vérifier que tous les fichiers audio et transcriptions sont présents.

## 4 Choix et inférence du modèle de base

Le modèle de base choisi est **facebook/wav2vec2-base** pour les raisons suivantes :

- Il a été entraîné sur des milliers d'heures d'audio multilingue.
- Il fonctionne très bien même avec peu de données.
- Compatible nativement avec PyTorch + Transformers.

- C'est celui utilisé dans presque tous les défis ASR low-resource.
- et surtout car Digital Umuganda utilise Wav2Vec2 dans ses projets Afrivoice. Donc c'est cohérent avec l'écosystème du challenge.

Dans le fichier `inference_base.py` on charge le modèle de base qui a été pré-entraîné, on applique ce modèle aux données de test de notre jeu de données, on en tire la transcription obtenu en inférant sur le modèle sans les adaptateurs que l'on stocke dans `base_transcriptions.txt`.

Ecrivez le code suivant dans le fichier `src/inference/inference_base.py` du répertoire `ASR_code` :

```

1 import torchaudio
2 from transformers import Wav2Vec2Processor, Wav2Vec2ForCTC
3 import os
4
5 def transcribe_audio(model, processor, path):
6     speech, sr = torchaudio.load(path)
7     speech = torchaudio.functional.resample(speech, sr, 16000).squeeze()
8     inputs = processor(speech, sampling_rate=16000, return_tensors="pt")
9     logits = model(inputs.input_values).logits
10    pred_ids = logits.argmax(-1)
11    return processor.decode(pred_ids[0])
12
13 def run_inference(test_dir, output_file):
14     processor = Wav2Vec2Processor.from_pretrained("facebook/wav2vec2-base")
15     model = Wav2Vec2ForCTC.from_pretrained("facebook/wav2vec2-base")
16
17     with open(output_file, "w", encoding="utf8") as f:
18         for file in sorted(os.listdir(test_dir)):
19             if file.endswith(".wav"):
20                 text = transcribe_audio(model, processor, os.path.join(test_dir,
21                             file))
22                 f.write(f"{file}\t{text}\n")
23
24 if __name__ == "__main__":
25     run_inference("data/test", "base_transcriptions.txt")

```

Puis l'exécuter.

## 5 Entraînement des adaptateurs

### 5.1 Définition du module adaptateur

Créez un fichier `src/models/adapters.py` et y écrire le code suivant :

```

1 import torch
2 import torch.nn as nn
3
4 class BottleneckAdapter(nn.Module):
5     def __init__(self, hidden_size=768, bottleneck=64):
6         super().__init__()
7         self.down = nn.Linear(hidden_size, bottleneck)
8         self.activation = nn.ReLU()
9         self.up = nn.Linear(bottleneck, hidden_size)
10
11     def forward(self, x):
12         return x + self.up(self.activation(self.down(x)))

```

C'est dans ce fichier que l'on défini la structure interne des adaptateurs ( nombre de noeuds, fonction d'activation, etc.).

## 5.2 Script d'entraînement

L'entraînement icic fait référence à l'entraînement des adaptateurs sur les données du jeu de données qui nous est fournis. On commence par charger le modèle, geler ses poids et y ajouter les adaptateurs. Une fois cela effectué on entraîne uniquement les adaptateurs puis on sauvegarde les poids entraînés des adaptateurs.

Créer le fichier `src/training/train.py` et recopier le code suivant :

```
1 import torch
2 from torch.utils.data import DataLoader
3 from transformers import Wav2Vec2ForCTC, Wav2Vec2Processor
4 from datasets import load_dataset
5 from src.models.adapters import BottleneckAdapter
6
7 def freeze_base_model(model):
8     for p in model.parameters():
9         p.requires_grad = False
10
11 def insert_adapters(model, adapter_class):
12     for layer in model.wav2vec2.encoder.layers:
13         layer.adapter = adapter_class()
14     return model
15
16 def train():
17     processor = Wav2Vec2Processor.from_pretrained("facebook/wav2vec2-base")
18     model = Wav2Vec2ForCTC.from_pretrained("facebook/wav2vec2-base")
19
20     freeze_base_model(model)
21     model = insert_adapters(model, BottleneckAdapter)
22
23     train_data = load_dataset("audiofolder", data_dir="data/train")
24
25     def preprocess(batch):
26         batch["input_values"] = processor(batch["audio"]["array"], sampling_rate
27                                         =16000).input_values[0]
28         batch["labels"] = processor.tokenizer(batch["text"]).input_ids
29         return batch
30
31     train_data = train_data.map(preprocess)
32
33     loader = DataLoader(train_data["train"], batch_size=4, shuffle=True)
34
35     optimizer = torch.optim.AdamW(
36         [p for p in model.parameters() if p.requires_grad], lr=1e-4
37     )
38
39     model.train()
40     for epoch in range(5):
41         for batch in loader:
42             optimizer.zero_grad()
43             out = model(input_values=torch.tensor(batch["input_values"]),
44                         unsqueeze(0),
45                         labels=torch.tensor(batch["labels"]).unsqueeze(0))
```

```

44         loss = out.loss
45         loss.backward()
46         optimizer.step()
47
48     torch.save(model.state_dict(), "weights/adapters/adapter_weights.pth")
49
50 if __name__ == "__main__":
51     train()

```

## 6 Inférence avec les adaptateurs entraînés

Il s'agit d'inférer le modèle fine-tuné obtenu après l'entraînement des adaptateurs.

Dans le répertoire `src/inference/`, créer un fichier `inference_finetuned.py` et y recopier le code suivant :

```

1 import torch
2 from transformers import Wav2Vec2Processor, Wav2Vec2ForCTC
3 from src.models.adapters import BottleneckAdapter
4 from src.inference.inference_base import transcribe_audio
5 import os
6
7 def load_finetuned_model():
8     processor = Wav2Vec2Processor.from_pretrained("facebook/wav2vec2-base")
9     model = Wav2Vec2ForCTC.from_pretrained("facebook/wav2vec2-base")
10
11    for layer in model.wav2vec2.encoder.layers:
12        layer.adapter = BottleneckAdapter()
13
14    state = torch.load("weights/adapters/adapter_weights.pth")
15    model.load_state_dict(state, strict=False)
16
17    return model, processor
18
19 def run_inference(test_dir, output_file):
20     model, processor = load_finetuned_model()
21
22     with open(output_file, "w", encoding="utf8") as f:
23         for file in sorted(os.listdir(test_dir)):
24             if file.endswith(".wav"):
25                 text = transcribe_audio(
26                     model,
27                     processor,
28                     os.path.join(test_dir, file)
29                 )
30                 f.write(f"{file}\t{text}\n")
31
32 if __name__ == "__main__":
33     run_inference("data/test", "finetuned_transcriptions.txt")

```

Ce script :

- recharge le modèle de base,
- charge les adaptateurs,
- charge les poids entraînés,
- effectue l'inférence,

- enregistre les résultats dans `finetuned_transcriptions.txt`.

## 7 Calcul du WER final

A cette étape on calcule le WER (word error rate) afin d'évaluer la qualité du modèle. Pour que le modèle fine-tuné soit plus performant il faudrait que la valeur obtenu soit assez grande.

Créer un fichier `src/evaluation/evaluate.py` et y insérer le code suivant puis l'exécuter :

```

1 from jiwer import wer
2
3 def load_lines(path):
4     return [line.strip().split("\t")[1] for line in open(path, encoding="utf8")]
5
6 def compute_wer(refs_path, hyp_path):
7     ref = load_lines(refs_path)
8     hyp = load_lines(hyp_path)
9     score = wer(ref, hyp)
10    print("WER =", score)
11
12 if __name__ == "__main__":
13     compute_wer("data/test/references.txt", "finetuned_transcriptions.txt")

```

Ce script :

- charge les transcriptions référence et prédictions,
- calcule le WER du modèle finement ajusté,

## 8 Orchestration et automatisation de la chaîne de relâsation du projet

Afin d'automatiser l'exécution des différents fichiers de notre méthodologie, on crée un fichier `main.py` à la racine du répertoire `ASR_code`. Ce fichier `main.py` contiendra alors :

```

1 from src.utils.data_utils import download_dataset
2 from src.inference.inference_base import run_inference as run_base
3 from src.training.train import train
4 from src.inference.inference_finetuned import run_inference as run_finetuned
5 from src.evaluation.evaluate import compute_wer
6
7 def main():
8
9     print("\n==== tape 1 : T l chargement du dataset ===")
10    download_dataset()
11
12    print("\n==== tape 2 : Inf rence du mod le de base ===")
13    run_base("data/test", "base_transcriptions.txt")
14
15    print("\n==== tape 3 : Entra nement des adaptateurs ===")
16    train()
17
18    print("\n==== tape 4 : Inf rence du mod le fine-tun ===")
19    run_finetuned("data/test", "finetuned_transcriptions.txt")
20
21    print("\n==== tape 5 : valuation (WER) ===")
22    compute_wer()

```

```
23     "data/test/references.txt",
24     "finetuned_transcriptions.txt"
25 )
26
27 print("\n==== Pipeline termin  avec succ s ===")
28
29 if __name__ == "__main__":
30     main()
```

## Résultats

Les résultats montrent une amélioration significative du modèle grâce à un fine-tuning léger ( $\approx 4.4\%$  des poids entraînés), démontrant l'efficacité des adaptateurs pour l'adaptation à une langue à faible ressource.

## Discussion

Avec le résultat et les performances obtenus en utilisant les adaptateurs on se rend compte que la précision du modèle une fois fine-tuné est s'améliore. Il n'y a plus de risque de perte d'informations, de contexte et de connaissances acquises au cours de l'entraînement car il conserve ses poids initiaux et donc sa capacité initiale à laquelle vient s'ajouter les capacités des adaptateurs ajoutés au modèle.

## Références bibliographiques

Geef een lijst met alle gebruikte bronnen.

- Whisper Turns Stronger : Augmenting Wav2Vec 2.0 for Superior ASR in Low-Resource Languages — Anidjar, Or Haim ; Marbel, Revital ; Yozevitch, Roi (2024).
- Magic dust for cross-lingual adaptation of monolingual wav2vec-2.0 — Khurana, Sameer ; Laurent, Antoine ; Glass, James (2021).