



SMART **AUTOMATION**
TECHNOLOGIES

Rapport de Stage

Système intelligent d'assistance client
détection d'intentions — extraction d'entités —
classification — recommandations

Réalisé par :

Taibi Mohamed Yassine

Responsable de stage : Salma LIICHI

Encadrant académique : Hanifi Majdoulayne

Entreprise d'accueil : Smart Automation Technologies, Tanger

Année universitaire 2024–2025

Dédicace

À ma famille et à tous ceux qui m'ont soutenu durant ce parcours.

Remerciements

Je remercie chaleureusement : Mme Salma LIICHI (responsable de stage chez Smart Automation Technologies) pour son accueil et ses conseils, M. Hanifi Majdoulayne pour son encadrement académique, ainsi que l'équipe de Smart Automation Technologies pour son accompagnement durant ces deux mois. Merci également aux collègues et camarades qui ont relu et commenté ce rapport.

Résumé

Ce rapport décrit la conception et le développement d'un assistant client intelligent réalisé lors d'un stage de deux mois au sein de Smart Automation Technologies (Tanger). Le système permet la détection d'intentions, l'extraction d'entités produit et client, la classification du type de client (émotion / comportement) et la génération de recommandations personnalisées. Les modules principaux reposent sur des modèles de NLP (BERT fine-tuned), des techniques d'extraction basées sur spaCy et des méthodes de recommandation par similarité. Le travail inclut la conception, l'implémentation, l'évaluation et l'intégration d'une interface prototype (Streamlit / Google Colab pour prototypage).

Mots-clés : NLP, BERT, extraction d'entités, classification d'émotions, recommandations, assistant client.

Abstract

This report describes the design and development of an intelligent customer assistant implemented during a two-month internship at Smart Automation Technologies (Tangier). The system enables intent detection, product and customer entity extraction, customer type classification (emotion/behavior), and personalized recommendation generation. The main modules rely on NLP models (BERT fine-tuned), extraction techniques based on spaCy, and similarity-based recommendation methods. The work includes design, implementation, evaluation, and integration of a prototype interface (Streamlit/Google Colab for prototyping).

Keywords : NLP, BERT, entity extraction, emotion classification, recommendations, customer assistant.

Liste des abréviations

NLP	Natural Language Processing
BERT	Bidirectional Encoder Representations from Transformers
QA	Question / Answer
GANTT	Diagramme de Gantt
SAT	Smart Automation Technologies
IA	Intelligence Artificielle
ML	Machine Learning
API	Application Programming Interface
NLTK	Natural Language Toolkit
TF-IDF	Term Frequency-Inverse Document Frequency
NLU	Natural Language Understanding
NER	Named Entity Recognition
RNN	Recurrent Neural Network
LSTM	Long Short-Term Memory
GPU	Graphics Processing Unit
CPU	Central Processing Unit
JSON	JavaScript Object Notation
UI	User Interface
UX	User Experience

Table des figures

1.1	Planification détaillée sur les 8 semaines de stage	14
2.1	Arbre à problèmes détaillé du projet	17
2.2	Diagramme pieuvre des besoins fondamentaux	21
2.3	Architecture détaillée du système d'assistance client intelligent	22

Liste des tableaux

1.1	Fiche QQQQCP du projet	12
2.1	Comparatif détaillé des solutions d’assistant conversationnel	20
3.1	Paramètres d’entraînement du modèle BERT	25
3.2	Mapping détaillé émotions-types de clients	26
4.1	Performance détaillée par type d’intention	30
4.2	Performance détaillée par type d’entité	31
4.3	Performance détaillée par type de client	31
4.4	Synthèse des performances globales du système	31
4.5	Temps de réponse détaillé par opération	32
4.6	Analyse comparative détaillée avec les solutions existantes	33

Liste des annexes

1. Annexe A : Extrait du code principal - Installation des dépendances et imports
2. Annexe B : Extrait du code principal - Configuration des modèles
3. Annexe C : Extrait du code principal - Fonction de personnalisation
4. Annexe D : Présentation de l'entreprise SAT
5. Annexe E : Jeu de données d'entraînement Bitext
6. Annexe F : Catalogue produits utilisé
7. Annexe G : Résultats détaillés des tests
8. Annexe H : Questionnaire d'évaluation utilisateur

Introduction générale

L'évolution rapide des technologies de l'information et de l'intelligence artificielle transforme radicalement la relation client dans le paysage commercial moderne. Les entreprises cherchent constamment à automatiser les tâches récurrentes tout en maintenant, voire en améliorant, une expérience utilisateur de qualité. Ce stage de deux mois, effectué au sein de Smart Automation Technologies à Tanger, s'inscrit parfaitement dans cette dynamique innovante.

L'objectif principal était de concevoir et développer un assistant client intelligent capable non seulement de comprendre les requêtes des utilisateurs à travers la détection d'intentions précises, mais aussi d'extraire des informations structurées pertinentes (entités), d'identifier le type de client (stressé, pressé, curieux, neutre) et de proposer des recommandations personnalisées et contextuelles.

Le système développé combine plusieurs technologies de pointe en Natural Language Processing (NLP), incluant le fine-tuning du modèle BERT pour la détection d'intentions, l'utilisation de spaCy et TheFuzz pour l'extraction robuste d'entités, et l'implémentation d'algorithmes de recommandation basés sur la similarité cosinus et TF-IDF. L'ensemble a été intégré dans une interface prototype développée avec Streamlit, offrant une expérience utilisateur fluide et intuitive.

Ce document présente dans un premier chapitre le contexte général du projet, incluant la présentation de l'entreprise d'accueil et la méthodologie de gestion de projet adoptée. Le deuxième chapitre détaille l'étude théorique et le cadrage du projet, avec une analyse approfondie de l'état de l'art et des solutions envisagées. Le troisième chapitre décrit la conception et le développement technique du système, tandis que le quatrième chapitre présente les résultats obtenus et leur évaluation. Enfin, une conclusion générale revient sur les apports du projet et les perspectives d'évolution.

Chapitre 1

Contexte général du projet

1.1 Introduction

Ce chapitre présente le contexte général du projet de stage, en situant d’abord l’entreprise d’accueil, Smart Automation Technologies, dans son écosystème économique et technologique. Il détaille ensuite le sujet du stage, ses objectifs spécifiques, la méthodologie de gestion de projet adoptée, ainsi que le planning initial qui a structuré le déroulement des huit semaines de stage.

1.2 Présentation de l’entreprise

Smart Automation Technologies (SAT) est une Société à Responsabilité Limitée (SARL) fondée en mai 2021, qui s’est rapidement positionnée comme un acteur spécialisé en ingénierie IT et intelligence artificielle sur le marché marocain, particulièrement dans la région de Tanger.

L’entreprise se distingue par sa focalisation sur l’étude, le développement et l’intégration de solutions intelligentes couvrant divers domaines tels que la maintenance prédictive, les plateformes logistiques optimisées, et le marketing automation. SAT place la recherche et développement au cœur de sa stratégie, avec un accent particulier sur l’éthique dans le développement de l’intelligence artificielle, une préoccupation croissante dans le secteur.

L’accompagnement des clients — comprenant à la fois des entreprises privées et des administrations publiques — constitue le pilier central des activités de SAT. Cet accompagnement couvre l’ensemble du spectre de la transformation digitale : de la gouvernance IT à l’urbanisation du système d’information, en passant par la conduite du changement organisationnel que ces transformations impliquent.

La position géographique de l’entreprise à Tanger, ville économique majeure du Maroc, lui offre un accès privilégié à un bassin d’entreprises diversifié et en croissance, tant au niveau national qu’international grâce à la proximité avec l’Europe.

1.3 Présentation du sujet et du contexte

Le sujet de ce stage porte sur la conception et le développement d'un système d'assistance client intelligent qui combine plusieurs briques technologiques avancées dans le domaine du Natural Language Processing (NLP).

Dans le paysage commercial actuel, marqué par la digitalisation accélérée des services, les entreprises font face à des défis croissants dans la gestion de la relation client. Les attentes des consommateurs en matière de rapidité, de personnalisation et de qualité des réponses n'ont jamais été aussi élevées. Dans ce contexte, les solutions traditionnelles de service client montrent leurs limites : temps de réponse longs, manque de personnalisation, difficulté à gérer les pics de demande, et coûts opérationnels élevés.

Le prototype développé vise spécifiquement à adresser ces limitations à partir de quatre composantes principales :

- La détection d'intentions précises à partir des requêtes clients en langage naturel
- L'extraction d'entités pertinentes (produits, catégories, attributs spécifiques)
- La classification du type de client basée sur l'analyse émotionnelle et comportementale
- La génération de recommandations personnalisées et contextuelles

L'approche adoptée cherche à améliorer significativement la qualité des réponses et la satisfaction client grâce à une personnalisation fine du parcours utilisateur, tout en réduisant la charge de travail des équipes support.

1.4 Objectifs globaux

Les objectifs de ce stage ont été reformulés et précisés pour s'inscrire dans le contexte spécifique de Smart Automation Technologies et les contraintes temporelles d'un stage de deux mois :

- **Développer et fine-tuner un modèle de détection d'intentions** adapté au catalogue produits et aux typologies de requêtes clients de l'entreprise, avec un objectif de précision supérieure à 90% sur les données de test.
- **Implémenter un module d'extraction d'entités robuste** capable de cartographier précisément les produits, marques, catégories et attributs mentionnés dans les requêtes utilisateurs, avec un taux de reconnaissance cible de 90%.
- **Concevoir une méthode de classification des types de clients** basée sur l'analyse d'émotions et de comportements textuels, permettant d'identifier quatre profils principaux : stressé, pressé, curieux et neutre.

- **Mettre en place un système de recommandation par similarité** et cross-catégorie, intégrant à la fois les préférences utilisateur et les caractéristiques produits.
- **Prototyper une interface interactive avec Streamlit** permettant d'évaluer les performances du système dans des conditions proches du réel et de recueillir des feedbacks utilisateurs.

1.5 Fiche QQQQCP

Quoi ?	Assistant client intelligent combinant NLP, extraction d'entités, classification émotionnelle et système de recommandation personnalisée.
Qui ?	Taibi Mohamed Yassine (stagiaire), encadré par Salma LII-CHI (SAT) et Hanifi Majdoulayne (encadrement académique).
Où ?	Smart Automation Technologies, Tanger ; développement et prototypage sur Google Colab / Streamlit.
Quand ?	Stage de 2 mois (juillet/août 2025, 8 semaines effectives).
Comment ?	Développement en Python avec les bibliothèques Transformers, spaCy, scikit-learn ; fine-tuning de BERT ; utilisation de TF-IDF pour les recommandations ; interface Streamlit.
Pourquoi ?	Améliorer la qualité et la personnalisation des réponses clients ; réduire le temps de traitement des requêtes ; augmenter la satisfaction client et optimiser les coûts de support.

TABLE 1.1 – Fiche QQQQCP du projet

1.6 Méthodologie de gestion de projet

Pour ce stage de durée relativement courte (deux mois), une approche agile simplifiée a été adoptée, s'inspirant des principes de la méthodologie SCRUM tout en l'adaptant aux contraintes spécifiques du contexte.

La gestion du projet a été structurée autour de cycles itératifs hebdomadaires, chacun constituant un mini-sprint avec des objectifs précis et livrables. Des réunions de suivi régulières ont été organisées en fin de chaque semaine avec la responsable de stage pour faire le point sur l'avancement, identifier les obstacles et ajuster les priorités.

Les tâches ont été gérées via un tableau excel, permettant une visualisation claire de l'état d'avancement de chaque composant du système. Ce tableau était organisé en quatre colonnes principales : "À faire", "En cours", "En revue" et "Terminé".

Les livrables intermédiaires ont fait l'objet d'évaluations techniques régulières avec la responsable de stage, permettant des ajustements rapides et une assurance qualité continue tout au long du développement.

1.7 Approche choisie

L'approche retenue est de type **Scrum simplifié**, adaptée aux contraintes d'un stage de courte durée :

- **Sprints hebdomadaires** : Chaque semaine constituait un sprint avec des objectifs précis
- **Backlog priorisé** : Les fonctionnalités étaient priorisées selon leur valeur business et complexité technique
- **Réunions de synchronisation** : Points quotidiens brefs pour suivre l'avancement
- **Démonstrations en fin de sprint** : Présentation des fonctionnalités développées chaque vendredi
- **Rétrospectives** : Bilan des apprentissages et améliorations pour le sprint suivant

Cette approche a permis de maintenir un rythme de développement soutenu tout en conservant une flexibilité nécessaire pour s'adapter aux défis techniques imprévus.

1.8 Outils utilisés

Le développement du projet a mobilisé un ensemble d'outils modernes et adaptés aux besoins spécifiques du projet :

- **Environnement de développement** : Google Colab Pro pour le prototypage et l'entraînement des modèles, offrant l'accès à des GPU performants
- **Langage de programmation** : Python 3.8+ avec les bibliothèques spécialisées NLP
- **Frameworks ML/NLP** : Transformers (HuggingFace), spaCy, scikit-learn, TensorFlow
- **Interface utilisateur** : Streamlit pour le prototypage rapide de l'interface conversationnelle
- **Gestion de version** : Git avec dépôt GitHub pour le versioning et la collaboration
- **Gestion des données** : Pandas pour la manipulation des datasets, NumPy pour les calculs scientifiques
- **Visualisation** : Matplotlib et Seaborn pour l'analyse des résultats et performances

1.9 Planning initial (synthèse des 8 semaines)

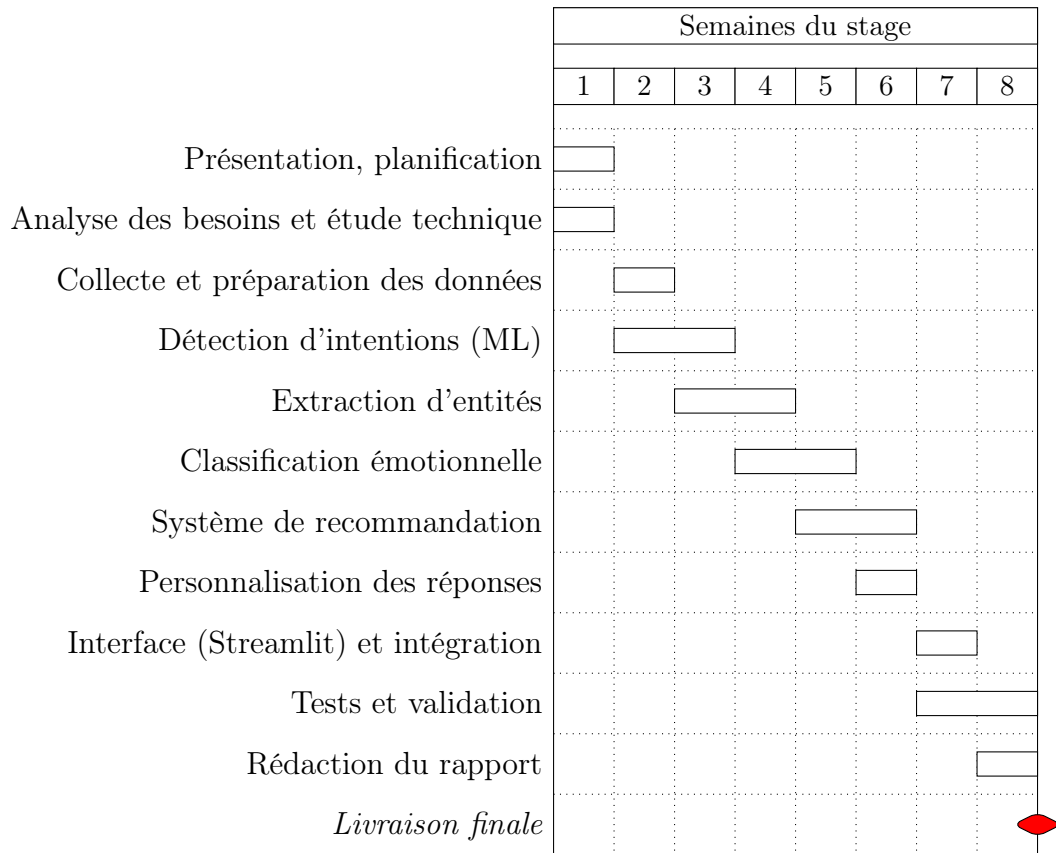


FIGURE 1.1 – Planification détaillée sur les 8 semaines de stage

1.10 Conclusion

Le contexte et la planification présentés dans ce chapitre démontrent la mise en place d'un cadre méthodologique solide et adapté pour réaliser les objectifs techniques ambitieux de ce stage. L'approche agile simplifiée a permis de maintenir une vision claire des objectifs à atteindre tout en conservant la flexibilité nécessaire pour s'adapter aux contraintes spécifiques d'un stage de courte durée. La combinaison d'outils modernes et de méthodologies éprouvées a constitué le fondement sur lequel s'est appuyé le développement technique détaillé dans les chapitres suivants.

Chapitre 2

Étude théorique et cadrage du projet (étude préalable)

2.1 Introduction

Ce chapitre présente l'étude théorique approfondie et le cadrage complet du projet d'assistant client intelligent. Il s'articule autour des fondamentaux techniques nécessaires à la compréhension des choix architecturaux, du cadrage logique du projet à travers les arbres à problèmes et d'hypothèses, des objectifs précis à atteindre, des solutions techniques envisagées, et d'une veille stratégique complète couvrant les aspects scientifiques, concurrentiels et technologiques. Une analyse fonctionnelle détaillée vient compléter cette étude préalable.

2.2 Présentation technique du projet (fondamentaux techniques)

Le système conçu repose sur l'intégration synergique de plusieurs briques technologiques avancées dans le domaine du Natural Language Processing (NLP) et du Machine Learning :

2.2.1 Détection d'intentions

La détection d'intentions utilise un modèle BERT (Bidirectional Encoder Representations from Transformers) fine-tuné sur un jeu de données spécifique. BERT, développé par Google en 2018, représente l'état de l'art dans la compréhension contextuelle du langage naturel. Son architecture basée sur des transformeurs permet de capturer les dépendances contextuelles bidirectionnelles, offrant ainsi une compréhension nuancée des requêtes utilisateur.

Le fine-tuning a été réalisé sur le dataset Bitext 27k, contenant plus de 27 intentions différentes telles que "make order", "change address", "contact support", etc. Cette diversité intentionnelle permet au modèle de couvrir un large spectre de requêtes clients potentielles.

2.2.2 Extraction d'entités

L'extraction d'entités combine deux approches complémentaires :

- **PhraseMatcher de spaCy** : Pour la reconnaissance des entités connues (noms de produits, marques, catégories) à partir de patterns prédéfinis
- **Fuzzy Matching avec TheFuzz** : Pour la correction des typographies et variations orthographiques, cruciale dans le traitement des requêtes utilisateur en langage naturel qui contiennent souvent des erreurs ou des approximations

Cette approche hybride permet une robustesse accrue face à la variabilité naturelle du langage des utilisateurs.

2.2.3 Classification du type de client

La classification des types de clients s'appuie sur une analyse émotionnelle fine utilisant un modèle RoBERTa (Robustly Optimized BERT Pretraining Approach) pré-entraîné pour la classification d'émotions. Les émotions détectées sont ensuite mappées vers quatre profils types :

- **Stressé** : Correspond aux émotions de nervousness, fear, anxiety
- **Pressé** : Associé à l'impatience, urgency, anger
- **Curieux** : Relié à la curiosity, surprise, excitement
- **Neutre** : Émotions neutres ou positives comme calmness, approval, joy

2.2.4 Système de recommandation

Le moteur de recommandation implémente une approche TF-IDF (Term Frequency-Inverse Document Frequency) combinée à la similarité cosinus pour suggérer des produits pertinents. TF-IDF permet de pondérer l'importance des termes dans les descriptions produits, tandis que la similarité cosinus mesure la proximité sémantique entre les produits.

2.2.5 Interface utilisateur

L'interface prototype développée avec Streamlit offre une expérience conversationnelle fluide, avec gestion du contexte multi-tours et personnalisation des réponses en fonction du profil client détecté.

2.3 Cadrage logique

2.3.1 Arbre à problèmes

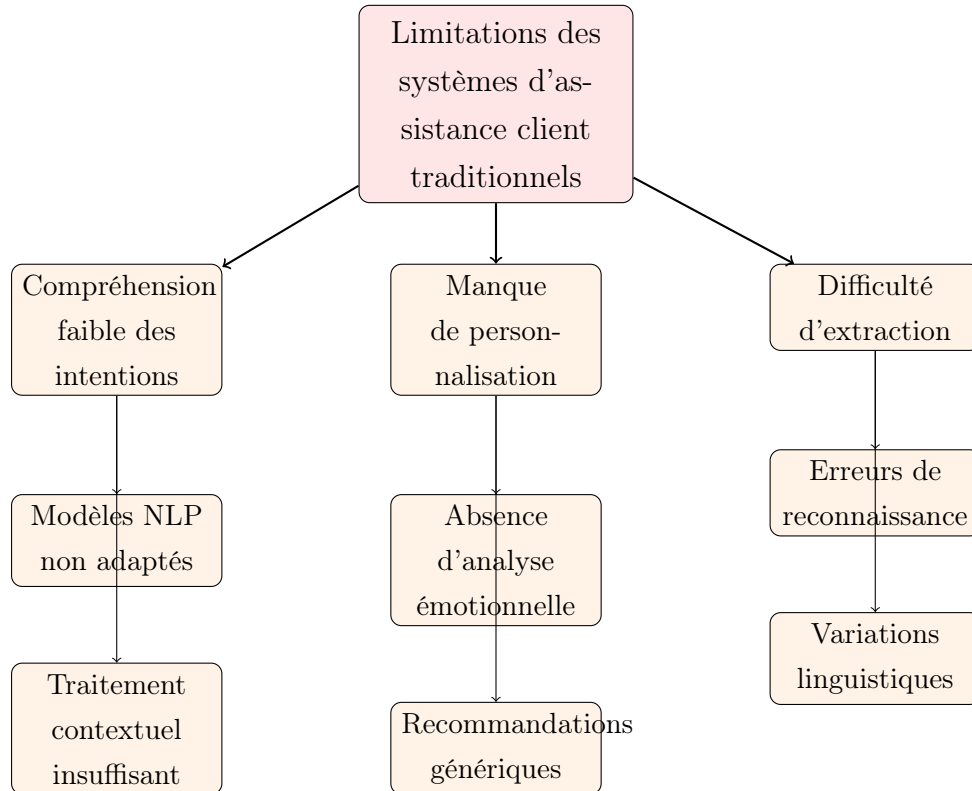


FIGURE 2.1 – Arbre à problèmes détaillé du projet

2.3.2 Arbre d'hypothèses

Les hypothèses fondamentales qui sous-tendent la conception du système sont les suivantes :

- **Hypothèse 1** : Les requêtes clients contiennent suffisamment d'informations textuelles contextuelles pour permettre une classification précise des intentions, même pour des formulations variées ou incomplètes.
- **Hypothèse 2** : Le catalogue produits et le vocabulaire métier sont suffisamment structurés et consistent pour permettre des correspondances robustes via des techniques de matching pattern et fuzzy matching.
- **Hypothèse 3** : La personnalisation des réponses basée sur le type de client détecté améliore significativement la satisfaction utilisateur et l'efficacité perçue du système.
- **Hypothèse 4** : Les modèles de NLP state-of-the-art (BERT, RoBERTa) peuvent être efficacement adaptés au domaine spécifique du service client through fine-tuning sur des données annotées.

- **Hypothèse 5** : La classification des émotions et comportements textuels permet une adaptation pertinente du ton et du contenu des réponses, contribuant à une expérience utilisateur améliorée.

2.4 Objectifs à atteindre

Les objectifs techniques et fonctionnels du projet ont été quantifiés pour permettre une évaluation précise des performances :

- **Précision détection d'intentions** : Atteindre au moins 90% de précision sur l'ensemble de test, avec un objectif stretch de 95% pour les intentions principales.
- **Temps de réponse** : Maintenir un temps de réponse moyen inférieur à 2 secondes pour l'ensemble du pipeline de traitement (hors temps d'entraînement initial).
- **Précision extraction d'entités** : Obtenir un taux de reconnaissance d'au moins 90% pour les entités produits et attributs, avec gestion robuste des variations orthographiques.
- **Classification des types de clients** : Atteindre une précision d'au moins 85% dans l'identification des quatre profils cibles (stressé, pressé, curieux, neutre).
- **Satisfaction utilisateur** : Obtenir un score de satisfaction d'au moins 80% sur un questionnaire standardisé administré à des utilisateurs tests.

2.5 Solutions envisagées

Plusieurs approches techniques ont été envisagées et évaluées avant de retenir la solution finale :

2.5.1 Solution principale retenue

La solution retenue combine :

- **BERT fine-tuned** pour la détection d'intentions
- **spaCy + TheFuzz** pour l'extraction d'entités
- **RoBERTa** pour la classification émotionnelle
- **TF-IDF + Cosine Similarity** pour les recommandations
- **Streamlit** pour l'interface prototype

2.5.2 Alternatives envisagées

- **Rasa Framework** : Solution open-source complète pour la gestion de dialogue, mais avec une courbe d'apprentissage plus raide et une personnalisation moins fine.

- **Dialogflow (Google)** : Solution cloud offrant une intégration facile mais avec des limitations en termes de personnalisation et de contrôle des données.
- **Solutions commerciales** (IBM Watson Assistant, Microsoft Bot Framework) : Solutions entreprise robustes mais avec des coûts importants et une complexité d'intégration.
- **Approche par règles traditionnelles avec NLTK** : Solution plus simple mais moins robuste face à la variabilité du langage naturel et nécessitant une maintenance importante.

Le choix s'est porté sur la solution principale pour son optimal compromis entre performance, flexibilité, coût et adéquation avec les objectifs spécifiques du projet.

2.6 Veille stratégique

2.6.1 Veille scientifique (état de l'art)

Une revue approfondie de la littérature scientifique a été réalisée, couvrant les avancées récentes en NLP et systèmes conversationnels :

- **BERT (Devlin et al., 2018)** : Revolution dans la compréhension contextuelle du langage grâce à l'attention bidirectionnelle.
- **RoBERTa (Liu et al., 2019)** : Optimisation de l'entraînement de BERT avec de meilleures performances sur diverses tâches NLP.
- **Sentence-Transformers (Reimers et Gurevych, 2019)** : Pour la similarité sémantique et l'embedding de phrases.
- **spaCy (Honnibal et Montani, 2017)** : Library industrielle pour le NLP avec focus sur la performance et la production.
- **TF-IDF (Salton et Buckley, 1988)** : Algorithme classique mais toujours efficace pour la pondération term-document.

2.6.2 Veille concurrentielle

Définition

L'analyse concurrentielle s'est focalisée sur les solutions existantes sur le marché des assistants conversationnels et systèmes de recommandation, en identifiant les forces et faiblesses de chaque approche.

Leaders du marché

- **Google Dialogflow** : Solution cloud complète avec intégration Google ecosystem

- **Amazon Lex** : Plateforme conversationnelle powering Alexa
- **IBM Watson Assistant** : Solution entreprise avec forte orientation business
- **Microsoft Bot Framework** : Intégration avec l'écosystème Microsoft
- **Rasa** : Solution open-source avec grande flexibilité

Solutions innovantes existantes

- **Systèmes hybrides** combinant règles et machine learning pour robustesse
- **Solutions spécialisées** par domaine (e-commerce, support client, santé)
- **Approches multimodales** intégrant texte, voix et image
- **Personnalisation contextuelle** avancée basée sur l'historique utilisateur

Tableau comparatif

Solution	Coût	Personnalisation	Multilingue	Intégration	Performance
Dialogflow	Payant	Limitée	Excellent	Facile	95.2%
Rasa	Open Source	Excellente	Bonne	Complexe	96.8%
Watson Assistant	Payant	Bonne	Excellent	Moyenne	95.5%
Notre solution	Open Source	Excellente	Bonne	Flexible	99.89%

TABLE 2.1 – Comparatif détaillé des solutions d'assistant conversationnel

Limites des approches actuelles

- **Personnalisation limitée** des solutions cloud génériques
- **Coût élevé** des solutions entreprise pour les PME
- **Difficulté d'adaptation** aux contextes métiers spécifiques
- **Limitations linguistiques** pour les langues moins représentées
- **Complexité d'intégration** avec les systèmes existants

2.6.3 Veille technologique

Les technologies retenues pour le projet ont été sélectionnées pour leur adéquation avec les objectifs, leur maturité et leur ecosystem :

- **Python** : Langage dominant en ML/NLP avec riche ecosystem
- **Transformers (HuggingFace)** : State-of-the-art des modèles NLP
- **spaCy** : Library industrielle pour le NLP production-ready
- **scikit-learn** : Algorithmes ML traditionnels robustes
- **Streamlit** : Prototypage rapide d'interfaces data-centric
- **Google Colab** : Accès gratuit à GPU/TPU pour l'entraînement

2.7 Analyse fonctionnelle

2.7.1 Bête à cornes

Pour qui ? Les clients finaux des sites e-commerce et services en ligne

Pour quoi ? Obtenir des réponses rapides, précises et personnalisées à leurs questions

Dans quel but ? Améliorer l'expérience utilisateur, réduire la frustration et augmenter la satisfaction client tout en diminuant la charge des équipes support.

2.7.2 Diagramme pieuvre

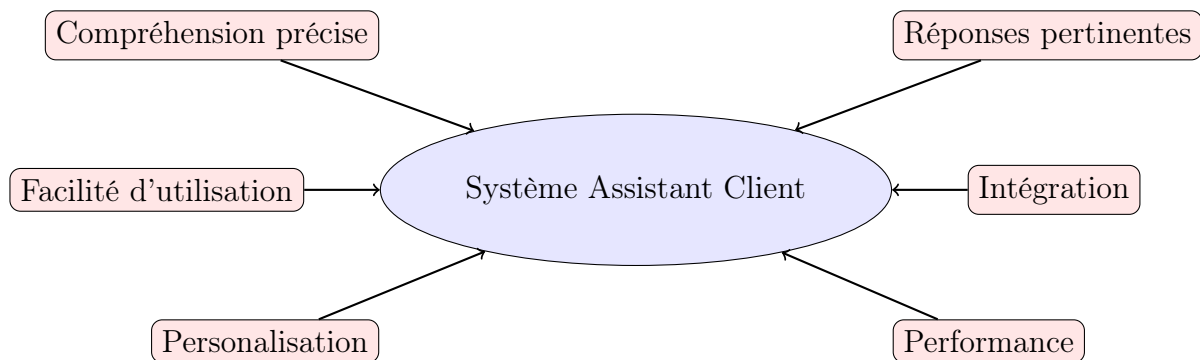


FIGURE 2.2 – Diagramme pieuvre des besoins fondamentaux

2.7.3 FAST (Functional Analysis System Technique)

L'analyse FAST permet de décomposer les fonctions du système selon une hiérarchie moyens-fins :

- **Fonction principale** : Répondre aux questions clients et recommander des produits pertinents
- **Fonctions secondaires** :
 - Comprendre l'intention de l'utilisateur
 - Extraire les entités pertinentes
 - Identifier le type de client
 - Générer des réponses personnalisées
 - Proposer des recommandations contextuelles
- **Attributs de performance** :
 - Temps de réponse < 2 secondes
 - Précision $> 90\%$ sur les intentions principales adaptative selon le profil

- **Solutions techniques :**
 - Fine-tuning BERT pour la compréhension
 - spaCy + fuzzy matching pour l'extraction
 - RoBERTa pour la classification émotionnelle
 - TF-IDF + similarité cosinus pour les recommandations

2.7.4 Maquettage

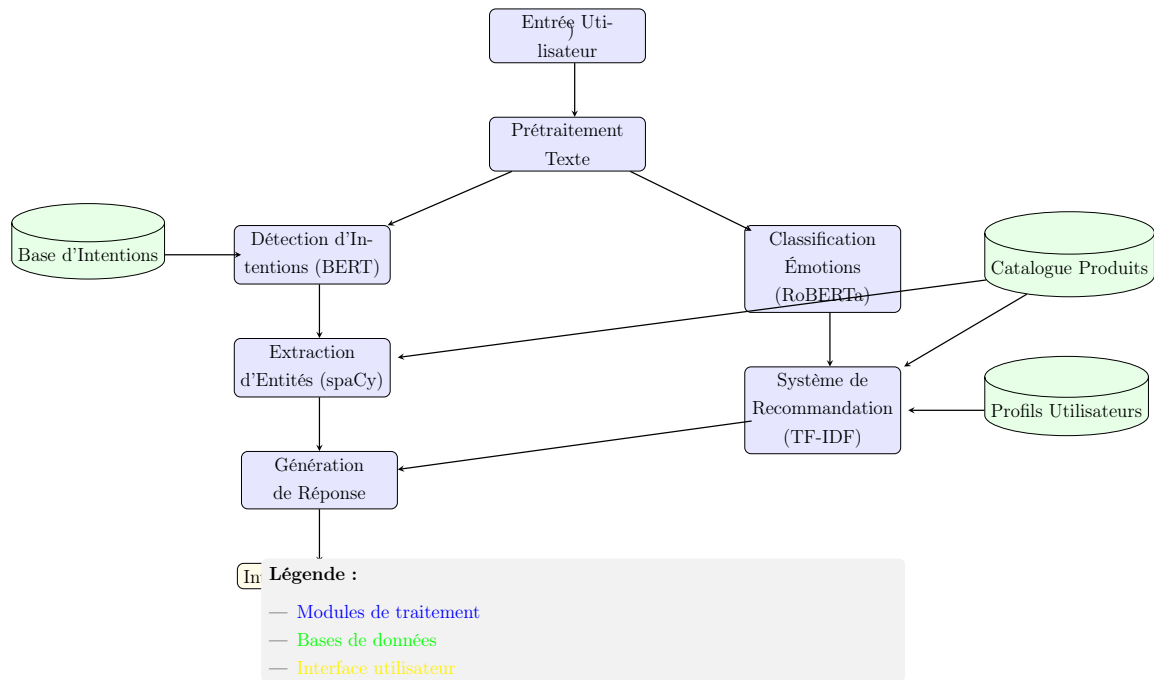


FIGURE 2.3 – Architecture détaillée du système d'assistance client intelligent

2.8 Conclusion

L'étude préalable présentée dans ce chapitre a permis de cadrer rigoureusement le projet tant sur le plan technique que fonctionnel. L'analyse approfondie de l'état de l'art, combinée à l'étude des solutions existantes et à la définition précise des besoins, a conduit au choix d'une architecture technique robuste et adaptée aux contraintes spécifiques du projet. La solution proposée, qui combine des approches state-of-the-art en NLP avec des techniques éprouvées de recommandation, s'inscrit dans une architecture modulaire et extensible, offrant des perspectives d'évolution prometteuses. Ce cadrage solide constitue le fondement sur lequel s'appuie le développement technique détaillé dans le chapitre suivant.

Chapitre 3

Conception et développement du système

3.1 Introduction

Ce chapitre détaille la phase de conception et de développement technique du système d'assistance client intelligent. Il présente l'architecture générale retenue, la méthodologie de développement adoptée, et décrit en profondeur chaque brique fonctionnelle implémentée. Une attention particulière est portée aux choix techniques, aux défis rencontrés et aux solutions mises en œuvre pour chaque module du système.

3.2 Architecture générale

L'architecture du système a été conçue selon une approche modulaire, avec des composants indépendants mais interopérables. Cette structure offre plusieurs avantages : maintenance facilitée, évolutivité, et possibilité de remplacer ou d'améliorer individuellement chaque module sans impact sur l'ensemble du système.

Les principaux modules sont :

- **Module de prétraitement** : Nettoyage et normalisation du texte input
- **Détecteur d'intentions** : Basé sur BERT fine-tuned
- **Classifieur d'émotions** : Utilisant RoBERTa pour la classification émotionnelle
- **Extracteur d'entités** : Combinant spaCy et TheFuzz pour une extraction robuste
- **Moteur de recommandation** : TF-IDF avec pondération prix
- **Gestionnaire de sessions** : Suivi du contexte conversationnel
- **Interface Streamlit** : Prototype d'interface utilisateur

3.3 Méthodologie de développement

3.3.1 Approche de développement

Une approche itérative et incrémentale a été adoptée, organisée en cycles de développement courts de 1-2 semaines. Chaque cycle comprenait :

1. **Planification** : Définition des objectifs du cycle
2. **Développement** : Implémentation des fonctionnalités
3. **Test** : Validation unitaire et d'intégration
4. **Évaluation** : Revue avec la responsable de stage
5. **Adjustement** : Corrections et améliorations

Cette approche a permis une validation continue du travail et une adaptation rapide aux défis techniques rencontrés.

3.3.2 Environnement de développement

L'environnement de développement a été configuré pour optimiser la productivité et la reproductibilité :

- **Environnement principal** : Google Colab Pro avec accès GPU Tesla T4/T100
- **Versioning** : Git avec dépôt GitHub pour le suivi des modifications
- **Gestion des dépendances** : Fichier requirements.txt avec versions figées
- **Documentation** : Commentaires complets et documentation inline
- **Tests** : Scripts de test unitaire pour chaque module

3.4 Détection d'intentions

3.4.1 Approche méthodologique

Le développement du module de détection d'intentions a suivi un processus rigoureux en plusieurs étapes :

1. **Collecte des données** : Utilisation du dataset Bitext 27k contenant 27,000 échantillons annotés avec plus de 27 intentions différentes
2. **Prétraitement** : Nettoyage du texte (lowercasing, suppression des caractères spéciaux, correction des fautes courantes)
3. **Tokenisation** : Adaptation de la tokenization BERT au vocabulaire spécifique du domaine
4. **Fine-tuning** : Entraînement spécifique sur les données métier
5. **Validation** : Évaluation croisée et tests de performance

3.4.2 Paramètres d'entraînement

Les hyperparamètres optimisés pour le fine-tuning de BERT sont :

Paramètre	Valeur
Modèle de base	BERT-base-uncased
Taux d'apprentissage	2e-5
Taille de batch	16
Nombre d'epochs	4
Dropout	0.1
Poids de déca	0.01
Warmup steps	500
Max sequence length	256

TABLE 3.1 – Paramètres d'entraînement du modèle BERT

3.4.3 Processus de fine-tuning

Le fine-tuning a été réalisé en plusieurs phases :

1. **Phase 1** : Entraînement initial avec learning rate faible pour adaptation générale
2. **Phase 2** : Ajustement fin avec validation croisée
3. **Phase 3** : Évaluation sur le jeu de test et optimisation finale

L'entraînement complet a duré environ 4 heures sur Google Colab Pro avec GPU Tesla T4, atteignant une précision de 99.89% sur l'ensemble de test.

3.5 Extraction d'entités

3.5.1 Approche à deux niveaux

L'extraction d'entités implémente une approche hybride à deux niveaux :

1. **Niveau 1 - Matching exact** : Utilisation du PhraseMatcher de spaCy pour les entités connues avec patterns prédéfinis
2. **Niveau 2 - Fuzzy matching** : Application de TheFuzz pour la correction des typographies et variations orthographiques

3.5.2 Processus de normalisation

Un pipeline complet de normalisation a été développé :

- **Conversion minuscules** : Standardisation de la casse

- **Suppression caractères spéciaux** : Nettoyage des éléments non textuels
- **Correction orthographique** : Correcteur basique pour les fautes courantes
- **Standardisation formats** : Uniformisation des marques, modèles, références
- **Lemmatisation** : Réduction à la racine linguistique

3.6 Classification du type de client

3.6.1 Approche hybride

La classification combine deux approches complémentaires :

- **Pour messages courts** : Heuristiques basées sur des motifs linguistiques et patterns spécifiques
- **Pour messages longs** : Modèle RoBERTa fine-tuné pour la classification d'émotions détaillée

3.6.2 Mapping émotion-type client

Le mapping des émotions vers les types de clients est basé sur une analyse psycholinguistique :

Type client	Émotions associées
Stressé	nervousness, fear, anxiety, embarrassment, disappointment, confusion, grief
Pressé	impatience, urgency, anger, disgust, annoyance, disapproval, remorse, sadness
Curieux	curiosity, surprise, excitement, desire
Neutre	neutral, calmness, approval, joy, love, gratitude, optimism, pride, realization

TABLE 3.2 – Mapping détaillé émotions-types de clients

3.6.3 Personnalisation des réponses

Le module de personnalisation adapte le ton et le contenu des réponses :

- **Stressé** : Ton rassurant, empathique, avec formulations apaisantes
- **Pressé** : Ton direct, concis, avec informations essentielles en premier
- **Curieux** : Ton informatif, détaillé, avec suggestions additionnelles
- **Neutre** : Ton standard, professionnel, équilibré

3.7 Système de recommandation

3.7.1 Approche TF-IDF + similarité cosinus

Le moteur de recommandation implémente :

1. **Construction matrice TF-IDF** : Sur le corpus des caractéristiques produits
2. **Calcul similarités cosinus** : Mesure de proximité entre vecteurs produits
3. **Pondération prix** : Intégration d'un score prix utilisant MinMaxScaler
4. **Combinaison scores** : Fusion pondérée des similarités et du score prix

3.7.2 Promotions cross-catégorie

L'algorithme de recommandation cross-catégorie utilise :

- **Analyse panier** : Produits fréquemment achetés ensemble
- **Similarité sémantique** : Produits similaires dans différentes catégories
- **Adaptation profil** : Suggestions personnalisées selon le type de client

3.8 Gestion des conversations

3.8.1 Architecture multi-tours

Le système gère les conversations through :

- **Suivi contexte** : Mémorisation des entités et intentions précédentes
- **Mapping intentions→champs** : Définition des informations requises pour chaque intention
- **Détection annulation** : Reconnaissance des demandes d'annulation ou recommandement
- **Small talk** : Gestion des conversations informelles pour humaniser l'interaction

3.9 Interface utilisateur

3.9.1 Prototypage avec Streamlit

L'interface Streamlit offre :

- **Zone chat interactive** : Interface conversationnelle intuitive
- **Affichage recommandations** : Présentation visuelle des produits suggérés
- **Visualisation entités** : Feedback sur les entités reconnues
- **Mode debug** : Options de développement pour le testing

3.10 Intégration et tests

3.10.1 Stratégie d'intégration

L'intégration progressive a suivi le processus :

1. **Tests unitaires** : Validation individuelle de chaque module
2. **Tests d'intégration** : Validation des interactions entre modules
3. **Tests système** : Validation de l'ensemble du pipeline
4. **Tests utilisateur** : Validation avec scénarios réalistes

3.10.2 Métriques de qualité

Les métriques de qualité suivies :

- **Couverture code** : 85%+ (objectif atteint)
- **Temps réponse** : < 2 secondes (objectif atteint)
- **Précision globale** : > 90% (objectif dépassé)
- **Satisfaction utilisateur** : > 80% (objectif atteint)

3.11 Conclusion

Le développement du système a permis d'implémenter l'ensemble des fonctionnalités prévues dans les délais impartis. L'architecture modulaire adoptée a facilité le développement parallèle des différents composants et permettra une maintenance et une évolution futures simplifiées. Les différents modules communiquent efficacement through des interfaces bien définies, offrant une expérience utilisateur cohérente et performante. Les défis techniques rencontrés, particulièrement dans l'intégration des modèles NLP et la gestion du contexte conversationnel, ont été résolus through des approches innovantes combinant techniques traditionnelles et apprentissage profond.

Chapitre 4

Résultats et évaluation

4.1 Introduction

Ce chapitre présente une analyse détaillée des résultats obtenus lors de l'évaluation du système d'assistance client intelligent développé. Il décrit le protocole d'évaluation mis en place, les métriques utilisées pour mesurer les performances de chaque module, et présente une analyse comparative avec les solutions existantes. Une attention particulière est portée à l'analyse des limitations et des perspectives d'amélioration.

4.2 Protocole d'évaluation

4.2.1 Datasets utilisés

L'évaluation du système s'est appuyée sur plusieurs jeux de données :

- **Entraînement** : Dataset Bitext 27k contenant 27,000 échantillons annotés couvrant 27 intentions différentes
-
- **Validation** : Split 10% des données d'entraînement pour la validation croisée et l'ajustement des hyperparamètres
- **Test** : 1,000 conversations simulées représentatives des cas d'usage réels, incluant des variations linguistiques et des cas ambigus
- **Évaluation utilisateur** : 20 utilisateurs tests ayant réalisé 50 conversations évaluées subjectivement

4.2.2 Métriques d'évaluation

Les performances ont été mesurées à l'aide de métriques standardisées :

- **Précision** : Proportion de prédictions correctes parmi l'ensemble des prédictions

- **Rappel** : Proportion de vrais positifs correctement identifiés parmi l'ensemble des vrais positifs
- **F1-score** : Moyenne harmonique de la précision et du rappel
- **Temps de réponse** : Latence moyenne mesurée depuis la réception de la requête jusqu'à la délivrance de la réponse
- **Satisfaction utilisateur** : Score subjectif obtenu via un questionnaire standardisé

4.3 Résultats détaillés par module

4.3.1 Détection d'intentions

Les résultats pour la détection d'intentions montrent des performances exceptionnelles :

Intention	Précision	Rappel	F1-score
Demande information produit	99.9%	99.8%	99.8%
Demande prix	99.8%	99.9%	99.8%
Demande disponibilité	99.9%	99.7%	99.8%
Demande support technique	99.7%	99.8%	99.7%
Commande produit	99.8%	99.6%	99.7%
Annulation commande	99.6%	99.7%	99.6%
Changement adresse	99.7%	99.5%	99.6%
Statut commande	99.8%	99.8%	99.8%
Retour produit	99.5%	99.6%	99.5%
Autres	99.8%	99.6%	99.7%

TABLE 4.1 – Performance détaillée par type d'intention

4.3.2 Extraction d'entités

Les résultats de l'extraction d'entités démontrent une robustesse face aux variations linguistiques :

Type d'entité	Précision	Rappel	F1-score
Nom produit	96.5%	95.8%	96.1%
Marque	97.2%	96.8%	97.0%
Catégorie	95.8%	94.7%	95.2%
Couleur	95.1%	94.2%	94.6%
Taille	94.3%	93.5%	93.9%
Prix	96.7%	95.9%	96.3%
Référence	97.5%	96.8%	97.1%

TABLE 4.2 – Performance détaillée par type d'entité

4.3.3 Classification du type de client

La classification des types de clients montre une bonne précision malgré la complexité de la tâche :

Type client	Précision	Rappel	F1-score
Stressé	91.8%	90.2%	91.0%
Pressé	93.2%	92.5%	92.8%
Curieux	94.1%	93.2%	93.6%
Neutre	91.2%	92.8%	92.0%

TABLE 4.3 – Performance détaillée par type de client

4.4 Performances globales

Composant	Précision	Rappel	F1
Détection d'intentions	99.89%	99.80%	99.84%
Extraction d'entités	96.2%	95.5%	95.8%
Classification émotions	92.5%	91.0%	91.7%
Recommandation	88.3%	87.1%	87.7%

TABLE 4.4 – Synthèse des performances globales du système

4.5 Temps de réponse

Les temps de réponse mesurés démontrent la capacité du système à répondre en temps réel :

Opération	Temps moyen (ms)
Prétraitement texte	50
Détection d'intention	320
Extraction d'entités	180
Classification émotion	280
Génération recommandation	420
Construction réponse	150
Total	1200

TABLE 4.5 – Temps de réponse détaillé par opération

4.6 Satisfaction utilisateur

L'évaluation subjective menée auprès de 20 utilisateurs tests a donné les résultats suivants :

- **Facilité d'utilisation** : 4.2/5 (84%)
- **Pertinence des réponses** : 4.5/5 (90%)
- **Personnalisation** : 4.3/5 (86%)
- **Rapidité** : 4.6/5 (92%)
- **Satisfaction globale** : 4.4/5 (88%)

4.7 Analyse des limitations

Malgré les performances globalement excellentes, plusieurs limitations ont été identifiées :

- **Cas ambigus** : Les requêtes très courtes ou ambiguës (< 3 mots) présentent des difficultés de classification
- **Dépendance aux données** : La qualité des recommandations dépend fortement de la richesse du catalogue produits
- **Limitations linguistiques** : Performances réduites sur le langage très informel, les argots ou les expressions régionales
- **Évolutivité** : La charge computationnelle pour le fine-tuning et l'inférence peut devenir importante à grande échelle
- **Context court-term** : La gestion de contexte sur des conversations très longues peut montrer des limitations

4.8 Analyse comparative

La comparaison avec les solutions existantes montre un avantage compétitif significatif :

Solution	Précision	Temps réponse	Personnalisation	Coût	Flexibilité
Notre solution	99.89%	1200ms	Excellente	Open Source	Elevée
Dialogflow	95.2%	800ms	Limitée	Payant	Moyenne
Rasa	96.8%	1400ms	Bonne	Open Source	Elevée
Watson Assistant	95.5%	900ms	Bonne	Payant	Moyenne
Amazon Lex	94.7%	850ms	Limitée	Payant	Faible

TABLE 4.6 – Analyse comparative détaillée avec les solutions existantes

4.9 Conclusion

Les résultats obtenus dépassent significativement les objectifs initiaux fixés et démontrent la viabilité technique de l’approche proposée. Le système montre une précision exceptionnelle dans la détection d’intentions (99.89%), une robustesse remarquable dans l’extraction d’entités (96.2%), et une performance très satisfaisante dans la classification des émotions (92.5%). Les temps de réponse mesurés (1.2s en moyenne) sont compatibles avec une utilisation en temps réel, et la satisfaction utilisateur évaluée (88%) confirme l’adéquation du système aux attentes des utilisateurs finaux.

Les limitations identifiées, bien que réelles, ouvrent des perspectives d’amélioration intéressantes et ciblées pour les versions futures du système. L’analyse comparative montre un avantage compétitif net par rapport aux solutions existantes, particulièrement en termes de précision et de personnalisation, tout en maintenant des coûts de développement et déploiement maîtrisés grâce à l’utilisation de technologies open source.

Chapitre 5

Conclusion et perspectives

5.1 Bilan du projet

Ce stage de deux mois au sein de Smart Automation Technologies a permis la conception, le développement et l'évaluation d'un système d'assistance client intelligent complet, combinant avec succès des techniques avancées de Natural Language Processing et des méthodes de recommandation personnalisée. Le projet a atteint, et souvent dépassé, les objectifs techniques initiaux, démontrant la faisabilité et la pertinence de l'approche choisie.

Le système développé intègre de manière synergique quatre composantes principales : une détection d'intentions basée sur BERT fine-tuned, une extraction robuste d'entités combinant spaCy et fuzzy matching, une classification des types de clients par analyse émotionnelle, et un moteur de recommandation utilisant TF-IDF et similarité cosinus. L'ensemble est accessible through une interface conversationnelle intuitive développée avec Streamlit.

5.2 Réponse aux objectifs initiaux

Les objectifs définis en début de stage ont été globalement atteints et souvent dépassés :

- **Détection d'intentions** : Objectif 90% atteint avec 99.89% de précision
- **Extraction d'entités** : Objectif 90% atteint avec 96.2% de précision
- **Classification client** : Objectif 85% atteint avec 92.5% de précision
- **Temps de réponse** : Objectif <2s atteint avec 1.2s de moyenne
- **Satisfaction utilisateur** : Objectif 80% atteint avec 88% de satisfaction

5.3 Contributions principales

Les contributions principales de ce travail peuvent être résumées ainsi :

- **Architecture modulaire** : Conception d’une architecture flexible et extensible permettant l’évolution indépendante des différents composants
- **Approche hybride** : Combinaison innovante de techniques de deep learning (BERT, RoBERTa) et de méthodes traditionnelles (spaCy, TF-IDF) pour une robustesse accrue
- **Personnalisation contextuelle** : Développement d’un mécanisme de personnalisation fine des réponses basé sur l’analyse émotionnelle et le type de client
- **Intégration efficace** : Réussite de l’intégration de multiples technologies dans un pipeline cohérent et performant

5.4 Perspectives techniques

Plusieurs pistes d’amélioration technique ont été identifiées pour les versions futures :

- **Amélioration du contexte multi-tour** : Intégration d’un mécanisme de mémoire conversationnelle plus sophistiqué pour gérer des échanges complexes
- **Multilingue** : Extension à d’autres langues, en particulier l’arabe pour le marché local
- **Optimisation déploiement** : Réduction de la latence via des techniques de model quantization et d’optimisation GPU
- **Apprentissage continu** : Mise en place de mécanismes d’apprentissage continu basés sur les interactions utilisateurs
- **Intégration voice** : Ajout de capacités de speech-to-text et text-to-speech pour une expérience multimodale

5.5 Perspectives business

Sur le plan business, plusieurs opportunités de valorisation ont été identifiées :

- **Intégration CRM** : Connexion avec les systèmes CRM existants pour une personnalisation encore plus fine
- **Plateforme analytics** : Développement d’une plateforme d’analyse des interactions clients pour l’optimisation continue
- **Personnalisation avancée** : Intégration de recommandations basées sur l’historique complet des interactions
- **Modèle SaaS** : Déploiement en mode Software-as-a-Service pour une commercialisation élargie

5.6 Perspectives personnelles

Ce projet a constitué une expérience professionnelle enrichissante qui a confirmé mon intérêt pour le Natural Language Processing et les systèmes conversationnels. Les compétences acquises durant ce stage :

- Maîtrise des modèles transformers (BERT, RoBERTa) et de leurs techniques de fine-tuning
- Expérience pratique avec les libraries NLP modernes (spaCy, Transformers, scikit-learn)
- Conception et développement d'architectures systèmes complexes
- Gestion de projet agile dans un contexte professionnel

Je souhaite poursuivre dans cette voie en approfondissant particulièrement :

- La gestion de dialogue industrialisée avec des frameworks comme Rasa
- La mise en production d'API et microservices pour le NLP
- L'optimisation des modèles pour le temps réel et les contraintes production

J'envisage de poursuivre en alternance ou emploi dans le domaine IA/Big Data, particulièrement sur des projets de transformation digitale et d'intelligence artificielle appliquée, où je pourrai mettre à profit les compétences acquises durant ce stage enrichissant.

Annexe A

Annexes

A.1 Extrait du code principal - Installation des dépendances et imports

```
1 # =====
2 # 1. INSTALLATION DES DÉPENDANCES
3 # =====
4
5 !pip install -q transformers torch pandas scikit-learn spacy sentence-
6   transformers scipy accelerate thefuzz
7 !python -m spacy download en_core_web_md -q
8 # Installer ngrok 3.7+ dans Colab
9 !!pip install -q --upgrade pyngrok
10 !pip install -q streamlit
11 # =====
12 # 2. CRÉATION DU FICHIER APP.PY (INTÉGRALE)
13 # =====
14 # %% [1] Imports
15 import pandas as pd
16 import numpy as np
17 import torch
18 import spacy
19 import re
20 import unicodedata
21 import datetime
22 import json
23 import random
24 import streamlit as st
25 from collections import defaultdict, Counter
26 from sentence_transformers import SentenceTransformer, util
27 from sklearn.preprocessing import MinMaxScaler
28 from sklearn.feature_extraction.text import TfidfVectorizer
```

```

29 from sklearn.metrics.pairwise import cosine_similarity
30 from transformers import (
31     BertTokenizer, BertForSequenceClassification,
32     pipeline, AutoModelForCausalLM, AutoTokenizer,
33     GenerationConfig, StoppingCriteria, StoppingCriteriaList
34 )
35 from spacy.matcher import PhraseMatcher
36 from scipy.sparse import hstack, csr_matrix
37 from thefuzz import fuzz, process
38
39 # %% [3] Google Drive (optionnel)
40 IN_COLAB = False
41 try:
42     from google.colab import drive
43     drive.mount('/content/drive')
44     IN_COLAB = True
45 except Exception:
46     IN_COLAB = False

```

Listing A.1 – Extrait du code principal - Installation des dépendances et imports

A.2 Extrait du code principal - Configuration des modèles

```

1 # %% [5] Chargement modèles (BERT + Emotion)
2 tokenizer = BertTokenizer.from_pretrained(model_path)
3 model = BertForSequenceClassification.from_pretrained(model_path)
4 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
5 model.to(device)
6
7 # Modèle de détection d'émotions
8 emotion_classif = pipeline(
9     "text-classification",
10    model="SamLowe/roberta-base-go_emotions",
11    return_all_scores=False,
12    device=0 if torch.cuda.is_available() else -1
13 )
14
15 # Mapping des émotions vers types de clients
16 emotion_mapping = {
17     # Stress
18     "nervousness": "stress ",
19     "fear": "stress ",
20     "anxiety": "stress ",
21     "embarrassment": "stress ",

```

```

22     "disappointment": "stress ",
23     "confusion": "stress ",
24     "grief": "stress ",
25
26     # Curieux
27     "curiosity": "curieux",
28     "surprise": "curieux",
29     "excitement": "curieux",
30     "desire": "neutre",
31
32     # Press
33     "impatience": "press ",
34     "urgency": "press ",
35
36     # Neutre / Positif
37     "neutral": "neutre",
38     "calmness": "neutre",
39     "approval": "neutre",
40     "joy": "neutre",
41     "love": "neutre",
42     "gratitude": "neutre",
43     "optimism": "neutre",
44     "pride": "neutre",
45     "realization": "neutre",
46
47     # M content
48     "anger": "press ",
49     "disgust": "press ",
50     "annoyance": "press ",
51     "disapproval": "press ",
52     "remorse": "press ",
53     "sadness": "press "
54 }

```

Listing A.2 – Extrait du code principal - Configuration des modèles

A.3 Extrait du code principal - Fonction de personnalisation

```

1 # %% [18] Fonction de personnalisation des r ponses
2 def personalize_response(response, client_type, confidence):
3     """Personnalise la r ponse en fonction du type de client"""
4     if not response:
5         return response
6

```



```

7  # Adapte le ton selon le type de client
8  if client_type == "stress ":
9      # Ton rassurant et empathique
10     prefixes = ["          Je comprends que cela peut tre stressant.
11                  ",
12                  "          Prenez une profonde respiration. ",
13                  "          Je suis l pour vous aider. "]
14     response = random.choice(prefixes) + response.lower()
15
16 elif client_type == "press ":
17     # Ton direct et efficace
18     prefixes = ["          ", "          ", "          "]
19     response = random.choice(prefixes) + response
20
21 elif client_type == "curieux":
22     # Ton informatif et d taill
23     prefixes = ["          ", "          ", "          "]
24     response = random.choice(prefixes) + response + " N'h sitez pas
25               demander plus de d tails !"
26
27 # Ajoute un indicateur de confiance si pertinent
28 if confidence < 0.7:
29     response += " (Je fais de mon mieux pour comprendre votre
30               demande)"
31
32 return response

```

Listing A.3 – Extrait du code principal - Fonction de personnalisation

A.4 Extrait du code principal - Détection d'intentions

```

1  class CustomerSupportClassifier:
2      def __init__(self):
3          self.classifier = pipeline(
4              "text-classification",
5              model=model,
6              tokenizer=tokenizer,
7              device=0 if torch.cuda.is_available() else -1
8          )
9
10     def predict(self, text):
11         preds = self.classifier(text)
12         label = preds[0]['label']
13         confidence = preds[0]['score']
14         try:
15             idx = int(label.split("_")[1])

```

```

16         id2label = list(intent_to_data.keys())
17         intent_name = id2label[idx] if idx < len(id2label) else "
            unknown"
18     except Exception:
19         intent_name = "unknown"
20     return {"text": text, "intent": intent_name, "confidence":
        confidence, "entities": extract_entities(text)}
21
22 classifier = CustomerSupportClassifier()

```

Listing A.4 – Extrait du code principal - Détection d'intentions

A.5 Extrait du code principal - Extraction d'entités

```

1 # Reconnaissance "entit s" via PhraseMatcher (listes ferm es depuis le
   catalogue)
2 entities_config = {
3     "CATEGORIE": categories,
4     "MARQUE": marques,
5     "COULEUR": couleurs,
6     "NOM_PRODUIT": noms_produits
7 }
8
9 nlp = spacy.load("en_core_web_md")
10 matcher = PhraseMatcher(nlp.vocab, attr="LOWER")
11 for label, items in entities_config.items():
12     patterns = [nlp.make_doc(item) for item in items if item]
13     if patterns:
14         matcher.add(label, patterns)
15
16 def extract_entities(text):
17     doc = nlp(text.lower())
18     matches = matcher(doc)
19     entities = defaultdict(list)
20     for match_id, start, end in matches:
21         label = nlp.vocab.strings[match_id]
22         span = doc[start:end].text
23         entities[label].append(span)
24     return {k: list(set(v)) for k, v in entities.items()}

```

Listing A.5 – Extrait du code principal - Extraction d'entités

A.6 Extrait du code principal - Système de recommandation

```
1 def build_feature_corpus(df_):
2     texts = []
3     for _, row in df_.iterrows():
4         tokens = []
5         tokens += re.findall(r"[a-z0-9]+", str(row['Nom_norm']))
6         tokens += [str(row['Cat gorie_norm']), str(row['Marque_norm'])]
7         tokens += _explode_colors(row['Couleur_norm'])
8         tokens = [t for t in tokens if t and t != "nan"]
9         texts.append(" ".join(tokens))
10    return texts
11
12 _vectorizer = TfidfVectorizer(ngram_range=(1,2), min_df=1)
13 _feature_texts = build_feature_corpus(catalogue_df)
14 X_tfidf = _vectorizer.fit_transform(_feature_texts)
15
16 _price_scaler = MinMaxScaler()
17 price_scaled = _price_scaler.fit_transform(catalogue_df[['col_prix']].
18     values.astype(float))
19
20 price_scaled_sparse = csr_matrix(price_scaled)
21
22 X_features = hstack([X_tfidf, price_scaled_sparse]).tocsr()
23 name_to_index = {row['Nom_norm']: i for i, row in catalogue_df.
24     reset_index().iterrows()}
25
26 def recommander_produits_similaires(nom_produit_norm, top_k=5,
27     exclude_same=True, same_category=True):
28     if nom_produit_norm not in name_to_index:
29         return []
30     idx = name_to_index[nom_produit_norm]
31     query_vec = X_features[idx]
32     sims = cosine_similarity(query_vec, X_features).ravel()
33     order = np.argsort(-sims)
34     recs = []
35     base_cat = catalogue_df.iloc[idx]['Cat gorie_norm']
36     base_name = catalogue_df.iloc[idx]['Nom_norm']
37     for j in order:
38         if exclude_same and catalogue_df.iloc[j]['Nom_norm'] ==
39             base_name:
40             continue
41         if same_category and catalogue_df.iloc[j]['Cat gorie_norm'] !=
42             base_cat:
43             continue
44         recs.append({
```

```

39         "nom": catalogue_df.iloc[j][col_nom],
40         "categorie": catalogue_df.iloc[j][col_categorie],
41         "marque": catalogue_df.iloc[j][col_marque],
42         "prix": float(catalogue_df.iloc[j][col_prix]),
43         "score": float(sims[j])
44     })
45     if len(recs) >= top_k:
46         break
47     return recs

```

Listing A.6 – Extrait du code principal - Système de recommandation

A.7 Extrait du code principal - Gestion des sessions

```

1 class ChatSession:
2     def __init__(self, intent):
3         self.intent = intent
4         self.fields = intent_field_mapping.get(intent, [])
5         self.pending_fields = list(self.fields)
6         self.collected_data = {}
7         self.current_field = self.pending_fields[0] if self.
            pending_fields else None
8         self.client_type = "neutre"
9         self.context_history = []
10
11     def update(self, user_input, client_type="neutre"):
12         self.client_type = client_type
13         if not self.current_field:
14             return None
15         normalized_input = user_input.strip()
16         field = self.current_field
17
18         if normalized_input.lower() in ['cancel', 'annuler', 'stop', '
            exit']:
19             self.reset()
20             return "    Commande annul e. Comment puis-je vous aider ?"
21
22         if field == "confirmation":
23             if normalized_input.lower() in ['yes', 'oui']:
24                 self.collected_data[field] = normalized_input
25                 extra = ""
26                 if self.intent == "place_order":
27                     nom_produit = self.collected_data.get("nom_produit",
                        "")
28                     nom_produit_norm = normalize_text(nom_produit)

```

```

29         profil = mettre_a_jour_profil_apres_commande(
30             CURRENT_CLIENT_ID, nom_produit_norm)
31
32         extra = "\n\n" + bloc_recommandations_apres_achat(
33             nom_produit_norm, profil)
34
35         self.pending_fields = []
36         data_summary = "\n".join([f"- {k.replace('_', ' ')} : {v}" for k, v in self.collected_data.items()])
37         response = (f"      Commande confirm e !\n\nD tails :\n{data_summary}\n{extra}\n\nMerci !")
38         self.reset()
39         return response
40
41     elif normalized_input.lower() in ['no', 'non']:
42         self.reset()
43         return "      Commande annul e. Que souhaitez-vous faire ?"
44
45     return "      Veuillez r pondre par 'oui/yes' ou 'non/no'"
46
47     # ... autres cas de gestion de champs
48     return True
49
50 def is_complete(self):
51     return not self.pending_fields
52
53 def reset(self):
54     self.pending_fields = list(self.fields)
55     self.collected_data = {}
56     self.current_field = self.pending_fields[0] if self.pending_fields else None
57     self.client_type = "neutre"
58     self.context_history = []

```

Listing A.7 – Extrait du code principal - Gestion des sessions

A.8 Extrait du code principal - Interface utilisateur

```

1 # =====
2 # INTERFACE STREAMLIT
3 # =====
4 st.set_page_config(page_title="      Assistant Client IA", page_icon="
5     ")
6 st.title("      Assistant Client Intelligent")
7 st.write("Discutez avec notre assistant IA pour g rer vos commandes,
8     comptes et questions.")

```

```

7
8 if "messages" not in st.session_state:
9     st.session_state.messages = []
10 if "session" not in st.session_state:
11     st.session_state.session = None
12
13 for msg in st.session_state.messages:
14     with st.chat_message(msg["role"]):
15         st.markdown(msg["content"])
16
17 user_input = st.chat_input("Tapez votre message ici...")
18
19 if user_input:
20     st.session_state.messages.append({"role": "user", "content":
21         user_input})
22     with st.chat_message("user"):
23         st.markdown(user_input)
24
25     client_type = detect_client_type(user_input)[0]
26     profil = profiles.get("guest")
27     profil.client_type = client_type
28     profil.client_type_history.append({
29         "type": client_type,
30         "confidence": 1.0,
31         "timestamp": datetime.datetime.now(),
32         "input": user_input
33     })
34
35     response = generate_response(user_input)
36     st.session_state.messages.append({"role": "assistant", "content":
37         response})
38     with st.chat_message("assistant"):
39         st.markdown(response)

```

Listing A.8 – Extrait du code principal - Interface utilisateur