

ENTITE DEPOSITAIRE

Nom de la structure dépositaire : ARDIAN

Type d'entité dépositaire : ☒ Entreprise ☐ Startup ☐ Collectivité ☐ Association

Nom et prénom du dépositaire : PATRICK PFUNDSTEIN

Fonction : BRM

E-mail : Patrick.Pfundstein@ardian.com

Téléphone :

Canal de communication à privilégier : ☒ e-mail ☐ téléphone ☒ visioconférence

☒ J'accepte d'être présent.e lors de la première journée, le **mardi 14 mai 2024 à l'Université de technologie de Troyes (10)** pour une séance de travail avec les étudiants et répondre à leurs questions.

Votre présence durant les 2 journées suivantes, mercredi 15 mai et jeudi 16 mai, seront appréciées selon vos impératifs mais ne sont pas indispensables.

☒ Je reconnais avoir lu et approuvé le règlement et les conditions tarifaires de l'UTT Innovation Crunch Time 2024, accessibles sur le site internet [ici](#)

☐ J'autorise l'Université de technologie de Troyes – UTT et ses associations étudiantes à utiliser le logo de la structure dépositaire ainsi que l'image de ses représentants telle que reproduite dans les séquences audiovisuelles (prises de vue photographique, interview filmée et audio, citations écrites : commentaires et noms/prénoms), réalisées par l'équipe/journalistes dans le cadre de l'UTT Innovation Crunch Time 2024 à l'Université de technologie de Troyes (10).

☐ J'envoie le logo de la structure dépositaire en haute définition format .eps ou .jpg à crunch-responsable@utt.fr

Dans le cadre de la gestion des contacts, l'Université de Technologie de Troyes recueille les données de types : identité, fonction, numéro de téléphone et adresse e-mail. Ces informations sont enregistrées dans un fichier informatisé par le service de la Direction de la Formation et de la Pédagogie (DFP), sont conservées pendant une durée de 1 an et pourront être transmises à la Direction des Relations Entreprise et la Direction de Communication. Conformément à la loi "informatique et libertés" modifiée en 2018 par le règlement européen général sur la protection des données (RGPD), vous pouvez exercer vos droits d'accès, de rectification et d'effacement des données vous concernant en vous adressant au service de la Direction de la Formation et de la Pédagogie (DFP) ou au délégué à la protection des données personnelles par les adresses mails suivantes : dfp@utt.fr ou DPO@utt.fr ou encore par courrier postal au : 12, Rue Marie Curie, CS 42060, 10004 Troyes cedex. Si vous estimez, après nous avoir contactés, que vos droits Informatique et Libertés ne sont pas respectés ou que le dispositif n'est pas conforme aux règles de protection des données, vous pouvez adresser une réclamation en ligne à la CNIL ou par voie postale.

SUJET D'INNOVATION

Pour toutes questions vous permettant de compléter cette fiche, nous vous invitons à contacter : crunch-responsable@utt.fr

Sujet confidentiel : (merci de rayer la mention inutile) Oui / Non

Thématique choisie : (une seule thématique à choisir)

Pour plus d'informations sur les thématiques, rendez-vous sur le site internet [ici](#)

- ☐ Thématique 1 - Systèmes automatisés de production et dispositifs embarqués
- ☐ Thématique 2 - Création, transmission et stockage responsable de la donnée
- ☐ Thématique 3 - Communication et Interactions Responsabilité Sociétale des Entreprises (RSE)
- ☒ Thématique 4 - Intelligence artificielle et technologies Green au service des SI
- ☐ Thématique 5 - Organisation responsable des Systèmes industriels
- ☐ Thématique 6 - Nouveaux matériaux pour une démarche responsable et innovante
- ☐ Thématique 7 - Environnement, recyclage et réduction des déchets
- ☐ Thématique 8 - Finance et climat
- ☐ Thématique 9 - Thématique transversale : sujets variés selon les enjeux des entreprises

participantes,
par exemple : conception mécanique, gestion de projets complexes, organisation du travail...

INTRODUCTION

Ardian est le leader européen du private equity. Fondée en 1996, l'entreprise a accumulé depuis sa création un très grand nombre de documents confidentiels. Dans une démarche de valorisation de ces documents, l'équipe IT est à la recherche d'une solution pour extraire le texte de ces Documents. Pour mieux appréhender cette problématique, l'équipe stagiaire 100% UTTienne du projet ArdianBrowser vous lance ce défi.

DEFI

Enoncé

« Extraire le maximum de texte brut de ce PDF, de la manière la mieux structurée possible. »

Explication

Le « texte brut » ciblé est l'ensemble des mots avec le maximum de ponctuation.

La récupération de la « structure » visée concerne la conservation de l'unité des paragraphes. A titre d'information, ceci est le 10^e paragraphe de ce document

REGLES

Contraintes

1- Format

Votre solution sera sous la forme d'un projet python qui contiendra la fonction suivante :

```
def extract_text(path):  
    """  
    Main function that extracts text from a file in a structured way to return the list of its paragraphs.  
    :param path: path to the targeted file  
    :return: list of paragraph strings  
    """  
  
    # TODO  
  
    return raw_text
```

2- Compliance

L'intégralité de la documentation interne d'Ardian est confidentielle, pour des raisons évidentes de compliance, votre solution ne doit en aucun cas utiliser un service tiers qui aurait accès au document. C'est-à-dire pas d'appel à une API d'OCR (ilovepdf, adobe, etc..). Le fichier doit être intégralement traité en local sur votre machine.

Vous êtes libres d'utiliser n'importe quel outil ou logiciel tiers à condition de prouver que la compliance a été respectée (L'utilisation de l'outil n'entraîne pas l'accès au fichier traité par un tiers). Si vous utilisez des outils tiers, vous êtes invités à décrire brièvement les aspects pertinents de leur installation et de leur configuration

3- Scalability

Votre solution devra être exécutable sur un corpus qui dépasse le million de fichiers, avec une évolution linéaire du temps d'exécution en fonction du nombre de fichiers à traiter. (3 points bonus si vous prouvez empiriquement la quasi-linéarité de votre solution sur une période d'une heure minimum)

Vous détaillerez les différents problèmes de mémoire qui pourraient potentiellement apparaître avec un corpus d'une telle taille, ainsi qu'une explication sur comment y remédier. Vous indiquerez également le nombre de fois que votre fonction `extract_text()` réussit à traiter ce fichier en 10 minutes sur la meilleure de vos machines. (En précisant les caractéristiques de celle-ci)

Barème

Rien de vous y oblige, mais si vous êtes friand de challenges, essayez de récolter le maximum de points de notre barème :

| Si vous parvenez à : | Récompense : |
|--|------------------|
| Récupérer le texte brut du paragraphe d'introduction | 1 point |
| Conserver la structure du paragraphe d'introduction | 4 points |
| Récupérer tout le texte de ce tableau | 3 points |
| Conserver la structure de ce tableau (i.e. 1 cellule = 1 paragraphe) | 4 points |
| Récupérer tous les pays européens de l'item 1 | 2 points |
| Pour chaque pays non-européen correctement récupéré dans l'item 1 | 2 points |
| Conserver la structure du paragraphe « <i>Market definition</i> » dans l'item 2 | 5 points |
| Récupérer le texte de toutes les cellules et de tous les libellés du tableau « <i>Market dynamics</i> » dans l'item 2 | 2 points |
| Conserver la structure de chaque cellule/libellé de ce même tableau (i.e. 1 cellule/libellé = 1 paragraphe) | 5 points |
| Récupérer la totalité du texte du schéma « <i>ADE addressed market overview</i> » dans l'item 2 | 3 points |
| Conserver la structure des paragraphes dans l'item 3 (plus de 13 paragraphes) | 7 points |
| Conserver la structure en 4 paragraphes dans la section « <i>SNAPSHOT</i> » de l'item 4 | 6 points |
| Conserver la structure en 3 paragraphes dans la section « <i>M&A STRATEGY</i> » de l'item 4 (6 paragraphes tolérés avec les dates) | 6 points |
| TOTAL | 50 points |

Bonus

Parce qu'on est sympa :)

Si vous parvenez à récupérer au moins une fois le filigrane « CONFIDENTIEL » → 5 points

Si vous parvenez à récupérer « ARDIAN » avec le logo du pied de page → 2 points

Si vous étendez votre solution pour qu'elle prenne en charge d'autres formats de fichiers contenus dans la liste suivante → 2 points par extension

['.doc', '.docx', '.ppt', '.pptx', '.png', '.jpeg', '.jpg']

(Il va sans dire que les fichiers Word & PowerPoint contiennent évidemment des images avec du texte à récupérer)

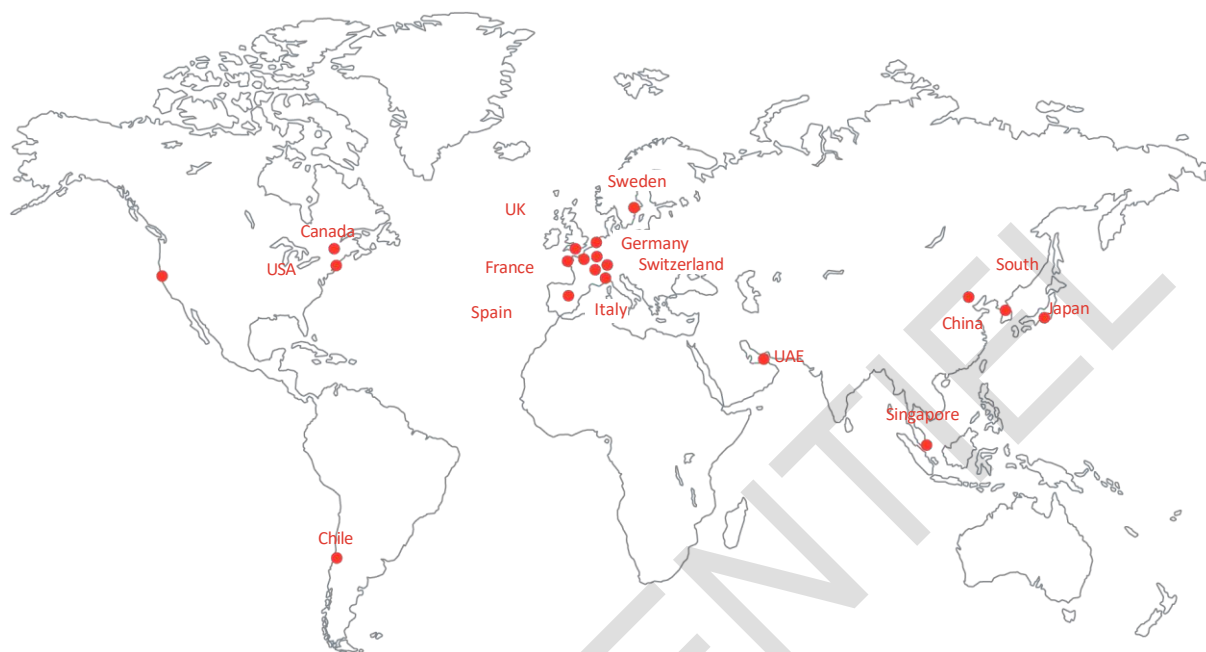
Si vous prouvez empiriquement la quasi-linéarité de votre solution sur une période d'une heure minimum (un graphe nb de traitements en fonction du temps suffira) → 3 points

Si vous n'utilisez que des outils open-source : 3 points

CONFIDENTIEL

Real data

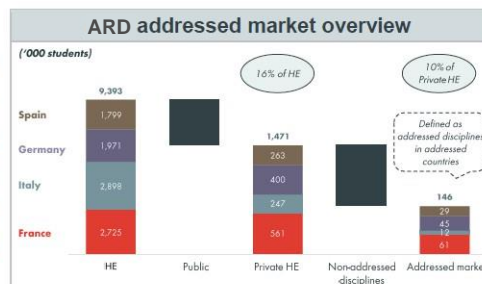
Item 1



Item 2

Market definition

The Higher Education ("HE") in Europe represents 23m students and ARD is present in France (63% of revenues), Italy (12% of revenues), Germany (12% of revenues) and Spain (13% of revenues), together representing 9.4m HE students (40% of the total HE market in Europe) and 1.5m students in private HE (50% of the private market). Out of these 1.5m students, 150k students (10%) are in ARD **addressed disciplines** of Creative Arts, namely: (i) Design & Graphical Arts as well as Digital in France, (ii) Design in Italy, (iii) Creative Arts and Media & Communications in Germany and (iv) Audiovisual (professional curriculum) in Spain.



Market dynamics

| | Demographics | | GER | | Share of discipline | | Share of private HE | | International students | | Market volume | | Average fee | | Market value | |
|--------------------------------------|--------------|-------|-------|-------|---------------------|-------|---------------------|-------|------------------------|-------|---------------|--------|-------------|-------|--------------|--------|
| | 16-20 | 20-25 | 16-20 | 20-25 | 16-20 | 20-25 | 16-20 | 20-25 | 16-20 | 20-25 | 16-20 | 20-25 | 16-20 | 20-25 | 16-20 | 20-25 |
| CAGR | | | | | | | | | | | | | | | | |
| France - Creative Arts | +0.3% | +0.7% | +0.7% | +0.7% | +1.4% | +1.3% | +1.8% | +0.9% | +0.5% | +0.2% | +4.6% | +3.8% | +1.5% | +1.5% | +6.2% | +5.3% |
| France - Digital | +0.3% | +0.7% | +0.7% | +0.7% | +9.6% | +4.2% | 0.0% | 0.0% | 0.0% | 0.0% | +10.5% | +5.6% | +1.5% | +1.5% | +12.2% | +7.2% |
| Italy - Design | -0.4% | +0.2% | +1.5% | +1.0% | +5.9% | +4.2% | +1.6% | +1.9% | +0.4% | +0.4% | +9.1% | +7.6% | +1.0% | +1.0% | +10.2% | +8.7% |
| Germany - Design & Media | 0.0% | -1.3% | +0.4% | +0.7% | -0.6% | -1.3% | +4.8% | +5.7% | +0.7% | +0.7% | +5.3% | +4.5% | +1.0% | +1.0% | +6.3% | +5.5% |
| Spain - Audiovisual | -0.8% | +0.7% | +4.0% | +3.6% | +1.6% | +1.5% | +7.1% | +4.7% | +0.2% | +0.2% | +12.1% | +10.4% | +1.0% | +1.0% | +13.3% | +11.5% |
| Total | -0.2% | +0.4% | +1.7% | +1.6% | +2.9% | +2.0% | +3.4% | +2.6% | +0.3% | +0.2% | +8.2% | +6.9% | | | +9.5% | +8.2% |
| Total (weighted avg. ADE enrolments) | -0.1% | +0.3% | +1.4% | +1.2% | +2.6% | +1.7% | +3.0% | +2.3% | -0.4% | +0.3% | +7.3% | +5.8% | | | +8.6% | +7.1% |

The addressed market for ARD has been growing in **volumes** at +8.2% p.a. over 2016-20 and in **value** at +9.5% p.a., mostly driven by increased share of private education in addressed disciplines (c.40% of vol. growth), increased share of

Item 3

- > We are in friendly territories with Kaminauseri and Hulumi acting as KG's advisors (governance and Manpack);
- > We have a strong value proposition with our geographical footprint in the Company's core geographies and potential proprietary cross-fertilization with Willow Group's sourcing network;
- > **We have further consolidated our set of experts and advisors:**
 - We have assembled an outstanding team of operating partners and advisors:
 - o **Manu Albani** ("MA"): Former CEO of Colitransiviks (2013-2019 – French platform of c.€115m sales). Killian Guoguenar knows him well and validated him to work alongside us;
 - o **Dorian Libouri** ("DL"): CEO of Firewelli, digital platform acquired by Sadisi group (spin-off from Imperial College);
 - o **Jean-Michel Courtalinay** ("JMC" – Laguardes Associés): Founder of Enodisu, Head of Laguardes Associés, working with all education players in France, specifically on accreditations, employability and placement topics;
 - Our team is also comprised of a remarkable group of third-party advisors, including:
 - o **Advaloremi** (M&A advisor – Guillaume Pierre & Rodrick Van Bure): they have worked successfully in precedents in education and are both former Beyond executives, with a direct line and close relationship with the Beyond team;
 - o **Athenadorlures** (Commercial DD): Colliseums is one of the very few consulting firms very knowledgeable on the education space (notably the Partner Anabelle Roger). They have performed a high-quality phase 1 strat. DD;
 - o **LUSCO** (Finance DD - Ghilem Anequinn): LUSCO has worked alongside us on Houras and has extensive knowledge and experience in the education sector;
 - o **KART** (Business intelligence): KART is performing a reputational and operational assessment of the management team, including school directors;
 - o **Private& Expense** (M&A, Structuring & Financing – Nicholett Ramirez & Beatrice Conssegretto): Nicholett Ramirez is close to the personal advisor of Guillaume Pierre (Rodrick Van Buren) as they were former colleagues (at Willow).

Item 4

Kinoa Firman

