



ក្រសួងអប់រំ យុវជន និងកីឡា

ពិភាក្សាលទ្ធផលថ្មីកម្ពុជា

ខេត្តកែវិច្ឆិក និងការអប់រំ និងកីឡា និងការអប់រំ



សញ្ញាណសញ្ញាណបញ្ជីកម្រិត

ប្រធានបទ	: បំណែលដែលអាចបង្កើតត្រួវនិចបាយស្ថីយប្រព័ន្ធ
ឯកសារ	: ក្រសួងអប់រំ យុវជន និងកីឡា និងការអប់រំ
ប្រធានបទ	: ក្រសួងអប់រំ យុវជន និងកីឡា និងការអប់រំ
ឆ្នាំសំខាន់	: ២០២៤-២០២៥

MINISTERE DE L'EDUCATION,
DE LA JEUNESSE ET DES SPORTS

INSTITUT DE TECHNOLOGIE DU CAMBODGE

DEPARTEMENT DE GENIE INFORMATIQUE ET
COMMUNICATION

MEMOIRE DE FIN D'ETUDES

Titre	: AI-Powered Email Sorting: Leveraging Text and Visual Cues for Intelligent Triage
Etudiant	: KAO Visal
Spécialité	: Génie Informatique et Communication
Tuteur de stage	: Mr. SOK Kimheng
Année scolaire	: 2024-2025



ក្រសួងអប់រំ យុវជន សិក្សា

និគ្ងាសាលបច្ចេកទិន្នន័យ



លេខាសីម៉ែត្រ ៩ ទៀតនៃការណានិភ័យ និគ្ងាសាលបច្ចេកទិន្នន័យ

សញ្ញាណព័ត៌មាន

ឈ្មោះ: កែវ ពិសោធន៍

ភោជន៍: ខេត្តកំពង់ចាម សង្កាត់ ៩ ក្នុង ២០២៤

ឈ្មោះ: កែវ ពិសោធន៍

ឈ្មោះ: _____

ថ្ងៃខែឆ្នាំ: ៩ ១ ២០២៤

ប្រធានបទ: ចំណោមដោយអាជីវកម្មបច្ចេកទិន្នន័យ

ចំណោមដោយអាជីវកម្មបច្ចេកទិន្នន័យ

ឈ្មោះ: _____

: ត្រូវបានបញ្ជាក់

ប្រធានលេខាសីម៉ែត្រ

: លេខាសីម៉ែត្រ _____

ស្ថាបនបានបញ្ជាក់

: លេខាសីម៉ែត្រ _____

ឯកចាយបញ្ជាក់: លេខាសីម៉ែត្រ _____

រាល់ការណានិភ័យ ៩ ក្នុង ២០២៤



MINISTERE DE L'EDUCATION,
DE LA JEUNESSE ET DES SPORTS



INSTITUT DE TECHNOLOGIE DU CAMBODGE
DEPARTEMENT DE GENIE INFORMATIQUE ET COMMUNICATION

MEMOIRE DE FIN D'ETUDES

DE M. KAO Visal

Date de soutenance : le 8 septembre 2025

« Autorise la soutenance du mémoire »

Directeur de l'Institut : _____

Phnom Penh, le 8 septembre 2025

Titre : AI-Powered Email Sorting: Leveraging Text and Visual Cues for Intelligent Triage

Etablissement du stage : Groupe CSW

Chef du département : M. LAY Heng _____

Tuteur de stage : M. SOK Kimheng _____

Responsable de l'établissement : M. Jérôme SICO _____

PHNOM PENH, 2025

ទេរចក្ចិតទេច

ការបែលដឹងកិច្ចការប្រើដឹងអាយុវត្ថុជាការងារអូតូម៉ាទិក ជាការប្រើប្រាស់បច្ចេកវិទ្យាបញ្ហាសិប្បនិមិត (AI) ដែលគេធ្វើប្រចាំថ្ងៃបំផុតក្នុងសម្រាប់បច្ចុប្បន្ន។ ក្រុមហ៊ុន និងអង្គភាពជាប្រើប្រាស់ទំនាក់ទំនងដោយអគិតិ កំពុងអនុវត្តន៍និងប្រើប្រាស់បញ្ហាសិប្បនិមិតដើម្បីដឹងក្នុងក្រុងខ្លួន និងធ្វើទំនាក់ទំនងដោយអគិតិ ជនដោយស្ម័យប្រភេទ។ និស់យម្ចាយក្នុងបំណែមិស្ទ័យដែលមានសត្ថិភាពលម្អិត និងការដោក់បញ្ហាល AI គឺ ប្រព័ន្ធបញ្ហាលដោយស្ម័យប្រភេទ។ នៅក្នុងអង្គភាពមួយចំនួន ប្រអប់សារូម (shared mailbox) ត្រូវបានប្រើប្រាស់ដោយបុគ្គលិកប្រើប្រាស់ក្នុងក្រុមហ៊ុន ដើម្បីធ្វើទំនាក់ទំនងដោយការគ្រប់គ្រងដោយបុគ្គលិកមួយរួចប្រើប្រាស់ដែលទទួលខុសត្រូវក្នុងការអាន វិភាគ និងបែកអុំមែលដែលបានធ្វើចូលមក អាយុវត្ថុអ្នកដែលត្រូវទទួលភ័យដែល ធ្វើយ និងទទួលសារនោះ។

គោលបំណងនៃនិភ័យបច្ចេនេះ គឺដើម្បីស្រាវជ្រាវ ពិសោធន៍យ និងដោក់អាយុវត្ថុ និងបុគ្គលិក (software) មួយដែលអាច អានយល់ពីអត្ថន័យបស់អុំមែល ពីវត្ថុបំណងរបស់អ្នកដៅ និងអាចបែកចាយអុំមែលនោះទៅ ឱ្យបុគ្គលិកណាមួយ ក្នុងក្រុមហ៊ុន ដែលអាចធ្វើយនូវតម្លៃការបែកចាយរបស់ក្រុមហ៊ុន។ ការដោនេះត្រូវបានអនុវត្តនៅ ក្រុមហ៊ុន Groupe CSW (Cosywee) ដោយប្រើប្រាស់បច្ចេកវិទ្យាបញ្ហាសិប្បនិមិត ដែលរួមមាន Machine Learning និង មួយចំនួន RoBERTa បន្ថែមដោមួយ Computer Vision Methods ដើម្បី មួយដែល RoBERTa (CamemBERT) ត្រូវបានដើរសិសមកប្រើប្រាស់ក្នុងគម្រោងនេះ ព្រមទាំងបណ្តុះបណ្តាល (Train) ដោយទិន្នន័យប្រើប្រាស់ជាការសាងគារការងារ របស់ក្រុមហ៊ុន ប្រើប្រាស់ self-attention រួម ទាំងមុខងារ Bidirectional ដែលអាចចាប់យកពាណិជ្ជកម្មនៃយោប់យោ ប្រើប្រាស់ក្នុងការងារ។

ដំណើរការបណ្តុះបណ្តាលមួយខែល បានប្រើប្រាស់បណ្តាលឲយ និងឧបករណ៍ជាប្រើប្រាស់ដើម្បីធានាចូល
ប្រសិទ្ធភាព និងមានភាពដាយស្រួល។ Weights & Biases (wandb) ត្រូវបានប្រើប្រាស់ម្រាប់តាមដានការ
សាកលវិធី កំណត់ និងប្រព័បែងចុះរបស់មួយខែលនីមួយៗ។ HuggingFace ត្រូវបានប្រើ
សម្រាប់ទាញយកមួយខែលដែលបានបណ្តុះបណ្តាលរួច និងធ្វើការប្រព័បែងតាម benchmark។ ហច្ខកវិញ្ញា
ទាំងនេះ ត្រូវបានប្រើប្រាស់ដើម្បីមានសមត្ថភាពធ្វើមាត្រាតាន (scalability) មានភាពបត់បែនឌួល
និងដាយស្រួលក្នុងការគ្រប់គ្រងមួយខែល។ នៅចុងបញ្ហាប់នៃគម្រោង បណ្តាលកម្មវិធីត្រូវបានអនុវត្តន៍ និងសាក
លុងក្នុងបរិបទពិតប្រាកដ ដោយប្រើកំពុទ្ធដែលមានកម្លាំងតណាតិចបំផុត។

លទ្ធផលនៃការស្រាវជ្រាវនេះ បានដឹងដូចមួយថាអ្នកបានប្រើប្រាស់ដំណើរការដាយដែលបានបង្កើតឡើង
90% ដែលបានដឹងថាគាត់បន្ទាយចំណាយ និងបង្កើនប្រសិទ្ធភាពនៅក្នុងក្រុងក្រោម។

Résumé

L'automatisation des tâches répétitives est l'une des applications les plus courantes de l'intelligence artificielle aujourd'hui. Les entreprises, qu'elles soient locales ou internationales, adoptent et utilisent de plus en plus l'intelligence artificielle pour traiter automatiquement des documents physiques, extraire des informations et interagir avec les clients. L'un des domaines prometteurs de l'intégration de l'IA est le routage automatique des e-mails. Dans certaines organisations, les boîtes mail partagées sont généralement gérées par une ou plusieurs personnes chargées de lire, d'analyser et de classer manuellement les e-mails afin de les transférer aux destinataires appropriés.

Le processus d'entraînement s'est appuyé sur plusieurs bibliothèques et outils afin d'assurer efficacité et reproductibilité. Weights & Biases (wandb) a été utilisé pour suivre les expériences, gérer les hyperparamètres et comparer les performances des différents modèles. HuggingFace a permis d'accéder à des modèles pré-entraînés et d'évaluer des benchmarks. Ces technologies ont été choisies pour leur scalabilité, leur flexibilité, et leur facilité d'intégration pour la gestion des modèles. À la fin du projet, la solution a été déployée et testée dans un environnement réel sur une machine disposant de ressources minimales.

L'objectif de ce rapport est de présenter le travail réalisé au sein du Groupe CSW (Cosy-wee) pour automatiser cette tâche, en utilisant une combinaison d'algorithmes traditionnels d'apprentissage automatique et de modèles basés sur RoBERTa, intégrés à des techniques de détection d'indices visuels sur les images. Le modèle CamemBERT (basé sur l'architecture RoBERTa) a été choisi pour son encodeur Transformer avec mécanisme d'attention, qui permet au modèle de comprendre le contexte des mots, quelle que soit leur distance. De plus, CamemBERT a été pré-entraîné sur un large corpus en français, ce qui le rend particulièrement adapté à ce projet.

Cela a permis de réduire de près de 90 % le temps de traitement manuel, réduisant ainsi considérablement les coûts et améliorant l'efficacité au travail.

Abstract

Automating repetitive tasks is one of the most common applications of Artificial Intelligence today. Companies, both local and international, have been increasingly adopting and using Artificial Intelligence to automatically process physical documents, extract information, and interact with clients. One promising area of AI integration is Automatic Email Routing systems. In some organizations, shared mailboxes are usually managed and controlled by one or more people who are responsible for manually reading, analyzing, and classifying those emails to forward them to appropriate recipients.

The objective of this report is to present the work carried out at Groupe CSW (Cosywee) to automate this task, using a combination of traditional machine learning algorithms and RoBERTa-based models, integrated with visual clue detection on images. The CamemBERT model (based on the RoBERTa architecture) was chosen due to its Transformer encoder with self-attention. This enables the model to capture the context of words regardless of their distance. Additionally, CamemBERT was pre-trained on a large French corpus, making it particularly well-suited for this project. The training process leveraged several libraries and tools to ensure efficiency and reproducibility. Weights & Biases (wandb) was used to monitor experiments, track hyperparameters, and compare the performance of different models. HuggingFace was employed to access pre-trained models and evaluate benchmarks. These technologies were chosen for their scalability, flexibility, and ease of integration for model management. At the end of the project, the solution was deployed and tested in a real-world environment on a machine with minimal resources.

This resulted in a reduction of almost 90 % in manual processing time, significantly cutting costs and improving workplace efficiency.

ACKNOWLEDGEMENTS

The success of this project is not just the result of my own effort, but a shared journey made possible by the support, encouragement, and kindness of many people. Before diving into the work itself, I want to take a moment to express my deep gratitude to those who have helped me and are still helping me along the way.

First, I wish to extend my sincere appreciation to **H.E. Dr. PO Kimtho**, Director of the Institute of Technology of Cambodia. His excellent leadership and the strong cooperation with partner universities both locally and internationally have significantly contributed to improving the quality of education for engineers in our country, and that cooperation is the reason why I get a chance to pursue my double degree in France as well.

I am also deeply appreciative of **Mr. LAY Heng**, head of GIC department, for his valuable support and guidance during my time as a student in his department.

I would like to acknowledge the faculty and administrative staff of the GIC department (Génie Information et Communication) at Institute of Technology of Cambodia for fostering an engaging and intellectually rich environment. I am especially thankful to **Mr. SOK Kimheng**, my academic supervisor, for his kindness and ongoing assistance throughout this internship. I am grateful to **Dr. VALY Donna**, the head of ViLa Research Lab, for introducing me into the field of Artificial Intelligence, providing me with materials and giving me a chance to set my foot in the gate of this domain.

I want to thank **Mr. Jérôme SICO**, my supervisor during this end-of-study project (Projet de Fin d'Études) at Cosywee, for his trust, guidance, and constant support. His sharp insights, thoughtful feedback, and professional approach helped shape this thesis and pushed me to grow both technically and personally.

Above all, I owe the deepest gratitude to my family, including my parents, my aunt, my siblings, and my partner, for their unconditional love, sacrifices, and motivation. Their belief in me has been the foundation of every success I've achieved.

ABBREVIATION LIST

AI	: Artificial Intelligence
DICE	: Department of Information and Communication Engineering
BDC	: Bon de Commande (Purchase Order)
ITC	: Institute of Technology Cambodia
CSW	: Groupe CSW (Cosywee)
BERT	: Bidirectional Encoder Representations from Transformers model
ML	: Machine Learning
NLP	: Nature Language Processing
TF-IDF	: Term Frequency - Inverse Document Frequency
RoBERTa	: A Robustly Optimized BERT Pretraining Approach
CamemBERT	: RoBERTa model pre-trained on French corpus
MP	: Marseille-Perpignan (shortcut for a label in classification model)
OCR	: Optical Character Recognition

TABLE OF CONTENTS

Résumé	4
Abstract	5
ACKNOWLEDGEMENTS.....	6
ABBREVIATION LIST	7
TABLE OF CONTENTS.....	8
LIST OF FIGURES	11
LIST OF TABLES	12
I - INTRODUCTION.....	13
1.1 General Presentation	13
1.1.1 Company Presentation.....	13
1.1.2 Company and Supervisor contact information.....	13
1.2 Project Presentation	13
1.2.1 Problem Statement	13
1.2.2 Mission Objective	14
1.2.3 Types of Incoming Emails	15
1.2.4 Trade-offs.....	16
1.3 Project Timeline.....	16
II - RELATED WORK.....	18
III - TECHNOLOGIES	19
3.1 Tools and Technologies	19
3.1.1 Python	19
3.1.2 Anaconda.....	19
3.1.3 Sourcetree.....	19
3.1.4. Bitbucket	20
3.1.5. Jira	20
3.1.6. Visual Studio Code.....	20
3.1.7. Transformers	21
3.1.8. Spacy	21
3.1.9. OpenCV.....	21
3.1.10. Pytesseract.....	21
3.1.11. EasyOCR.....	21
3.1.12. nltk.....	22
3.1.13. Microsoft Graph API.....	22

3.1.14. HuggingFace	22
3.1.15. Weights and Bias (WandB).....	22
3.2. Coding Guideline	23
IV - MODEL DEVELOPMENT.....	24
 4.1. Data Collection	24
4.1.1. Problem	24
4.1.2. Solution	26
 4.2. Data Cleaning.....	26
4.2.1. Email body cleaning.....	26
4.2.2. Attachment cleaning.....	27
4.2.3. Final Data.....	30
 4.3. Model Experimentation.....	30
 4.4. Evaluation Metrics.....	34
Confusion Matrix and Terminology.....	35
Evaluation Methods.....	35
 4.5. Results	36
4.5.1. Intent classification model	36
4.5.2. Alès, Lyon, MP classification	37
4.5.3. Marseille, Perpignan, Roller Shutter (<i>VR</i>) classification.....	39
 4.6. Final Pipeline.....	41
V - IMPLEMENTATION & PROTOTYPE BUILDING.....	43
 5.1. Project Development.....	43
5.1.1. Environment configuration	43
5.1.2. Project Structure.....	44
 5.2. Heuristics methods.....	45
5.2.1. Original Sender of the email	45
5.2.2. Document Creator	45
5.2.3. Product extraction from Purchase Order	48
5.2.4. Intent extraction from Purchase Order	50
 5.3. Helper Functions with Microsoft Graph API	52
5.3.1. Read Emails.....	53
5.3.2. Get Folder ID	54
5.3.3. Move Email.....	54
 5.4. Final Classification Pipeline Implementation	55
 5.5. Optimization and Improvement.....	58
5.5.1. Version v1.0	58

5.5.2. Version v2.0	59
VI - RESULTS AND DISCUSSIONS	61
6.1. Results	61
6.2. Future improvements	63
6.2.1. Optimizing handwriting text recognition model	63
6.2.2. Integration of human feedback loop.....	63
VII - CONCLUSION	65
References	66
APPENDIX	69
Types of Emails	69
Log files	71

LIST OF FIGURES

Figure 1 - Email Flows Within The Service.Clients Mailbox.....	15
Figure 2 - Project Timeline.....	16
Figure 3 – Log Files	23
Figure 4 - Dataset	30
Figure 5 - Overview Of Model Pipeline.....	31
Figure 6 - Confusion Matrix Of Intents Classification.....	37
Figure 7 - Fully Integrated Ai Pipeline.....	41
Figure 8 - Project Structure	44
Figure 9 - Employee's Name To Company.....	46
Figure 10 - Employee's Name Extraction From Purchase Order Flowchart	47
Figure 11 - Purchase Order For Mosquito Net.....	48
Figure 12 - Matching Product's Name/Code To Company.....	49
Figure 13 - Product Name Extraction From Purchase Order Flowchart	50
Figure 14 - Orb Features Matching	51
Figure 15 - Orientation Corrected Purchase	52
Figure 16 - Works Completed Flowchart	64
Figure 17 - Emails With No Attachments, Containing Only Plain Text Order Or Price Quotation Requests	69
Figure 18 - Emails With Handwritten Attachments Indicating Order Details	69
Figure 19 - Emails With Computer-Generated Attachments (Printed Text) Indicating Order Details.....	69
Figure 20 - Emails Replying To A Cosywee Employee's Email	70
Figure 21 - Emails Inquiring About Order Progress (With Only A Reference Number And No Product Information)	70
Figure 22 - Order Emails With Cosywee's Purchase Order Form (Bon De Commande) Attached	70
Figure 23 - Emails Containing Cosywee-Generated Documents As References.....	70
Figure 24 - Time Records Log: Logs The Displacement Or Movement Of Emails, Which Is Useful For Debugging Cases Where Emails Appear To Be Missing.	71
Figure 25 - Prediction Log: Stores Predictions And Confidence Scores From Both The Camembert Model And The Random Forest Model.)	71
Figure 26 - Complete Data Log: Stores All Extracted Data Along With Their Predictions. This Log Is Useful For Fine-Tuning New Models In The Future.....	71

LIST OF TABLES

Table 1 - Comparison Of Email Processing Libraries.....	25
Table 2 - Performance Of Intents Classification Model.....	37
Table 3 - Confusion Matrix Of Ales, Lyon, Mp Classification.....	38
Table 4 - Performance Of Alès, Lyon, Mp Classification Model	38
Table 5 - Performance Comparison Of Algorithms - Alès, Lyon, Mp Classification.....	39
Table 6 - Confusion Matrix Of Marseille, Perpignan, Roller Shutter Classification	40
Table 7 - Performance Of Marseille, Perpignan, Roller Shutter Classification Model.....	40
Table 8 - Performance Comparison Of Algorithms - Marseille, Perpignan, Roller Shutter Classification.....	40
Table 10 - Processing Time Vs Computing Power On A High Computing Power Machine ..	58
Table 11 - Processing Time Vs Computing Power On A Low Computing Power Machine ...	59
Table 12 - Processing Time Vs Computing Power Of Version V2.0 On A Low Computing Power Machine.....	60
Table 13 - Ai System Capability (With 500 Emails Sorted)	62

I - INTRODUCTION

1.1 General Presentation

1.1.1 Company Presentation

The CSW Group (Cosywee), specialized in closure and solar protection products, has been building and maintaining a relationship of trust with professionals since 2003. The CSW Group comprises several companies at the south of France including Cosywee Alès, Cosywee Perpignan, Cosywee Nîmes, Cosywee Lyon, Cosywee Marseille, and Cosywee Bourg-en-Bresse. Each of these entities contributes to the Group's specialization in closure and solar protection products, ensuring a constant level of expertise and service in the different regions.

With a nationwide presence, the CSW Group has developed over the years on the basis of top-of-the-range products assembled in France. The group's 130 employees work at 6 production sites and 1 design office located near Perpignan, Nîmes, Marseille and Lyon.

1.1.2 Company and Supervisor contact information

COSYWEE Alès (bureau de l'équipe informatique) is at Alès 30100, France :

- Google Map: <https://maps.app.goo.gl/Uitary9pKzYfv5bJLA>
- Adresse : 22b Avenue du général De Gaulle, 30100 ALES
- Tel : +33 4 48 40 99 89 • Email : jerome.sico@cosywee.com
- Website: <https://www.cosywee.com>

1.2 Project Presentation

1.2.1 Problem Statement

Within the company, there is a shared mailbox where all employees in Service Clients Department communicate with clients answering product inquiry, customer survey, account opening (New clients registration), price quotation, orders, after-sales services, order confirmation, payments, accounting related problems, etc. As all employees shared the same mailbox, all emails sent and received flow into one email address, "service.clients@cosywee". To make things even more complicated, each product has a unique code, a different Bon de Commande (Purchase order) template, and is produced by different company (factory) of the Group.

Usually, there is a group of employees, at each company of the group (or each factory), who are responsible for reading emails, including attachments, responding to them and manually

forwarding those emails to the corresponding people or factory. As there are more and more mails incoming, it can easily get out of hand and can be one of the problems for the productivity and efficiency of the company.

In order to solve this problem, there was an attempt to use keyword matching to parse the email's body and subject to redirect those emails, but as the patterns are too complicated, with too many categories and purpose involves, and missing information from attachments, this option is out of consideration.

1.2.2 Mission Objective

The main objective of this internship is to implement artificial intelligence to improve productivity in the company. My scope (responsibility) includes problem definition, requirement gathering, data collection, finetuning a pre-trained model or training one from scratch, up to model implementation, pipeline building, and model/program hosting as well.

Basically, this main objective is to build an end-to-end system improving efficiency, reducing and to eliminate manual, repetitive tasks in the company. In this case, the aim of the internship to build an end-to-end program to sort or redirect incoming email from Service Clients mailbox to appropriate recipients, within a reasonable delay and high accuracy to reduce manual tasks. This includes problem formulation, requirement gathering, data collection, data cleaning, model experimentation, complete program pipeline, and hosting. This might include data extraction, data scraping, API integration, intention classification, Optical Character Recognition (OCR), keywords matching, etc. In the diagram below (Figure 1), you can see the flow of emails in the Service.Clients mailbox. When a client sends an email to the Service.Clients address, employee from Client Service department reads the message, analyzes its intent, checks whether any products are mentioned, identifies which company or factory is responsible for the product, and then re-routes the email to the appropriate recipients.

Furthermore, since all emails sent by employees from the Service Clients Department use the "service.clients@cosywee.com" address, any replies to those emails must also be read and re-routed to the original sender.

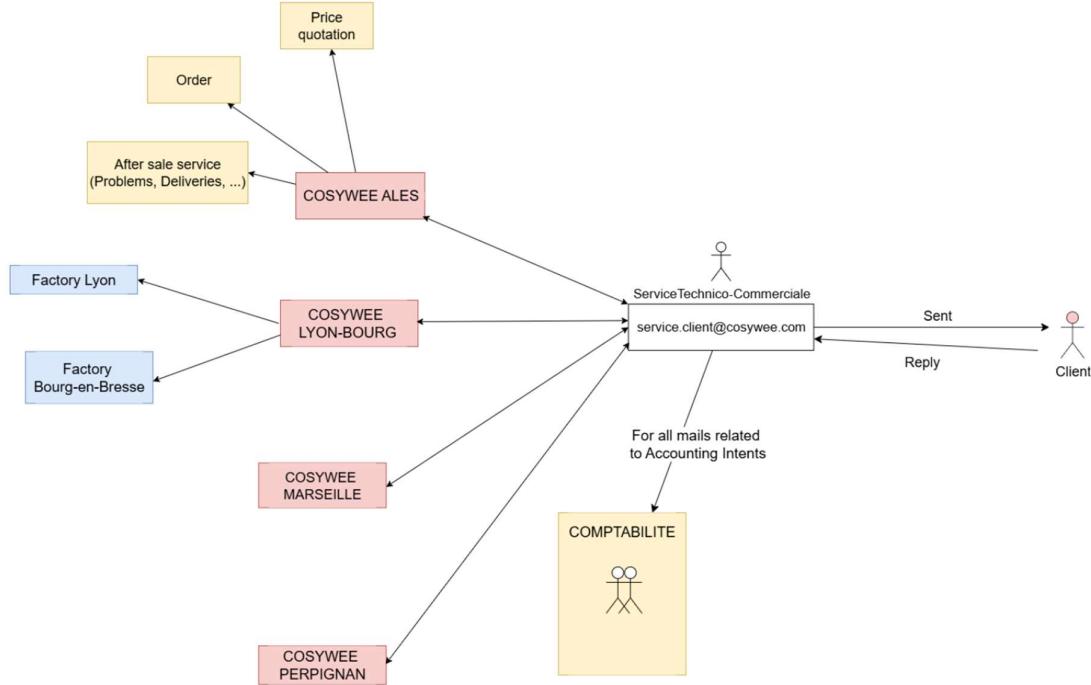


Figure 1 - Email flows within the service.clients mailbox

Each email can contain multiple types of attachments, including images, PDF files, .txt files, .csv files, .xlsb files, .docx files, etc. These files come in different formats, related to different products and serving unique purposes.

1.2.3 Types of Incoming Emails

There are many types of emails that arrive in service.client@cosywee.com mailbox.

Below are some common types of emails that we would get daily:

- Emails with no attachments, containing only plain text order or price quotation requests
- Emails with handwritten attachments indicating order details
- Emails with computer-generated attachments (printed text) indicating order details
- Emails replying to a Cosywee employee's email
- Emails inquiring about order progress (with only a reference number and no product information)
- Order emails with Cosywee's Purchase Order form (Bon de Commande) attached
- Emails containing Cosywee-generated documents as references
- Emails announcing unrelated information (vacation dates, congratulations, etc.)

You can find examples of each type of email in [Types of Emails](#).

1.2.4 Trade-offs

In this section, I will talk about the constraints, problems, and trade-offs considered prior to proposing a solution.

First of all, **COSYWE MARSEILLE** and **COSYWE PERPIGNAN** share overlapping product lines. These products include all items related to **Roller shutters (Volet Roulant)** (Traditional Roller shutters, Bloc-Blae roller shutters, Roller Shutter Renovation, Solar Roller Shutters, Special Roller Shutters, etc.). Due to this overlap, all products related to Roller shutters must be classified manually by a human (excluded from the classification pipeline).

Secondly, the processing time of the program must be minimized as it must deal with multiple incoming emails in one second. In this sense, the company is willing to accept a small trade-off in accuracy in favor of the program's processing time. Similarly, the computational resources required by the program must be kept as low as possible, since it will be hosted on a shared company server that also supports other critical applications.

Thirdly, the company prioritizes explainability over small gains in accuracy. This indicates that if two algorithms perform quite similarly in terms of accuracy and precision, we tend to choose the one with more explainability regardless of its performance.

Finally, security is a critical pillar of the system. Since the corporate mailbox contains confidential business information, it represents one of the most sensitive components of the company's infrastructure. As such, security protocols must be implemented with the highest priority and rigor.

1.3 Project Timeline

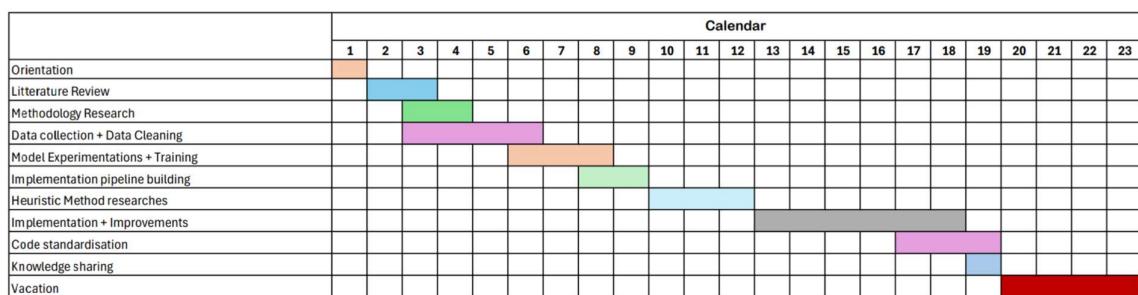


Figure 2 - Project Timeline

The project started in early March (week 1 in Figure 2) when I was introduced to the problem. I then spent some time on literature review and researching potential methodologies. After that, I began collecting data, finding ways to extract and clean it, spending a significant amount of time on the OCR part, as it provides the most critical information from attachments.

Next, I experimented with different machine learning algorithms and deep learning architecture before starting to build the pipeline and API endpoints needed to read and move emails for implementing the AI system. I spent about two weeks setting up a secure API inference using Microsoft Graph API and configuring the API to access emails in a shared mailbox.

Note that, to implement this application, we had to create a new rule, a new user, and a dedicated account for the entire automation system. After that, I focused on developing heuristic methods to improve the system's accuracy and efficiency before deploying it to production in the 13th week. From that point onward, I began releasing new versions of the program every few days, addressing reported bugs and adding new features.

And since the 17th week, I started to standardize my code and write detailed reports on how to fine-tune the model, how to retrain it from scratch, and how the company can maintain and scale the program in the future.

It is worth noting that the project was completed by week 19. Between weeks 20 and 23, Groupe CSW, as well as many other companies in France, observed the annual summer holiday break.

II - RELATED WORK

Email classification and text classification have long been active research areas in the field of Natural Language Processing (NLP). Previous research and published works in these areas have laid a strong foundation for building a good, reliable and robust multi-modal systems.

Historically, email classification, or text classification in general, relied on traditional machine learning methods like Naive Bayes and Random Forest, while leveraging features selection techniques such as TF-IDF or bag-of-words vectors, etc. These ancient approaches are considered efficient and interpretable but they all assume words independence, making them less suitable for real-world tasks like semantic or intent classification.

Newer methods introduced a few years back include transformer-based language models like BERT (Al., 2018) and its variances. One well-known variance for French Natural Language Processing is CamemBERT (Martin, et al., 2020), a RoBERTa-based (Liu, et al., 2019) transformer that was pretrained on a large French corpus by Facebook researchers and researchers from Sorbonne University. It has shown a promising context-awareness ability and is often used as a benchmark, a state-of-the-art model for French language across multiple tasks including classification, named entity recognition, and semantic similarity, etc.

In this project, a combination of Camembert and Random Forest were used to make predictions regarding products and entities, to reduce bias and prevent overconfident outputs. Additionally, the CamemBERT model alone is used to predict the intent of emails as well.

III - TECHNOLOGIES

3.1 Tools and Technologies

In the course of this project, a variety of tools, libraries, and platforms were used to ensure efficient development, collaboration, and deployment. Below is an overview of the key technologies:

3.1.1 Python

Python is the first programming language choice for Artificial Intelligence tasks thanks to its simplicity and feasibility to implement. One of the main advantages of using Python in this project is its diverse ecosystem and libraries available such as torch (Pytorch), Tensorflow and Scikit-learn, etc. These tools have many features and functionalities that can facilitate the process of data processing, cleaning, model experimentation and deployment as well.



3.1.2 Anaconda

Anaconda is a platform of virtual environments and packages management that is mostly used in data science and AI development. In this project, I used Anaconda to manage different environments (Testing environment, development environment, and deployment environment) since I had to test different models with diverse dependency, which would cause a problem if I only used global environment.



3.1.3 Sourcetree

Sourcetree is a Git client that allows simplification of version control process for all software development projects. It provides engineers and developers with a graphic interface that

allows us to effectively review visually of changes we made on code and to manage other git manipulation operations as well.



3.1.4. Bitbucket

Bitbucket is a web platform to host git projects, essentially for version control and collaboration in a software development team. In this project, Bitbucket is used to manage the code repository and version control.



3.1.5. Jira

Jira, developed by Atlassian, is a tool of project management used to follow tasks and workflows, mostly in the software development life cycle (Development, improvement, bug fixing, maintenance, etc). In this project, Jira is used to measure the completion of each step, each task. Jira facilitated task management and progress monitoring, ensuring efficient project coordination and on-time delivery.



3.1.6. Visual Studio Code

Visual Studio Code, developed by Microsoft, is one of the most popular code editors used for software development today. It offers many functionalities such as syntax highlighting, debugging, version control integration, and more.

In this project, Visual Studio code is used to write and test code, to run and to debug the project as well.



3.1.7. Transformers

This library was essential for using and fine-tuning state-of-the-art pre-trained language models for NLP tasks such as text classification and sentiment analysis.

3.1.8. Spacy

Spacy (spaCy French Language Model: fr_core_news_lg, 2024) is utilized for fast and efficient natural language processing, especially for tasks such as noise removal, named entity recognition, and part-of-speech tagging in this project.

3.1.9. OpenCV

OpenCV is an image-processing-focused library, utilized for image pre-processing and manipulation in this project (particularly in document layout analysis tasks).



3.1.10. Pytesseract

Pytesseract (pytesseract: Python wrapper for Google Tesseract-OCR, 2024) is a Python wrapper for Google's Tesseract-OCR Engine, widely used for Optical Character Recognition (OCR) mostly on printed text. In this project, it is leveraged for text detection and extraction from scanned documents.



3.1.11. EasyOCR

EasyOCR (JaideAI, 2024) is a Python-based OCR library available on Github that supports many languages and is also known for handling both printed and handwritten text. Just like Pytesseract, in this project, EasyOCR is primarily used for text extraction (both handwritten and printed text) from scanned documents and images.



3.1.12. nltk

Natural Language Toolkit (NLTK) is a leading library in Python, mainly used for Natural Language Processing tasks. It provides tools for text processing, stopword removal, and more



3.1.13. Microsoft Graph API

Microsoft Graph API is a unified REST API that allows us, developers, to access Microsoft 365 as well as its services like Outlook, OneDrive in a programmatic way.



In this project, Microsoft Graph API is used to connect from the developed application to the company's Outlook inbox to fetch and read emails to classify.

3.1.14. HuggingFace

Hugging Face is a platform that provides state-of-the-art pretrained models, which can be used in the development of this project via its Transformers library and Model Hub. In this project, Hugging Face is used to load pretrained models such as **almanach/camembert-base** (Martin, et al., 2020) and **Jean-Baptiste/camembert-ner** (CamemBERT NER - HuggingFace model, 2021).



Figure 3.L: HuggingFace Logo

3.1.15. Weights and Bias (WandB)

Weights and Biases (WandB) is an experiment tracking tool widely used in machine learning for logging training runs, visualizing metrics, and managing hyperparameter tuning.



In this project, WandB is used to monitor model training in real-time as well as to compare different experiments to choose the best model for implementation.

3.2. Coding Guideline

In this project, strict coding guideline is implemented to ensure a high-quality and more readable codebase. The rules I followed include:

- Variable names must follow the *snake_case* (like `variable_name`) format (not camelCase).
- There must be no unnecessary spaces between lines of code.
- The code must not contain unnecessary comments.
- Each function must be documented, explaining its purpose, parameters, and outputs.
- Variable names must be meaningful to facilitate code review and debugging by other teammates.
- Four log files are created to record predictions, email routing, processing time, and extracted data for debugging purposes.
- Log files are zipped every day to reduce disk space and zip files that are older than ten days will be deleted.

These guidelines help maintain a clean and organized codebase, facilitating collaboration within the team.

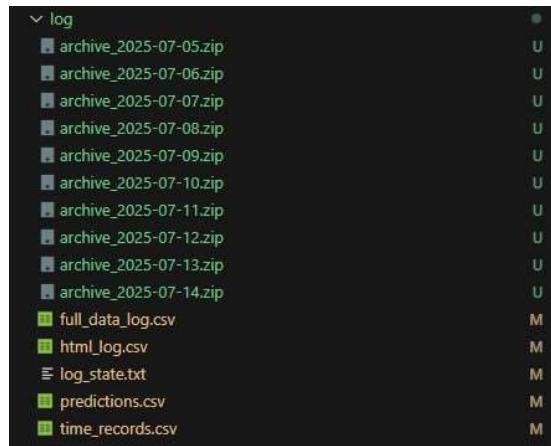


Figure 3 – Log files

IV - MODEL DEVELOPMENT

This chapter will present the process of data collection, data cleaning, literature reviews, model experimentations, evaluation metrics, and model implementation.

4.1. Data Collection

Since we are working with a custom intention classification model, the data that our model will be trained on is fully from the **service.clients@cosywee.com** itself. Due to security reasons, I did not have access to the full dataset from past years, and I only had my hands on the email samples starting from January 2025 onward.

The data is divided into 4 classes for 4 different sociétés/companies:

- Cosywee Alès
- Cosywee Lyon
- Cosywee Marseille
- Cosywee Perpignan

4.1.1. Problem

The first problem that I had was data collection. The data that I first received was in .PST file format. And in order to extract email content and attachments from .PST file, I experimented with some Python libraries including:

- Aspose.Email (Aspose.Email for Python via .NET, 2024)
- extract_msg (extract_msg: Python library to extract emails from MSG files, 2024)
- pypff (libpff/pypff: Library and Python bindings to access Personal Folder File (PST, OST), 2024)
- Outlook COM Automation (win32com)

Each one of these libraries has its own limitations and I evaluate them based on these important categories:

- Format Supported
- Metadata extraction (Subject, sender, recipients, body, attachments)
- Ability to identify inline vs Outline attachments (To remove signatures and email footer images)

- Performance and Setup easibility
- Platform supported
- Pricing and Limitation

Below is a table of performance of each library and its limitations.

Library	Format Support	Metadata Extraction	Inline vs Outline Attachment	Performance & Setup	Platform	Licensing
Aspose.Email	.PST, others	Yes	Yes	Fast & Easy	Crossplatform	Free for 50 mails, subscription required
Extract_msg	.MSG only	Yes	No	Fast & Easy	Crossplatform	Open-source, format limited
pypff	.PST	Yes	No	Slow, Complex Setup	Crossplatform	Open-source, hard to automate
Outlook COM	.MSG, .PST	Yes	Yes	Good but Platform bound	Windows only	Free, not portable

Table 1 - Comparison of Email Processing Libraries

- **Aspose.Email** is capable of extracting attachments, body text, subject lines, and other metadata accurately. It can also distinguish between inline and regular attachments, which is useful for filtering out elements like logos and email signatures. However, it comes with a major limitation: it's subscription-based. The free version only allows extraction from the first 50 emails of each .PST file, after which a paid license is required.
- **Extract_msg** can accurately extract metadata such as sender, recipients, subject, and body text. However, it only supports .MSG files, making it incompatible with our current dataset. Additionally, it cannot differentiate between inline and regular attachments—an essential feature for removing embedded images and signature blocks.
- **pypff** offers similar functionality to extract_msg but is notably slower and more complex to set up. This makes it unsuitable for our goal of building a repeatable and

automated system, especially one that can be easily re-executed via shell scripts when retraining the model.

- **Outlook COM Automation** is another option that provides access to both .MSG and .PST formats and supports inline attachment detection. However, it only works on Windows, which limits its portability and makes it less suitable for a cross-platform or server-based deployment.

4.1.2. Solution

Since each library has its own limitation and there is none of them fits it all, I had to work around them. The solution I implemented was to combine the **Aspose.Email**'s ability to process .PST file and **extract_msg** library to extract information.

The script to collect dataset is divided into 3 parts:

- Converting .PST file to multiple .MSG files (one .MSG for one email) with **Aspose.Email**
- Utilizing **extract_email** library to extract information and attachment from .MSG files.
- Implementing safeguards and filters to remove signature and footer images and attachments.

4.2. Data Cleaning

4.2.1. Email body cleaning

To extract clean and usable data from customer emails, I developed a multi-step text preprocessing pipeline. The primary goal was to isolate the main message content and remove non-informative elements like email signatures, automatic footers, greetings, and previous replies.

I begin by parsing the HTML content of the email using the BeautifulSoup(8) library. Tags such as
, <p>, and <div> were converted into newline placeholders to ensure the structure is preserved when extracting plain text. Special attention was given to <div> elements with class or ID values containing the word "signature", these were removed entirely to avoid capturing footers or company branding.

After converting the HTML to text, I inserted line breaks before certain key phrases like "Bonjour (*Translation: Hello*)", "Envoyé (*Translation: Sent - Indicating sent time*)", or "De : (*Translation: From*)" to improve the segmentation of email components. I also

eliminated hidden or corrupt Unicode characters and split the text into individual lines for easier analysis.

Then each line was then cleaned further:

Lines containing metadata (e.g., sender, recipient, subject) or unwanted content (URLs, phone numbers, formal closings like "Cordialement" (*Translation: Best Regards*)) were ignored.

If the system detected the beginning of a reply to Cosywee's customer service team, it stopped processing further, as these are often auto-generated or duplicate content.

The remaining lines, assumed to be part of the user's message, were concatenated into a unified body.

Lastly, we performed post-processing to remove repeated symbols (e.g., dashes or underscores), and polite phrases that do not add analytical value (e.g., "Merci beaucoup (*Translation: Thank you*)", "Bonne journée (*Translation: Hope you have a good day*)"). The resulting cleaned message was then stored in a structured format.

4.2.2. Attachment cleaning

Since I was using **Tesseract** (Pytesseract) (pytesseract: Python wrapper for Google Tesseract-OCR, 2024) and **EasyOCR** (JaideAI, 2024) to extract texts and OCR images and PDF files, I had to clean those extracted text again to remove low-quality texts and noises. Here is a function where I clean text, remove noises and OOV (out of vocabulary) words, except custom words used in the company.

Here is the function to extract texts from images and pdf files:

Algorithm 1 Extract Forwarded Custom Email Data

```
1: procedure EXTRACTEMAILDATA(html)
2:   Parse html with BeautifulSoup into soup
3:   Replace <br>, <p>, <div> with [[NEWLINE]]
4:   Remove <div>s with "signature" in class or id
5:   Convert soup to plain text
6:   Replace [[NEWLINE]] with real line breaks
7:   Insert line breaks before keywords (e.g. "Envoyé", "De :," "Bonjour")
8:   Remove non-UTF characters
9:   Save result to log file
10:  Initialize data dictionary with empty fields
11:  Initialize body_content as empty string
12:  for each line in the cleaned text do
13:    Normalize and lowercase line
14:    if line starts with metadata (e.g., "de :", "cc :", "objet :") then
15:      continue
16:    else if line is a known Cosywee reply signature then
17:      break
18:    else if line contains contact info or common closing phrases then
19:      continue
20:    else
21:      Append cleaned line to body_content
22:    end if
23:  end for
24:  Clean body_content (remove dashes, signatures, greetings)
25:  Set data.body and data.body_detached
26:  return data
27: end procedure
```

Algorithm 2 Extract Text Word by Word with EasyOCR, Tesseract and Trained TrOCR

```
1: procedure EXTRACTTEXT(image, header)
2:   Correct image orientation
3:   Initialize empty text and seen words set
4:   Crop image from header to bottom as img_full
5:   Convert img_full to grayscale
6:   Determine orientation (vertical or horizontal)
7:   Divide image dimension into 3 sections
8:   for i ∈ [0, 1, 2] do
9:     Extract section slice from grayscale image
10:    Try
11:      Run EasyOCR on section
12:    Catch
13:      Set EasyOCR results to empty
14:      Extract bounding boxes from EasyOCR results with confidence threshold
15:      Run Tesseract OCR on section, extract bounding boxes with confidence threshold
16:      Merge and sort all bounding boxes from both OCRs
17:      for each bounding box do
18:        Pad and crop bounding box area
19:        if crop valid then
20:          Run OCR on cropped image to get word (trained TrOCR)
21:          if word not empty and unseen then
22:            Add word to text and seen set
23:            Draw bounding boxes on debug images
24:          end if
25:        end if
26:      end for
27:    end for
28:    Save debug images
29:    clean accumulated text with clean_text() function
30:  return accumulated text
31: end procedure
```

Algorithm 2 above shows the process of extracting text from an image. First, the scanned document(image) orientation is corrected (by rotating the document four times with different angles and selecting the angle that results in the most text written from left to right). Then, the scanned document is cropped into three different smaller images (or

sections), and each section is passed to EasyOCR and Pytesseract OCR engines to detect text.

The reason why we use both OCR engines is because EasyOCR tends to work better with handwritten text, while Pytesseract performs well on printed text.

Next, for each detected text, the program extracts their bounding boxes (coordinates of the detected text). Using those coordinates, smaller images are cropped, one image per word or phrase in a line. Each image is then passed through OCR again, and the extracted text is added to a list to build the dataset. This is how the program combines text detection capabilities from EasyOCR and Pytesseract with the text extraction capabilities of TrOCR.

And here is my function to clean the text after performing OCR as described above. In this case, I used two additional resources for the French language, including:

- spacy's '**fr_core_news_lg**' (spaCy French Language Model: fr_core_news_lg, 2024) model
- **Jean-Baptiste/camembert-ner** (CamemBERT NER - HuggingFace model, 2021) model from HuggingFace

Algorithm 3 above shows the process of cleaning text using spaCy and Named Entity Recognition (NER) filtering. First, the input string is cleaned by removing commas and semicolons. Then, spaces are inserted around punctuation that may be attached to words. The cleaned string is then processed with spaCy's NLP pipeline, which converts the text into lowercase and tokenizes it.

For each token, the algorithm checks whether the word is found in the company's custom vocabulary. If so, it is preserved and directly added to the list of cleaned tokens. Otherwise, it continues the filtering steps. Tokens that are punctuation marks, spaces, numbers, or contain any digits are skipped. It also checks whether the token's base form is a color or a month, etc. If it actually is, the token will be removed too. After that, some tokens will be filtered out if they are recognized with high confidence by the NER engine as locations, organizations, or people, etc. Finally, a token is only kept if it is not a stopword, not out-of-vocabulary, has at least two characters, and is not a duplicate.

4.2.3. Final Data

In this subsection, the final data is presented. The training and testing data are shuffled and selected randomly from a dataset that is saved in a .CSV (Comma-separated Values) file. This file consists of multiple columns containing different information including sender, recipients, subject (Object of the email), body, contents extracted from attachments, etc. You may find an example of the dataset in Figure 4 - Dataset.

```
"sep="",
label;senders;recipients;subjects;body;body_detached;full_content;attachments
PERPIGNAN;SARL ALPLAST <alplast@wanadoo.fr>;;commande;ci-joint une commande cdt 2 av de l'hérault 11110 coursan;joint commande cdt av;L'ALES;BRUGUIER MENUISERIE <bruguier.menuiserie@orange.fr>;;commande usseglio moustiquaires;;;l'objectif de cet e-mail est commande usseg;
ALES;contact@garciafermetures.fr <contact@garciafermetures.fr>;;commande ref laignelot - moustiquaires;veuillez trouver ci-joint une nou
ALES;Jean-Baptiste <griveljh@yahoo.fr>;;bon de commande pour vénitiens ;grivel stores et fenêtres 35 route de tully 74200 thonon les bai
ALES;GRANCHETTE, Pauline <pauline.granchette@lapeyre.fr>;;re: demande faisabilité ;oui c'est possible, je vous fais un devis dès que pos
LYON;Marianne <service.logistique@etg-enseigne.fr>;;re: commande - buffalo grill choleterrait-il possible d'ajouter des capuchons extre
ALES;yoann BETSCH <yoann@map-betsch.com>;;re: demande de prix ref ccrlp;manoeuvre par chainette ccrlp ; nous ne faisons pas de vénitien
MARSEILLE;DAVID PRUD'HOMME <dpm.prudhomme@gmail.com>;;commande; vous trouverez ci-joint une commande .;trouverez joint commande;l'objec
ALES;2B MENUISERIE - Yannick JOURDAIN <yannick.jourdain@2bmenuiserie.com>;;2b menuiserie commande yj/jourdain ;en rapport à la commande
ALES;Esteban Saunier - SECODOIS <esteban@secodois.fr>;;devis moustiquaire ezia;m. laurenti, est-ce que vous pouvez me faire un devis po
ALES;patrick sevenier <patrick-sevenier@orange.fr>;;bonne commande ammathieu;mathias sevenier annule et remplace par celle ci commande
```

Figure 4 - Dataset

This data was then split into training, testing, and validation set. Since two types of AI models/architectures were used (to be detailed in the next chapters), the data is split differently as well.

- **CamemBERT**: Mainly 70% for training, 15% for validation and 15% for testing *
- **Traditional Machine Learning Algorithms (Random Forest)**: Mainly 80% for training, and 20% for testing **

* *This is default partition. Some models are trained on different partition depending on availability of data.*

** *This is partition is for training. For model experimentation, only 70% of the data were used to reduce workload.*

4.3. Model Experimentation

After all of the data are cleaned and written into a .CSV file, we finally get dataset with 7 columns: **Label, Sender, Recipients, SentTime, Subject, Body text, and Text extracted From Attachments**.

This data were used to trained multiple algorithms and architectures of models including Naive Bayes, Random Forest, Logistic Regression, Support Vector Machine, etc. Moreover, it was also used to finetune existing models like CamemBERT (French Large Language Model) as well.

As a reminder, we have to classify between four different factories, COSYWEE ALES, COSYWEE LYON, COSYWEE MARSEILLE, COSYWEE PERPIGNAN. At the same time, we also have to classify email's intents (whether it is an order, an invoice, a price quotation request, an After-Sale service request, or anything related to Accounting department, etc). Moreover, we have to take into account the fact that two factories produces same type of products, meaning if the mail is related to the shared products, the program have to sort them to a sub-folder for human to manually classify them.

In figure 4.B, you may find the proposed pipeline.

- An AI model for intent classification
- An AI model for classification between **COSYWEE ALES, COSYWEE LYON, and MP** (short of Marseille and Perpignan)
- Another AI model for classification between **COSYWEE MARSEILLE, COSYWEE PERPIGNAN** and **Roller Shutters** products.

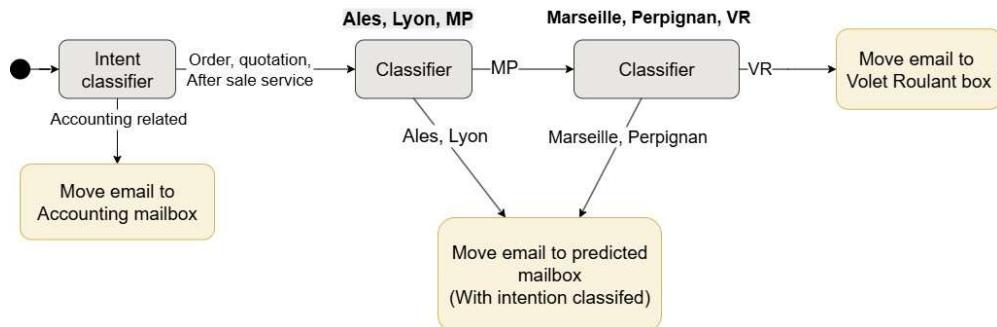


Figure 5 - Overview of model pipeline

As for the model experimentation and development, **Optuna** (Optuna: A Next-generation Hyperparameter Optimization Framework, 2019), and Scikit-learn (Scikit-learn: Machine Learning in Python, 2011) were used to find the best performing model to implement into our pipeline. Each algorithm (Traditional Machine Learning algorithms only) was experimented with hundreds of combinations of parameters with number of trials = 250.

```

1 nltk.download('stopwords')
2 french_stopwords = stopwords.words('french')
3
4 def get_common_preprocessing(trial):
5     max_features = trial.suggest_categorical('vectorizer__max_features', [200,
6         400, 500, 750, 1000, 2000, 3000, 4000, 5000])
7     ngram_str = trial.suggest_categorical('vectorizer__ngram_range', ['1,1', '1,2',
8         ', 1,3', '2,2'])
9     ngram_range = tuple(map(int, ngram_str.split(',')))
10    stop_words_option = trial.suggest_categorical('vectorizer__stop_words', [
11        'none', 'french'])
12    stop_words = None if stop_words_option == 'none' else french_stopwords
13    binary = trial.suggest_categorical('vectorizer__binary', [True, False])
14    k = trial.suggest_categorical('selector__k', [70, 80, 100, 150, 200, 300,
15        500, 700])
16    k = min(k, max_features)
17    score_func_name = trial.suggest_categorical('selector__score_func', ['chi2',
18        'f_classif'])
19    score_func = chi2 if score_func_name == 'chi2' else f_classif
20    vectorizer = CountVectorizer(max_features=max_features, ngram_range=
21        ngram_range, stop_words=stop_words, binary=binary, preprocessor=
22        clean_text_for_company_words)
23    selector = SelectKBest(score_func=score_func, k=k)
24
25    return [('vectorizer', vectorizer), ('selector', selector)]

```

The algorithms that I experimented on include:

- Logistic Regression (LR)

```

1 def lr_objective(trial):
2     common_steps = get_common_preprocessing(trial)
3     C = trial.suggest_float('C', 0.1, 5.0, log=True)
4     penalty = trial.suggest_categorical('penalty', ['l1', 'l2'])
5     solver = 'liblinear'
6     classifier = LogisticRegression(
7         C=C,
8         penalty=penalty,
9         solver=solver,
10        max_iter=1000
11    )
12     pipeline = Pipeline(common_steps + [('classifier', classifier)])
13     scores = cross_val_score(pipeline, data["full_with_att"], data[""
14         "label_encoded"], cv=5, scoring='accuracy')
15
16     return scores.mean()

```

- Support Vector Classifier (SVC)

```

1 def svm_objective(trial):
2     common_steps = get_common_preprocessing(trial)
3     C = trial.suggest_float('C', 0.1, 10.0, log=True)
4     kernel = trial.suggest_categorical('kernel', ['linear', 'rbf'])
5     classifier = SVC(C=C, kernel=kernel)
6     pipeline = Pipeline(common_steps + [('classifier', classifier)])
7     scores = cross_val_score(pipeline, data["full_with_att"], data["label_encoded
8         "], cv=5, scoring='accuracy')
9
9     return scores.mean()

```

- eXtreme Gradient Boosting (XGB) (14)

```

1 def xgb_objective(trial):
2     steps = get_common_preprocessing(trial)
3     n_estimators = trial.suggest_categorical('classifier__n_estimators',
4                                              [100, 300, 400, 500])
5     max_depth = trial.suggest_int('classifier__max_depth', 3, 7)
6     learning_rate = trial.suggest_float('classifier__learning_rate', 0.01,
7                                         0.2, log=True)
8     subsample = trial.suggest_categorical('classifier__subsample', [0.7,
9                                         1.0])
10    model = XGBClassifier(
11        n_estimators=n_estimators,
12        max_depth=max_depth,
13        learning_rate=learning_rate,
14        subsample=subsample,
15        use_label_encoder=False,
16        eval_metric='logloss'
17    )
18    pipeline = Pipeline(steps + [('classifier', model)])
19    scores = cross_val_score(pipeline, data["full_with_att"], data[""
20        "label_encoded"].astype(int), cv=5, scoring='accuracy')
21
22    return scores.mean()

```

- Random Forest (RF)

```

1 def rf_objective(trial):
2     steps = get_common_preprocessing(trial)
3     n_estimators = trial.suggest_categorical('classifier__n_estimators',
4                                              [20, 30, 50, 70, 85, 90, 100, 150, 200, 250, 300, 400, 500])
5     max_depth = trial.suggest_categorical('classifier__max_depth', [None, 2,
6                                         3, 5, 7, 10, 20, 30])
7     min_samples_split = trial.suggest_categorical('
8         classifier__min_samples_split', [2, 3, 4, 5, 7, 8, 10])
9     model = RandomForestClassifier(
10        n_estimators=n_estimators,
11        max_depth=max_depth,
12        min_samples_split=min_samples_split
13    )
14    pipeline = Pipeline(steps + [('classifier', model)])
15    scores = cross_val_score(pipeline, data["full_with_att"], data[""
16        "label_encoded"], cv=5, scoring='accuracy')
17
18    return scores.mean()

```

- K-Nearest Neighbors (KNN)

```

1 def knn_objective(trial):
2     steps = get_common_preprocessing(trial)
3     n_neighbors = trial.suggest_categorical('classifier_n_neighbors', [3, 5, 7,
4         10])
5     weights = trial.suggest_categorical('classifier_weights', ['uniform',
6         'distance'])
7     algorithm = trial.suggest_categorical('classifier_algorithm', ['auto',
8         'ball_tree', 'kd_tree', 'brute'])
9     p = trial.suggest_categorical('classifier_p', [1, 2]) # 1 = Manhattan, 2 =
10        Euclidean
11    model = KNeighborsClassifier(
12        n_neighbors=n_neighbors,
13        weights=weights,
14        algorithm=algorithm,
15        p=p
16    )
17    pipeline = Pipeline(steps + [('classifier', model)])
18    scores = cross_val_score(pipeline, data["full_with_att"], data["label_encoded"],
19        cv=5, scoring='accuracy')
20
21    return scores.mean()

```

- Naive-Bayes (NB)

```

1 def nb_objective(trial):
2     steps = get_common_preprocessing(trial)
3     nb_type = trial.suggest_categorical("classifier_type", ["multinomial",
4         "bernoulli"])
5     if nb_type == "multinomial":
6         alpha = trial.suggest_float("classifier_alpha", 1e-3, 1.0, log=True)
7         classifier = MultinomialNB(alpha=alpha)
8     else:
9         alpha = trial.suggest_float("classifier_alpha", 1e-3, 1.0, log=True)
10        binarize = trial.suggest_float("classifier_binarize", 0.0, 1.0)
11        classifier = BernoulliNB(alpha=alpha, binarize=binarize)
12    pipeline = Pipeline(steps + [('classifier', classifier)])
13    scores = cross_val_score(pipeline, data["full_with_att"], data["label_encoded"],
14        cv=5, scoring='accuracy')
15
16    return scores.mean()

```

- CamemBERT (French Large Language Model)

As for the CamemBERT model, which is a transformer-based model, I experimented on multiple hyperparameters including : **Learning rates, train and test batch sizes, Train epoches, Regularization techniques**, etc.

4.4. Evaluation Metrics

In this section, we discuss the performance of our models by evaluating them on the three metrics including Accuracy, Precision and F1-score.

Confusion Matrix and Terminology

These metrics are based on four possible outcomes in a binary classification context:

- **TP (True Positives)**: Instances correctly predicted as positive.
- **TN (True Negatives)**: Instances correctly predicted as negative.
- **FP (False Positives)**: Instances incorrectly predicted as positive (also called Type I error).
- **FN (False Negatives)**: Instances incorrectly predicted as negative (also called Type II error).

These values are typically organized in a confusion matrix as follows:

Actual / Predicted	Positive	Negative
Positive	True Positive (TP)	False Negative (FN)
Negative	False Positive (FP)	True Negative (TN)

Evaluation Methods

- **Accuracy:**

Accuracy measures the proportion of correct predictions out of the total number of predictions made.

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad (4.1)$$

- **Precision:**

Precision measures the proportion of positive predictions that were actually correct.

$$\text{Precision} = \frac{TP}{TP + FP} \quad (4.2)$$

- **F1 Score:**

The F1 Score is the harmonic meaning of Precision and Recall. It balances both metrics, which are especially useful when dealing with imbalanced datasets.

$$\text{F1 Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (4.3)$$

4.5. Results

As a reminder, in our pipeline, we need 5 distinct trained models to better capture the context, intent, and products mentioned in each email, while also helping to reduce model complexity. Please check on figure 4.B for the pipeline architecture.

4.5.1. Intent classification model

Because this model focuses on sentiment analysis instead of relying on keywords (product codes or names), I chose to implement BERT architecture model instead of traditional machine learning methods or other architectures.

After multiple experimentations and hyper-parameters testing, these following hyperparameters yield the best results:

```
learning_rate=3e-6
per_device_train_batch_size=8,
per_device_eval_batch_size=6,
num_train_epochs=70,
weight_decay=0.01,
metric_for_best_model="eval_loss",
greater_is_better=False,
early_stopping_patience=10
```

With these parameters, the statistics were:

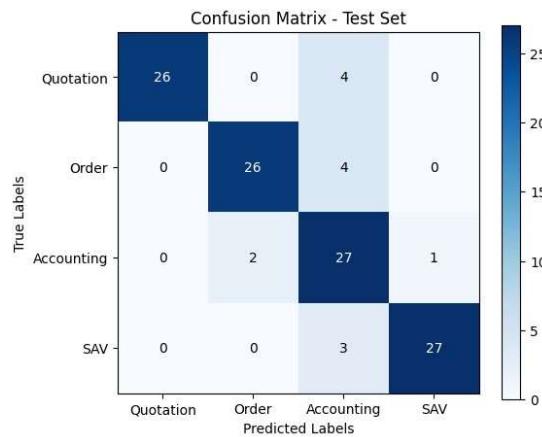


Figure 6 - Confusion matrix of intents classification

As shown in Table 2 and Figure 6, approximately 9 out of 10 emails are classified correctly with this model. Some emails were wrongly classified in Class 2 (Accounting) because of the nature of the emails themselves. Some emails contain multiple intents, which led to confusion. For example, an email is sent to place an order and also to request bank account information (in the same email), etc.

Data	Accuracy	Precision	F1
Validation Set	92.30%	92.41%	93.07%
Test Set	88.33%	90.08%	88.76%

Table 2 - Performance of intents classification model

4.5.2. Alès, Lyon, MP classification

4.5.2.1. CamemBERT Model

After extensive experimentation and hyperparameter tuning, the following settings yielded the best results:

```

learning_rate=6e-6
per_device_train_batch_size=8,
per_device_eval_batch_size=6,
num_train_epochs=70,
weight_decay=0.01,
metric_for_best_model="eval_loss",
greater_is_better=False,
early_stopping_patience=12

```

With these settings, the performance metrics were:

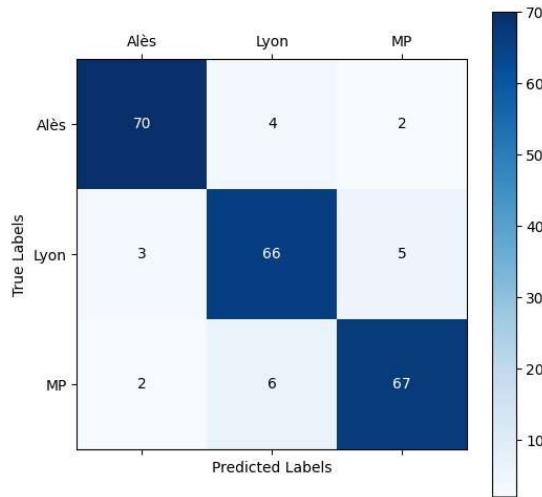


Table 3 - Confusion Matrix of Alès, Lyon, MP classification

As shown in Figure 4.E and Table4.3, around 90% of emails are correctly classified (in the test set). The reason why Class 0 (Alès) has the most correctly classified samples is that the majority of the products in this class are different from the other two classes.

Data	Accuracy	Precision	Recall
Validation Set	92.48%	92.77%	92.51%
Test Set	90.22%	90.24%	90.22%

Table 4 - Performance of Alès, Lyon, MP classification model

4.5.2.2 Traditional Machine Learning Algorithms

Utilizing Optuna and Scikit-learn, I have experimented with seven different models such as XGB, Random Forest,...

I used only 70% of the data for this experimentation, and below are each model's accuracy on the balanced dataset:

Type of Model	Hyperparameters	Accuracy
Logistic Regression	C = 4.2676, penalty = L2	90.28%
Support Vector Classifier	k = 200, Kernel = rbf, C = 8.9056	90.48%
XGB	learning_rate = 0.0469, sub_sample = 1	91.01%
Random Forest	n_estimators = 85, max_depth = None = 1, min_samples_split = 2	90.79%
K-Nearest Neighbors	n_neighbors = 10, algorithm = Brute, p = 1	84.46%
Naive-Bayes	type = multinomial, alpha = 0.0090	87.69%

Table 5 - Performance Comparison of Algorithms - Alès, Lyon, MP classification

Based on the results above, the XGBoost algorithm achieves the highest performance (91.01% accuracy). However, considering the importance of model interpretability in operational environments, for traceability, auditing, and maintainability, the Random Forest algorithm (90.79% accuracy) was selected. This model offers a strong balance between performance and transparency, making it more suitable for practical deployment in this context.

After selecting Random Forest as our base model, I split the dataset into training and testing sets using an 80/20 split. Then, I retrained the model, achieving an accuracy of **91.08%** on test set.

4.5.3. Marseille, Perpignan, Roller Shutter (*VR*) classification

4.5.3.1 CamemBERT Model

After extensive experimentation and hyperparameter tuning, the following settings yielded the best results:

```
learning_rate=5e-6
per_device_train_batch_size=8,
per_device_eval_batch_size=6,
num_train_epochs=70,
weight_decay=0.01,
metric_for_best_model="eval_loss",
greater_is_better=False,
early_stopping_patience=7
```

With these settings, the performance metrics were:

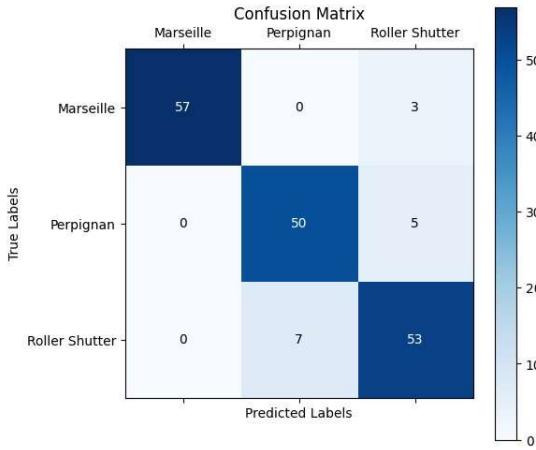


Table 6 - Confusion Matrix of Marseille, Perpignan, Roller Shutter classification

In Table 7 and Table 6, the statistics of the results of the Marseille, Perpignan, Roller Shutter classification model are shown. There is a small confusion between Class 1 and Class 2 (Perpignan and Roller Shutter), because sometimes, the characteristics used to describe the product lines of the 2 factories are identical, making it harder for the model to predict without a definite keyword.

Data	Accuracy	Precision	F1
Validation Set	92.03%	92.29%	92.13%
Test Set	91.44%	91.53%	91.44%

Table 7 - Performance of Marseille, Perpignan, Roller Shutter classification model

4.5.3.2 Traditional Machine Learning Algorithms

Utilizing **Optuna** and Scikit-learn, I have experimented with seven different models such as XGB, Random Forest, ...

I used only 70% of the data for this experimentation and below is each model's accuracy on balanced dataset:

Type of Model	Hyperparameters	Accuracy
Logistic Regression	C = 1.5059, penalty = L1	89.55%
Support Vector Classifier	k = 200, Kernel = rbf, C = 8.26	89.55%
XGB	Learning_rate = 0.0373, sub_sample = 0.7	90.17%
Random Forest	n_estimators = 150, max_depth = None = 1, min_samples_split = 3	90.29%
K-Nearest Neighbors	n_neighbors = 3, algorithm = kd_tree, weights = distance, p = 1	83.01%
Naive-Bayes	type = multinomial, alpha = 0.05657	80.91%

Table 8 - Performance Comparison of Algorithms - Marseille, Perpignan, Roller Shutter classification

Based on the results above, the Random Forest algorithms achieve the highest performance metrics (90.29%).

With an 80/20 split, I retrained the model, achieving an accuracy of approximately **92.63%**. This might seem low but considering the circumstances and the nature of the class (overlapping with no clear boundary), these results satisfied the company's expectations.

4.6. Final Pipeline

After the experimentation described above and multiple iterations of trial and error, I implemented a pipeline that combines several models, including CamemBERT-based models and a Random Forest classifier.

As shown in 4.F, the CamemBERT model is used for intent classification, while a combination of CamemBERT and Random Forest is employed to classify between COSYWEE ALES, COSYWEE LYON, COSYWEE MARSEILLE, COSYWEE PERPIGNAN, and Roller shutter (*Volet Roulant - VR*).

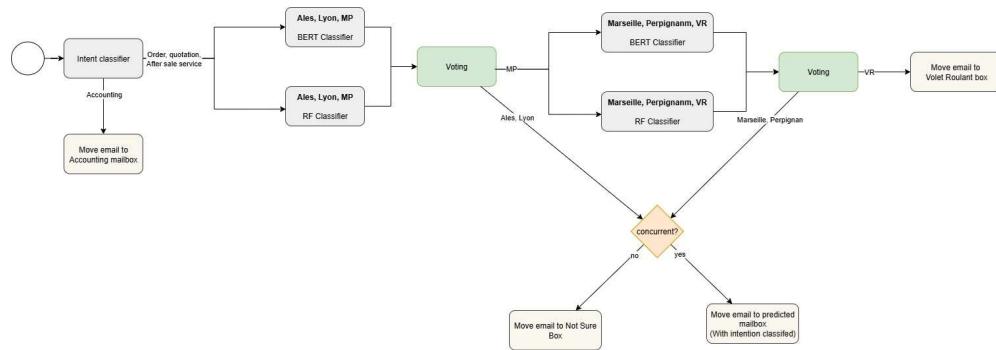


Figure 7 - Fully integrated AI pipeline

As illustrated in the flowchart above, I employed a combination of CamemBERT-based models and Random Forest models to classify its class using a hard voting strategy. This generally means that if both models agree on a specific class, the email would be classified automatically. But if their votes contradict, we will move them to a folder called "**A CLASSER (Not Sure)**", which will be sorted manually later.

This approach has proven effective as it captures spam emails, and emails with insufficient information to the "**A CLASSER (Not Sure)**" folder for further inspection, ensuring all emails are correctly classified.

V - IMPLEMENTATION & PROTOTYPE BUILDING

After a detailed experimentation and literature review have been conducted, I began to implement the pipeline into production.

Although until now this project may seem like a normal text-classification problem, in order to maximize accuracy, taking advantage of existing clues, and keys I found on project analysis, I decided to implement the trained text classification AI pipeline with other heuristics as well as other methods as well.

Other methods and heuristics include Image classification, OCR, Regular Expression, etc.

5.1. Project Development

5.1.1. Environment configuration

Before diving into the implementation, it is crucial to configure the development and testing environment correctly first. This project needs multiple libraries and tools to function efficiently and correct without error like Pytorch(15), Pytesseract(9), OpenCV(16), beautifulsoup4(8), EasyOCR(10), Pandas(17), Transformers, PyMuPDF(18), etc. for the text classification process.

Those libraries include:

```
opencv-python == 4.11.0.86
pytorch      == 2.7.0
pytesseract  == 0.3.13
beautifulsoup4 == 4.13.3
PyMuPDF     == 0.0.1.dev2
EasyOCR      == 1.7.2
transformers  == 4.50.1
docx2txt==0.9 == 0.9
nltk         == 3.9.1
paddleocr    == 2.4
extract_msg   == 0.54.1
spacy        == 3.7.5
msal         == 1.32.0
pandas       == 1.5.3
pillow       == 10.4.0
```

ftfy == 6.3.1

These specific library versions are required to ensure the program runs correctly. In particular, compatibility issues exist between newer versions of opencv-python and paddleocr. Starting from version 4.12.0.0 of opencv-python, certain internal dependencies and APIs used by paddleocr are no longer supported or behave differently, leading to runtime errors such as segmentation faults, unexpected crashes, or broken image preprocessing pipelines. So we have to install the libraries version as shown above to ensure the compatibility between those libraries.

5.1.2. Project Structure

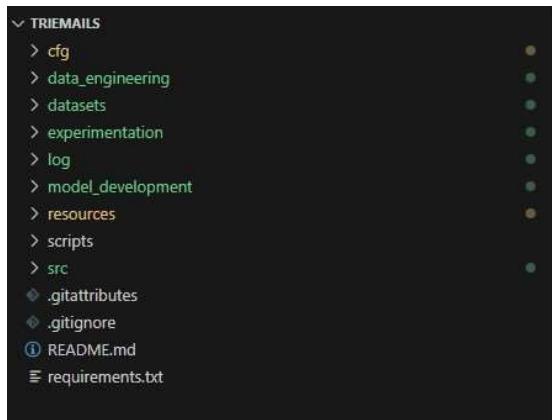


Figure 8 - Project Structure

Figure 8 shows the structure of the entire project:

- **cfg**: Contains configuration files such as access tokens and path settings.
- **data_engineering**: Contains scripts to scrape and clean data from .PST and .MSG files.
- **datasets**: Contains all scraped and cleaned datasets.
- **experimentation**: Contains all .ipynb and .py scripts used to experiment with different libraries and machine learning model architectures.
- **log**: Contains log files (mainly for debugging).
- **model_development**: Contains trained models and code (with documentation) to retrain or fine-tune them if necessary.
- **resources**: Contains important files such as mappings between products and companies, or employees and companies.
- **scripts**: Contains shell scripts (.sh) used to automate parts of the project.

- **src**: Contains all script files needed to deploy the project.
- **requirements.txt**: A text file listing all packages and libraries required to run the program.

5.2. Heuristics methods

To quote one of my professors in class during his first lecture, Monsieur Sébastien Harispe, "**Avoid Machine Learning as much as possible**", as the system built with it would need maintenance over time once deployed.

With that said, I tried to integrate other heuristic approaches to achieve the same tasks as much as possible before using our trained model as the last resort.

That being said, I analyzed the flow of the emails, the attachments, and the senders themselves and found some things that I could take advantage of using some heuristic methods.

Those characteristics include:

5.2.1. Original Sender of the email

As previously mentioned above, the email flows into **service.client@cosywee.com** is bidirectional. Some emails are initiated by clients (Inquiry about new products, orders, or price quotation etc). Others are replies from client responding to the company employees. Upon careful analysis, a consistent pattern is identified. All emails originally sent by company employee starts with this pattern in the subject (*objet*) :

"(Cosywee-XXX) Commande n°" or "(Cosywee-XXX) Devis n°"

The program checks if the email was originally sent by an employee by parsing its subject with Regular Expression to find the specific patterns above. If the Regular Expression is matched, it filters out all lines of the email except the signature section to find the name of employee who sent the email . This then matches this name to the company associated with it, classifying emails correctly.

5.2.2. Document Creator

Another heuristic that can be leveraged involves identifying the document creator. Cosywee uses its own in-house software platform, developed by the company's IT team,

called Nexus. This platform is used to create and manage business documents such as Quotation, Purchase Order, Order Confirmation, among others.

Particularly, each document generated by an employee using this internal software includes the creator's name embedded in the document itself. As I did some research and analysis on customer behavior, I found that customers sometimes attach these documents to their emails as a reference for their inquiries.

Based on the information above, this feature of **Nexus** actually presents an opportunity for another heuristic approach with the help of OCR and the pdf_extraction library.



Figure 9 - Employee's name to Company

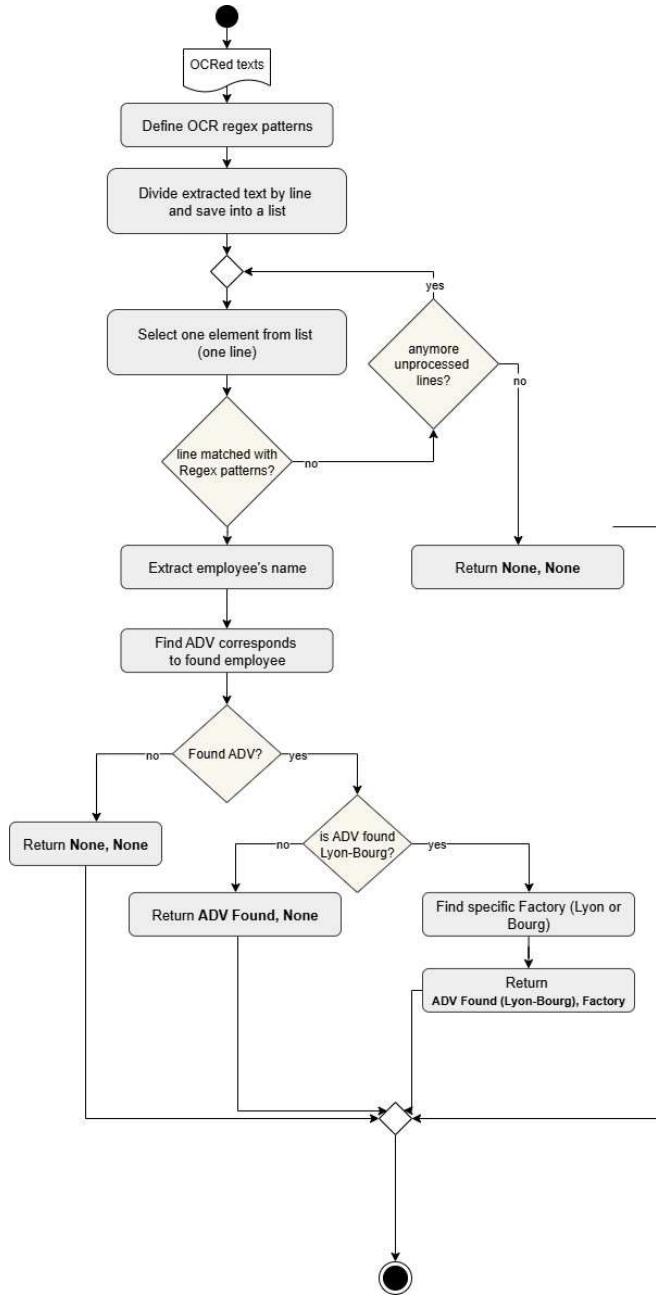


Figure 10 - Employee's name extraction from Purchase Order Flowchart

Therefore, by parsing extracted/OCRed text with a regular expression that searches for the pattern '**Document établi par XXXXXXXX**' (English translation: Document created by XXXXXXXX), we can then extract the name of the employee. This name can then be mapped and matched to the associated company, allowing a more reliable and accurate email classification feature Figure 9. Figure 10 shows the flow of the process to find the employee's name from the document generated.

5.2.3. Product extraction from Purchase Order

For each product produced by Cosywee, there is a Purchase Order that allows clients to complete the relevant product details such as dimensions, color, etc. This purchase order is created by Marketing team and follows the format as shown in Figure 11.

BON DE COMMANDE MOUSTIQUAIRES

DEMANDE DE DEVIS
 COMMANDE Selon devis n°:
 DATE: 07/02/25
 RÉF. CHANTIER: SEGARD
 LIVRAISON SOUHAITÉE SEMAINE: ...
 (sous réserve du planning des travaux)

RAISON SOCIALE: **Menuiseries Stores Méditerranéens**
 ADRESSE: ZA LA MASSANE Rue des Bauxites
 19210 St Rémy-de-Provence
 Tél: 06 22 70 27 51 - Fax: 04 90 87 46 51
 www.entreprise-msm.com
 Siret: 462 941 742 00023 - ARF 4112 B

Cachet et signature

Menuiseries Stores Méditerranéens
 ZA LA MASSANE Rue des Bauxites
 19210 St Rémy-de-Provence
 Tél: 06 22 70 27 51 - Fax: 04 90 87 46 51
 www.entreprise-msm.com
 Siret: 462 941 742 00023 - ARF 4112 B

COSYWE
 PÔLE PROTECTION SOLAIRE
 tél: 04 74 10 79 01
 service.clients@cosywee.com

BRC/2025/MOIS/GENERAL																
1 Optima Enroul. Verticale	2 Esta Enroul. Verticale	3 Evolution Enroul. Verticale	4 Lussa Enroul. Latérale	5 Classique Enroul. Latérale	6 Golfe Enroul. Latérale	7 Flexeo Enroul. Latérale	8 Planès Enroul. Latérale	9 Finée Plissée pour fenêtre	10 Plissée 2 Plissée pour porte	11 Listelle Cadre fixe	12 Slalom Coulissante	13 Altezza 2 Battante	14 Altezza 1 Battante	15 Altezza 3 Battante		
1A - Auto Clic and Push 1B - Auto Clic clac 1C - Châssis	3A - Auto régulée 3B - Auto non réglée 3C - Châssis	5A - 1 vantail droit 5B - 1 vantail gauche (vue intérieur) 5C - 2 vantails	6A - 1 vantail 6B - 2 vantails	7A - 1 vantail 7B - 2 vantails	8A - 1 vantail 8B - 2 vantails	9A - 1 vantail 9B - 1 vantail réversible 9C - 2 vantails	10A - 1 vantail 10B - 2 vantails centraux 10C - 2 vantails latéraux	11A - En applique 11B - Tableau	12A - 1 vantail 12B - 2 vantails	13A - 1 vantail poussant droit 13B - 2 vantails poussant droit 13C - 2 vantails poussant gauche	13D - 2 vantails poussant gauche					
COTES																
Repère	Type 1 à 13	Qté	Largeur de fabrication (mm)	Hauteur de fabrication (mm)	COLORIS (Selon modèle)											
					Blanc RAL 9010 (Base noir/rouge)	Marbre RAL 8017 (Pierre/rouge)	Beige RAL 9015 (Terrelé/bleu)	Vert RAL 6005 (Terrelé/vert)	Gris RAL 7015 (Terrelé/bleu)	Noir RAL 9005 (Terrelé/noir)	Fumé RAL 9003 (Terrelé/fumé)	Fluo Rose RAL 3000 (Terrelé/rose)	Autre RAL N°			
					Tube Noir	Tube Gris	Tube Beige	Tube Vert	Tube Gris	Tube Beige	Tube Fumé	Tube Rose				
					Tube Supra (1/2, 3/4, 1, 1 1/2, 2)	Forme A (1/2, 3/4, 1, 1 1/2, 2)	Forme B (1/2, 3/4, 1, 1 1/2, 2)	Posé de Rev (1/2, 3/4, 1, 1 1/2, 2)	Joint Bulé (1/2, 3/4, 1, 1 1/2, 2)	Enroulé (1/2, 3/4, 1, 1 1/2, 2)	Côté paravape châssis (1/2, 3/4, 1, 1 1/2, 2)	No. traverse renforcée (1/2, 3/4, 1, 1 1/2, 2)	Non traverse renforcée (1/2, 3/4, 1, 1 1/2, 2)	Non traverse supplémentaire (1/2, 3/4, 1, 1 1/2, 2)	Non plinthe supplémentaire (1/2, 3/4, 1, 1 1/2, 2)	Fermé joint renfor
														OBSERVATIONS		
bc-1 1690 2030 X																

Figure 11 - Purchase order for mosquito net

As you can see in Purchase order Figure 12, you can find the name of the product on the top left of the document (**BON DE COMMANDE MOUSTIQUAIRES**). By using the algorithm presented in Figure 12, we can extract the name of the product from the Purchase order.

The flow of this program starts by calculating the dimension of the input image. Then it crops the top-left section of the image (assuming the name of the product is in that section). Then it performs OCR with Pytesseract to extract texts from the image. If it is unsuccessful, this mostly means the image is too small, meaning it's likely that the image is not a Purchase Order. Else, we split the extracted text line by line, comparing them to find the pattern (seeing if the line starts with "**BON DE COMMANDE**" (Translation: PURCHASE ORDER ...)). If it does, we extract the texts after the matched words, which is the name of the product in the purchase order. Then we search the name of the product in our lists to match them with its corresponding manufacturer Figure 13Figure 12; then we can simply forward that message to the correct manufacturer.



Figure 12 - Matching Product's name/code to Company

Sometimes the rotated scan does crop out the header of the Purchase Order, which leads to missing product names as well. To counter this scenario, a few lines below the header are also OCRed to find the identification code as well. In Figure 11 the code is in this format "**BDC/2025/MOUST/GENERAL**". With this code extracted, we can map it to the name of the product Figure 12, as well as to the company/factory that produces those products as well.

To sum up, with the name of the product on attachments, we can map it to the producer, leading to direct email classification without using a trained AI model.

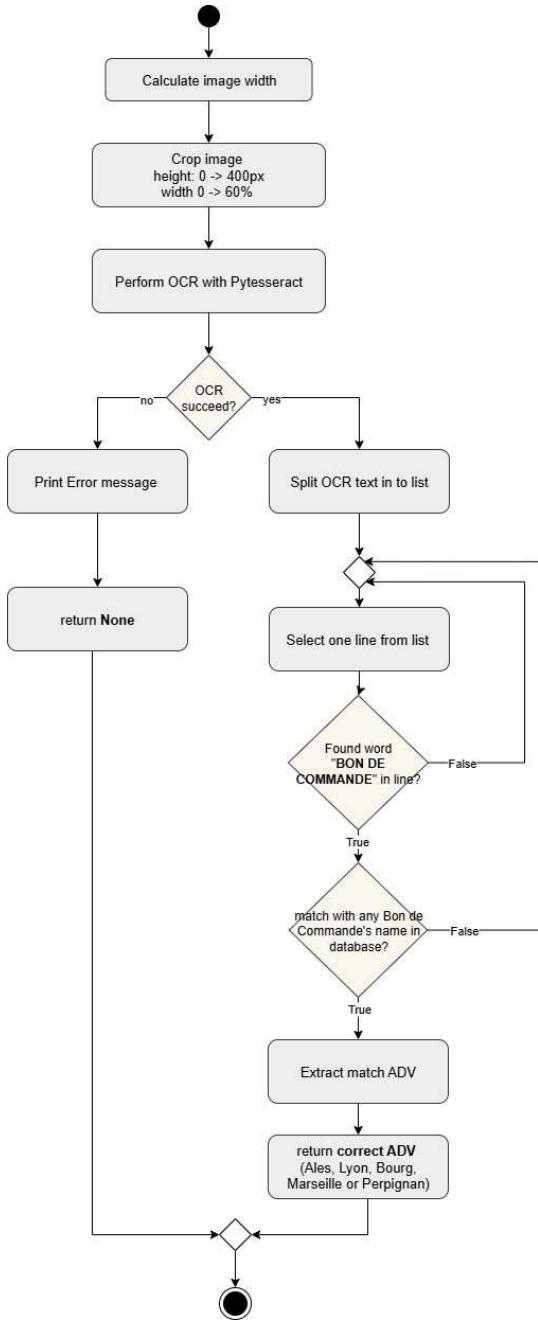


Figure 13 - Product name extraction from Purchase Order Flowchart

5.2.4. Intent extraction from Purchase Order

One more thing we can heuristically extract in the Purchase order is the intent of emails. You can see below the product name (Title of the document), there are two checkboxes, one for *Devis* (Quotation) and another one for *Commande* (Order). Clients may express their intent by selecting one of these checkboxes, and I took this into account by developing a function that uses Image Processing techniques to correct orientation and detect the marked checkbox.

Taking into account the nature of the scanned document as an input, the orientation, placement, scale, and the quality of the scanned *Purchase Order* are not always perfect. For example, as shown in Figure 15, the scanned document is at an angle. With these issues, it becomes almost impossible to detect the checkbox marks directly without applying preprocessing and cleaning steps first.

To solve this problem, a Feature Matching algorithm called Oriented FAST and Rotated BRIEF (ORB) was implemented. This algorithm was introduced in the paper **ORB: an efficient alternative to SIFT or SURF** by Rublee et al (ORB: An efficient alternative to SIFT or SURF, 2011). in 2011, proposing a new and more efficient way to extract interesting features from images compared to other previously introduced methods such as Scale-Invariant Feature Transform (SIFT) (Distinctive Image Features from Scale-Invariant Keypoints, 2004), and Speeded Up Robust Feature (SURF) (Machine learning for high-speed corner detection, 2006).

```

1 orb = cv2.ORB_create(8500)
2 keypoints_scanned, descriptors_scanned = orb.detectAndCompute(scanned_image, None
   )
3 keypoints_model, descriptors_model =orb.detectAndCompute(model_image, None)
4 # Use BFMatcher to match descriptors
5 bf = cv2.BFMatcher(cv2.NORM_HAMMING, crossCheck=True)
6 matches = bf.match(descriptors_scanned, descriptors_model)
7 # Sort matches by distance (best matches first)
8 matches = sorted(matches, key=lambda x: x.distance)
9 # Filter out low quality points
10 points_scanned = np.zeros((len(matches), 2), dtype=np.float32)
11 points_model = np.zeros((len(matches), 2), dtype=np.float32)
12 for i, match in enumerate(matches):
13     points_scanned[i, :] = keypoints_scanned[match.queryIdx].pt
14     points_model[i, :] = keypoints_model[match.trainIdx].pt
15 # Find the homography matrix
16 homography_matrix, mask =cv2.findHomography(points_scanned, points_model, cv2.
   RANSAC)
17 # Warp the scanned image to align with the reference model image
18 height, width = scanned_image.shape
19 corrected_scanned_image = cv2.warpPerspective(
20     scanned_image,
21     homography_matrix,
22     (width, height),
23     borderValue=(255,255, 255)
24 )

```

My method is to implement ORB feature matching between the scanned document (scanned Purchase Order), and a reference model image (a saved image of perfectly scanned document) matching those descriptors using Brute Force Hamming to correct the orientation and scale of the scanned Purchase Order.



Figure 14 - ORB Features matching

EXAMPLE: Figure 15 (left side) show a rotated scan of a Purchase Order. If we were to crop out the checkbox for Demande de devis (Price Quotation request) and the checkbox for Commande (Order) with pre-defined coordinates, we wouldn't get the checkbox itself.

Figure 14, you can find the matched descriptors between the scanned image and the model image (while we are using ORB detection). Figure 15presents the result of the scanned image after rotation and translation correction.

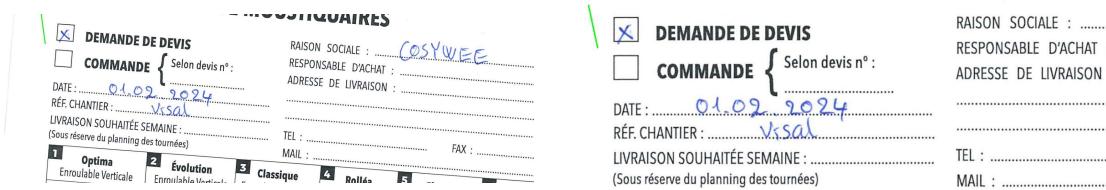


Figure 15 - Orientation corrected Purchase

After this, we can detect objects in each checkbox of the Purchase Order with the Contours detection feature from the OpenCV Library (OpenCV: Open Source Computer Vision Library, 2000).

5.3. Helper Functions with Microsoft Graph API

The main objective of this project is to build an automated system utilizing Artificial Intelligence to sort incoming emails to designated recipients. One main task of the program is to be able to read and move emails within the mailbox itself.

This subsection presents the methodology and steps involved in setting up the API to access company shared mailbox emails securely. This task requires a lot of research and careful attention due to strict security protocols and risk assessment aimed at preventing attacks in the company mailbox.

Email inboxes are among most sensitive components of a company's infrastructure. They usually contain confidential and secret client information, personal data (of employee and clients), internal conversations, all of which are highly susceptible to cyberattacks and data breaches.

To integrate this application on Microsoft Azure, first we had to choose an appropriate login method. Considering security protocols and importance, I opted for Access Token method which expires and is renewable every two hours using a refresh token and client secret (Client secret expires and requires manual reset every 18 months). This method offers a good balance between security and accessibility, ensuring a secure connection and program efficiency. The API permissions are **Mail.Send**, **Mail.ReadWrite** and **Mail.ReadWrite.Shared**. This allows the application to read and move emails in a shared mailbox.

5.3.1. Read Emails

In this subsection, I will present the flow to extract the metadata of the email such as the sender, subject, body, and any attachments. This function sends a request to `f"https://graph.microsoft.com/v1.0/users/shared_mailbox/mailFolders/folderID"` to read all emails within the mailbox.

```

1 def read_emails():
2     global access_token, headers
3     if inbox_folder_id is None or results_id is None:
4         print("Error while fetching folder IDs")
5         return
6     messages_endpoint = f"[MAILFOLDER]/{inbox_folder_id}/messages?$top=100"
7     start_time_program = datetime.datetime.now()
8     while True:
9         messages_response = requests.get(messages_endpoint, headers=headers)
10        if messages_response.status_code != 200:
11            print(f"Error retrieving messages: {messages_response.status_code}, {messages_response.json()}")
12
13        access_token = authentication.get_access_token() # Assume
14        authentication module exists
15        if not access_token:
16            print("Access token not available")
17            return
18        print("Acquired New Token")
19        headers = {AUTHORIZATION: f"Bearer {access_token}"}
20        messages_response = requests.get(messages_endpoint, headers=headers)
21        if messages_response.status_code != 200:
22            return
23
24        for index_message, message in enumerate(messages_response.json().get("value", [])):
25            message_id = message["id"]
26            message_body_html = message.get("body", {}).get("content", "")
27            sender = (
28                f"{message['sender']['emailAddress']['name']} <{message['sender']
29                    ['emailAddress']['address']}>")
30            subject = message.get("subject") or ""
31            # .... Other processing code ...

```

5.3.2. Get Folder ID

This function returns folder ID, based on the folder path. In Microsoft graph API, we cannot move email with the folder name explicitly. In order to move email, we have to pass three parameters:

- **Header** : the authentication token, containing access token to access the mailbox
- **MessageID**: the identification of the message to move
- **DestinationFolderID**: the identification of the folder we want to move message to

The `get_folder_id` function above allows us to get identification of the destination folder, with the display name and path of the folder itself.

```
1 def get_folder_id(headers, folder_name, parent_folder_id=None):
2     if parent_folder_id:
3         endpoint = f"{MAILFOLDER}/{parent_folder_id}/childFolders"
4     else:
5         endpoint = MAILFOLDER
6     response = requests.get(endpoint, headers=headers)
7     if response.status_code == 200:
8         folders = response.json().get('value', [])
9         for folder in folders:
10             if folder['displayName'] == folder_name:
11                 return folder['id']
12     print(f"{folder_name} folder not found.")
13 return None
```

5.3.3. Move Email

To move email from one folder to another folder, we must send a request to **graph.microsoft.com** with a **/move** request.

The function below shows how a message is moved from one folder to another. It takes header (which is basically `access_token`), ID of the message to move, and the destination folder ID as parameters.

```

1 def move_email(headers, message_id, destination_folder_id):
2     move_endpoint = f"https://graph.microsoft.com/v1.0/me/messages/{message_id}/
3         move"
4     payload = json.dumps({"destinationId": destination_folder_id})
5     move_headers = {
6         AUTHORIZATION: headers[AUTHORIZATION],
7         CONTENT_TYPE: APPLICATION_JSON
8     }
9     move_response = requests.post(move_endpoint, headers=move_headers, data=
10        payload)
11    if move_response.status_code == 201:
12        print(f"Moved message {message_id} to folder {destination_folder_id}.")
13        new_message_id = move_response.json().get('id')
14        mark_as_unread(headers, new_message_id)
15        return True
16    else:
17        print(f"Error moving message {message_id}: {move_response.status_code}, {
18            move_response.json()}")
19        return False

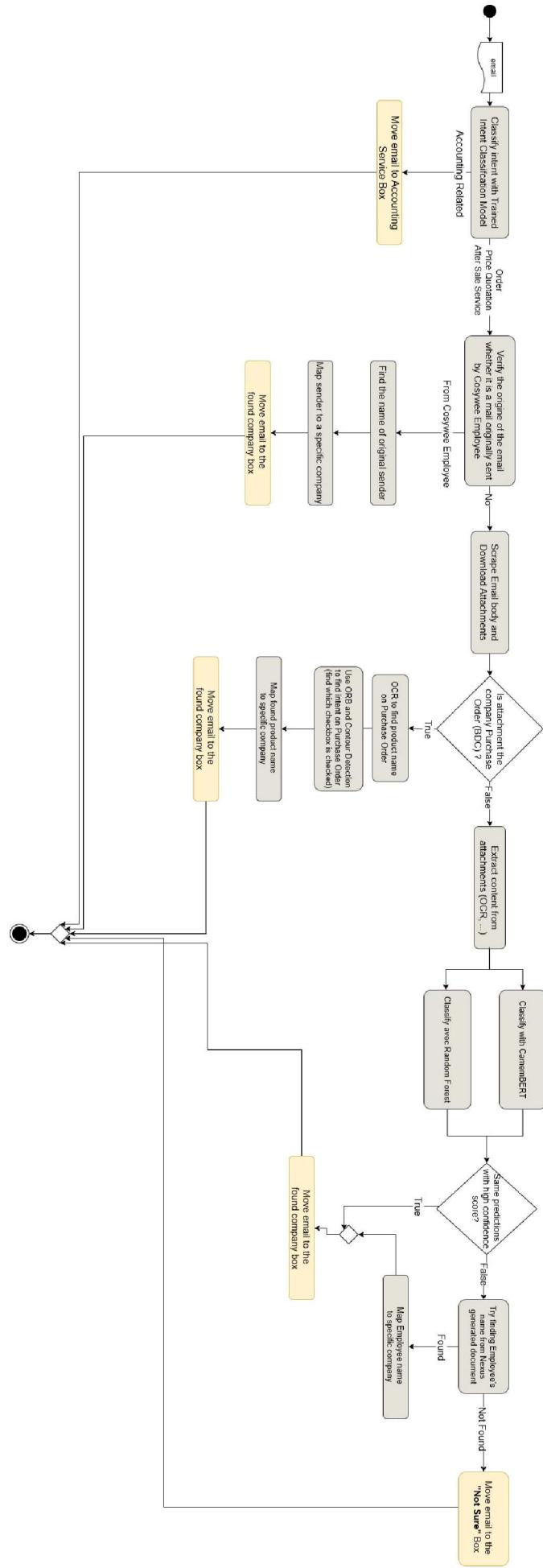
```

5.4. Final Classification Pipeline Implementation

The section presents the whole pipeline of the fully developed solution.

1. When a new email arrives in the inbox, our system will paste its subject and body text to our trained intent classification model.
2. If it predicts that the intent of email is related to **Accounting purpose** (Service Comptabilité), the email will be sorted directly to Accounting mailbox and the classification pipeline ends there. Other than that, it extracts other metadata like sender, recipients, attachments, etc before checking if the email is a reply to a message originally sent by a Cosywee employee. If so, the system will automatically forward the message to the original sender.
3. If the email is not a reply, the system will parse the HTML body text and combine it with information extracted from any attachments. This includes text from .docx, .csv, .txt, and .pdf files, as well as OCR-extracted text from images or scanned PDFs.
4. In this stage, we apply another heuristic, which is to determine whether the attachment is actually a Purchase Order or not. If it is a Purchase order, we simply extract product name, extract intent from the checkbox, before mapping the product name to a specific company that produces it. Once the company responsible for the product is found, the system forwards the email directly to the corresponding department or recipient responsible for that specific Purchase Order type.

5. In the event the attachment is not a Purchase Order, all the gathered content above will be cleaned and passed through our trained models, which will predict and vote on the most likely recipient. The email will then be forwarded to the predicted recipient.
6. In cases where the predictions between the models are nonconcurrent, we pass them through another heuristic method, which is finding employee name from Nexus generated document. Just like in step 2 if the employee's name is found, it is mapped to a specific company the employee works at, then we can simply move email to their mailbox. If none of conditions above are matched, we simply move that email to **A CLASSEUR** (Not sure) box for manual inspection.



5.5. Optimization and Improvement

As at Cosywee, projects are in **Agile** environment, we implemented the first prototype since early June to test our program. The goal is to develop the first program, receive feedback, improve, release a new version, receive more feedback, so on and so forth.

This environment really did help me in terms of agility, productivity and efficiency as I repair crucial bugs as I go instead of focusing on aspect that didn't matter as much.

5.5.1. Version v1.0

5.5.1.1 Pros

The first prototype of the AI system employs a custom-built OCR system and PaddleOCR (PaddleOCR: Optical Character Recognition tools based on PaddlePaddle, 2024) to extract texts from images and PDF files as we want to capture as much context as possible (PaddleOCR is used as a text detector, and the text images were then cropped and passed to a finetuned TrOCR model for text extraction). The goal was to OCR all texts, printed or handwritten, structured or non-structured data.

Figure below shows the capability of the system to extract non-structure and handwritten text from scanned document despite poor quality of document itself.

5.5.1.2 Cons

Although the program works well in all attachment conditions, it has its own perks as well.

As the AI system employs PaddleOCR and custom-trained OCR model (finetuned model from TrOCR), it takes a lot of computational power and time compared to other OCR engines like Pytesseract and EasyOCR.

The table below shows the differences between computational power, processing time and system capabilities between different setups.

	GPU Available?	RAM	CPU cores	Processing Time
With no attachment	Yes	16 GB	32	3 seconds per email
With attachment	Yes	16 GB	32	45 seconds per email

Table 9 - Processing Time vs Computing Power on a high computing power machine

	GPU Available?	RAM	CPU cores	Processing Time
With no attachment	No	8 GB	16	15 seconds per email
With attachment	No	8 GB	16	55 minutes per email

Table 10 - Processing Time vs Computing Power on a low computing power machine

As shown in Table 10, the processing time for emails with attachments on a low-resource computer takes too long (at least 55 minutes per email), which doesn't meet the company's criteria for low computing power and fast processing time.

After analyzing the attachment statistics, it was discovered that less than 6% of the attachments were handwritten. Considering the small proportion of handwritten attachments and the significant processing time they require, and after careful consideration, experimentation, and discussion with my project supervisor, we decided to release a new version of the AI email routing system with reduced OCR capability but with much more efficient processing time and lower computational cost.

5.5.1.3 Improvement aspects

This program performs quite well in terms of accuracy, but if we were to consider the processing time, it is not optimized yet. It requires a lot of processing power and takes too much time to process one email, sometimes even longer than manual classification. In this reason, this version is not deployed in production and was later improved in Version v2.0.

5.5.2. Version v2.0

This version is an improvement from last version (v1.0) as it drastically improves processing time, reduces computational costs and is more suitable for a server-based sub-process that runs on the background.

5.5.2.1 Pros

In this newer version, Pytesseract OCR engine is used to replace the combination of PaddleOCR and custom finetuned TrOCR to improve processing time as well as to reduce computational cost.

Table below shows the statistical comparison between our AI system with (PaddleOCR + TrOCR) and the AI system that employs the Pytesseract engine on a small and less powerful machine.

GPU required?	RAM	CPU cores	Processing Time
----------------------	------------	------------------	------------------------

With no attachment	No	4 GB	4	2 seconds per email
With attachment	No	4 GB	4	15 seconds per email

Table 11 - Processing Time vs Computing Power of version v2.0 on a low computing power machine

5.5.2.2 Cons

The newer version of the AI system performs significantly better in terms of computing power and processing time compared to the older version. However, there is one downside: it cannot extract handwritten text from attachments, as Pytesseract is mainly built for printed and uniform text.

Although this limitation exists, if we were to consider the differences in processing time and computational efficiency, it becomes clear that the newer version is better suited for production deployment compared to version v1.0.

VI - RESULTS AND DISCUSSIONS

For the last few months, I have been making my way up the ladder to develop an AI system that is capable of classifying and routing emails automatically to designated recipients from a shared mailbox.

I am pleased to say I developed a reliable and production-ready system for email classification, taking into account the body content, subject, and attachments' content.

In Fig 6.A, the process diagram of my work during this project is presented. In the first week, I was presented with the project, then I performed research and literature reviews while starting to collect and clean data to train the models. After all data were collected, multiple experiments were performed, and a few final models were trained to be used in production. I set up the API, project routes, project structures, and mail pipeline, then did some research on available heuristics to improve the accuracy of the AI system before implementing them and deploying it in production. Since then, I have been receiving feedback and reports from people in each factory and released multiple newer versions until a reliable and sustainable system is in place. At the end, I standardized my code (commenting, renaming variables to improve coherence, etc). In the last few weeks, I created a script to finetune models and clean data for the company to make changes when new problems/errors occur in the future. I also gave a presentation to my team's developers to share knowledge on topics such as model fine-tuning, data cleaning, and evaluating trained models.

Compared to rule-based email classification features proposed in Microsoft Outlook and other email provider services, this method can analyze the context, intent, and product mentioned in emails, making the classification more reliable and sentiment-based instead of keyword matching. The essence of this method is to implement as many heuristics as possible to ensure the reliability of the program before passing it to our trained classification model.

6.1. Results

Moreover, this project is run on a lower resource and computing power computer:

- **CPU:** 8 Cores, 12 threads
- **RAM:** 4GB
- No Graphics Processing Unit (GPU)

Other than this project, my role at the company also involves other tasks as well, such as:

- Setting up, organizing and lead weekly meetings every other week to talk about strength, weaknesses and interesting things every IT team member has faced in the prior week, as well as to evoke problems that everyone faces so that the team can help with their own experience.
- Reviewing, testing code, documentation written by other teammates
- Documenting processes and ensuring project scalability.

The final implementation is described as follows:

- Models Used: A hard voting ensemble composed of 2 Random Forest classifiers and 2 CamemBERT-based transformer models.
- Performance Optimization: In alignment with company requirements, handwriting recognition was deprioritized in favor of faster processing. As a result, the system uses Tesseract OCR (PyTesseract), which improves processing speed by over 80% compared to handwriting-capable OCR models. This approach also significantly reduces resource requirements, operating efficiently without GPU acceleration and with less than 8GB of RAM.
- Processing Time: The system processes emails without attachments in under 2 seconds. For emails with attachments requiring OCR and text extraction, the average processing time is approximately 15 seconds per attachment.
- Error Rate: The overall error rate remains below 2%. Approximately 12% of attachments require manual review due to factors such as handwritten content or insufficient information for confident classification.

Table 6.1 indicates how many emails (in percentage) are sorted with each heuristic method and with our trained model. It also shows the amount of email that the AI system failed to classify due to lack of information and/or other problems that will be discussed below.

	Average Processing Time	Error rate	Email to classify manually*
With 500 emails	15 s / email	< 2%	< 12%

Table 12 - AI System capability (With 500 emails sorted)

***Note:** Emails that require manual classification don't necessarily indicate a failure of the AI pipeline. Rather, it's usually due to the nature of those emails being inherently difficult to classify automatically. These cases include:

- Emails with handwritten and unstructured texts
- Emails that don't contain enough information (e.g., only a reference number without mentioning a product or service)
- Spam emails
- Emails describing very generic information (like dimensions or price) that could apply to multiple products from different companies
- Emails with unusual purposes (e.g., registering new clients, announcements about vacation days from partner companies, etc.)

6.2. Future improvements

While current system archives a reliable and satisfactory performance in classify and routing emails, several improvements can be made to enhance accuracy, and robustness of the program.

6.2.1. Optimizing handwriting text recognition model

Currently, **Cosywee** has its own finetuned OCR model for handwriting text recognition. The base model of this OCR model is TrOCR, which is only capable of text recognition but lacks the strength of text detection. Moreover, it can only process text line by line, word by word and it takes a lot of time to predict/recognize each input. As the document has a lot of words/phrases, it would take even longer to process them, that is why in this system, we prioritize the speed of the program over accuracy of handwritten text. If I were to have more time, my priority would be to improve processing time of Cosywee's custom OCR model, which leads to reliable and automated handwriting-based attachment classification as well.

6.2.2. Integration of human feedback loop

Incorporating a real time feedback loop from human validation can help document any error the program has made along the way, which would be a great help for finetuning, and this problem can also be framed as a Reinforcement Learning project as well.

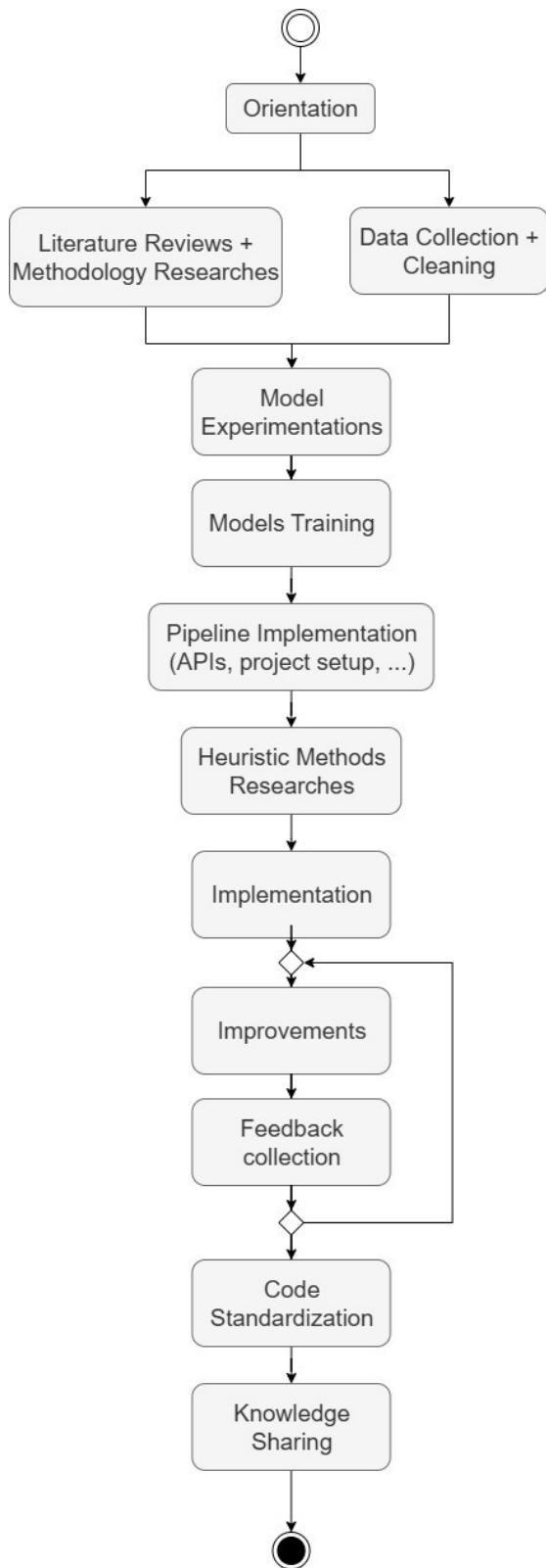


Figure 16 - Works completed flowchart

VII - CONCLUSION

To summarize, this report presented an end-to-end process of building an automatic email routing system from requirement scoping to data collect, to model experimentation, up to finally, deployment. By developing a custom AI models pipeline, the system is now fully capable of reading and understanding email content, analyzing its intent then forwarding them to appropriate recipients. This significantly reduces manual workload, eliminates repetitive tasks, and most importantly, boosts workplace efficiency.

Despite certain limitations, such as occasional misclassification among some emails, the system itself demonstrates robust and reliable performance across different types of emails.

Overall, this project marks a first step in Cosywee's milestone as it is the first homemade (in-house built) AI system that is implemented in the company. With the improvements in the previous section, this system can evolve and be improved into an even more dependable and intelligent solution in the future.

References

- Al., D. e. (2018). *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding.*
- Aspose.Email for Python via .NET. (2024). Retrieved from
<https://products.aspose.com/email/python-net/>
- Bay, H., Tuytelaars, T., & Gool, L. V. (2006). SURF: Speeded Up Robust Features.
- Beautiful Soup: HTML and XML parser for Python. (2024). Retrieved from
<https://www.crummy.com/software/BeautifulSoup/>
- CamemBERT NER - HuggingFace model. (2021). Retrieved from
<https://huggingface.co/Jean-Baptiste/camembert-ner>
- Chen, T., & Guestrin, C. (2016). XGBoost: A Scalable Tree Boosting System. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 785-794.
- Distinctive Image Features from Scale-Invariant Keypoints. (2004). *International Journal of Computer Vision*, 60, 91-110.
- Experimentation in Software Engineering - An Introduction.* (2012). Kluwer Academic Publishers.
- extract_msg: Python library to extract emails from MSG files. (2024). Retrieved from
<https://github.com/mattgwwalker/msg-extractor>
- JaidedAI. (2024). *EasyOCR*. Retrieved from <https://github.com/JaidedAI/EasyOCR>
- Li, M., Lv, T., Chen, J., Cui, L., Lu, Y., Dinei, F., . . . Wei, F. (2021). TrOCR: Transformer-based OCR with pre-trained Vision and Language Models. Retrieved from
<https://arxiv.org/abs/2109.10282>
- libpff/pypff: Library and Python bindings to access Personal Folder File (PST, OST). (2024). Retrieved from <https://github.com/libyal/libpff>
- Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., . . . Stoyanov, V. (2019). *RoBERTa: A Robustly Optimized BERT Pretraining Approach.*
- Machine learning for high-speed corner detection. (2006)., (pp. 430-443).
- Martin, L., Muller, B., Suàrez, P. J., Dupont, Y., Laurent, R., Clergerie, É. V., . . . Sagot, B. (2020). CamemBERT: a Tasty French Language Model. *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics.*
- OpenCV: Open Source Computer Vision Library. (2000). *Dr. Dobb's Journal of Software Tools.*

Optuna: A Next-generation Hyperparameter Optimization Framework. (2019)., (pp. 2623-2631).

ORB: An efficient alternative to SIFT or SURF. (2011)., (pp. 2564-2571).

PaddleOCR: Optical Character Recognition tools based on PaddlePaddle. (2024). Retrieved from <https://github.com/PaddlePaddle/PaddleOCR>

pandas-dev/pandas: Pandas. (2020). *Zenodo*. doi:10.5281/zenodo.3509134

PyMuPDF (fitz): A Python binding for MuPDF PDF and document rendering library. (2024). Retrieved from <https://github.com/pymupdf/PyMuPDF>

pytesseract: Python wrapper for Google Tesseract-OCR. (2024). Retrieved from <https://github.com/madmaze/pytesseract>

PyTorch: An Imperative Style, High-Performance Deep Learning Library. (2019). *Advances in Neural Information Processing Systems (NeurIPS)*, 32, 8026-8037.

Scikit-learn: Machine Learning in Python. (2011). *Journal of Machine Learning Research*, 12, 2825-2830.

spaCy French Language Model: fr_core_news_lg. (2024). Retrieved from https://spacy.io/models/fr#fr_core_news_lg

Term-weighting approaches in automatic text retrieval. (1988). *Information Processing & Management*, 24, 513-523.

APPENDIX

Types of Emails

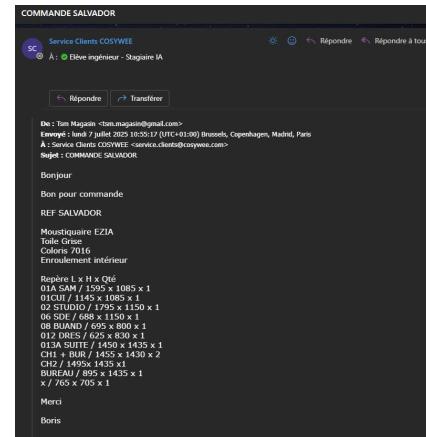


Figure 17 - Emails with no attachments, containing only plain text order or price quotation requests

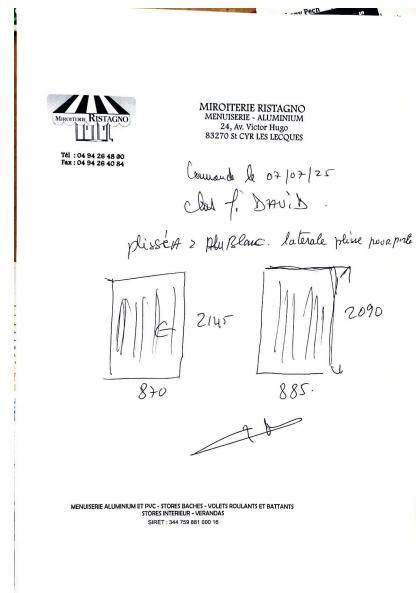


Figure 18 - Emails with handwritten attachments indicating order details



Figure 19 - Emails with computer-generated attachments (printed text) indicating order details

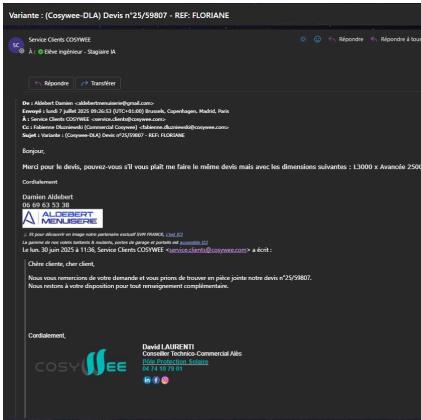


Figure 20 - Emails replying to a Cosywee employee's email

Figure 21 - Emails inquiring about order progress (with only a reference number and no product information)

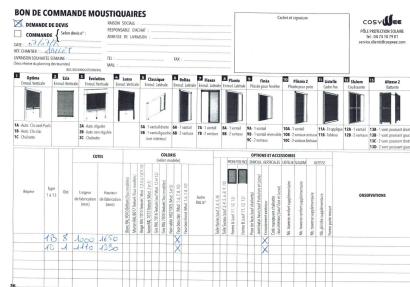
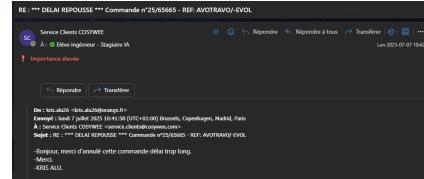


Figure 22 - Order emails with Cosywee's Purchase Order form (Bon de Commande) attached

Figure 23 - Emails containing Cosywee-generated documents as references



DANS L'ATTENTE DE VOTRE ACOMpte DE 50% SOIT 157.27 EUROS POUR POUVOIR lANCER LA COMMANDE***
Mode de règlement : ACpte 50% S/CDR PAR VRH + SOLDE LCR 30 J DF
*** ATTENTION ! Toute demande de modification/annulation doit intervenir
sois 24h après réception de l'ARS de commande. ***

Total HT	242,04 €
Taxe de transport client	20,04 €
TVA 20% France métrop. hors Corse	48,48 €
Total TTC	214,56 €

GIE COSYWEE, QUAI DES BARDINES, 30000 LES MAGES
service.client@cosywee.com - www.cosywee.com
1000 € - RCS Nîmes - Siret: 622 39 452 001 0 42542 - TVA FR220994463

