



Ecole Polytechnique de l'Université de Tours
Département Informatique
64 avenue Jean Portalis
37200 Tours, France
Tél. +33 (0)2 47 36 14 14
polytech.univ-tours.fr

Research & Development Project 2022-2023

Deep-Agora

Incremental segmentation of images of old documents



Company
Centre d'études supérieures de la Renaissance



Industrial supervisor
Rémi JIMENES

Student
Théo BOISSEAU (DI5)

Academic supervisor
Jean-Yves RAMEL

30th November 2022

List of contributors

Company

Centre d'études supérieures de la Renaissance
59, rue Néricault Destouches
37020 Tours, France
cesr.univ-tours.fr



Name	Email	Quality
Théo BOISSEAU	theo.boisseau@etu.univ-tours.fr	Student DI5
Jean-Yves RAMEL	jean-yves.ramel@univ-tours.fr	Academic supervisor, Computer Science Department
Rémi JIMENES	remi.jimenes@univ-tours.fr	Industrial supervisor



Warning

This document was written by Théo BOISSEAU hereinafter referred to as the author.

The company Centre d'études supérieures de la Renaissance is represented by Rémi Jimenes hereinafter referred to as the industrial supervisor.

The Polytechnic School of the University of Tours is represented by Jean-Yves RAMEL hereinafter referred to as the academic supervisor.

By the use of this document model, all parties of the project accept the conditions defined hereunder.

The author acknowledges taking full responsibility for the content of the document, as well as any judicial consequences that might ensue as a result of non-compliance with laws or copyright.

The author certifies that the statements of the document are truthful and takes full responsibility for the veracity of the statements.

The author certifies not to appropriate the work of others and that the document contains no plagiarism.

The author certifies that the document contains no statements that are defamatory or reprehensible under the law.

The author acknowledges that they may not publish this document in part or in whole in any form without prior approval of the academic supervisor and the industrial supervisor.

The author authorizes the Polytechnic School of the University of Tours to publish this document in whole or in part in any form, including after modification with due acknowledgement of the source. These publications must be free of charge and come along with this warning.



To cite this document

Théo BOISSEAU, *Deep-Agora: Incremental segmentation of images of old documents*,
Research & Development Project, Ecole Polytechnique de l'Université de Tours, Tours,
France, 2022-2023.

```
@mastersthesis{
  author={BOISSEAU, Théo},
  title={Deep-Agora: Incremental segmentation of images of old documents},
  type={Research & Development Project},
  school={Ecole Polytechnique de l'Université de Tours},
  address={Tours, France},
  year={2022-2023}
}
```

Contents

List of contributors	a
Warning	b
To cite this document	c
Contents	i
List of Figures	iv
1 Introduction	1
1 Actors, issues and context	1
2 Objectives.....	2
3 Hypotheses	2
4 Methodological bases	2
2 General description	4
1 Project environment	4
2 User characteristics.....	4
3 System features	4
4 General structure of the system	5
3 State of the art / Technology watch	6
1 Section 1	6
Paragraph 1	6
Paragraph 2	6
2 Section 2	6

4	Analysis and design	7
1	Analysis.....	7
1.1	Assumptions used	7
1.2	Specifications.....	7
2	Proposed modelling	7
5	Implementation	8
1	Tools and library used	8
2	Implementation elements, technical choices	8
3	Analysis of results, evaluation, quality.....	9
4	Main HMIs	9
4.1	HMI 1	9
6	Assessment and conclusion	10
1	Semester 9 review	10
2	Review of semester 10.....	10
3	Quality assessment.....	10
4	Self-critical review	10
	Appendices	11
A	Planning, project management	12
1	Evolution of the project	12
2	Job description	12
	Task 1: Speaking Heading	13
	Task 2: Speaking Heading	13
B	Description of the interfaces	14
1	Hardware/software interfaces	14
2	Human/machine interfaces.....	14
C	Specification booklet	15
1	Functional specifications	15
1.1	Features to be developed	15
1.2	Definition of function 1: speaking title	15
	Presentation of function 1 :	15
1.3	Definition of function 2: speaking title	15
	Presentation of function 2:.....	15
	Description de la fonction 2 :	15
2	Non-functional specifications.....	16

2.1	Development constraints and design.....	16
2.2	Functional and operational constraints.....	16
2.2.1	Performance	16
2.2.2	Capabilities	16
2.2.3	Controllability	16
2.2.4	Security	16
D	Developer's Workbook	17
1	Introduction	17
2	Architectural diagrams and UML.....	17
3	Detailed descriptions of data used.....	17
4	Detailed descriptions of classes, modules, achievements	17
E	Installation document	18
F	User document	19
G	Test booklet	20
1	Unit testing	20
2	Integration testing	20



List of Figures

2	General description	
2.1	Use cases diagram.....	5
2.2	Component diagram	5
3	State of the art / Technology watch	
3.1	Fouille de données et visualisation.....	6
A	Planning, project management	
A.1	Le diagramme de Gantt Final	12
A.2	Le diagramme de Gantt Final	12

1

Introduction

1 Actors, issues and context

The **Research & Development Project (R&D project)** is the final work that the student engineer must complete to obtain his diploma. It places the future engineer in a project situation by making him/her produce personal work and invites him/her to show initiative and maturity regarding a specific high-level problem. The R&D project, which lasts at least two days a week throughout the fifth year, i.e. 26 weeks, is the subject, each semester, of a dissertation and an oral presentation to a jury.

This report aims to provide both the main document that everyone can read and all the technical and methodological elements. It consists mainly of complete sections of the different documents produced, with the technical sections in the appendix.

The actors of this project are:

- the client, which here are **Centre for Advanced Renaissance Studies (fr. Centre d'études supérieures de la Renaissance) (CESR)**, for which a contact is Rémi Jimenes, lecturer and researcher.
- the **Project/Product Owner (fr. Maître d'ouvrage) (fr. MOA)**, who is Jean-Yves Ramel, professor of computer science, director of **Laboratory of Fundamental and Applied Computer Science of Tours (fr. Laboratoire d'Informatique Fondamentale et Appliquée de Tours) (LIFAT)** and academic tutor for this project. He is responsible for representing the client by ensuring that the deadlines are met and that the product conforms.
- the **Project Manager / Scrum Master (fr. Maître d'œuvre) (fr. MOE)**, who is me, Théo Boisseau, an engineering student in his final year of study. I decide on the technical means used to design the product by what was defined by the product owner.

To convert these historical books into accessible digital libraries, LIFAT is developing image processing software that participates in a complete processing chain, including layout analysis, text/illustration separation (i.e. segmentation of content elements), optical character recognition (i.e. OCR) and text transcription. This project focuses on layout analysis and segmentation of content elements of historical documents.

The client expressed the need for easy-to-use interactive software so that its users, historians, could create their own scenarios for extracting **elements of content (EOCs)** from images of historical documents. These historical documents are mainly Renaissance corpora, accessible

from the CESR database, and contain mainly printed or manuscript text, illustrations and page ornaments.

The simplicity of creating extraction scenarios, their reuse and their adaptation to different documents are essential dimensions of the requirement.

However, this simplicity should not unduly compromise the reliability and performance of the software. Image processing of historical documents is a particularly difficult task notably because of broken characters, stains, and poor paper quality.

In recent years, the performance of some deep learning techniques has surpassed that of shallow methods established by experts on various image processing tasks. As this progress has made many computer vision tools available, it now seems possible to meet this need with a completely new approach.

2 Objectives

This project aims to propose a new approach based on deep learning neural networks to solve this image processing problem.

To this end, the Deep-Agora R&D project aims to build a prototype of an optimisation software capable of extracting textual and decorative elements of content from images of historical documents.

The user should not be responsible for training the models. Therefore, several deep learning models can be created and trained to extract the content elements required in the different use cases of the software.

Due to its nature as a prototype, the system will need to be composed of computational documents combining scripts and good documentation. It must also provide access to training datasets and parameter storage files to reproduce the deep learning models created.

If the objective is achieved, the project can be continued and a scenario creation subsystem can be implemented to deploy the models created within it.

3 Hypotheses

Blablabla.

4 Methodological bases

An Agile project management method will be used to create learning loops to quickly gather and integrate feedback. Therefore, the Scrum method should be preferred in which ideology is to:

- learn from experience
- to self-organise and prioritise
- to reflect on gains and losses to continuously improve

Therefore, contact with the product owner should be maintained as much as possible, as it will help me to improve and learn considerably as the project progresses.

To this end, we set sprints with a fixed duration of 2 weeks. At least one deliverable, containing an e-mail, should be sent to the product owner at least every two weeks and preferably once a week. During the implementation phase, a meeting to get feedback about the product should be scheduled at the end of each sprint.

I use GitHub for configuration management, by creating two different repositories:

- Deep-Agora, which contains the source code of the project
- Deep-Agora_DOC, which contains all the deliverables of the specification, analysis and modelling part of the first semester

As a project management tool, I will also use GitHub. As of this year, it offers a similar feature to Trello called Projects, an adaptable spreadsheet that can also integrate with my issues and pull requests on GitHub to help me plan and track my work efficiently.

? programming rules, references to supporting documents such as the quality assurance and/or test plan?

2

General description

1 Project environment

This project is part of a larger research project between CESR and LIFAT. It is currently being carried out as part of a programme for the regional valorisation of old books (mainly dating from the Renaissance), namely the *Humanist Virtual Libraries* controlled by the CESR.

CESR does not have powerful computing machines capable of training deep neural networks, but it has several machines and a large amount of remote and on-premises storage.

Agora, the software developed and published ten years ago by LIFAT to process images of historical documents, is undergoing a complete overhaul in this project. Its technologies need to be updated and, above all, its overhaul should meet the previously unattainable need for simplicity in scenario creation.

Therefore, no takeover of the existing system is planned, as it has to be completely redesigned.

The developer will train deep neural networks, whose task is not intended for the end users of Deep-Agora. This part of the project is to be carried out outside the software system, but within the environment, as an engineer's system.

2 User characteristics

End users of Deep-Agora are all historians of CESR.

They have a sufficient but moderate command of computer tools. They often use them but need extensive training or solid documentation to use them in the case of advanced tools with complex functions. They did not have a satisfactory experience with Agora, as its interface was too complex. They do not need user access rights to use Agora.

3 System features

blablabla....

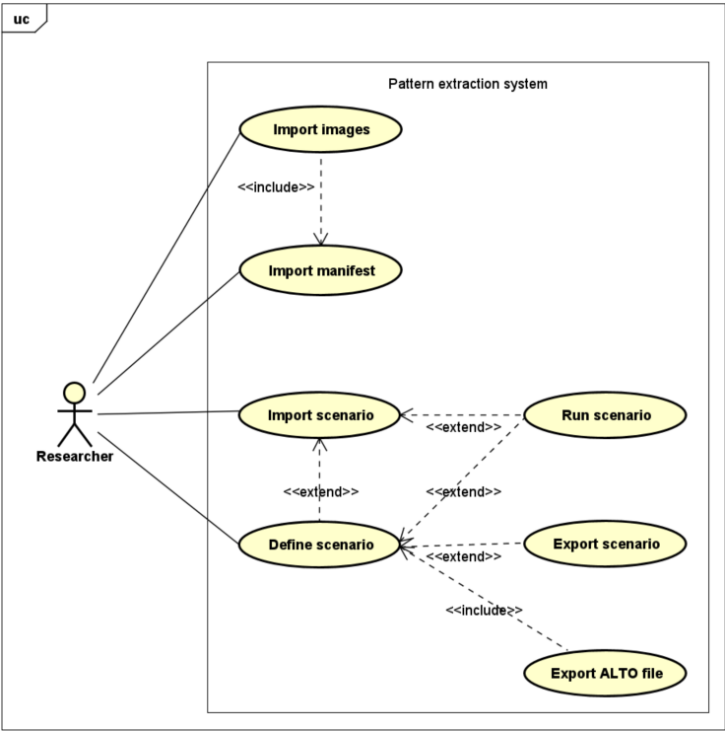


Figure 2.1: Use cases diagram

4 General structure of the system

blablabla....

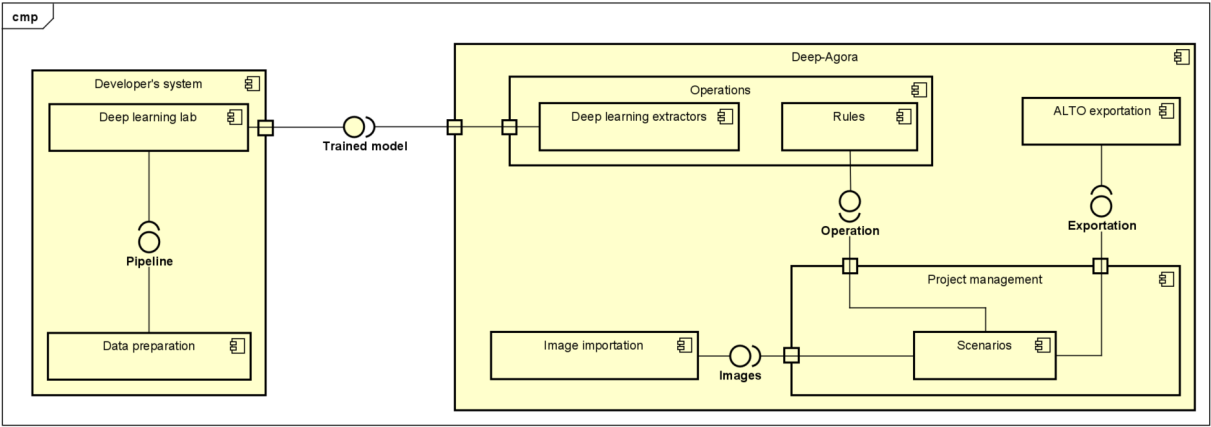


Figure 2.2: Component diagram

blablabla....

3

State of the art / Technology watch

Sujet a définir en concertation avec votre encadrant Polytech

1 Section 1

PENSEZ à bien insérer TOUTES les références bibliographiques utilisées dans votre bibliographie. Voici un exemple de citation: [DBLP:journals/corr/abs-1804-02767]. Et une autre: [DBLP:journals/corr/RedmonDGF15].

Paragraphe 1

blablabla

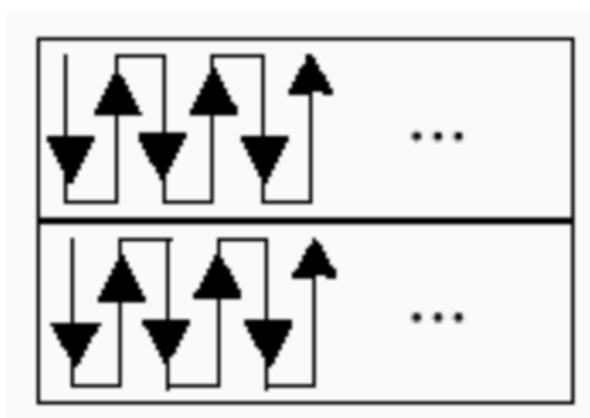


Figure 3.1: Fouille de données et visualisation

Paragraphe 2

blablabla

2 Section 2

blablabla

4

Analysis and design

1 Analysis

1.1 Assumptions used

blablabla

1.2 Specifications

Inclure ici un résumé du cahier de spécification qui sera inséré en ANNEXE

2 Proposed modelling

Inclure ici une description du système à développer pouvant notamment inclure les principaux diagrammes UML non détaillés et démontrant sa faisabilité durant la phase de mise en œuvre.

Les modes de validation prévus pour les différents éléments à produire pourront être précisés ici.

5

Implementation

Description de vos productions et de leurs modes de réalisation.

(résumé du cahier de développement inséré en ANNEXE)

blablabla

1 Tools and library used

blablabla

2 Implementation elements, technical choices

```
1 #include <iostream>
2 using namespace std;
3
4 int main () {
5     cout << "Hello , world !";
6     return 0;
7 }
```

Un exemple de PHP:

```
1 class pdfOrder extends FPDF
2 {
3     function _check($x,$y,$width,$checked) {
4         if ($checked)
5             $this->rect($x,$y,$width,$width,'F');
6         else
7             $this->rect($x,$y,$width,$width);
8     }
9     function LI($sansFrais = false) {
10         $LI = 'LI';
11         $coord = 'Laboratoire informatique'
```



```

12 64, avenue Jean Portalis
13 37200 Tours
14 Tél. : 02 47 36 14 42
15 Fax. : 02 47 36 14 22';
16 $this->Image(dirname(__FILE__) . '/li.jpg',10,2,20);
17 $this->SetFont('Times','B',20);
18 $this->SetFont('Times','',9);
19 $this->setXY(35,3);
20 $this->Multicell(80,4,utf8_decode($coord),0,'LT');
21 }

```

3 Analysis of results, evaluation, quality

blablabla

4 Main HMIs

4.1 HMI 1

Résumé des principaux éléments présent dans le Guide de l'utilisateur avec d'éventuels compléments d'information sur leur mode de mise en œuvre.

6

Assessment and conclusion

1 Semester 9 review

Liste des taches faites, en cours, à faire CF Planning S9 et S10 à fournir en annexe

2 Review of semester 10

Bilan global \Rightarrow respect du cahier des charges (fait / à faire)

3 Quality assessment

blablabla

4 Self-critical review

blablabla

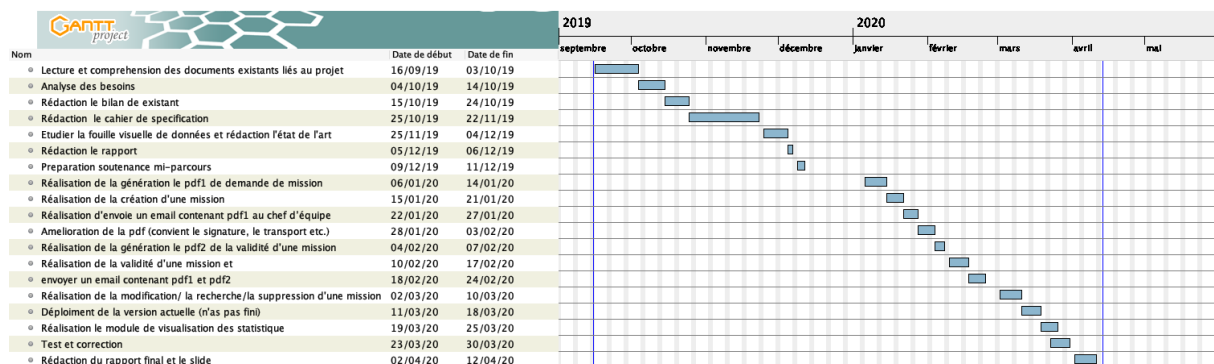
Appendices

A

Planning, project management

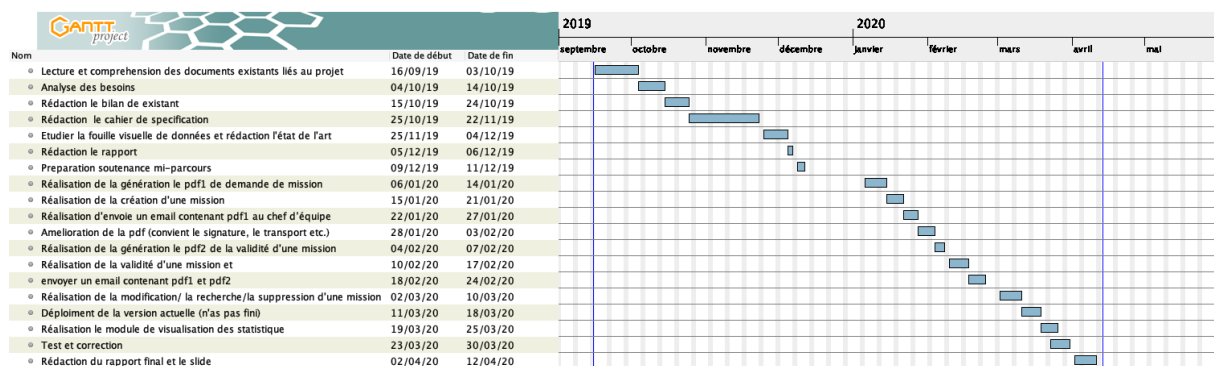
1 Evolution of the project

Le diagramme de Gantt Initial pour la planification de ce projet



?figurename? A.1: Le diagramme de Gantt Final

Le diagramme de Gantt Final de ce projet est comme Figure A.2.



?figurename? A.2: Le diagramme de Gantt Final

2 Job description

Task 1: Speaking Heading

- Date de début: 16/09/2019
- Date de fin: 03/10/2019
- Durée: 17 jours
- Description: éléments à faire, liste des entrées (pré-requis) et sorties (livrables)s de la tache.

Task 2: Speaking Heading

- Date de début: 16/09/2019
- Date de fin: 03/10/2019
- Durée: 17 jours
- Description: éléments à faire, liste des entrées (pré-requis) et sorties (livrables)s de la tache.

B

Description of the interfaces

1 Hardware/software interfaces

blablabla

2 Human/machine interfaces

blablabla



Specification booklet

1 Functional specifications

1.1 Features to be developed

1.2 Definition of function 1: speaking title

Presentation of function 1 :

Élément 1

Entrée : ????

Sortie : ????

Préconditions : ????

Postconditions : ????

Élément 2

Entrée : ????

Sortie : ????

Préconditions : ????

Postconditions : ????

1.3 Definition of function 2: speaking title

Presentation of function 2:

- Nom de la fonction : Visualisation des statistiques
- blablabla
- Primordiale

Description de la fonction 2 :

blablabla

2 Non-functional specifications

2.1 Development constraints and design

blablabla

2.2 Functional and operational constraints

2.2.1 Performance

blablabla

2.2.2 Capabilities

blablabla

2.2.3 Controllability

blablabla

2.2.4 Security

blablabla

D

Developer's Workbook

1 Introduction

blablabla

2 Architectural diagrams and UML

blablabla

3 Detailed descriptions of data used

blablabla

4 Detailed descriptions of classes, modules, achievements

blablabla

E

Installation document

Ce document regroupe toutes les informations nécessaires pour l'installation du projet sur les machines, ainsi que pour sa mise en production.



F

User document

blablabla



Test booklet

Les tests visent à garantir l'exactitude, l'intégrité, la sécurité et les performances du logiciel.

1 Unit testing

blablabla

IDENTIFICATION OF COMPONENT
Afficher toutes les missions de l'utilisateur identifié
DESCRIPTION OF THE TEST (granularity, scenario, values, actions)
Action :blablabla. blablabla.
EXPECTED RESULTS
Cas 1: blablablaa. Cas 2: blablabla.
OBTAINED RESULTS
blablabla

2 Integration testing

blablabla

Deep-Agora : Incremental segmentation of images of old documents

Théo BOISSEAU

Supervisor : Jean-Yves RAMEL



In collaboration with Centre d'études
supérieures de la Renaissance

Objectifs

- point 1
- point 2
- point 3



LABORATOIRE D'INFORMATIQUE FONDAMENTALE ET APPLIQUÉE DE TOURS

Mise en œuvre

1. point 1
2. point 2
3. point 3



LABORATOIRE D'INFORMATIQUE FONDAMENTALE ET APPLIQUÉE DE TOURS

Résultats attendus

Voici du texte. Voici du texte. Voici du texte.
Voici du texte. Voici du texte. Voici du texte.



LABORATOIRE D'INFORMATIQUE FONDAMENTALE ET APPLIQUÉE DE TOURS



Deep-Agora : Incremental segmentation of images of old documents

Théo BOISSEAU

Supervisor : Jean-Yves RAMEL

Objectifs

- point 1
- point 2
- point 3

Mise en œuvre

1. point 1
2. point 2
3. point 3



In collaboration with Centre d'études
supérieures de la Renaissance

Résultats attendus

Voici du texte. Voici du texte. Voici du
texte. Voici du texte. Voici du texte. Voici
du texte.



LABORATOIRE D'INFORMATIQUE FONDAMENTALE ET APPLIQUÉE DE TOURS

LABORATOIRE D'INFORMATIQUE FONDAMENTALE ET APPLIQUÉE DE TOURS

LABORATOIRE D'INFORMATIQUE FONDAMENTALE ET APPLIQUÉE DE TOURS

Ecole Polytechnique de l'Université de Tours
Département Informatique
64 avenue Jean Portalis, 37200 Tours, France
polytech.univ-tours.fr



POLYTECH
TOURS
Computer Science

Deep-Agora

Incremental segmentation of images of old documents

Résumé

Voici le résumé de ce PRD. Voici le résumé de ce PRD. Voici le résumé de ce PRD. Voici le résumé de ce PRD. Voici le résumé de ce PRD. Voici le résumé de ce PRD. Voici le résumé de ce PRD. Voici le résumé de ce PRD.

Mots-clés

motcle1, motcle2, etc.

Abstract

Here is the abstract of this project. Here is the abstract of this project. Here is the abstract of this project. Here is the abstract of this project. Here is the abstract of this project. Here is the abstract of this project. Here is the abstract of this project. Here is the abstract of this project.

Keywords

word1, word2, etc.

Company

Centre d'études supérieures de la Renaissance



Industrial supervisor

Rémi JIMENES

Student

Théo BOISSEAU (DI5)

Academic supervisor

Jean-Yves RAMEL