

---

## **CYB6043 - Atelier pratique en cybersécurité**

Entraînement de modèles à partir de la base de données CVE fixes pour détecter du code vulnérable

BERIOT Ehouarn, BILLY Théo, LEROY Jules, PUECH  
Émilien, YAKOVLEV Mikhael



15 décembre 2024

## Table des matières

<b>Introduction</b>	<b>3</b>
Contexte et enjeux . . . . .	3
Objectifs du projet . . . . .	3
<b>CVE fixes</b>	<b>3</b>
Origine et pertinence des données . . . . .	3
Structure et spécificité du jeu de donnée . . . . .	3
Structure des données . . . . .	4
<b>Prétraitement des données</b>	<b>5</b>
Objectifs du prétraitement . . . . .	5
Étapes du prétraitement . . . . .	5
Lecture du fichier de base de données . . . . .	5
Filtrage des CVE récentes . . . . .	5
Requête SQL pour extraire le code . . . . .	5
Filtrage des données par taille des tokens . . . . .	6
Résultats du prétraitement . . . . .	8
Exemple de données . . . . .	8
<b>Entraînement du modèle</b>	<b>9</b>
Outils . . . . .	9
Github . . . . .	9
Kaggle . . . . .	9
Hugging-Face . . . . .	9
Entraînement . . . . .	10
Importation des modules . . . . .	10
Lecture du jeu de donnée et vérification du contenu . . . . .	10
Séparation des données de test et d'entraînement . . . . .	10
Tokenization . . . . .	11
Fine tuning du modèle . . . . .	11
<b>Résultats et analyses</b>	<b>12</b>
<b>Perspectives</b>	<b>16</b>
Modèle . . . . .	16
Fine-Tuning . . . . .	16
Jeu de données . . . . .	16

**Sources**

**16**

## Introduction

### Contexte et enjeux

Dans le monde du développement de logiciels, et plus particulièrement les logiciels open-source, la détection proactive des vulnérabilités dans les codes source est devenue une priorité pour renforcer la sécurité du code. Avec l'augmentation exponentielle des vulnérabilités répertoriées, des solutions basées sur l'intelligence artificielle (IA) se révèlent être des réponses efficaces pour automatiser cette tâche complexe. Les vulnérabilités logicielles ont des impacts considérables, non seulement sur la sécurité des systèmes, mais aussi sur l'intégrité, la confidentialité, la disponibilité des données et la confiance des utilisateurs. Face à cela, il est crucial de détecter et de corriger rapidement les failles avant qu'elles ne soient exploitées.

### Objectifs du projet

Afin de répondre à ce défi, nous avons décidé d'entraîner un modèle d'IA en utilisant le jeu de données CVE fixes, qui contient des informations détaillées sur les vulnérabilités logicielles et leurs correctifs. Ce jeu de données nous offre une base solide pour entraîner un modèle en analysant les différences entre les versions avant et après correction.

## CVE fixes

### Origine et pertinence des données

Le jeu de données CVE fixes provient de la base Common Vulnerabilities and Exposures (CVE), une source standardisée et fiable qui garantit que les données sont bien documentées et représentent les vulnérabilités les plus récentes et pertinentes. En plus de fournir des informations sur les failles de sécurité, le jeu de données offre également le code avant et après correction, avec une colonne diff qui met en évidence les différences spécifiques entre ces versions. Ici un réseau de neurones va pouvoir faire un apprentissage de notre jeu de données pour classer une nouvelle donnée et détecter des vulnérabilités dans un code.

### Structure et spécificité du jeu de donnée

La base de données CVE fixes pèse environ 60 Go (format sqlite) et regroupe un grand nombre d'exemples. Elle récolte des informations directement depuis la National Vulnerability Database (NVD),

jusqu'à une date donnée (celle-ci est mise à jour plutôt régulièrement). En outre, les auteurs du jeu de données ont mis à disposition un script permettant de mettre à jour la base de données avec les nouvelles CVE, garantissant ainsi une actualisation continue des données. Par exemple, la dernière CVE a été publiée le **22 juillet 2024**. Jusqu'à cette date, la base répertorie **12 923 fixes**, fournissant une couverture exhaustive des vulnérabilités logicielles récentes.

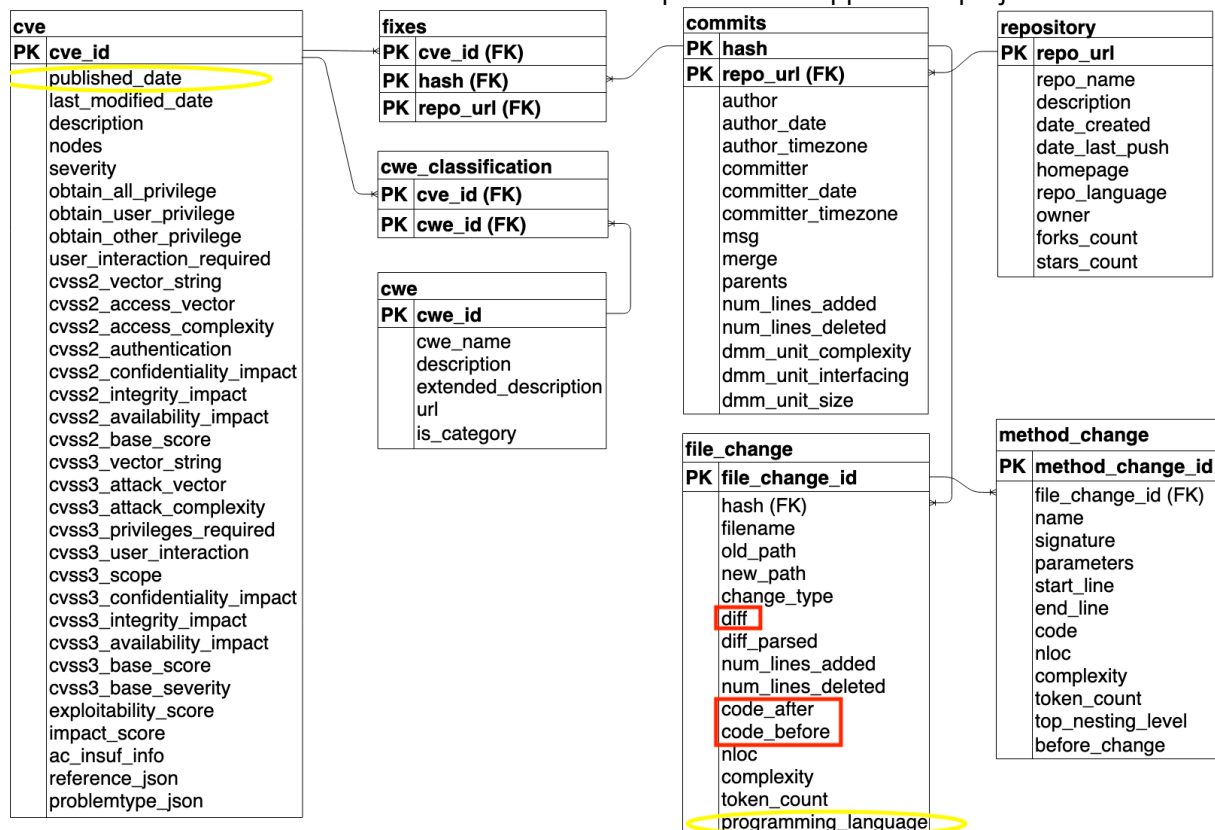
## Structure des données

Le jeu de données est structuré avec plusieurs colonnes clés, telles que `code_before`, `code_after`, et `diff` dans la table `file_change`.

- `code_before` contient le code supposé vulnérable avant correction.
- `code_after` contient le code après correction.
- `diff` représente les différences entre les deux versions, c'est-à-dire les lignes modifiées ou ajoutées pour corriger la vulnérabilité (au format git).

Ces informations sont essentielles pour l'entraînement du modèle d'IA, qui devra apprendre à identifier des motifs récurrents dans les modifications apportées au code afin de détecter des vulnérabilités similaires dans des applications nouvelles.

Voici la structure de la base de données CVE fixes fournie par les développeurs du projet :



## Prétraitement des données

### Objectifs du prétraitement

L'objectif principal de cette étape est d'extraire, de filtrer et de transformer les données issues de la base de données fournie afin de créer un fichier CSV prêt à être utilisé par la suite pour l'entraînement de notre modèle.

### Étapes du prétraitement

#### Lecture du fichier de base de données

La première étape consiste à ouvrir le fichier de base de données SQLite contenant les informations sur les vulnérabilités et leurs correctifs. Cela permet d'exécuter des requêtes SQL pour extraire les données nécessaires.

```
1 import config
2 import sqlite3
3
4 db_path = 'config.DB_PATH' # Variable configurable dans un fichier de
    configuration
5 conn = sqlite3.connect(db_path)
6 cursor = conn.cursor()
```

#### Filtrage des CVE récentes

Pour garantir que les données utilisées soient actuelles et pertinentes, seules les CVE publiées au cours des trois dernières années sont extraites. La date limite est calculée dynamiquement en fonction de la date actuelle.

```
1 from datetime import datetime, timedelta
2
3 date_limite = datetime.now() - timedelta(days=3*365)
4 date_limite_str = date_limite.strftime('%Y-%m-%d')
5 print(f"[INFO] Extraction des CVE publiés après le {date_limite_str}")
```

#### Requête SQL pour extraire le code

Deux requêtes SQL distinctes sont utilisées pour récupérer :

- Les codes corrigés (`code_after`) avec le label 0.
- Les codes vulnérables (`code_before`) avec le label 1.

Ces requêtes incluent des conditions pour : - Ne sélectionner que les fichiers en PHP. - Exclure les enregistrements sans code avant correction.

```
1 query = f"""
2 SELECT
3     file_change.code_after code,
4     file_change.diff,
5     file_change_id,
6     '0' label
7 FROM cve
8 JOIN fixes ON cve.cve_id = fixes.cve_id
9 JOIN commits ON fixes.hash = commits.hash
10 JOIN file_change ON commits.hash = file_change.hash
11 WHERE cve.published_date >= '{date_limite_str}'
12        AND file_change.programming_language = 'PHP'
13        AND file_change.code_before IS NOT "None"
14 ORDER BY cve.published_date DESC;
15 """
16
17 query2 = f"""
18 SELECT
19     file_change.code_before code,
20     file_change.diff,
21     file_change_id,
22     '1' label
23 FROM cve
24 JOIN fixes ON cve.cve_id = fixes.cve_id
25 JOIN commits ON fixes.hash = commits.hash
26 JOIN file_change ON commits.hash = file_change.hash
27 WHERE cve.published_date >= '{date_limite_str}'
28        AND file_change.programming_language = 'PHP'
29        AND file_change.code_before IS NOT "None"
30 ORDER BY cve.published_date DESC;
31 """
```

### Filtrage des données par taille des tokens

Étant donné que nous utilisons le tokenizer [microsoft/codebert-base](#) nous devons respecter une contrainte qui est de ne pas dépasser **512 tokens** une fois le code tokenizer. Chaque entrée est donc analysée pour extraire les modifications spécifiques (diff) et s'assurer que le nombre total de tokens générés ne dépasse pas 512. Si une entrée dépasse cette limite, elle est supprimée.

Étapes :

- Identifier les plages de lignes affectées ([start\\_line](#), [end\\_line](#)) à partir du diff.
- Calculer le nombre de tokens dans les modifications.
- Si le total dépasse 512 tokens, ignorer l'entrée.
- Ajouter un contexte avant et après les modifications pour maximiser l'information.

```
1 # Parcours des lignes du DataFrame
2 for n, row in combined_df.iterrows():
3     if n % 100 == 0:
4         print(f"[INFO] Traitement de l'entrée {n}/{len(combined_df)}")
5
6     try:
7         code = row['code']
8         diff = row['diff']
9         label = row['label']
10
11         # Identifier les plages de lignes affectées
12         if label == '1': # Code vulnérable
13             range_info = diff.split(' ')[1]
14         elif label == '0': # Code corrigé
15             range_info = diff.split(' ')[2]
16         else:
17             print(f"[WARNING] Label inconnu pour l'entrée {n}. Ignorée.")
18             to_remove.append(n)
19             continue
20
21         if ',' in range_info:
22             start_line, nb_lines = map(int, range_info.split(','))
23         else:
24             start_line = int(range_info)
25             nb_lines = 1
26
27         # Extraction des tokens de `diff` et calcul du nombre total de
28         # tokens
29         start_line = abs(start_line)
30         end_line = start_line + nb_lines
31         splitted_code = code.split("\n")
32         code_in_diff = "\n".join(splitted_code[start_line:end_line])
33         code_in_diff_tokenized = tokenizer.tokenize(code_in_diff)
34         diff_token_count = len(code_in_diff_tokenized)
35
36         # Si le diff dépasse 512 tokens, supprimer cette entrée
37         if diff_token_count > 512:
38             print(f"[INFO] Entrée {n} : Diff trop long ({diff_token_count} tokens). Ignorée.")
39             to_remove.append(n)
40             continue
41
42         # Calculer le contexte en fonction des tokens restants
43         remaining_tokens = 512 - diff_token_count
44         context_tokens = remaining_tokens // 2 # Égalité avant et après
45
46         # Tokenisation complète du code pour manipuler les tokens
47         splitted_code_tokenized = tokenizer.tokenize("\n".join(splitted_code))
```



```

        splitted_code))
47
48     # Identifier les indices des tokens
49     start_token_index = len(tokenizer.tokenize("\n".join(
        splitted_code[:start_line])))
50     end_token_index = start_token_index + diff_token_count
51
52     # Calcul des indices pour le contexte
53     start_context_index = max(0, start_token_index - context_tokens
        )
54     end_context_index = min(len(splitted_code_tokenized),
        end_token_index + context_tokens)
55
56     # Extraction des tokens avant, dans et après le diff
57     context_before_tokens = splitted_code_tokenized[
        start_context_index:start_token_index]
58     context_after_tokens = splitted_code_tokenized[end_token_index:
        end_context_index]
59
60     # Reconstruction des tokens finaux
61     final_code_tokens = context_before_tokens +
        splitted_code_tokenized[start_token_index:end_token_index] +
        context_after_tokens
62
63     # Conversion des tokens en texte
64     final_code = tokenizer.convert_tokens_to_string(
        final_code_tokens)
65     combined_df.at[n, 'code'] = final_code
66
67     except Exception as e:
68         print(f"[ERROR] Erreur lors du traitement de l'entrée {n} : {e}
        ")
69     to_remove.append(n)

```

## Résultats du prétraitement

Statistiques des données prétraitées :

- Nombre total d'instances initiales : 125572
- Nombre final d'instance dans notre jeu de données : à 10285

## Exemple de données

```

1 head path_to_csv_file
2
3
4 0          }\n          }\n\n          code  label
   ...      0

```

5	1	} \n	} \n	return \$output; \n	...	0
6	2	<?php \n \n namespace SilverStripe\Reports\Tests;	...			0
7	3	<?php \n \n namespace SilverStripe\Reports\Tests;	...			0
8	4	int with deprecated class Webauthn \n \n \n PublicKe...				0

## Entraînement du modèle

### Outils

#### Github

Pour le projet, nous avons créé un dépôt github, celui-ci contient :

- Le script de prétraitement complet
- Un script de validation des données (vérification de la longueur du code une fois tokenizer avec codebert-base)
- Notre jeu de données (fichier csv)
- Une sauvegarde du notebook jupyter pour entraîner notre modèle

La technologie git permet de créer un environnement collaboratif très simplement et de garder un versioning. C'est un outil indispensable pour les projets de groupe qui contiennent du code comme le notre.

#### Kaggle

Kaggle nous a permis d'exécuter notre notebook d'entraînement en utilisant des ressources plus importantes que nos machines locales. Cet outil propose également une fonctionnalité de sauvegarde des différentes versions d'un notebook, et d'importer des jeux de données de manière persistante. L'utilisation d'un notebook par opposition à un script entier permet d'isoler chaque étape de l'entraînement et de les lancer indépendamment.

#### Hugging-Face

Hugging Face propose une multitude de modèles préentraînés. Nous avons d'abord utilisé le modèle "bert-base-cased", spécialisé dans le Natural Language Processing, puis nous avons évolué vers le modèle "codebert-base", spécialisé dans le traitement de code.

Le modèle est importé en quelques lignes dans notre notebook d'entraînement.

```
1 from transformers import AutoTokenizer
2 tokenizer = AutoTokenizer.from_pretrained('bert-base-cased')
```

## Entraînement

Nous avons fine-tuné notre modèle dans notre notebook. Celui-ci est organisé en bloc.

### Importation des modules

```
1 import numpy as np
2 import pandas as pd
3 import os
4 import matplotlib.pyplot as plt
5
6
7 from sklearn.model_selection import train_test_split
8 from transformers import RobertaTokenizer, RobertaConfig, RobertaModel,
   pipeline
9 from datasets import Dataset
10 from transformers import AutoModelForSequenceClassification, Trainer,
   TrainingArguments, AutoTokenizer
11 from sklearn.metrics import accuracy_score, confusion_matrix,
   ConfusionMatrixDisplay, precision_score
```

### Lecture du jeu de donnée et vérification du contenu

```
1 data = pd.read_csv('/kaggle/input/cvefixes-php-under-512-token/
   cvefixes_diff_under_512_centered_v3.csv')
2 data.head()
3 data.count()
```

### Séparation des données de test et d'entraînement

```
1 X = data['code']
2 Y = data['label']
3
4 # X, Y
5
6 X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size
   =0.2, shuffle=True, random_state=0)
```

## Tokenization

C'est à cette étape que nous commençons à utiliser le modèle "microsoft/codebert-base". Le nombre maximum de tokens (512) demandé par BERT est bien spécifié, pour lever un avertissement en cas de données non valide.

```
1 tokenizer = RobertaTokenizer.from_pretrained("microsoft/codebert-base")
2
3 train_data = Dataset.from_dict({"code": X_train, "label": Y_train})
4 test_data = Dataset.from_dict({"code": X_test, "label": Y_test})
5
6 def preprocess_function(examples):
7     return tokenizer(
8         examples['code'],
9         truncation=True,
10        padding='max_length',
11        max_length=512
12    )
13
14 tokenized_datasets = {
15     "train": train_data.map(preprocess_function, batched=True,
16                             batch_size=32, load_from_cache_file=False),
17     "validation": test_data.map(preprocess_function, batched=True,
18                                 batch_size=32, load_from_cache_file=False)
19 }
```

## Fine tuning du modèle

C'est ici que commence le fine-tuning du modèle pré-entraîné. La variable `training_args` contient les hyperparamètres. Nous avons testé plusieurs combinaisons de paramètres différentes en cherchant à optimiser le modèle.

```
1 model = AutoModelForSequenceClassification.from_pretrained("microsoft/
   codebert-base", num_labels=2)
2
3 training_args = TrainingArguments(
4     per_device_train_batch_size=32, # 16 ok, 32 ça passe
5     learning_rate=2e-5, # valeur à tester
6     warmup_ratio=0.1, # valeur à tester
7     weight_decay=0.1, # pour régulariser le modèle et d'éviter l'
   overfitting
8     output_dir="./results",
9     evaluation_strategy="epoch", # epoch vs steps
10    save_strategy="epoch", # epoch vs steps
11    report_to="none",
12    logging_steps=4,
13    num_train_epochs=5,
```

```
14     logging_dir="./logs"
15 )
16
17 def compute_metrics(eval_pred):
18     logits, labels = eval_pred
19     predictions = np.argmax(logits, axis=-1)
20     accuracy = accuracy_score(labels, predictions)
21     return {"accuracy": accuracy}
22
23 trainer = Trainer(
24     model=model,
25     args=training_args,
26     train_dataset=tokenized_datasets["train"],
27     eval_dataset=tokenized_datasets["validation"],
28     compute_metrics=compute_metrics
29 )
30
31 trainer.train()
```

## Résultats et analyses

Une fois notre modèle fine-tuner nous devons évaluer ces performances et comprendre quels paramètres modifier pour optimiser ces performances, et faire ce processus de manière itératif.

Le modèle que nous évaluons ci-dessous est le modèles codebert-base de Microsoft qui est fine-tuner avec notre jeu de données spécifique au code php de moins de 3 ans, sur les données spécifiques aux modifications avec moins de 512 tokens. Nous l'avons fine-tuner avec seulement 5 epoch, ce qui qui est trop peu mais par manque de temps et de ressource matérielle nous n'avons pas peu malheureusement aller plus loin.

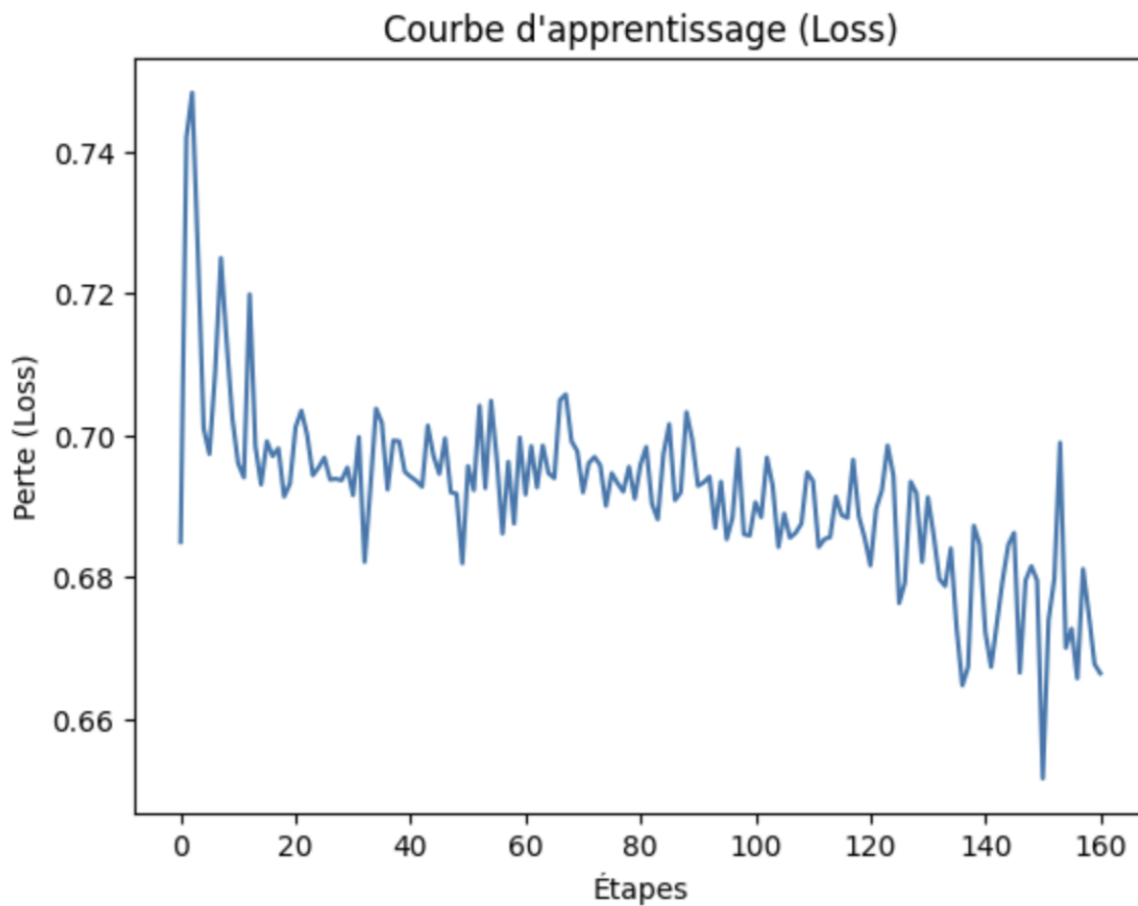
Pour évaluer notre modèle, nous avons utilisé :

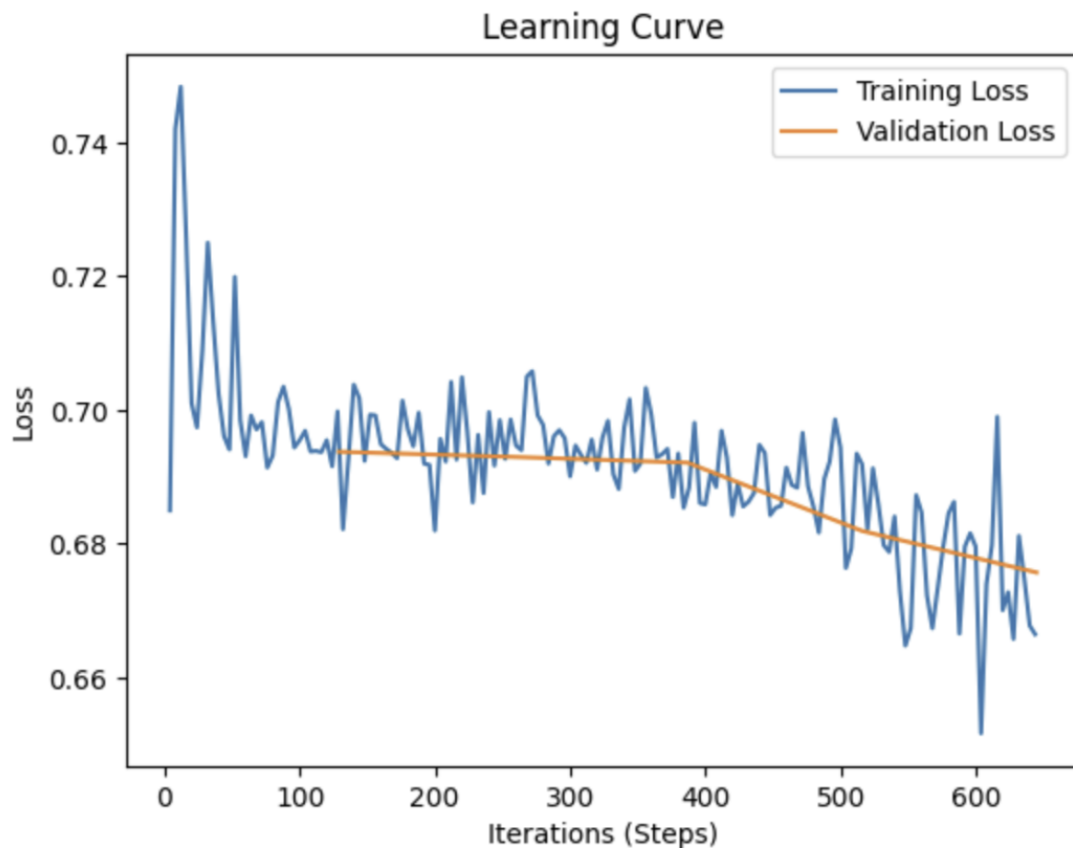
- Un **test de précision** : il mesure la capacité du modèle à éviter les faux positifs dans les prédictions de vulnérabilités. En termes techniques, la précision est le rapport entre les vrais positifs et la somme des vrais positifs et des faux positifs. Une précision proche de 1 indique que le modèle est très précis, tandis qu'une précision de 0,5 signifie une performance médiocre, équivalente à une classification aléatoire (pile ou face).

```
Validation Metrics: {'eval_loss': 0.6757381558418274, 'eval_accuracy': 0.5545851528384279, 'eval_runtime': 36.633, 'eval_samples_per_second': 56.261, 'eval_steps_per_second': 3.521, 'epoch': 5.0}
```

Le test de précision révèle que notre modèle atteint une précision de 0.55, ce qui est légèrement au-dessus d'une classification aléatoire (0.5). Cela indique que le modèle est capable de prédire correctement certaines vulnérabilités, mais il génère encore un nombre significatif de faux positifs. Cette précision pourrait être améliorée en augmentant le nombre d'époques ou en diversifiant le jeu de données d'entraînement.

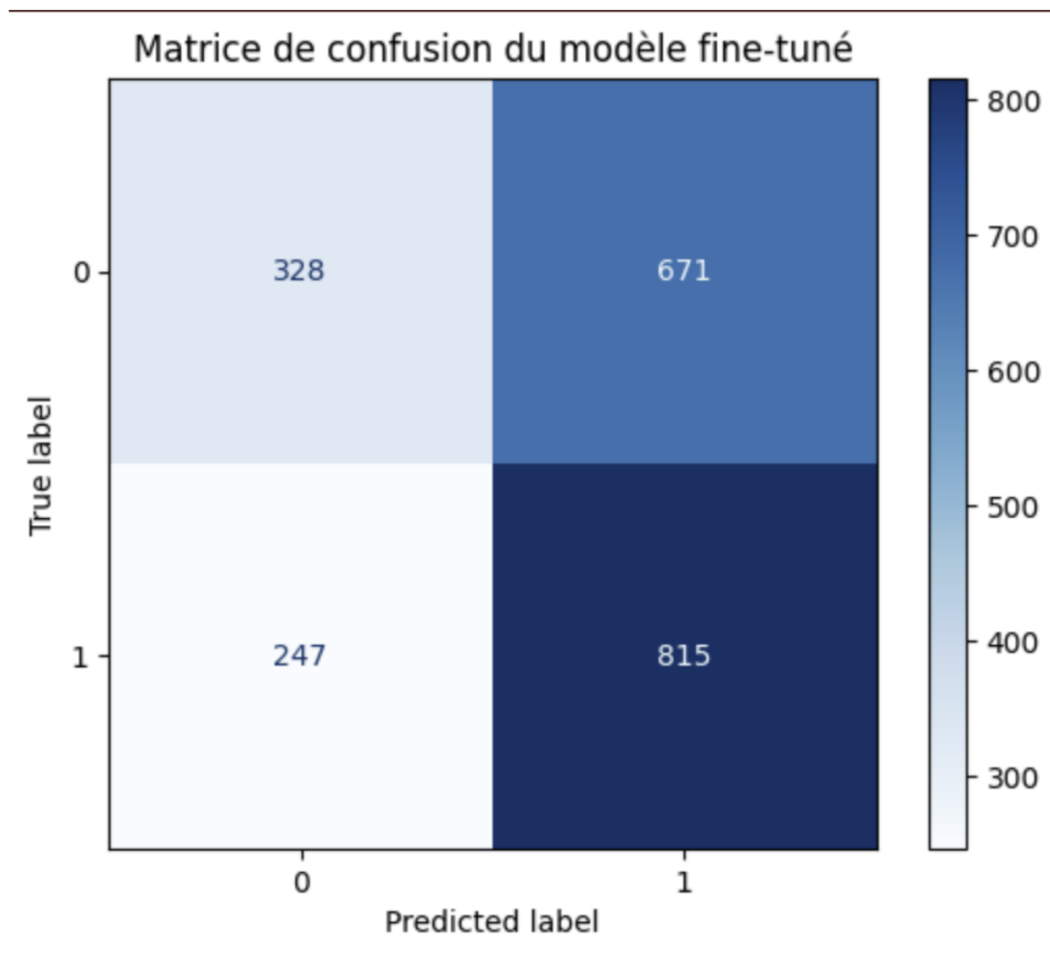
- La **courbe d'apprentissage** : Elle permet de visualiser la perte (loss) du modèle au fil des étapes d'entraînement. Elle permet de détecter des problèmes comme un surapprentissage (overfitting) ou un sous-apprentissage (underfitting). Une courbe stable et décroissante indique que le modèle apprend efficacement.





- Visualiser matrice de confusion : La matrice de confusion est un outil visuel permettant d'évaluer les performances de classification du modèle en comparant les prédictions avec les vraies étiquettes. Elle est composée de quatre éléments principaux:
  - **Vrais positifs (TP)** : cas où le modèle a correctement prédit une vulnérabilité.
  - **Faux positifs (FP)** : cas où le modèle a prédit une vulnérabilité inexistante.
  - **Vrais négatifs (TN)** : cas où le modèle a correctement prédit qu'il n'y avait pas de vulnérabilité.
  - **Faux négatifs (FN)** : cas où le modèle a manqué une vulnérabilité réelle.

En visualisant la matrice de confusion, nous pouvons identifier les biais du modèle, comme une tendance à surestimer ou sous-estimer les vulnérabilités. Cet outil est essentiel pour comprendre où se situent les erreurs du modèle et ainsi orienter les ajustements nécessaires pour améliorer ses performances.



La matrice de confusion met en évidence les erreurs du modèle avec une précision de 0.55. Le modèle nous donne un nombre de faux positifs (**FP**) élevé par rapport aux vrais positifs (**TP**). Cela indique une tendance du modèle à signaler des vulnérabilités inexistantes.

Ce comportement pourrait être corrigé en équilibrant le jeu de données ou en introduisant des techniques comme la régularisation ou le fine-tuning avec un jeu de validation mieux adapté. De plus, un nombre élevé de faux négatifs (**FN**) compromettrait la capacité du modèle à détecter les vulnérabilités réelles, ce qui nécessiterait une amélioration du rappel.

Bien que ces conclusions soient pertinentes elles ne sont pas forcément correctes car le modèle n'a pas été assez fine tuner, 5 époques n'est pas suffisent.

Une fois que nous avons compris les résultats de notre modèle, nous devons le réentraîner en optimisant les hyperparamètres pour avoir des résultats satisfaisant, c'est un processus itératif.

Afin d'optimiser le fine-tuning de notre modèle, nous pouvons modifier les hyperparamètres, qui sont les paramètres d'apprentissage de notre modèle, voici les différents paramètres :



- `per_device_train_batch_size` : Taille du batch, c'est-à-dire le nombre d'échantillons traités simultanément par le modèle à chaque étape d'entraînement. Une valeur plus élevée peut accélérer l'entraînement, mais nécessite plus de mémoire.
- `num_train_epochs` : Nombre d'époques, soit le nombre de fois que le modèle parcourt l'intégralité des données d'entraînement, ce paramètre a un impact important sur le temps d'apprentissage.
- `learning_rate` : Taux d'apprentissage, qui contrôle la vitesse à laquelle le modèle ajuste ses paramètres pendant l'entraînement. Une valeur trop élevée peut empêcher la convergence, tandis qu'une valeur trop faible ralentit l'apprentissage.

## Perspectives

### Modèle

Dans notre projet nous nous sommes concentrés sur le modèle pré-entraîné codebert-base, il serait intéressant de comparer nos résultats avec d'autres modèles pré-entraînés comme code T5.

### Fine-Tuning

Étant donné que nos données sont plutôt conséquentes (beaucoup de ligne de code), nous avons entraîné notre modèle avec des réglages précis, que nous comparons dans la partie précédente. Nous pouvons imaginer beaucoup plus de tests avec des réglages différents pour évaluer lesquels sont les plus optimisés.

### Jeu de données

En effet, nous avons choisi le jeu de données CVE fixes, mais nous avons pris qu'une seule partie de celle-ci, bien que le langage PHP représente un gros pourcentage dans la base de données CVE fixes, nous pouvons imaginer que nous ne nous limitons pas à un seul langage de programmation.

## Sources