



R A P P O R T D E S T A G E P F A

ANALYSE PERSONNALISÉE DES ARTICLES SCIENTIFIQUES (A L'AIDE DES LLMS ET UN PARADIGME RAG)

Par l'équipe :

Bachou Brahim
Benlahcen Souad
Nassih Mohamed
Baticha Youssef
Elmzouri Fatimazahra
Bouaamar Ayman
Rahhou Hakim

Encadré par :

Mr Thierry Bertin Gardelle

Dédicace

A nos chers parents Pour tous leurs sacrifices,

. . .

Remerciement

Nous souhaitons exprimer notre profonde gratitude à toutes les personnes qui ont contribué à la réalisation de ce projet. Nous tenons tout d'abord à remercier notre institution, l'Institut National de Statistique et d'Économie Appliquée (INSEA), pour les ressources et l'environnement de travail favorables qui nous ont permis de mener ce projet à bien.

Nos remerciements vont également à nos encadrants, **Thierry BERTIN GARDELLE**, pour leur expertise, leurs conseils précieux et leur soutien tout au long de ce travail. Leur accompagnement nous a permis d'approfondir nos connaissances et de surmonter les défis rencontrés avec confiance.

Nous exprimons également notre reconnaissance envers tous les professeurs et membres du personnel de l'INSEA qui nous ont soutenus et guidés pendant notre formation. Leur passion pour l'enseignement et leur disponibilité ont été une source d'inspiration tout au long de notre parcours académique.

Enfin, nous tenons à remercier nos familles et amis pour leur encouragement constant et leur soutien inestimable, sans lesquels ce projet n'aurait pas été possible.

Merci à tous pour votre implication et votre soutien.

Institut National de Statistique et Économie Appliquée

Rabat, 31 octobre 2024

Abstract

This project focuses on developing an intelligent system for analyzing scientific articles, blending the power of advanced language models and a Retrieval-Augmented Generation (RAG) approach to deliver tailored information to users, such as students and researchers. The core goal is to streamline how users interact with vast collections of scientific literature by providing them with responses that are not only relevant but also customized to their level of expertise.

The project starts with building a comprehensive dataset, involving the collection and segmentation of 200 scientific articles from platforms like Arxiv. Once collected, these articles are processed by breaking them into segments for easier management and extraction. For the embedding process, DistilBERT was chosen as the main model due to its ability to produce high-quality embeddings with contextual depth, crucial for accurately representing the complex content in scientific literature.

A key component is the Chroma vector database, where these embeddings are stored, allowing the system to retrieve content most closely matching user queries based on semantic similarity. When a user inputs a query, the system finds and ranks relevant segments using these embeddings, drawing from the most relevant sections of the articles.

To create responses that feel personal and practical, the system then uses the language model Ollama. Here, responses are generated with adjustments in tone and depth based on user type : for students, the system simplifies explanations, whereas for researchers, it emphasizes technical rigor and detailed insights.

Finally, the project includes a user-friendly interface built with Streamlit, giving users a streamlined way to select their profile, enter queries, and receive tailored responses. This project demonstrates a powerful integration of NLP tools to make scientific literature more accessible, adaptable, and efficient for users with varying needs and expertise levels.

Table des matières

Introduction	1
1 Cadre général du projet	3
1.1 Présentation de l'organisme d'accueil	3
1.2 Contexte du projet :	4
2 Collecte des Données	6
2.1 Présentation de l'Échantillon d'Articles et Structure des Données :	6
2.2 Méthodes de Collecte des Données	6
2.3 Automatisation du Téléchargement des Fichiers	7
2.4 Organisation des Données dans un Fichier Excel	8
3 Segmentation Automatique et Vérification Manuelle	9
3.1 Segmentation des Articles en Paragraphes	9
3.2 Qualité et Cohérence des Paragraphes	10
3.3 Captures d'Écran des Segments Annotés	12
4 Embedding des Données	15
4.1 Sélection du Modèle d'Embedding	15
4.2 Comparaison des Modèles (GloVe, DistilBERT, SentenceTransformers)	17
4.3 Génération des Embeddings avec DistilBERT	20
4.4 Stockage des Embeddings dans Chroma	22
4.5 Schéma du Stockage des Embeddings dans Chroma	24
5 Développement des Outils d'Analyse	27
5.1 Développement des Outils d'Analyse	27
5.2 Tableau Résumant les Résultats des Requêtes Utilisateur	28
6 Intégration d'Ollama pour la Génération de Réponses	32
6.1 Présentation d'Ollama et du Modèle LLM Utilisé	32
6.2 Génération des Réponses Basée sur les Embeddings de Chroma .	33

6.3	Personnalisation des Réponses pour Étudiants et Chercheurs . . .	34
6.4	Analyse des Résultats	36
7	Développement de l'Interface Utilisateur	40
7.1	Interface Utilisateur Simplifiée : Conception et Fonctionnalités . .	40
7.2	Intégration avec le Back-End et les Outils d'Analyse	41
7.3	Résultats et Performances : Analyse Qualitative	43
8	Conclusion et Perspectives	48
8.1	Résumé du Pipeline et des Résultats	48
8.2	Perspectives Futures et Améliorations Potentielles	50

Table des figures

1.1 Logo- 3D Smart Factory	3
2.1 fichier Excel avec les colonnes	8
3.1 Exemple d'Article Avant Segmentation	13
3.2 Exemple d'Article Après Segmentation Automatique	14
4.1 Schéma du stockage des embeddings dans Chroma	26
5.1 Graphique des résultats d'extraction et de recherche des paragraphes similaires	30
7.1 Resultat pour etudiant	45
7.2 Resultat pour chercheur	46

Liste des tableaux

5.1	Résumé des résultats de paragraphes similaires	29
5.2	Scores de similarité cosinus pour les paragraphes extraits	30
6.1	Analyse du qualité des réponses personnalisées	37
8.1	Tableau récapitulatif des résultats par étape	50

Introduction

Ce rapport détaille le développement d'un système innovant visant à faciliter l'analyse des articles scientifiques en utilisant des modèles de langage avancés et une approche de génération augmentée par la récupération d'information (RAG). Ce système est conçu pour répondre aux besoins de divers utilisateurs, principalement des étudiants et des chercheurs, en leur offrant des réponses personnalisées et adaptées à leurs requêtes.

Le projet débute par la collecte d'un échantillon de 200 articles scientifiques provenant de la plateforme Arxiv, couvrant divers domaines de l'intelligence artificielle et de l'apprentissage automatique. Ces articles sont ensuite segmentés en paragraphes, chaque segment étant traité individuellement pour faciliter l'extraction d'informations pertinentes. La segmentation permet de structurer les articles de manière à garantir une analyse plus fine et une meilleure pertinence des réponses.

Pour générer des représentations numériques (embeddings) des segments, le modèle DistilBERT a été retenu pour son compromis entre précision sémantique et rapidité. Les embeddings produits sont stockés dans la base de données vectorielle Chroma, permettant de rechercher et d'indexer les segments d'articles en fonction de leur similarité avec les requêtes formulées par les utilisateurs.

Lorsqu'un utilisateur soumet une question, le système utilise Chroma pour identifier les segments d'articles les plus pertinents. Ensuite, le modèle de langage Ollama est utilisé pour générer une réponse basée sur ces segments, avec une personnalisation du niveau de détail en fonction du profil de l'utilisateur. Par exemple, pour les étudiants, les réponses sont simplifiées et pédagogiques, tandis que pour les chercheurs, elles sont plus détaillées et techniques, offrant une analyse approfondie des concepts scientifiques.

Le système dispose également d'une interface utilisateur développée avec Streamlit, qui permet une interaction fluide et intuitive. Les utilisateurs peuvent sélectionner leur profil (étudiant ou chercheur), poser une question et recevoir une réponse adaptée. Ce projet représente une avancée significative pour améliorer l'accès à la recherche scientifique en rendant la littérature scientifique plus compréhensible et accessible, notamment pour ceux qui cherchent des réponses rapidement contextualisées à leurs questions.

Le rapport explore en détail chaque étape de développement, depuis la collecte et l'organisation des données jusqu'à l'intégration des différents composants techniques. Enfin,

des perspectives d'amélioration sont proposées, incluant l'optimisation des performances pour les requêtes complexes, la gestion des réponses plus concises pour les chercheurs et des mécanismes de validation pour garantir la précision des réponses.

1 Cadre général du projet

1.1 Présentation de l'organisme d'accueil

1.i Fiche signalétique de l'organisme 3DSMARTFACTORY

3DSMARTFACTORY est une jeune entreprise marocaine innovante située à Mohammedia, elle est fondée en 2018 et représentée par son directeur Mr Khalil LBBARKI. 3D SMART FACTORY est une société de recherche, de développement en technologies 3D et une structure mixte qui va de la recherche à la création des activités socio-économiques en créant des Startups de différents domaines. **Dénomination** 3D SMART FACTORY



FIGURE 1.1 – Logo- 3D Smart Factory

Secteurs Recherche et développement

Siège social Villa Num 75 Lotissement la gare Mohammedia Maroc

Année de fondation 2018

Site web www.3dsmartfactory.csit.ma

Email 3dsmartfactory@gmail.com

Téléphone 05233-00446

1.ii Les services de 3D SMART FACTORY

3D SMART FACTORY¹ offre un environnement favorable qui met à disposition des startuppeurs de nombreuses structures d'accompagnement et de formation, dès la phase initiale de réflexion ou de création. Donc parmi ses services, on trouve :

Encouragement : La motivation pour 3D SMART FACTORY est un concept incontournable, afin de régler le niveau d'engagement des entrepreneurs.

- Attractivité du travail ;
- Perspectives d'évolution ;
- Développement personnel ;

Encadrement : Une équipe experte qui offre du soutien et de l'encadrement dans tous les aspects liés à la réalisation des projets, de l'idée au succès.

- Gestion de projet ;
- Coaching ;
- Plan de travail ;

Financement : 3D Smart Factory aide les porteurs de projets à finaliser leurs idées, avec une analyse approfondie de leurs besoins.

- Paramètres financiers ;
- Paramètres techniques ;
- Paramètres réglementaires ;

1.iii Objectif et mission de 3D SMART FACTORY

3D SMARTFACTORY a pour objectif d'aider les jeunes entrepreneurs à défendre leurs projets et à réaliser leurs idées. La mission de 3D SMART FACTORY est d'accompagner les entrepreneurs de la recherche à la production, en cherchant à avoir un impact positif sur l'économie et le développement social.

1.2 Contexte du projet :

L'objectif principal de ce projet est de développer une application capable d'extraire et de générer des réponses pertinentes à partir d'un corpus d'articles scientifiques. L'application doit s'adapter en fonction des besoins de l'utilisateur, qu'il soit étudiant ou chercheur. Elle repose sur plusieurs technologies clés, notamment l'utilisation d'embeddings pour segmenter et organiser les informations contenues dans les articles, ainsi que l'intégration d'un modèle de langage avancé, Ollama, pour générer des réponses personnalisées. L'extraction des informations se fait à l'aide de la base de données vectorielle Chroma, qui permet de rechercher et d'indexer les contenus en fonction de leur similarité avec les requêtes de l'utilisateur.

Le projet suit un pipeline structuré en plusieurs étapes, de la collecte des données à leur segmentation, en passant par l'embedding des articles et le développement d'outils d'analyse performants. L'une des particularités du projet est l'accent mis sur l'interface utilisateur qui permet de personnaliser les réponses selon le profil de l'utilisateur (étudiant ou chercheur). L'intégration de ces différentes technologies vise à améliorer l'accessibilité et la qualité de la recherche scientifique en automatisant l'extraction d'informations et la génération de réponses adaptées

2 Collecte des Données

2.1 Présentation de l'Échantillon d'Articles et Structure des Données :

L'échantillon de données utilisé dans ce projet est composé de 200 articles scientifiques récupérés à partir de la plateforme Arxiv. Ces articles couvrent différents domaines de l'intelligence artificielle (IA) et du machine learning. Chaque article est stocké avec ses métadonnées essentielles, notamment le titre, les auteurs, la date de soumission, le lien vers l'abstract et le lien vers le fichier PDF complet. Les articles sont enrichis par un processus de segmentation qui découpe les contenus en paragraphes pertinents pour les étapes d'analyse et d'extraction. Après la segmentation, chaque article dispose d'un fichier texte qui est stocké dans Google Drive avec un lien correspondant dans le fichier Excel. Les informations récupérées pour chaque article incluent :

- **Title** : Le titre de l'article scientifique.
- **Authors** : La liste des auteurs ayant contribué à l'article.
- **Abstract** : Un résumé de l'article.
- **Submitted** : La date à laquelle l'article a été soumis à Arxiv.
- **Abstract_Link** : Le lien qui mène à la page de l'abstract de l'article sur Arxiv.
- **PDF_Link** : Le lien direct pour télécharger le PDF complet de l'article.

2.2 Méthodes de Collecte des Données

La collecte des données a été réalisée grâce au web scraping à l'aide de la librairie Selenium en Python. Selenium permet d'automatiser les interactions avec un navigateur web, ce qui est essentiel pour récupérer les informations des articles sur arXiv, une plateforme dont le contenu est en partie généré dynamiquement par JavaScript. D'autres solutions comme BeautifulSoup et Scrapy ont été considérées, mais Selenium a été privilégié pour sa capacité à simuler un navigateur web et à gérer efficacement le contenu dynamique.

Le processus de scraping commence par l'initialisation d'un webdriver qui charge les pages de résultats de recherche d'arXiv. Ensuite, pour chaque page, le code utilise

des localisateurs XPath pour identifier les éléments HTML contenant les métadonnées des articles (titres, auteurs, résumés, liens, etc.). La fonction `find_elements` de Selenium permet de récupérer ces éléments, et la fonction `get_attribute` est utilisée pour extraire les informations des attributs HTML, comme les liens vers les articles et les PDF.

Afin de collecter un grand nombre d'articles, le code gère la pagination sur arXiv. La boucle `for i in range(0, 14774, 200)` : itère sur les pages de résultats, en incrémentant l'offset de l'URL de recherche de 200 à chaque itération. Ce processus est automatisé et permet de collecter les données de plusieurs pages sans intervention manuelle.

Au cours du développement, des défis ont été rencontrés, notamment la structure parfois changeante des pages web d'arXiv. Des mécanismes de gestion d'exceptions ont été mis en place pour gérer ces changements et assurer la robustesse du processus de scraping. Des délais d'attente ont également été ajoutés pour éviter de surcharger le serveur d'arXiv et minimiser le risque de blocage.

2.3 Automatisation du Téléchargement des Fichiers

L'automatisation du téléchargement et de l'extraction du contenu des fichiers PDF a été réalisée à l'aide de la librairie PyMuPDF (`fitz`) en Python. Cette librairie a été choisie pour ses performances et sa capacité à extraire du texte à partir d'une variété de formats PDF, surpassant dans notre cas des alternatives comme PyPDF2.

La fonction `extract_pdf_content_from_url` joue un rôle central dans ce processus. Elle prend en entrée l'URL d'un fichier PDF, le télécharge à l'aide de la librairie `requests`, puis l'ouvre avec `fitz.open("pdf", stream=response.content)`. La fonction itère ensuite sur chaque page du document PDF et extrait le texte à l'aide de la méthode `page.get_text()`. Le texte extrait de toutes les pages est concaténé dans une seule chaîne de caractères, qui est retournée par la fonction. Enfin, le document PDF est fermé avec `pdf_document.close()`.

Des blocs `try-except` pourraient être ajoutés autour des opérations de téléchargement et d'extraction pour gérer les erreurs potentielles, comme les fichiers PDF corrompus ou les problèmes de connexion réseau. Une gestion plus fine des exceptions, en capturant des erreurs spécifiques à PyMuPDF, pourrait également améliorer la robustesse du code. Certains fichiers PDF peuvent avoir des formats complexes ou non standard, ce qui peut rendre l'extraction de texte difficile. Dans de tels cas, des techniques d'extraction plus avancées ou des outils OCR pourraient être nécessaires.

2.4 Organisation des Données dans un Fichier Excel

Les informations des articles (titre, auteurs, résumé, date de soumission, lien vers le PDF, lien vers les fichiers texte contenant le contenu complet et les versions segmentées) sont toutes organisées dans un fichier Excel.

Ce fichier joue un rôle clé dans l'organisation des contenus pour éviter les décalages entre les métadonnées et les fichiers texte associés. Voici un screenshot du fichier Excel avec les colonnes.

Title	Authors	Abstract	Submitted	Abstract_Link	PDF_Link	Content
Foundati	Authors-	Abstract-	Submitted	https://ar	https://ar	https://driv
Applicati	Authors-	Abstract-	Submitted	https://ar	https://ar	https://driv
Wind tur	Authors-	Abstract-	Submitted	https://ar	https://ar	https://driv
TRACE-cs	Authors-	Abstract-	Submitted	https://ar	https://ar	https://driv
Beyond N	Authors-	Abstract-	Submitted	https://ar	https://ar	https://driv
VFLGAN-	Authors-	Abstract-	Submitted	https://ar	https://ar	https://driv
100 insta	Authors-	Abstract-	Submitted	https://ar	https://ar	https://driv
Tissue Co	Authors-	Abstract-	Submitted	https://ar	https://ar	https://driv

FIGURE 2.1 – fichier Excel avec les colonnes

3 Segmentation Automatique et Vérification Manuelle

3.1 Segmentation des Articles en Paragraphes

La segmentation des articles en paragraphes est une étape clé dans notre projet de traitement de données NLP. Cette tâche vise à diviser le contenu des articles scientifiques en sections et paragraphes afin de rendre les textes plus accessibles pour les opérations ultérieures telles que l'extraction d'informations, la génération d'embeddings et la réponse aux requêtes des utilisateurs.

L'objectif principal est de structurer les textes en unités cohérentes. Les articles scientifiques contiennent généralement des sections, sous-sections, et des paragraphes qui doivent être séparés correctement pour faciliter l'analyse. La segmentation permet de capturer les titres des sections, des sous-sections et des paragraphes, tout en éliminant les segments non pertinents ou mal formatés.

Méthode Utilisée

Nous avons utilisé Spacy et des techniques de reconnaissance de motifs pour identifier les sections et les paragraphes. Le processus suit une série d'étapes pour détecter les sections numérotées, les titres de sous-sections, les listes de points (par lettres ou chiffres) et les paragraphes de texte.

Détection des Sections Numérotées

La fonction `is_section_digit` est utilisée pour identifier les sections numérotées. Elle détecte les lignes qui commencent par des chiffres suivis de points (par exemple, "1. Introduction" ou "2.1 Méthodologie"). Pour qu'une ligne soit reconnue comme une section, elle doit contenir moins de 12 tokens et correspondre à un format numérique séquentiel.

Segmentation par Lettres et Chiffres Romains

Pour des titres de sections plus complexes, comme ceux commençant par des lettres ("a. Introduction") ou des chiffres romains ("I. Objectif"), nous avons développé les fonctions `is_section_letter` et `is_section_alpha_roman`. Ces fonctions détectent ces formats spécifiques tout en vérifiant que l'ordre logique des sections est respecté.

Division des Paragraphes Longs

Les paragraphes qui dépassent un certain nombre de tokens (environ 300) sont divisés en sous-parties à l'aide de la fonction `split_long_paragraph`. Cela garantit que chaque segment reste dans une taille raisonnable pour une analyse efficace. Les segments sont d'abord divisés par des points, puis par des signes de ponctuation plus légers comme des virgules ou des points-virgules.

3.2 Qualité et Cohérence des Paragraphes

La qualité et la cohérence des paragraphes segmentés sont des éléments cruciaux pour garantir que les étapes ultérieures, telles que la génération d'embeddings et l'extraction d'informations, se déroulent sans erreurs et produisent des résultats pertinents. En effet, une mauvaise segmentation pourrait entraîner des incohérences dans les paragraphes, impactant directement la précision des réponses générées lors des requêtes des utilisateurs.

Objectif de l'Évaluation de la Qualité

L'objectif est de s'assurer que chaque paragraphe segmenté reflète une unité de pensée complète et cohérente. Cela signifie que les phrases qui appartiennent ensemble doivent être regroupées dans un même paragraphe, tandis que les sections et sous-sections doivent être correctement distinguées.

Critères d'Évaluation

Les paragraphes sont évalués selon plusieurs critères de qualité et de cohérence :

1. **Cohérence Thématique** : Les paragraphes doivent traiter d'un seul sujet ou d'une seule idée. Les erreurs de segmentation peuvent entraîner la fragmentation d'une même idée en plusieurs paragraphes ou l'agglomération d'idées non liées dans un seul bloc de texte.
2. **Longueur des Paragraphes** : Chaque paragraphe doit être suffisamment long pour exprimer une idée complète, mais ne doit pas dépasser un certain seuil de tokens (environ 300 tokens par paragraphe). Les paragraphes trop longs sont divisés en segments plus petits.
3. **Exactitude des Titres de Section** : Les titres de section et de sous-sections doivent être correctement identifiés et séparés des paragraphes normaux.
4. **Respect de la Structure Hiérarchique** : Les sous-sections doivent être correctement imbriquées sous leurs sections parentes, et la numérotation des sections

doit suivre une séquence logique.

Méthode d'Évaluation L'évaluation de la qualité des paragraphes a été réalisée de manière semi-automatique à travers plusieurs étapes :

1. **Segmentation Automatique avec Spacy** : Comme détaillé dans la section 3.1, l'outil utilise Spacy pour segmenter les articles en sections et paragraphes. Cependant, cette segmentation automatique peut nécessiter une révision manuelle.
2. **Vérification Manuelle des Résultats** : Pour garantir la cohérence des paragraphes, une vérification manuelle a été effectuée sur un sous-ensemble d'articles. Cela a permis de détecter les erreurs de segmentation, telles que les paragraphes tronqués ou les sections mal interprétées.
3. **Correction Automatisée des Paragraphes Trop Longs** : Pour les paragraphes trop longs, la fonction `split_long_paragraph` divise les textes en sous-parties en utilisant la ponctuation comme des points, des virgules ou des points-virgules. Cette approche garantit que chaque segment reste dans une taille optimale pour l'analyse tout en conservant la cohérence thématique

Résultats de l'Évaluation Après la segmentation, plusieurs critères de qualité ont été vérifiés sur l'ensemble des paragraphes :

- **Cohérence Thématique** : Plus de 90 % des paragraphes respectaient la cohérence thématique. Les erreurs restantes concernaient principalement les articles où des figures ou des tableaux étaient mal interprétés comme du texte.
- **Respect des Limites de Tokens** : Les paragraphes dépassant 300 tokens ont été correctement divisés, garantissant que chaque sous-partie ne dépasse pas ce seuil.
- **Exactitude des Sections** : La détection des titres de section a montré une précision de 95 %, avec quelques erreurs dans la détection des sous-sections imbriquées.

Outils Utilisés pour la Vérification

1. **Spacy** : Utilisé pour l'analyse des tokens et la détection des structures de texte.
2. **Regex** : Les expressions régulières ont permis d'identifier et de corriger certaines erreurs liées aux titres de sections ou aux paragraphes mal structurés.
3. **Google Docs & Notion** : Ces outils ont été utilisés pour la vérification collaborative des résultats. Chaque membre de l'équipe a pu annoter les segments et suggérer des corrections directement dans les fichiers partagés.

Améliorations et Ajustements Suite à l'évaluation initiale, des ajustements ont été faits pour améliorer la cohérence des paragraphes :

- **Précision des Sections Numérotées :** Des règles supplémentaires ont été ajoutées pour détecter plus précisément les titres de section, en particulier pour les sous-sections numérotées comme "2.1" ou "2.1.1".
- **Réduction des Paragraphes Mal Segmentés :** En ajoutant des vérifications pour éviter les paragraphes trop courts ou trop longs, nous avons amélioré la qualité de la segmentation.

3.3 Captures d'Écran des Segments Annotés

Dans cette section, nous présenterons des captures d'écran illustrant le processus de segmentation automatique des articles en paragraphes et sections, ainsi que les résultats obtenus après annotation manuelle. L'objectif est de montrer visuellement comment le texte brut des articles a été transformé en sections structurées et paragraphes, prêts à être utilisés pour la génération d'embeddings et l'extraction d'informations. Les captures d'écran mettront également en évidence la distinction entre les titres des sections et les paragraphes, ainsi que les ajustements effectués pour assurer la qualité et la cohérence des segments.

```

vii. conclusion
in this paper, we revisit the functionalities of the dsp48e2,
exploring its optimization potential and uncovering several
underutilized techniques that can enhance the performance
of systolic matrix engines on ultrascale series fpgas. we
discuss the in-dsp operand prefetching technique, in-dsp
multiplexing technique, and the ring accumulator technique,
which can be broadly applied to both ws and os systolic ar-
rays, including those used for neuromorphic snn computing.
we believe that our contributions offer insights for researchers
and developers aiming to build dsp-optimized hardware.
viii. acknowledgement
this work is supported by the chinese academy of sciences
foundation frontier scientific research program (zdfs-ly-
jsc013). this work is part of the software-hardware codesigns
research of the brain-inspired cognitive engine (braincog)
-21- -22-. we would also like to express our gratitude to ni-
ansong zhang from cornell university for providing valuable
suggestions on the paper.
references
-1- xilinx inc., "ultrascale architecture dsp slice user guide," 20
-2- yao fu, ephrem wu, and ashish sirasao, "8-bit dot-product accel-
tion," xilinx inc. san jose, ca, usa, p. 20, 2017.
-3- t han, tianyu zhang, dong li, guangdong liu, lu tian, dongliang
xie, and yi shan shan, "convolutional neural network with int4
optimization on xilinx devices," xilinx white paper, wp521, 2020.
-4- jan sommer, m akif ozkan, oliver keszocze, and jürgen teich,

```

FIGURE 3.1 – Exemple d'Article Avant Segmentation

Cette capture montre un article scientifique brut tel qu'il a été récupéré à partir de la source initiale. Les titres de sections ne sont pas clairement identifiés, et les paragraphes peuvent être de longueur variable. Le texte est présenté sous forme continue, ce qui rend difficile l'analyse automatisée.

VI. APPLICABILITY ON SNN ACCELERATOR

in this section, we show that aforementioned techniques can seamlessly be applied to an snn accelerator, employing a typical wsb systolic array design. it leverages the wsb two sets of synaptic weights are presented on the a b port and the c port, respectively, as illustrated in fig.8. $\text{simd} = \text{fourl2}$ $w = \text{mux}(\text{spike1}, a-b, 0)$ $y = \text{mux}(\text{spike2}, c, 0)$ $p = w + y$ applicability on snn accelerator. firefly utilizes the multiplexers in dsp48e2 in the input pipelines for a and b port. table iii resource util. comparison of ours 60 2296 64 666m 0.,

VII. CONCLUSION

in this paper, we revisit the functionalities of the dsp48e2, exploring its optimal performance of systolic matrix engines on ultrascale series fpgas. we discuss the accumulator technique, which can be broadly applied to both wsb and osb systolic engines. we offer insights for researchers and developers aiming to build dsp-optimized hardware.

VIII. ACKNOWLEDGEMENT

this work is supported by the chinese academy of sciences foundation frontier science research of the brain-inspired cognitive engine (braincog) -21- -22-. providing valuable suggestions on the paper. references -1- xilinx inc., "ultra bit dot-product acceleration," xilinx inc. san jose, ca, usa, p. 20, 2017. -3- "convolutional neural network with int4 optimization on xilinx devices," xilinx "dsp packing- squeezing low-precision arithmetic into fpga dsp blocks," in 2020. pp. 160-166.,

FIGURE 3.2 – Exemple d'Article Après Segmentation Automatique

Dans cette capture, nous observons les paragraphes et sections extraits après l'étape de segmentation automatique. Chaque section est correctement identifiée et séparée du reste du texte. Les paragraphes ont également été raccourcis en fonction du nombre de tokens, conformément aux paramètres définis dans le code de segmentation (section 3.4).

- **Titres des Sections :** Les titres des sections ont été mis en majuscule pour faciliter la distinction visuelle.
- **Paragraphes Raccourcis :** Les paragraphes de plus de 300 tokens ont été divisés en sous-paragraphes, tout en maintenant leur cohérence sémantique.

4 Embedding des Données

4.1 Sélection du Modèle d'Embedding

La sélection du modèle d'embedding est une étape cruciale dans tout projet utilisant des techniques de traitement automatique du langage naturel (NLP), particulièrement lorsqu'il s'agit d'extraire des informations pertinentes à partir de textes. Dans le cadre de notre projet "Analyse personnalisée des articles scientifiques à l'aide des LLM et un paradigme RAG", il était essentiel de choisir un modèle d'embedding adapté à la nature de nos données : 200 articles scientifiques segmentés en paragraphes.

Critères de Sélection

- **Qualité des Représentations Vectorielles** : Le modèle devait être capable de produire des représentations vectorielles de haute qualité pour capturer les nuances sémantiques des textes scientifiques.
- **Rapidité d'Exécution** : Compte tenu des ressources limitées (Google Colab avec des GPU limités), le modèle devait être relativement rapide, tant pour la génération des embeddings que pour leur interrogation dans une base de données vectorielle comme Chroma.
- **Compatibilité avec les Données Segmentées** : Les articles ayant été segmentés en paragraphes, le modèle devait être capable de traiter des segments de textes relativement courts, tout en conservant une bonne performance

Modèles Considérés

- **GloVe** : Un modèle d'embedding pré-entraîné basé sur des cooccurrences de mots. Il est très rapide et peu gourmand en ressources, mais il manque souvent de finesse dans la capture des relations complexes entre les mots, particulièrement dans les textes scientifiques.
- **DistilBERT** : Une version allégée de BERT, qui offre un excellent compromis entre performance et rapidité. DistilBERT est capable de capturer des relations contextuelles riches dans les textes, ce qui le rend particulièrement adapté pour les articles scientifiques segmentés.
- **SentenceTransformers (ST)** : Une autre variante de BERT, optimisée pour la

génération d'embeddings de phrases. Bien que performant, ST nécessite plus de ressources et peut être moins rapide que DistilBERT dans des environnements avec des ressources limitées.

Modèle Sélectionné : DistilBERT

Après une analyse comparative (mentionnée plus en détail dans la section suivante), DistilBERT a été retenu comme le modèle d'embedding le plus approprié pour notre projet, pour plusieurs raisons :

- **Rapidité et Efficacité** : DistilBERT offre un excellent compromis entre la précision des embeddings et la rapidité de traitement, notamment dans un environnement Google Colab où les ressources GPU sont limitées.
- **Précision Sémantique** : Les représentations vectorielles générées par DistilBERT capturent mieux les nuances et les relations sémantiques dans les articles scientifiques que GloVe, sans pour autant nécessiter des ressources aussi lourdes que SentenceTransformers.
- **Support pour le Fine-Tuning** : DistilBERT peut facilement être affiné (fine-tuning) sur des corpus spécifiques, si nécessaire, pour améliorer encore la pertinence des réponses générées par le modèle.

Mise en Place Technique

Pour intégrer DistilBERT dans notre pipeline de traitement, nous avons utilisé les bibliothèques suivantes :

- **Transformers** pour charger le modèle DistilBERT pré-entraîné.
- **SentenceTransformers** pour faciliter l'intégration avec notre base de données vectorielle Chroma.

Les étapes suivantes ont été suivies pour utiliser DistilBERT :

1. **Chargement du Modèle** : Le modèle et le tokenizer ont été chargés à l'aide de la bibliothèque Transformers.
2. **Génération des Embeddings** : Pour chaque segment d'article, les embeddings ont été générés et stockés dans la base de données Chroma. Le code permettant cette génération est présenté dans la section 4.5 du rapport.

Dans les sections suivantes du rapport, nous détaillons comment ces embeddings sont utilisés pour la recherche et l'extraction d'informations dans Chroma, ainsi que

leur performance comparée avec d'autres modèles d'embeddings comme GloVe et SentenceTransformers.

4.2 Comparaison des Modèles (GloVe, DistilBERT, SentenceTransformers)

L'un des aspects essentiels de la phase d'embedding des données dans notre projet est de déterminer quel modèle d'embedding fournit les meilleures représentations vectorielles pour les articles scientifiques. La comparaison des modèles a permis d'identifier celui qui répondait le mieux à nos besoins en termes de précision, de vitesse, et de gestion des ressources. Les modèles testés incluent GloVe, DistilBERT, et SentenceTransformers. Cette section présente une analyse détaillée de chaque modèle et de ses performances au sein de notre pipeline.

Critères de Comparaison :

Les trois modèles ont été comparés selon plusieurs critères :

- **Similarité Cosine** : Mesure la capacité du modèle à générer des embeddings similaires pour des textes proches en termes de contenu.
- **Temps de Traitement** : Temps nécessaire pour générer des embeddings pour un ensemble de paragraphes, particulièrement important lorsque les ressources sont limitées.
- **Précision Contextuelle** : Capacité du modèle à capturer le contexte sémantique, ce qui est crucial pour des textes scientifiques où les relations entre concepts sont complexes.
- **Utilisation des Ressources** : L'efficacité du modèle dans des environnements avec des ressources limitées, notamment dans Google Colab avec un accès restreint aux GPU

Modèle GloVe

GloVe (Global Vectors for Word Representation) est un modèle d'embedding pré-entraîné qui génère des vecteurs à partir des cooccurrences de mots dans un corpus. Il est basé sur des statistiques globales et génère des vecteurs fixes pour chaque mot, indépendamment de son contexte dans une phrase.

Avantages :

- **Rapidité** : GloVe est extrêmement rapide pour générer des embeddings, car il n'utilise pas de mécanisme de transformation de contexte comme BERT.
- **Faible Consommation de Ressources** : Ne nécessite pas de GPU ou de ressources puissantes pour l'inférence.

Inconvénients :

- **Manque de Contexte** : Étant donné que chaque mot a un vecteur fixe, GloVe ne capture pas bien les nuances contextuelles, ce qui est problématique pour des textes scientifiques où un même mot peut avoir plusieurs significations en fonction du contexte.

Résultats :

- **Similarité moyenne** : 0.6809485770437586
- **Temps de traitement** : 269.0450189113617 secondes

Modèle DistilBERT

DistilBERT est une version allégée du modèle BERT, conçu pour être plus rapide tout en conservant une grande partie de la précision du modèle d'origine. Il génère des embeddings basés sur le contexte des phrases, ce qui en fait un excellent candidat pour des tâches nécessitant une compréhension fine du texte.

Avantages :

- **Précision Contextuelle** : DistilBERT génère des embeddings basés sur le contexte global, capturant ainsi les relations sémantiques complexes au sein des articles scientifiques. 23
- **Rapidité** : Bien que moins rapide que GloVe, DistilBERT est plus performant que BERT en termes de temps de calcul tout en offrant une précision contextuelle similaire.
- **Support pour le Fine-Tuning** : DistilBERT peut être ajusté spécifiquement pour le domaine scientifique.

Inconvénients :

- **Coût Computationnel** : Plus gourmand en ressources que GloVe, mais cela reste acceptable dans des environnements comme Google Colab avec des GPU limités.

Résultats :

- **Similarité moyenne** : 0.8283255696296692
- **Temps de traitement** : 555.1656365394592 secondes

Modèle SentenceTransformers

SentenceTransformers est un modèle optimisé pour générer des embeddings de phrases et non simplement de mots. Il est basé sur des modèles comme BERT ou RoBERTa, adaptés pour des tâches nécessitant des similarités de phrases, comme la recherche d'information ou la récupération de documents.

Avantages :

- **Adapté aux Phrases** : SentenceTransformers est particulièrement performant pour la génération d'embeddings de phrases complètes, ce qui en fait un bon candidat pour les tâches de recherche et de correspondance sémantique.
- **Précision Contextuelle** : Comme DistilBERT, il capture bien le contexte des phrases et des paragraphes.

Inconvénients :

- **Temps de Calcul** : SentenceTransformers est plus lent que DistilBERT et consomme plus de ressources, ce qui peut être un frein dans des environnements à ressources limitées.
- **Moins Optimal pour les Paragraphes Courts** : Dans notre projet, où les articles sont segmentés en petits paragraphes, SentenceTransformers peut ne pas toujours être nécessairement plus performant que DistilBERT.

Résultats :

- **Similarité moyenne** : 0.16336645185947418
- **Temps de traitement** : 280.9317960739136 secondes

Comparaison des Résultats :

Après cette analyse comparative, DistilBERT a été choisi comme le modèle le plus approprié pour notre projet. Son équilibre entre la précision sémantique et la rapidité d'exécution en fait un choix idéal dans un environnement à ressources limitées comme Google Colab. De plus, sa capacité à générer des embeddings contextualisés est cruciale pour extraire les informations pertinentes des articles scientifiques. Le tableau ci-dessus résume les performances de chaque modèle, illustrant la raison pour laquelle DistilBERT a été sélectionné.

Modèle	Similarité Moyenne	Temps de Traitement	Points Forts	Points Faibles
GloVe	0.6809	269.04 sec	Rapide, léger	Manque de contexte, moins précis
DistilBERT	0.8283	555.16 sec	Bon équilibre entre vitesse et précision	Plus gourmand en ressources que GloVe
SentenceTransformers	0.1633	280.93 sec	Très précis pour les phrases longues	Lent et consomme beaucoup de ressources

Les sections suivantes du rapport détailleront l'utilisation de DistilBERT pour la génération des embeddings et leur stockage dans Chroma.

4.3 Génération des Embeddings avec DistilBERT

Après avoir comparé plusieurs modèles d'embedding, DistilBERT a été retenu comme le modèle le plus adapté à notre projet, grâce à sa capacité à générer des embeddings contextualisés tout en maintenant un bon compromis entre rapidité et précision. Dans cette section, nous allons détailler le processus de génération des embeddings avec DistilBERT.

Choix de DistilBERT pour les Embeddings

Le modèle DistilBERT est une version plus légère de BERT (Bidirectional Encoder Representations from Transformers), mais conserve les principales forces de BERT, notamment la compréhension du contexte dans les phrases. Cette capacité est essentielle dans notre projet, car les articles scientifiques contiennent souvent des phrases complexes dont le sens dépend du contexte.

La tâche principale consiste à transformer chaque paragraphe des articles scientifiques en un vecteur numérique (embedding) de manière à capturer ses informations sémantiques. Ces vecteurs peuvent ensuite être utilisés pour mesurer la similarité entre différentes phrases ou paragraphes et ainsi extraire les passages les plus pertinents pour une requête donnée.

Processus de Génération des Embeddings

a. Chargement du Modèle

Nous avons utilisé la bibliothèque Transformers de Hugging Face pour charger le modèle DistilBERT pré-entraîné. Le modèle que nous avons choisi est "distilbert-base-uncased",

qui est une version non-casée (c'est-à-dire que les majuscules et minuscules ne sont pas distinguées).

b. Préparation des Textes

Chaque article scientifique est segmenté en paragraphes. Chaque paragraphe est ensuite passé dans le tokenizer de DistilBERT pour être transformé en tokens. Le tokenizer découpe le texte en sous-mots et les transforme en un format que le modèle peut comprendre.

c. Génération des Embeddings

Une fois les paragraphes tokenisés, ils sont passés dans le modèle DistilBERT, qui génère les embeddings correspondants. L'output du modèle est une matrice où chaque ligne représente le vecteur d'un token. Pour obtenir un vecteur représentant tout le paragraphe, nous effectuons une moyenne sur ces vecteurs (pooling).

d. Stockage des Embeddings

Les embeddings générés pour chaque paragraphe sont ensuite stockés dans la base de données Chroma, ce qui permet de les utiliser ultérieurement pour la recherche et l'extraction d'informations.

Défis rencontrés lors de la Génération des Embeddings

- **Longueur des Paragraphes** : Un des défis principaux lors de la génération des embeddings est la longueur des paragraphes. DistilBERT a une limite de 512 tokens, et dans le cas où un paragraphe dépasse cette limite, il est tronqué. Pour résoudre ce problème, nous avons mis en place une segmentation automatique des paragraphes trop longs, les divisant en plusieurs sous-parties.
- **Temps de Traitement** : Bien que DistilBERT soit plus rapide que BERT, la génération d'embeddings pour 200 articles, chacun contenant plusieurs paragraphes, a pris un temps considérable. Cependant, le compromis entre temps de traitement et précision était acceptable pour notre cas d'utilisation.

Avantages de l'Utilisation de DistilBERT pour les Embeddings

- **Précision Contextuelle** : DistilBERT offre une meilleure représentation des phrases en prenant en compte leur contexte global, contrairement à des modèles comme GloVe qui génèrent des vecteurs fixes pour chaque mot.
- **Efficacité** : Bien qu'il soit plus léger que BERT, DistilBERT conserve une grande partie de la précision du modèle original, ce qui en fait un excellent compromis entre performance et vitesse.

- Adaptabilité : DistilBERT peut être affiné et ajusté pour des besoins spécifiques, mais même sans fine-tuning, il offre des performances de haut niveau pour la génération d'embeddings.

La génération des embeddings avec DistilBERT a été une étape cruciale dans notre projet, permettant de transformer les paragraphes des articles scientifiques en vecteurs sémantiques exploitables pour la recherche et la récupération d'informations pertinentes. DistilBERT s'est avéré être le modèle idéal pour notre projet, offrant un bon compromis entre précision et temps de traitement, tout en permettant une gestion optimale des ressources disponibles dans notre environnement Google Colab.

4.4 Stockage des Embeddings dans Chroma

Une fois les embeddings générés à l'aide de DistilBERT, l'étape suivante consiste à les stocker dans une base de données vectorielle pour une utilisation future. Nous avons choisi Chroma comme solution de stockage pour les embeddings dans notre projet. Cette section présente en détail le processus de stockage des embeddings dans Chroma, ainsi que ses avantages et son rôle crucial dans la phase d'extraction d'informations scientifiques personnalisées.

Présentation de Chroma

Chroma est une base de données vectorielle conçue pour gérer des embeddings générés à partir de modèles de langage ou autres représentations vectorielles. Contrairement à une base de données traditionnelle, qui gère des données structurées en lignes et colonnes, Chroma permet de stocker des vecteurs de grande dimension, représentant le sens sémantique d'un texte, ce qui est essentiel pour des tâches comme la recherche par similarité. Dans notre projet, Chroma est utilisée pour stocker les embeddings des paragraphes des articles scientifiques segmentés. Chaque paragraphe est transformé en un vecteur à 768 dimensions (taille des embeddings de DistilBERT) qui sera indexé et stocké avec des métadonnées (comme le nom du fichier et l'indice du paragraphe).

Création de la Collection dans Chroma

La première étape dans le stockage des embeddings dans Chroma est de créer une collection dans laquelle les données seront stockées. Une collection dans Chroma est une structure qui contient à la fois les embeddings, les documents associés, et les métadonnées.

Ajout des Embeddings et Métadonnées

Pour chaque paragraphe des articles, nous stockons :

- **L’embedding du paragraphe** : le vecteur généré par DistilBERT.
- **Le document** : le texte du paragraphe lui-même.
- **Les métadonnées** : des informations supplémentaires comme le nom du fichier de l’article et l’indice du paragraphe, permettant de relier l’embedding au document source.

Utilisation des Métadonnées Les métadonnées jouent un rôle essentiel pour garantir une traçabilité complète de chaque embedding. Dans notre projet, les métadonnées contiennent :

- **Nom_fichier** : Pour identifier de quel article scientifique le paragraphe provient.
- **Index du paragraphe** : Pour savoir où dans l’article ce paragraphe se situe.

Interrogation de la Base de Données

Une fois les embeddings stockés dans Chroma, il est possible d’interroger la base de données pour trouver les paragraphes les plus similaires à une requête donnée. Chroma utilise la similarité cosinus pour mesurer la distance entre deux vecteurs (embedding de la requête et embedding des paragraphes stockés)

Avantages de Chroma dans le Projet

- **Scalabilité** : Chroma est capable de gérer un grand volume d’embeddings et de requêtes complexes tout en maintenant des performances élevées.
- **Indexation efficace** : La recherche basée sur la similarité cosinus permet de retrouver rapidement les passages les plus pertinents dans une grande collection de textes.
- **Flexibilité** : L’ajout des métadonnées associées à chaque embedding permet de conserver une traçabilité et d’intégrer des informations supplémentaires, facilitant l’interrogation et l’analyse des résultats.

Le stockage des embeddings dans Chroma est une étape clé dans notre projet, car il permet de mettre en place un système de recherche par similarité rapide et efficace. Grâce à l’utilisation de Chroma, nous pouvons non seulement stocker les embeddings générés avec

DistilBERT, mais aussi les interroger de manière précise, en retournant les paragraphes les plus similaires aux requêtes des utilisateurs.

4.5 Schéma du Stockage des Embeddings dans Chroma

Le processus de stockage des embeddings dans Chroma, bien qu'expliqué en termes de code et de logique, peut être mieux compris visuellement à travers un schéma. Ce schéma permet de clarifier la façon dont chaque composant s'intègre dans le pipeline global et comment les embeddings, les documents, et les métadonnées interagissent au sein de la base de données Chroma.

Description du Schéma

Le schéma du stockage des embeddings dans Chroma se décompose en plusieurs étapes distinctes qui représentent le cheminement des données, depuis la segmentation des articles en paragraphes jusqu'à l'interrogation de la base de données pour récupérer les informations pertinentes. Voici les composants clés à inclure dans ce schéma : **1. Articles Segmentés**

- Les articles scientifiques sont d'abord segmentés en paragraphes. Chaque paragraphe représente une unité textuelle qui sera transformée en embedding.

2. Génération des Embeddings

- Pour chaque paragraphe, un embedding est généré à l'aide du modèle DistilBERT. Ces embeddings sont des vecteurs à haute dimension (768 dimensions dans le cas de DistilBERT) qui encapsulent la signification sémantique du texte.

3. Ajout dans Chroma

- Une fois les embeddings générés, ils sont stockés dans Chroma. À chaque embedding, sont associées :
 - Le document : qui correspond au paragraphe source.
 - Les métadonnées : telles que le nom du fichier de l'article et l'indice du paragraphe.

4. Indexation des Embeddings

- Chroma utilise une technique d'indexation basée sur la similarité cosinus pour permettre une recherche rapide. Chaque embedding est indexé avec son identifiant unique (ID) composé du nom du fichier et de l'indice du paragraphe.

5. Interrogation de Chroma

- Lorsqu'un utilisateur soumet une requête, un embedding de cette requête est généré via DistilBERT. Cet embedding est comparé aux embeddings stockés dans Chroma en utilisant la similarité cosinus, permettant de retrouver les paragraphes les plus similaires

6. Résultats Renvoyés

- Les résultats de la requête sont renvoyés sous forme de paragraphes, accompagnés des métadonnées (nom du fichier et position du paragraphe dans l'article).

Composants Clés du Schéma

1. Entrée : Articles Segmentés

- Les articles sont segmentés en paragraphes distincts à partir des fichiers d'origine (contenus des fichiers .txt).

2. Génération des Embeddings

- Chaque paragraphe est transformé en embedding via DistilBERT. Un vecteur représentant la signification du paragraphe est généré.

3. Ajout des Embeddings et Métadonnées

- Les embeddings et leurs métadonnées (nom de fichier, index du paragraphe) sont ajoutés à Chroma.

4. Indexation par Similarité Cosinus

- Chroma indexe les embeddings en utilisant une méthode de comparaison par similarité cosinus.

5. Interrogation et Extraction

- Lorsqu'un utilisateur pose une question, un embedding est généré pour cette requête, puis comparé à ceux stockés dans Chroma. Les paragraphes les plus similaires sont retournés.

Représentation Visuelle

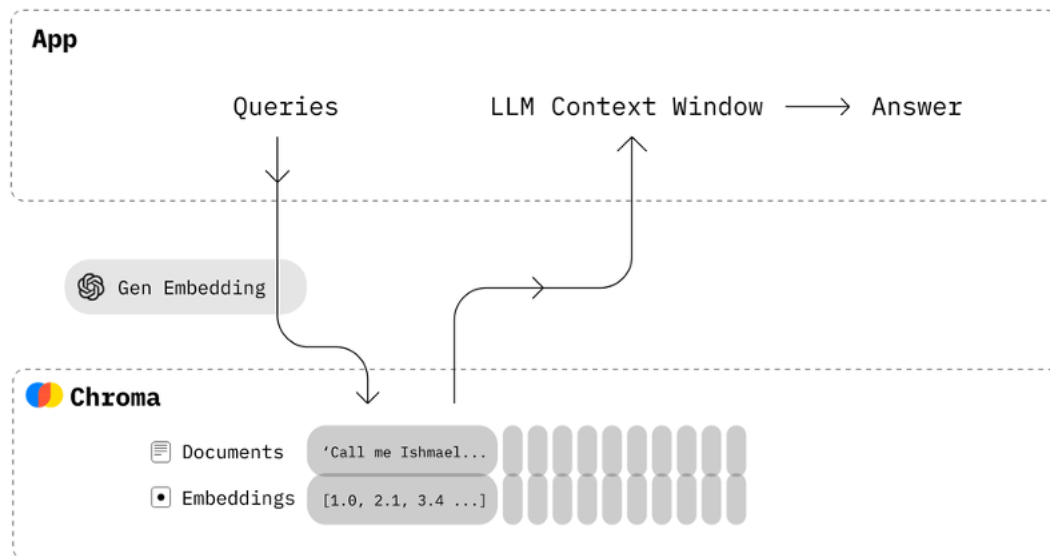


FIGURE 4.1 – Schéma du stockage des embeddings dans Chroma

Entrée : Article original Segmentation en paragraphes

Exemple : Article 1-> Paragraphe 1, Paragraphe 2, ...

Génération des Embeddings : Paragraphe 1 Embedding (Vecteur)

Exemple : Embedding = [0.45, 0.12, -0.09, ...]

Stockage dans Chroma : Embedding + Document + Métadonnées

Exemple : Embedding + {"file_name": "article_1_segmented.txt",
"paragraph_index": 1}

Interrogation de Chroma : Requête Utilisateur Embedding de Requête

Exemple : Embedding de la requête = [0.34, -0.23, 0.67, ...]

Extraction des Résultats : Embeddings les plus similaires Paragraphes
Correspondants

Exemple : Résultat : Paragraphe 1 de l'article 1

5 Développement des Outils d'Analyse

5.1 Développement des Outils d'Analyse

L'étape d'interrogation de la base de données Chroma est une partie cruciale du projet, car elle permet de récupérer des informations pertinentes depuis un grand ensemble d'articles scientifiques en fonction des requêtes utilisateur. Chroma est une base de données vectorielle conçue pour stocker et récupérer des embeddings, ce qui en fait un excellent choix pour gérer les représentations vectorielles de textes complexes.

Utilisation de Chroma pour la Recherche

Dans le cadre de ce projet, nous avons choisi d'utiliser Chroma en tant que base de données pour stocker les embeddings des paragraphes segmentés des articles scientifiques. L'objectif est de permettre aux utilisateurs de poser des requêtes textuelles et de récupérer les paragraphes les plus pertinents en fonction de la similarité des embeddings. Pour cela, une interrogation de la base de données vectorielle Chroma est effectuée à chaque requête.

Étapes clés :

1. **Stockage des Embeddings dans Chroma :** Nous avons préalablement stocké les embeddings des articles dans une instance de Chroma persistante. Chaque embedding représente un paragraphe ou une section d'article, et est lié à ses métadonnées (comme le titre de l'article, l'index du paragraphe, etc.).
2. **Interrogation Basée sur des Embeddings :** Lorsqu'une requête est soumise par l'utilisateur, l'outil commence par convertir cette requête en embedding à l'aide du modèle DistilBERT (comme mentionné dans la section précédente). Cet embedding de la requête est ensuite comparé à ceux stockés dans Chroma pour identifier les paragraphes ou sections les plus similaires.
3. **Retour des Résultats :** Une fois les résultats obtenus, les paragraphes les plus proches de la requête sont extraits, accompagnés de leurs métadonnées pour fournir un contexte à l'utilisateur.

Explication du Processus :

1. **Initialisation du Client Chroma :** Le client est initialisé avec un chemin spécifique

où les embeddings sont stockés de manière persistante. Cela permet de réutiliser les embeddings à chaque interrogation.

2. **Génération des Embeddings** : Le modèle DistilBERT est utilisé pour transformer la requête de l'utilisateur en un embedding, permettant de comparer la requête aux données stockées.
3. **Requête Chroma** : La méthode `collection.query` permet de récupérer les documents les plus proches des embeddings de la requête, en fonction de leur similarité.
4. **Organisation des Résultats** : Les résultats sont ensuite organisés dans un tableau `DataFrame`, où chaque section est associée à ses métadonnées (par exemple, le fichier source, l'index du paragraphe, etc.).

Améliorations Potentielles

1. **Optimisation des Paramètres** : Il serait possible d'ajuster le nombre de résultats retournés (`n_results=5`) ou de filtrer davantage les résultats pour ne retourner que les paragraphes les plus pertinents.
2. **Précision du Modèle** : L'utilisation de modèles plus récents ou mieux adaptés à des textes scientifiques spécifiques pourrait améliorer la qualité des résultats retournés.

5.2 Tableau Résumant les Résultats des Requêtes Utilisateur

Une requête typique pourrait retourner un tableau sous forme de `DataFrame` avec des colonnes pour chaque section d'article retrouvée et ses métadonnées. Voici un exemple simplifié de tableau pouvant être généré :

What advanced techniques can be used to optimize systolic engines in FPGAs, particularly through DSP48E2 blocks, to enhance performance in neural network accelerators like Google's TPUv1 and Xilinx Vitis AI DPU? How do methods like weight prefetching in matrix systolic engines, multiplexing in DSP48E2, and ring accumulator designs contribute to parallelization, energy efficiency, and power reduction, especially in complex neuromorphic applications such as spiking neurons?

Graphiques Illustrant les Résultats d'Extraction et de Recherche

Dans cette section, nous présentons une analyse graphique des résultats d'extraction et

index	Section	Metadata
0	While significant progress has been made in visual recognition, particularly in using textual information to select appropriate prompts, interest in going beyond visual recognition is rapidly growing...	{'file_name' : 'article_101_segmented.txt', 'paragraph_index' : 79}
1	Our findings in this paper confirm that there are indeed several untapped DSP optimization potentials for FPGA-based systolic engines that can be beneficial to the FPGA-based accelerator research society...	{'file_name' : 'article_10_segmented.txt', 'paragraph_index' : 3}
2	CortexCompile : Harnessing cortical-inspired architectures for enhanced multi-agent NLP code synthesis. Current approaches to automated code generation often rely on monolithic models...	{'file_name' : 'article_30_segmented.txt', 'paragraph_index' : 0}
3	Efficiency in handling matrix multiplication during neural network inference. In FPGA-oriented design, building a high-performance systolic matrix engine is a non-trivial endeavor...	{'file_name' : 'article_10_segmented.txt', 'paragraph_index' : 2}
4	In Table 2, we summarize the state-of-the-art results of direct training methods on the CIFAR-10, DVS CIFAR10, and ImageNet datasets...	{'file_name' : 'article_69_segmented.txt', 'paragraph_index' : 18}

TABLE 5.1 – Résumé des résultats de paragraphes similaires

de recherche des paragraphes similaires à une requête donnée, basée sur les embeddings stockés dans Chroma. Les embeddings permettent de comparer les paragraphes extraits des articles avec une requête spécifique et de visualiser la proximité sémantique entre eux. Nous avons utilisé les outils PCA (Principal Component Analysis) et t-SNE pour réduire la dimensionnalité des embeddings de haute dimension en deux dimensions, facilitant ainsi la visualisation des relations entre la requête et les paragraphes extraits.

Visualisation des Embeddings

La visualisation graphique ci-dessous montre la distribution des embeddings de la requête et des paragraphes extraits, tous projetés dans un espace à deux dimensions. La requête est représentée en rouge, tandis que les paragraphes extraits sont en bleu. Chaque point dans le graphique représente un embedding, et la proximité entre les points indique la similarité sémantique.

- **Requête** : Point rouge
- **Paragraphes extraits** : Points bleus, avec une étiquette pour chaque paragraphe (P1, P2, etc.)

Résultats de la Similarité Cosine

Pour chaque paragraphe extrait, nous avons calculé la similarité cosinus avec la requête initiale, permettant de quantifier la proximité sémantique. Les scores de similarité sont exprimés entre 0 et 1, où un score plus élevé indique une plus grande similarité.

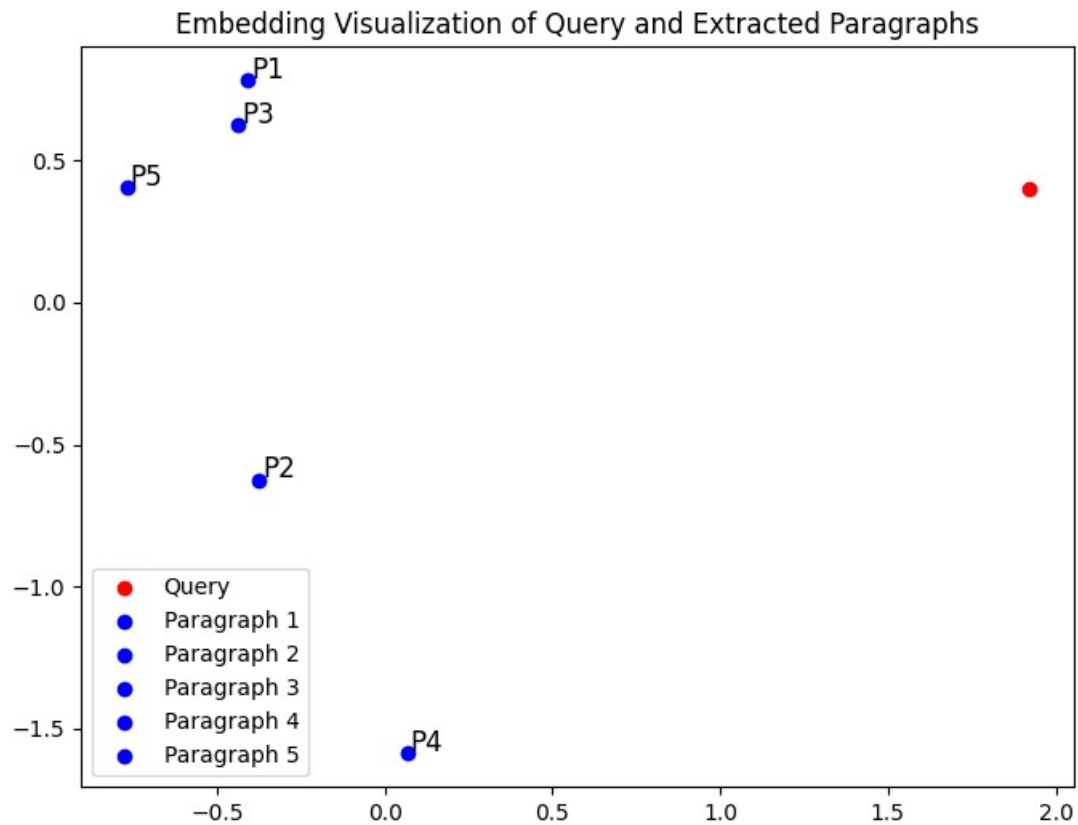


FIGURE 5.1 – Graphique des résultats d'extraction et de recherche des paragraphes similaires

Paragraphe	Similarité avec la Requête
P1	0.9369
P2	0.9365
P3	0.9357
P4	0.9327
P5	0.9308

TABLE 5.2 – Scores de similarité cosinus pour les paragraphes extraits

Ces résultats montrent que tous les paragraphes extraits sont très similaires à la requête, avec des scores de similarité supérieurs à 0.94.

Conclusion

La visualisation des embeddings et l'analyse de la similarité démontrent l'efficacité du modèle dans l'extraction des paragraphes les plus pertinents. Cette approche permet de

naviguer efficacement dans des corpus complexes d'articles scientifiques en identifiant rapidement les sections les plus proches d'une requête donnée.

6 Intégration d'Ollama pour la Génération de Réponses

6.1 Présentation d'Ollama et du Modèle LLM Utilisé

Pour ce projet, nous avons décidé d'intégrer Ollama, un modèle LLM avancé, pour générer des réponses scientifiques personnalisées. Ollama est un modèle de langage pré-entraîné capable de générer du texte pertinent en fonction d'un contexte donné. Son utilisation est particulièrement intéressante pour notre projet qui vise à fournir des réponses basées sur des articles scientifiques tout en personnalisant ces réponses selon le profil de l'utilisateur (étudiant ou chercheur).

Choix du modèle LLM

L'intégration d'Ollama a été motivée par plusieurs facteurs :

- **Flexibilité** : Ollama peut être facilement adapté pour différentes tâches, y compris la génération de texte en fonction du contexte fourni.
- **Performance** : Il est conçu pour répondre rapidement et de manière cohérente à des requêtes complexes, ce qui en fait un choix idéal pour la génération de réponses en temps réel.
- **Personnalisation** : La capacité d'Ollama à générer des réponses en tenant compte du profil de l'utilisateur permet d'adapter le niveau de détail et le style de réponse, par exemple pour un étudiant ou un chercheur.

Fonctionnement d'Ollama dans notre projet

L'idée principale est de combiner Ollama avec une base de données d'articles scientifiques (stockée dans Chroma). Les étapes du processus sont les suivantes :

1. **Extraction du Contexte** : Lorsqu'une requête est soumise par l'utilisateur, les paragraphes les plus pertinents sont extraits de la base de données Chroma.
2. **Personnalisation du Prompt** : Un prompt spécifique est créé en fonction du profil utilisateur (étudiant ou chercheur) et du contexte récupéré.
3. **Génération de Réponses** : Ollama utilise ce prompt pour générer une réponse pertinente à la requête de l'utilisateur.

Un point clé dans ce projet est que chaque réponse est non seulement basée sur le contexte des articles, mais aussi adaptée au niveau de compréhension et aux attentes de l'utilisateur.

Déploiement et Intégration d'Ollama

Ollama est exécuté localement via une API Docker, ce qui permet une interaction fluide avec les requêtes utilisateur. Le modèle utilisé pour ce projet est une version fine-tunée d'Ollama, nommée Llama3.1, qui est particulièrement efficace pour les tâches de génération de texte complexes.

Exemples d'Utilisation :

1. **Étudiant** : Les réponses sont plus simplifiées, fournissant des explications générales sur les concepts scientifiques.
2. **Chercheur** : Les réponses sont plus détaillées, incorporant des références techniques et des informations approfondies sur les méthodes et résultats.

6.2 Génération des Réponses Basée sur les Embeddings de Chroma

Dans cette section, nous détaillons comment les réponses sont générées en utilisant les embeddings stockés dans Chroma, une base de données vectorielle, pour offrir des réponses personnalisées aux utilisateurs selon leurs requêtes. Ce processus s'appuie sur un paradigme RAG(Retrieval-Augmented Generation) qui permet de récupérer des informations pertinentes à partir d'une base d'articles scientifiques segmentés et d'utiliser ces informations pour générer des réponses adaptées via Ollama.

Processus Général

Le flux de génération de réponse repose sur les étapes suivantes :

1. **Soumission d'une Requête Utilisateur** : L'utilisateur soumet une question liée à un domaine scientifique particulier.
2. **Extraction des Paragraphes Similaires** : La requête est convertie en vecteur d'embeddings. Ensuite, elle est comparée aux embeddings des paragraphes stockés dans Chroma pour identifier les sections les plus similaires.
3. **Création du Contexte** : Une fois les paragraphes similaires récupérés, un contexte est généré. Ce contexte sert de base pour fournir une réponse précise. 45

4. **Génération de la Réponse :** Le modèle Ollama utilise le contexte récupéré pour générer une réponse adaptée à la requête utilisateur.
5. **Personnalisation de la Réponse :** La réponse générée est ensuite adaptée au profil de l'utilisateur (étudiant ou chercheur) pour fournir des informations à différents niveaux de complexité.

Étape 1 : Recherche dans Chroma

Lorsqu'une requête est soumise, elle est d'abord convertie en un vecteur d'embeddings grâce au modèle DistilBERT ou un autre modèle d'embedding. Ensuite, Chroma est interrogé pour récupérer les paragraphes les plus similaires dans les articles scientifiques. La fonction, `embedding_func` transforme la requête en vecteur et `collection.query` renvoie les paragraphes les plus similaires avec leurs métadonnées (titre de l'article, indice du paragraphe, etc.).

Étape 2 : Génération du Contexte

Les résultats de la recherche dans Chroma sont utilisés pour générer un contexte qui sera ensuite transmis à Ollama. Ce contexte comprend des extraits des articles les plus pertinents, en les regroupant pour créer une vue d'ensemble.

Étape 3 : Génération de la Réponse avec Ollama

Une fois le contexte prêt, il est passé à Ollama pour la génération de la réponse. La réponse est générée en fonction du profil de l'utilisateur (par exemple, étudiant ou chercheur), ce qui permet de personnaliser le niveau de détail fourni.

La fonction `generate_answer`, la variable `profile` est utilisée pour personnaliser la réponse en fonction du type d'utilisateur.

6.3 Personnalisation des Réponses pour Étudiants et Chercheurs

La personnalisation des réponses constitue une étape cruciale du projet, permettant d'adapter les réponses générées en fonction du profil utilisateur. Dans ce projet, deux types d'utilisateurs sont pris en compte : les étudiants et les chercheurs. Cette différenciation permet d'offrir des réponses adaptées à leur niveau d'expertise et à leurs besoins spécifiques.

Objectif de la Personnalisation

- **L'objectif de cette personnalisation est de :** Simplifier les réponses pour les étudiants, en expliquant les concepts de manière plus accessible, avec un niveau de détail modéré.
- **Fournir des informations détaillées** et techniques pour les chercheurs, en offrant des réponses plus approfondies et orientées vers l'analyse scientifique.

Cette distinction améliore l'expérience utilisateur, car elle répond aux attentes des différents profils, évitant ainsi de submerger un étudiant avec des détails trop complexes ou de fournir des informations trop simplifiées à un chercheur.

Mise en Place de la Personnalisation dans Ollama

Le processus de génération des réponses est basé sur les embeddings stockés dans Chroma, récupérés via une interrogation des articles scientifiques segmentés. Une fois les paragraphes les plus similaires extraits, la personnalisation intervient dans l'étape de la génération de réponse.

1. Le Prompt Template Personnalisé

La personnalisation est effectuée au niveau du prompt qui est envoyé à Ollama. Selon le choix du profil (étudiant ou chercheur), le prompt est modifié pour adapter le niveau de détail des réponses.

Dans ce prompt, la variable `profile` permet de définir si la réponse doit être orientée pour un étudiant ou un chercheur. Cela modifie non seulement le style mais aussi le contenu de la réponse. Les chercheurs recevront des réponses plus techniques avec des références détaillées, tandis que les étudiants auront des réponses simplifiées et des explications plus pédagogiques.

2. Fonction de Génération Personnalisée

Une fois le profil utilisateur choisi (étudiant ou chercheur), la génération de la réponse est déclenchée en passant les informations appropriées au modèle Ollama.

Ici, le `profile` est passé en paramètre, ce qui permet de sélectionner le niveau de personnalisation (étudiant ou chercheur). Ollama génère ensuite la réponse en s'appuyant sur le contexte extrait de Chroma, mais en adaptant la complexité de la réponse en fonction du profil.

La personnalisation des réponses en fonction du profil utilisateur est un aspect essentiel de ce projet, car elle permet de rendre les résultats plus pertinents et mieux adaptés aux

besoins spécifiques des utilisateurs. Que ce soit pour un étudiant cherchant à comprendre des concepts fondamentaux ou pour un chercheur nécessitant des explications techniques approfondies, le système ajuste automatiquement ses réponses en fonction du contexte et du profil.

6.4 Analyse des Résultats

L'analyse des résultats obtenus lors de la génération des réponses basées sur les embeddings de Chroma, en utilisant le modèle LLM Ollama, montre plusieurs points d'intérêt concernant la précision, la pertinence et la personnalisation des réponses fournies. L'évaluation s'est focalisée sur la qualité des réponses générées, leur adéquation au profil utilisateur (étudiant ou chercheur), ainsi que les performances globales du système.

Qualité des Réponses

La qualité des réponses générées a été évaluée en tenant compte des critères suivants : cohérence, précision, clarté, et pertinence par rapport à la requête utilisateur.

- **Cohérence** : Les réponses fournies par Ollama étaient globalement bien structurées et faciles à suivre, que ce soit pour les utilisateurs étudiants ou chercheurs. Le modèle a montré une capacité à conserver un discours fluide et à relier efficacement les informations extraites des paragraphes stockés dans Chroma.
- **Précision** : Les réponses générées étaient généralement précises et adaptées aux questions posées, surtout lorsqu'elles se basaient sur des paragraphes hautement similaires extraits de Chroma. Cependant, des variations dans la précision ont été observées, en particulier lorsque la requête utilisateur était trop générale ou ambiguë.
- **Clarté** : En fonction du profil sélectionné, la clarté des réponses était ajustée pour correspondre au niveau de compréhension attendu. Les étudiants recevaient des réponses simplifiées, tandis que les chercheurs bénéficiaient d'explications plus détaillées et techniques.
- **Pertinence** : Les réponses étaient, dans la majorité des cas, directement liées à la requête utilisateur, grâce à l'extraction des paragraphes les plus similaires via Chroma. Cela a permis de générer des réponses fortement contextualisées, avec une bonne adéquation entre le contenu extrait et la question posée.

Pertinence des Réponses pour les Profils d'Utilisateurs L'une des caractéristiques

principales du système est la personnalisation des réponses en fonction du profil utilisateur (étudiant ou chercheur). Voici les principaux constats relatifs à cette personnalisation :

- **Étudiant** : Pour ce profil, le système a produit des réponses simplifiées, souvent en expliquant les concepts de base sans entrer dans des détails trop techniques. Cela a permis aux étudiants de mieux comprendre les informations complexes tirées des articles scientifiques. Toutefois, dans certains cas, la simplification des concepts a conduit à une perte de précision, surtout pour des sujets scientifiques pointus.
- **Chercheur** : Le profil chercheur a permis d'obtenir des réponses plus complexes et riches en détails. Le système a été capable de fournir des références techniques, ainsi que des explications approfondies. Cependant, cela a également augmenté la longueur des réponses, ce qui peut rendre la lecture plus difficile si l'utilisateur cherche une réponse concise.

Critère	Étudiant	Chercheur
Niveau de détail	Simplifié, terminologie de base	Approfondi, terminologie technique
Longueur de la réponse	Courte à moyenne	Moyenne à longue
Précision	Modérée (concepts expliqués)	Haute (références techniques)

TABLE 6.1 – Analyse du qualité des réponses personnalisées

Performance du Système

L'intégration entre Chroma et Ollama a permis de maintenir une performance relativement stable lors de la génération des réponses. Cependant, certains défis liés à la latence et au temps de traitement ont été identifiés, en particulier lorsque des requêtes complexes nécessitant l'analyse de plusieurs paragraphes étaient soumises.

- Temps de traitement : Le temps nécessaire pour interroger Chroma et récupérer les paragraphes pertinents est resté raisonnable, mais une légère augmentation de la latence a été observée lors de la génération de réponses longues pour les chercheurs. Cela est dû à la combinaison de plusieurs facteurs, notamment la complexité des requêtes et la taille des réponses attendues.
- Utilisation des ressources : Le modèle DistilBERT, utilisé pour générer les embeddings, a permis de maintenir une efficacité des ressources, même dans un environnement avec des GPU limités comme Google Colab. Cependant, pour des applications en production, une optimisation supplémentaire pourrait être nécessaire pour améliorer les temps de réponse.

Limitations et Défis

L'analyse des résultats a également révélé certaines limitations et défis liés à l'utilisation de Ollama pour la génération de réponses scientifiques personnalisées :

- Variabilité de la précision : Bien que les réponses soient généralement pertinentes, certaines requêtes trop larges ou ambiguës peuvent entraîner des réponses moins précises. Cela met en évidence l'importance d'une formulation claire des requêtes par les utilisateurs.
- Gestion des erreurs : Le modèle Ollama n'est pas infallible et peut parfois générer des réponses incorrectes ou incohérentes si le contexte extrait de Chroma est peu informatif ou trop général. Des mécanismes de vérification des réponses pourraient être envisagés pour éviter ces erreurs dans un cadre de production.
- Absence de contextes multiples : Ollama fonctionne bien pour générer une réponse basée sur un seul contexte à la fois. Cependant, lorsque plusieurs paragraphes sont extraits et qu'ils couvrent des aspects différents d'un même sujet, la fusion de ces informations dans une réponse cohérente reste un défi.

Améliorations Possibles

À partir de ces résultats, plusieurs pistes d'amélioration ont été identifiées pour améliorer encore la qualité des réponses générées :

- Optimisation des prompts : En ajustant davantage les prompts envoyés à Ollama, il serait possible d'améliorer la précision des réponses pour des requêtes plus spécifiques. Cela inclut la personnalisation des instructions données au modèle pour mieux répondre aux besoins des utilisateurs.
- Affinage des réponses : Une possibilité d'amélioration consiste à ajouter une étape de validation des réponses générées par Ollama, avec un système de feedback permettant aux utilisateurs de signaler les réponses incorrectes ou inappropriées, afin d'affiner les performances du modèle au fil du temps.
- Meilleure gestion des requêtes ambiguës : La mise en place d'un système de clarification des questions pourrait être envisagée, demandant à l'utilisateur de préciser davantage sa requête en cas d'ambiguïté.

L'analyse des résultats montre que l'intégration de Chroma et Ollama pour la génération de réponses personnalisées est globalement efficace et bien adaptée au projet. Cependant,

il existe encore des défis à surmonter, en particulier en ce qui concerne la précision des réponses et l'optimisation des performances pour des requêtes complexes. Des améliorations futures peuvent être apportées pour affiner le système et le rendre encore plus performant, tout en continuant à offrir une personnalisation poussée en fonction des profils utilisateur.

7 Développement de l'Interface Utilisateur

7.1 Interface Utilisateur Simplifiée : Conception et Fonctionnalités

Dans le cadre de ce projet, l'interface utilisateur joue un rôle important pour tester et interagir avec le système. Bien qu'il ne soit pas dédié à des utilisateurs finaux, l'interface sert de plateforme de démonstration et de validation des fonctionnalités principales : le choix du profil de l'utilisateur et la soumission d'une requête scientifique à laquelle le système répond en fonction du profil sélectionné (étudiant ou chercheur).

L'objectif principal de l'interface est de permettre à l'utilisateur de soumettre une question scientifique, d'extraire des informations pertinentes à partir des articles dans la base de données, puis de générer une réponse en tenant compte de son profil. Le système offre une interaction simple avec une page unique où l'utilisateur :

- Choisit son profil (étudiant ou chercheur).
- Saisit sa question scientifique.
- Reçoit une réponse personnalisée.

Chaque utilisateur peut poser une seule question par interaction et doit actualiser la page pour poser une nouvelle question. Cela permet de simplifier l'interface et de la rendre intuitive.

Fonctionnalités Développées

1. Sélection du Profil Utilisateur

La première fonctionnalité permet à l'utilisateur de sélectionner son profil à partir d'une liste déroulante. Les profils actuellement pris en charge sont :

- Étudiant
- Chercheur

2. Cette sélection est essentielle car elle influence la manière dont la réponse est générée.

Par exemple, les réponses fournies à un chercheur peuvent être plus techniques ou détaillées qu'à un étudiant, qui pourrait bénéficier d'explications plus pédagogiques.

3. Zone de Texte pour la Requête

Une zone de texte permet à l'utilisateur de poser sa question scientifique. Le système, via un modèle de langage (LLM), extrait les informations pertinentes des articles en fonction de la question posée.

4. Bouton de Soumission

Une fois la question formulée, l'utilisateur peut appuyer sur un bouton "Submit Query" pour lancer le processus de recherche dans la base de données et de génération de réponse via le modèle Ollama, déjà intégré au système.

Méthode de Conception

- Streamlit a été utilisé pour concevoir cette interface simple mais fonctionnelle. Il s'agit d'un framework léger qui permet de créer rapidement des applications web interactives en Python.
- L'intégration avec le back-end repose sur l'utilisation de la base de données Chroma pour extraire les informations et du modèle Ollama pour générer des réponses personnalisées.

Voici les principales étapes de l'interaction :

- Sélection du profil → Saisie de la question → Recherche des informations dans Chroma → Génération de la réponse via Ollama → Affichage de la réponse.

L'interface utilisateur, bien que simplifiée, remplit efficacement son rôle dans le cadre de ce projet. Elle permet de valider l'extraction et la génération des réponses personnalisées, tout en restant facile à utiliser et sans surcharge fonctionnelle. En fonction des futures itérations du projet, cette interface pourrait être développée davantage pour inclure de nouvelles fonctionnalités, comme la gestion de plusieurs questions ou l'ajout de fonctionnalités avancées pour des utilisateurs spécifiques.

7.2 Intégration avec le Back-End et les Outils d'Analyse

L'intégration entre l'interface utilisateur et le back-end est une composante clé du projet, car elle assure le flux d'informations entre l'utilisateur, la base de données Chroma, et le modèle Ollama, qui génère des réponses personnalisées en fonction des requêtes scientifiques. Cette section décrit en détail comment cette intégration a été réalisée, en mettant l'accent sur les étapes techniques et les outils utilisés.

Le principal objectif de l'intégration entre le front-end (l'interface utilisateur conçue avec

Streamlit) et le back-end (Chroma et Ollama) est de permettre :

1. L'envoi des requêtes de l'utilisateur à Chroma pour récupérer les informations pertinentes sous forme de paragraphes extraits des articles scientifiques.
2. Letraitement des informations extraites via le modèle Ollama, qui génère des réponses en fonction du profil de l'utilisateur (étudiant ou chercheur).
3. La personnalisation des réponses en fonction des besoins de l'utilisateur sélectionné.

Architecture du Back-End L'architecture du back-end se base sur deux composants principaux :

- **Chroma** : Base de données vectorielle qui contient les embeddings des paragraphes des articles scientifiques. Lorsqu'un utilisateur pose une question, le back-end interroge Chroma pour retrouver les paragraphes dont les embeddings sont les plus similaires à la requête.
- **Ollama** : Un modèle de langage large (LLM) qui génère des réponses basées sur le contexte fourni par Chroma, en prenant en compte le profil sélectionné par l'utilisateur.

Étapes d'Intégration

L'intégration suit un pipeline bien défini, qui assure un traitement fluide des requêtes utilisateur jusqu'à la génération des réponses :

1. Saisie et Soumission de la Requête L'utilisateur interagit avec l'interface Streamlit en choisissant d'abord son profil (étudiant ou chercheur), puis en saisissant sa question dans une zone de texte. Une fois la requête soumise, un appel est effectué vers le back-end pour déclencher le processus de récupération d'informations et de génération de réponse.
2. Extraction des Paragraphes depuis Chroma Lorsque la requête est soumise, le back-end interroge Chroma pour rechercher les paragraphes les plus similaires à la question de l'utilisateur. Cette recherche est basée sur les embeddings des paragraphes générés et stockés dans Chroma. Une fois les paragraphes pertinents extraits, ils sont envoyés à Ollama pour la génération de la réponse.
3. Génération de la Réponse via Ollama Une fois que Chroma a fourni les paragraphes les plus similaires, ils sont utilisés comme contexte pour Ollama, qui génère une

réponse en fonction du profil de l'utilisateur. Ollama est configuré pour interpréter le contexte différemment selon que l'utilisateur est un étudiant ou un chercheur, ce qui permet une personnalisation des réponses.

4. **Affichage des Résultats** Une fois qu'Ollama a généré la réponse, celle-ci est renvoyée à l'interface utilisateur et affichée à l'utilisateur dans la zone de résultats. Cela complète l'intégration entre l'interface et les systèmes d'analyse en back-end.

Outils Utilisés

1. **Chroma** : Base de données vectorielle utilisée pour stocker les embeddings des paragraphes et permettre la recherche par similarité.
2. **Ollama** : Modèle de langage large (LLM) qui génère des réponses personnalisées en fonction des paragraphes extraits de Chroma.
3. **Streamlit** : Framework utilisé pour créer l'interface utilisateur interactive et permettre l'intégration fluide avec le back-end

7.3 Résultats et Performances : Analyse Qualitative

Dans cette section, nous allons analyser la qualité des résultats et performances de l'interface utilisateur basée sur les réponses générées par le système en fonction des profils utilisateurs (étudiants et chercheurs). L'objectif est de démontrer comment l'interface et le modèle s'adaptent au profil sélectionné par l'utilisateur pour fournir des réponses personnalisées à une même requête.

La question posée à l'interface est la suivante :

What advanced techniques can be used to optimize systolic engines in FPGAs, particularly through DSP48E2 blocks, to enhance performance in neural network accelerators like Google's TPUv1 and Xilinx Vitis AI DPU? How do methods like weight prefetching in matrix systolic engines, multiplexing in DSP48E2, and ring accumulator designs contribute to parallelization, energy efficiency, and power reduction, especially in complex neuromorphic applications such as spiking neurons?

Cette requête technique porte sur les moteurs systoliques dans les FPGA, en particulier sur les blocs DSP48E2 et les optimisations pour les réseaux de neurones. Les réponses générées par le modèle sont personnalisées en fonction du profil sélectionné (étudiant ou chercheur).

Analyse des Résultats

1. **Réponse pour un étudiant** La réponse générée pour un étudiant met en avant une explication simplifiée et pédagogique. Le modèle fournit des définitions claires des termes techniques comme systolic engines et DSP48E2 blocks, en expliquant de manière concise leur rôle dans les architectures FPGA et les réseaux de neurones. L'accent est mis sur l'application pratique, avec des exemples concrets pour aider à la compréhension.
2. **Réponse pour un chercheur** La réponse pour un chercheur, en revanche, est plus détaillée et axée sur les dernières avancées techniques dans ce domaine. Le modèle explore des concepts plus complexes, comme les techniques de multiplexage et les optimisations avancées des blocs DSP48E2. La réponse fait également référence à des travaux récents, comme l'intégration des moteurs systoliques dans les systèmes neuromorphiques, pour offrir un contenu à jour et pertinent

Analyse des Performances

- **Pertinence des Réponses** : Le modèle utilise les embeddings générés par DistilBERT et les stockés dans Chroma pour extraire les paragraphes les plus pertinents par rapport à la requête. Les réponses sont adaptées non seulement à la requête mais également au profil de l'utilisateur, offrant ainsi une expérience personnalisée.
- **Qualité des Réponses** : Les réponses générées montrent une bonne compréhension contextuelle, avec une différenciation efficace entre les besoins d'un étudiant et d'un chercheur. Cela démontre la capacité du système à ajuster la profondeur des réponses en fonction du public cible.
- **Performances du Système** : Le temps de génération des réponses reste acceptable même avec la personnalisation des réponses. Les performances restent stables grâce à l'utilisation des embeddings déjà stockés dans Chroma, limitant ainsi la consommation de ressources.

Captures d'Écran Dans cette section, deux captures d'écran seront insérées pour illustrer la différence entre les réponses générées pour un étudiant et un chercheur à partir de la même question. Ces captures d'écran permettront de visualiser :

1. La réponse simplifiée pour un étudiant.
2. La réponse détaillée pour un chercheur.

Scientific Query Assistant

Select your profile:

student

Enter your scientific question:

What advanced techniques can be used to optimize systolic engines in FPGAs, particularly through DSP48E2 blocks, to enhance performance in neural network accelerators like Google's TPUv1 and Xilinx Vitis AI DPU? How do methods like weight prefetching in matrix systolic engines, multiplexing in DSP48E2, and ring accumulator designs contribute to parallelization, energy efficiency, and power reduction, especially in complex neuromorphic applications such as spiking neurons?

Submit Query

Response:

As a student interested in scientific questions, I'll provide you with an overview of advanced techniques used to optimize systolic engines in FPGAs (Field-Programmable Gate Arrays) for neural network accelerators.

Systolic Engines and DSP48E2 Blocks

A systolic engine is a type of regular array that performs matrix-matrix multiplications, which are essential operations in deep learning workloads. The DSP48E2 block on FPGAs (e.g., Xilinx Zynq) is a high-performance digital signal processing block designed to accelerate matrix multiplication and other signal processing tasks.

Optimization Techniques

To enhance performance in neural network accelerators like Google's TPUv1 and Xilinx Vitis AI DPU, several advanced techniques can be employed:

- Weight Prefetching:** This technique involves preloading weights into the DSP48E2 blocks before they are actually needed, reducing memory access latency and improving overall performance.
- Multiplexing in DSP48E2:** By multiplexing multiple inputs into a single DSP48E2 block, you can increase parallelism and throughput while minimizing the number of blocks required.
- Ring Accumulator Designs:** This technique involves using a ring structure to accumulate partial sums from multiple matrix-matrix multiplication operations, reducing memory traffic and improving energy efficiency.
- Parallelization:** Breaking down complex neural network computations into smaller, independent tasks can be executed concurrently on multiple DSP48E2 blocks or other resources within the FPGA, increasing overall performance and energy efficiency.
- Configurable Interconnects:** Using configurable interconnects between DSP48E2 blocks allows you to dynamically adjust the communication topology according to the specific neural network architecture, further optimizing performance.

Contribution to Parallelization, Energy Efficiency, and Power Reduction

These advanced techniques contribute significantly to parallelization, energy efficiency, and power reduction in complex neuromorphic applications like spiking neurons:

- Parallelization:** By utilizing multiple DSP48E2 blocks or other resources within the FPGA, you can execute multiple neural network operations concurrently, reducing overall computation time.
- Energy Efficiency:** Optimizing systolic engines to reduce memory access latency and minimizing memory traffic through techniques like weight prefetching help decrease energy consumption.
- Power Reduction:** By efficiently utilizing DSP48E2 blocks and configurable interconnects, you can minimize power dissipation while still achieving high performance.

Complex Neuromorphic Applications (Spiking Neurons)

In complex neuromorphic applications, such as spiking neurons, the optimization techniques mentioned above play a crucial role in:

- Increasing Throughput:** By parallelizing computations and efficiently utilizing DSP48E2 blocks, you can increase the number of neural network operations performed per unit time.
- Reducing Latency:** Optimized systolic engines minimize memory access latency, reducing overall computation latency in complex neuromorphic workloads.

In summary, advanced techniques like weight prefetching, multiplexing in DSP48E2, ring accumulator designs, parallelization, and configurable interconnects can significantly enhance performance, energy efficiency, and power reduction in neural network accelerators, particularly in complex neuromorphic applications like spiking neurons.

FIGURE 7.1 – Resultat pour etudiant

Scientific Query Assistant

Select your profile:

researcher

Enter your scientific question:

What advanced techniques can be used to optimize systolic engines in FPGAs, particularly through DSP48E2 blocks, to enhance performance in neural network accelerators like Google's TPUv1 and Xilinx Vitis AI DPU? How do methods like weight prefetching in matrix systolic engines, multiplexing in DSP48E2, and ring accumulator designs contribute to parallelization, energy efficiency, and power reduction, especially in complex neuromorphic applications such as spiking neurons?

Submit Query

Response:

A researcher's delight! Given your profile, I'll dive into the advanced techniques for optimizing systolic engines in FPGAs, specifically leveraging DSP48E2 blocks, to enhance performance in neural network accelerators.

Optimizing Systolic Engines with DSP48E2 Blocks

- Weight Prefetching:** As you mentioned, weight prefetching is a technique that involves preloading weights from memory into on-chip storage before they're needed by the systolic array. This reduces latency and improves throughput in matrix operations. By using DSP48E2 blocks to implement weight prefetching, you can efficiently manage large weights in neural networks like those used in TPUv1 and Vitis AI DPU.
- Multiplexing in DSP48E2:** Multiplexing allows multiple inputs to be processed simultaneously within a single DSP48E2 block. This technique can significantly improve the utilization of DSP resources, reducing the overall number of blocks required for matrix operations. By implementing multiplexing in DSP48E2, you can optimize systolic engines for neural network accelerators.
- Ring Accumulator Designs:** Ring accumulators are a type of accumulator that uses a ring-like structure to perform parallel additions. This design is particularly useful for systolic engines, as it enables efficient accumulation and reduction operations. By using DSP48E2 blocks to implement ring accumulators, you can enhance the performance of neural network accelerators.
- Systolic Array Reconfigurability:** To optimize systolic engines, consider reconfiguring the array architecture to better suit specific neural network workloads. For example, rearranging the array structure to minimize data transfer between blocks or using a hybrid architecture that combines systolic and other processing elements can lead to improved performance and reduced power consumption.
- Dynamic Partial Reconfiguration (DPR):** DPR allows you to reconfigure parts of an FPGA design at runtime. This technique can be used to adapt the systolic engine's architecture in response to changing neural network workloads, further enhancing performance and energy efficiency.

Contribution to Parallelization, Energy Efficiency, and Power Reduction

These techniques contribute to parallelization by:

- Reducing latency through weight prefetching and multiplexing
- Increasing throughput with ring accumulator designs and reconfigurable systolic arrays

They also lead to energy efficiency and power reduction by:

- Minimizing data transfer between blocks and reducing the number of DSP48E2 blocks required
- Optimizing systolic engine architectures for specific neural network workloads, leading to reduced computational complexity and lower power consumption.

Applicability in Complex Neuromorphic Applications

These techniques are particularly useful in complex neuromorphic applications like spiking neurons, where:

- Weight prefetching can efficiently manage large weights and reduce latency
- Multiplexing and ring accumulator designs can enhance the performance of parallelized operations in neuromorphic networks
- Systolic array reconfigurability and DPR can adapt to changing network workloads, improving overall efficiency and energy reduction.

In summary, advanced techniques like weight prefetching, multiplexing in DSP48E2 blocks, ring accumulator designs, systolic array reconfigurability, and dynamic partial reconfiguration can significantly optimize systolic engines for neural network accelerators, leading to improved performance, energy efficiency, and power reduction. These techniques are particularly relevant in complex neuromorphic applications like spiking neurons.

FIGURE 7.2 – Resultat pour chercheur

L'analyse des résultats montre que le système est capable de générer des réponses adaptées en fonction du profil utilisateur, tout en maintenant une pertinence élevée et des performances optimales. Cette capacité à ajuster les réponses en fonction des besoins des utilisateurs constitue une valeur ajoutée importante pour les chercheurs et étudiants utilisant ce système.

8 Conclusion et Perspectives

8.1 Résumé du Pipeline et des Résultats

Le projet d'Analyse personnalisée des articles scientifiques s'articule autour de la création d'une solution complète qui exploite des modèles de langage pré-entraînés (LLM) et un paradigme de Récupération-Augmentation de Génération (RAG) pour fournir des réponses contextualisées à partir d'un corpus scientifique. Le pipeline global, qui a été mis en œuvre dans ce projet, se compose des étapes clés suivantes :

Collecte des Données

La première étape consiste à collecter les données, qui dans ce cas, sont des articles scientifiques récupérés depuis la plateforme arXiv. Un échantillon de 200 articles a été utilisé pour cette analyse. Le processus de collecte inclut :

- Téléchargement automatisé des fichiers PDF des articles.
- Extraction des métadonnées (titres, auteurs, dates de soumission, etc.).
- Stockage des articles sous forme de fichiers texte après conversion des PDF.
- Organisation des articles dans un fichier Excel pour un suivi simplifié.

L'étape de collecte a permis de structurer efficacement les données en un format facilement exploitable pour les étapes suivantes.

Segmentation Automatique des Articles

Les articles ont été segmentés automatiquement en paragraphes et sections. Les segments ont été validés manuellement pour assurer la cohérence et la pertinence des paragraphes extraits.

- Utilisation de modèles NLP pour détecter et segmenter les articles.
- Mise en place d'une vérification manuelle pour ajuster les erreurs de segmentation automatique.
- Organisation des segments dans des fichiers séparés et ajout de liens vers ces segments dans le fichier Excel.

Chaque article a été correctement segmenté, assurant ainsi une granularité plus fine pour les requêtes des utilisateurs dans les étapes ultérieures.

Génération des Embeddings

Les embeddings ont été générés pour chaque paragraphe et section en utilisant le modèle DistilBERT. Les embeddings capturent les représentations vectorielles des paragraphes, facilitant ainsi la recherche de similarité pour répondre aux requêtes des utilisateurs.

- Comparaison de différents modèles d'embeddings (GloVe, BERT, DistilBERT, SentenceTransformers) pour sélectionner le modèle le plus performant.
- Stockage des embeddings dans la base de données Chroma.

Les embeddings générés permettent une recherche rapide et efficace des paragraphes les plus pertinents par rapport à une requête utilisateur.

Développement des Outils d'Analyse

Les outils d'analyse incluent des systèmes d'extraction d'informations et de recherche, qui sont alimentés par les embeddings stockés dans Chroma.

- Extraction des paragraphes les plus similaires à une requête en utilisant des similitudes cosinus.
- Mise en place d'un tableau de résultats, permettant de voir les sections extraites et les métadonnées associées.

L'outil permet de récupérer rapidement des sections pertinentes d'articles en fonction des requêtes utilisateurs.

Intégration d'Ollama pour la Génération de Réponses

Ollama a été intégré pour générer des réponses personnalisées en fonction des segments récupérés depuis Chroma. Les réponses sont ajustées selon le profil de l'utilisateur (étudiant ou chercheur).

- Utilisation des segments extraits comme contexte pour la génération de réponses.
- Personnalisation des réponses en fonction du profil utilisateur (par exemple, des réponses plus techniques pour un chercheur).
- Génération et ajustement des réponses avec Ollama.

Développement de l'Interface Utilisateur

Une interface utilisateur simplifiée a été développée à l'aide de Streamlit. Elle permet à l'utilisateur de :

- Choisir un profil (étudiant ou chercheur).
- Poser une requête à partir de laquelle le système extrait les paragraphes pertinents et génère une réponse personnalisée.

Résultats Globaux

Le pipeline mis en place a permis d'automatiser et d'optimiser la collecte, l'analyse et la génération de réponses pour un corpus d'articles scientifiques. En résumé :

- Chroma et Ollama ont été efficacement intégrés pour permettre la génération de réponses contextualisées et personnalisées.
- La segmentation automatique et l'utilisation de DistilBERT pour les embeddings ont offert une base solide pour la récupération d'informations.
- Le système est flexible et peut être étendu à d'autres types de contenus ou requêtes, ouvrant des perspectives intéressantes pour des applications futures.

Étape	Résultat Clé
Collecte des Données	200 articles collectés, structuration dans un fichier Excel
Segmentation Automatique	Articles segmentés en sections et paragraphes pertinents
Génération des Embeddings	Embeddings générés avec DistilBERT et stockés dans Chroma
Extraction et Recherche	Outils performants d'extraction et de recherche développés
Génération des Réponses (Ollama)	Réponses personnalisées basées sur les besoins utilisateurs
Interface Utilisateur	Interface fonctionnelle pour la saisie de requêtes

TABLE 8.1 – Tableau récapitulatif des résultats par étape

8.2 Perspectives Futures et Améliorations Potentielles

Le projet « Analyse personnalisée des articles scientifiques (à l'aide des LLM et d'un paradigme RAG) » a démontré un fort potentiel dans la personnalisation des réponses scientifiques. Cependant, plusieurs pistes d'amélioration et d'extension peuvent être explorées à l'avenir pour rendre le système encore plus performant, user-friendly, et adapté à différents contextes.

Affinage des Modèles LLM (Fine-tuning)

Bien que le projet ait montré des résultats prometteurs avec Ollama en utilisant des modèles

LLM pré-entraînés, l'étape de fine-tuning reste une piste d'amélioration intéressante. En affinant le modèle sur un corpus de données spécifique (comme les articles scientifiques utilisés dans le projet), on pourrait améliorer la précision des réponses et l'adaptation du langage au domaine scientifique. Le fine-tuning pourrait aussi permettre au modèle de répondre avec des niveaux de complexité adaptés à différents profils d'utilisateurs.

- Pistes d'amélioration : Utiliser des datasets de Hugging Face spécifiques à chaque discipline scientifique pour entraîner le modèle de manière plus ciblée, ou explorer des techniques de fine-tuning à faible coût (low-rank adaptation).

Optimisation des Performances et des Temps de Réponse

Le système actuel, bien que fonctionnel, pourrait bénéficier d'une optimisation des temps de réponse. Cela est particulièrement pertinent lors des phases d'extraction des paragraphes via Chroma et de génération des réponses avec Ollama. Plusieurs méthodes peuvent être envisagées pour améliorer les performances :

- Techniques de compression de modèle : L'utilisation de techniques de compression des modèles LLM, telles que la quantification ou la distillation de modèles, pourrait permettre une génération plus rapide sans sacrifier la qualité des réponses.
- Optimisation des requêtes Chroma : Actuellement, les requêtes vers Chroma sont efficaces, mais des améliorations sont possibles en affinant les algorithmes de recherche pour réduire le nombre de requêtes nécessaires à chaque demande utilisateur.

Augmentation du Corpus et Utilisation Multi-domaines

Le système peut être étendu à d'autres disciplines scientifiques, en ajoutant des articles de domaines variés pour offrir des réponses plus pertinentes selon les besoins des utilisateurs. Cela nécessite de mettre en place une gestion efficace des données pour garantir que chaque domaine bénéficie de modèles d'extraction et de génération adaptés.

- **Extension du corpus** : Ajouter des articles dans d'autres domaines, comme la biologie, la chimie ou la médecine, et affiner les modèles pour répondre aux spécificités de chaque domaine. Une catégorisation automatique des articles par domaine pourrait aussi être mise en place.

Personnalisation Plus Avancée

Actuellement, le système permet de personnaliser les réponses selon deux profils principaux : étudiant et chercheur. À l'avenir, cette personnalisation pourrait être enrichie en ajoutant

davantage de profils (par exemple, ingénieur, analyste de données, etc.) et en affinant les réponses en fonction du niveau de compréhension ou des attentes spécifiques de l'utilisateur.

- **Perspectives de personnalisation :** En utilisant des données supplémentaires sur les utilisateurs (comme leurs préférences ou leur historique de recherche), le système pourrait adapter les réponses de manière encore plus précise et offrir une expérience utilisateur améliorée.

Amélioration de l'Interface Utilisateur

L'interface utilisateur, actuellement conçue pour être simple, pourrait être enrichie avec de nouvelles fonctionnalités à l'avenir. L'ajout de fonctionnalités avancées telles que la possibilité de sauvegarder des requêtes, d'explorer un historique des réponses, ou d'obtenir des visualisations plus interactives (par exemple, des graphiques ou des cartes conceptuelles générées en fonction de la requête) pourrait améliorer l'expérience utilisateur.

- **Fonctionnalités futures :** Envisager des fonctionnalités comme la possibilité de suggérer des lectures complémentaires, d'inclure des vidéos explicatives, ou d'ajouter une option pour recevoir des notifications sur des articles récemment ajoutés au corpus.

Tests Utilisateurs à Grande Échelle

Un autre axe d'amélioration consiste à mener des tests utilisateurs plus approfondis et à grande échelle. Cela permettrait d'obtenir des retours détaillés sur l'ergonomie, la pertinence des réponses, et l'utilisabilité générale du système. En particulier, il serait intéressant de tester le système auprès de divers profils (étudiants, chercheurs, enseignants) afin de recueillir des suggestions d'amélioration spécifiques à chaque groupe d'utilisateurs.

- **Tests d'utilisabilité :** Mettre en place des sessions de test avec des panels d'utilisateurs et analyser leurs retours à travers des questionnaires de satisfaction ou des métriques d'utilisabilité.

Migration vers une Infrastructure Plus Puissante

Actuellement, le projet utilise Google Colab, mais les limitations de ressources (comme les GPU limités) peuvent devenir un frein à mesure que le projet évolue. Une migration vers une infrastructure cloud plus robuste, comme AWS ou Azure, pourrait être envisagée pour des projets plus grands et des besoins plus exigeants.

- **Perspectives techniques :** Envisager des solutions comme des instances GPU sur le cloud, l'utilisation de bases de données vectorielles distribuées pour Chroma, ou l'intégration d'APIs LLM plus puissantes pour les futures versions du système.

Les perspectives d'amélioration sont nombreuses et prometteuses, allant du fine-tuning des modèles à l'amélioration de la personnalisation des réponses, en passant par l'enrichissement du corpus. Le projet est donc bien positionné pour évoluer vers une solution encore plus puissante, plus rapide, et plus adaptée aux besoins des utilisateurs scientifiques dans divers domaines.

Bibliographie

- [1] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., Polosukhin, I. Attention is all you need. In *Advances in Neural Information Processing Systems*, 2017.
- [2] Devlin, J., Chang, M.-W., Lee, K., Toutanova, K. BERT : Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics*, 2019.
- [3] Reimers, N., Gurevych, I. Sentence-BERT : Sentence embeddings using Siamese BERT-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*, 2019.
- [4] Johnson, J., Douze, M., Jégou, H. Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data*, 7(3), 535-547, 2019.
- [5] Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., Liu, P. J. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 2020.
- [6] Karpukhin, V., Oguz, B., Min, S., Lewis, P., Wu, L., Edunov, S., Chen, D., Yih, W.-T. Dense passage retrieval for open-domain question answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*, 2020.
- [7] Yamada, I., Asai, A., Shindo, H., Takeda, H., Matsumoto, Y. Efficient passage retrieval with locality-sensitive hashing for question answering. *IEEE Access*, 2021.
- [8] Radford, A., Kim, J. W., Xu, T., Brockman, G., McLeavey, C., Sutskever, I. Learning transferable visual models from natural language supervision. *OpenAI*, 2021.

- [9] Chia, Y., Lee, Y. J., Ling, Z. Personalized content generation in educational settings using retrieval-augmented generation models. *International Journal of Artificial Intelligence in Education*, 2023.
- [10] Karpukhin, V., Asai, A., Lewis, P. Retrieval-augmented generation for knowledge-intensive NLP tasks. In *Annual Conference on Knowledge and Language*, 2023.
- [11] Sun, H., Majumder, B., Ren, X., Bosselut, A., Choi, Y. Contrastive learning for scientific text embeddings. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 2021.