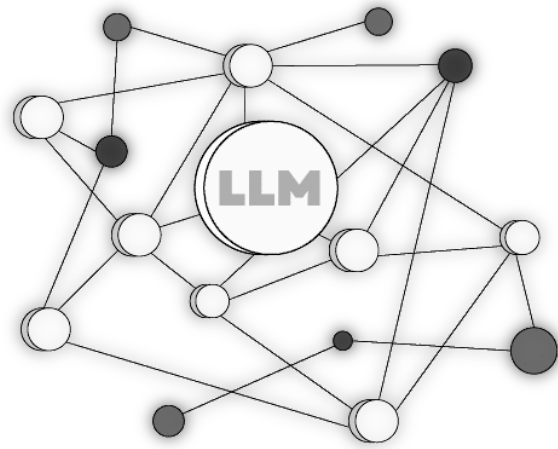


# Utilisation de LLM(large language models)

dans le Jeu de Role.



# Table des matières

<b>Introduction</b>	<b>4</b>
<b>Objets et Definitions</b>	<b>5</b>
Objets	5
Espace De Travail	5
<b>Applications sur notre espace</b>	<b>7</b>
<b>Encodeur</b>	<b>7</b>
définitions . . . . .	9
clarifications . . . . .	10
clarifications . . . . .	11
<b>Theorie et Modelisation</b>	<b>13</b>
<b>Mots et Phrases</b>	<b>13</b>
<b>Parlons sémantique</b>	<b>14</b>
fonctions des mots . . . . .	14
implémentation . . . . .	15
<b>Parlons Aquisition de Données</b>	<b>16</b>
Préambule . . . . .	16
Architècture de notre Modèl de Web Scraping . . . . .	17
Web Scrollers . . . . .	18
Rôle de nos Automates . . . . .	19
Jeu de données . . . . .	20
<b>Du code et encore du code</b>	<b>22</b>
séquenceurs . . . . .	22
activateurs . . . . .	25

<b>Biaisons Les Données</b>	<b>48</b>
<b>Modèle</b>	<b>48</b>
<b>Moteur de Jeu</b>	<b>50</b>
moteur de jeu . . . . .	51
<b>Générateur de texte</b>	<b>60</b>
<b>Chaine de markov</b>	<b>60</b>
implémentation . . . . .	60
performances . . . . .	61
<b>Réseau LSTM</b>	<b>61</b>
implémentation . . . . .	61
performances . . . . .	64
<b>Conclusions</b>	<b>64</b>
<b>Bibliographie Commentée</b>	<b>64</b>

# Introduction

## Introduction

Le récent essor des LLM, nous laisse espérer : une implémentation numérique satisfaisante, d'un jeu de rôle quelconque, similaire au jeu de rôles sur table : Dungeon et dragon.

Celui-ci (le jeu de rôle), entretient une relation étroite avec la littérature fantastique. Il en puise directement ses inspirations et en reprend les codes.<sup>1</sup>

Nous vient alors la question légitime : «Peut-on entraîner un modèle génératif à concevoir des scenarios de jeu de rôle, à partir d'ouvrages (fantastiques ou non), trouvés en libre accès ? »

Depuis la parution de Dungeon et dragon (DND) en 1974, nombreuses ont été les tentatives de portages numériques, du célèbre jeu de rôle. Ces bien heureuses tentatives, de par leur contraintes techniques, ont donné naissance aux mécaniques de RPG (qui font très largement consensus dans l'industrie contemporaine du jeu vidéo)<sup>2</sup>.

Cependant, après tant d'années d'expérimentations, il nous faut faire l'amer constat, de l'échec au moins partiel, de ces portages<sup>3</sup>. Pour cause, le CRPG (pour computer role playing game), perd en flexibilité, dans ses tentatives de réduire le jeu de rôle à son aspect statistique.

C'est là : une contrainte que le jeu de rôle, partage très largement avec «le traitement naturel du langage», et son tout récent essor, nous laisse espérer, que l'on puisse trouver une solution à ce premier problème en la résolution du second.

---

1. Nous pouvons citer en outre des auteurs comme Tolkien , Lovecraft , Glen Cook , Gorge R.R Martin , Robert E Howard , voir des auteurs de science-fiction tel que Philippe K. Dick (notoirement connu pour blade runner)

2. Connecting Worlds.Fantasy Role-Playing Games Ritual Acts and the Magic Circle 2005

3. wikipedia jeu de role

# Objets et Definitions

## Objets

**MAGMA :** Soit  $E$  un ensemble et  $*$  une *loi de composition interne* sur  $E$ ,  $(E, *)$  forme alors un magma.

**PARTIE STABLE  $L$  :** Soit  $E$  un magma,  $*$  sa loi de composition interne, et  $L \subseteq E$ ,  $L$  est alors une *partie stable* de  $E$  si

$$\forall (x, y) \in L^2 \quad x * y \in L$$

**ELEMENT NEUTRE :** Soit  $E$  un magma et  $*$  sa loi de composition interne, alors  $(E, *)$  admet un *element neutre*  $\varepsilon$  ; si

$$\forall a \in E \quad \varepsilon * a = a * \varepsilon = a$$

**LOI ASSOCIATIVE :** Soit  $(E, *)$  un magma.  $*$  est associative sur  $E$  si :

$$\forall (a, b, c) \in E^3 \quad (a * b) * c = a * (b * c)$$

## Espace De Travail

**ALPHABET  $\Sigma$  :** Soit  $\Sigma$  un ensemble de lettres, tel que  $\Sigma \neq \emptyset$  et que  $\text{Card}\Sigma < +\infty$ .

**MOT  $l \in \Sigma^*$  :** On appelle *mot* sur l'alphabet  $\Sigma$  toute suite  $l = (l_1, l_2, \dots, l_n) \in \Sigma^n$  avec  $n \in \mathbb{N}$ , que nous noterons  $l = l_1 l_2 \dots l_n$ .

**ALPHABET  $\Sigma^*$  :** Soit  $\Sigma$  un alphabet, alors  $\Sigma^*$  est l'ensemble des mots de  $\Sigma$ .

**LONGUEUR || D'UN MOT  $l \in \Sigma^*$  :** Soit  $\Sigma$  un alphabet, et  $l = l_1 l_2 \dots l_n \in \Sigma^*$ , alors  $|l| = n$  est la *longueur* de  $l$ .

**PHRASE  $s$  :** On appelle *phrase* sur l'alphabet  $\Sigma$  toute suite  $s = (m_1, m_2, \dots, m_n) \in (\Sigma^*)^n$  avec  $n \in \mathbb{N}$ .

### Proposition

Soit  $\Sigma$  un alphabet.

Soit  $*$  une loi interne associative sur  $\Sigma^*$ , tel que  $\forall (n, \eta) \in N^2$ ,  $\forall (m^{[1]}, m^{[2]}) \in \Sigma^n \times \Sigma^\eta$  avec  $m^{[1]} = (m_1^{[1]}, m_2^{[1]}, \dots, m_n^{[1]})$  et  $m^{[2]} = (m_1^{[2]}, m_2^{[2]}, \dots, m_\eta^{[2]})$

$$\text{alors } m^{[1]} * m^{[2]} = m_1^{[1]} m_2^{[1]} \dots m_n^{[1]} . m_1^{[2]} . m_2^{[2]} \dots m_\eta^{[2]}$$

$(\Sigma^*, *)$  est donc une structure de *magma*, et nous notons  $\varepsilon$  son *element neutre*, tel que  $|\varepsilon| = 0$ .

Itération

$$\begin{cases} \forall n \in N \ \forall x \in \Sigma^* \ x^{n+1} = x * x^n \\ x^0 = \varepsilon \end{cases} \quad (1)$$

Nous pouvons etre ammenes a note  $nx = x^n$

### Definition

Nous noterons simplement qu'un langage  $L$  sur l'alphabet  $\Sigma$ , est un sous ensemble de  $(\Sigma^*)^*$ .

### Language a Base Disjointe

Soit  $n \in N$ ,  $\forall i \in [1, n]$   $L_i \subseteq \Sigma^*$ , et tel que  $\forall (i, j) \in [1, n]^2$  avec  $i \neq j$ ,  $L_i \cap L_j = \{\varepsilon\}$ .

On note alors  $B = (L_1, L_2, \dots, L_n)$ , que l'on appellera *BaseDisjointe*.

Soit  $L$  un langage.

Soit  $Q \subseteq \{x | \exists k \in N \text{ tel que } x \in [1, n]^k\}$  avec  $Card Q < +\infty$ .

$$\text{Si } L = \sum_{\forall x \in Q} \sum_{\forall j \in x} L_j, \text{ alors } B \text{ est une base de } L.$$

$L$  est donc un langage a base disjointe.

# Applications sur notre espace

## Encodeur

L'idée d'utiliser nos ALGORITHMES CORRECTEURS sur du texte brut, nous vient de plusieurs constats.

En effet, extraire du texte de livres, pose de nombreux problèmes, dont le plus immédiat est l'encodage. Les règles de typographie<sup>1</sup>, les erreurs dans la conversion de PDF<sup>2</sup>, complexifient le problème, si bien, que l'on doit en permanence procéder à un nettoyage minutieux du texte. Celui-ci doit être nettoyé afin d'être analyser et correctement traiter<sup>3</sup>. Nous voulons alors prétraiter le texte, pour en simplifier son traitement, et rendre ainsi la donnée plus digeste<sup>4</sup>.

Une des solutions possibles, pourrait consister à encoder directement notre texte (en utf-8 pour l'exemple), puis le décoder dans la seconde qui suit. Il faudra en outre remplacer certains patterns récurrents qui complexifie inutilement notre donnée, par la fonction `replace` de python<sup>5</sup>. Seulement, cette méthode montre rapidement ses limites. Nous sommes amenés à traiter de grandes quantités de texte, et ces méthodes, coûteuses en temps, paraissent inélégante et par-dessus tout non-pertinentes. Pour peu que l'on ait une liste de patterns à interchanger, la fonction `replace` se verrait d'en l'obligation de parcourir notre chaîne de caractères à plusieurs reprises.

---

1. Il peut s'agir de "-" avant un dialogue, ou d'un "-" coupant un mot en 2 lors d'un retour à la ligne.

2. La bibliothèque que nous utilisons, laisse parfois certains artefacts lorsqu'elle extrait le texte de certains PDF, comme des erreurs d'encodage ou des mots dont les lettres sont espacées (exemple : e x e m p l e).

3. Le nettoyage ne doit pas être le même, pour l'analyse du texte que pour celle de son

<p>1 "âĀŞ 41 âĀŞ \nâĀŦ Bien,  mon sieur, trĀĀs bien !  rĀŦpliqua l'inspecteur des  fo-\nr ĀŦ t s , v o u s d  e m a n d e z , m a f i  l l e e n m a r i a g e ,  d ' a u t r e s f o n t  d e \nmĀĤme ; mon devoir  de pĀĤre m'oblige ĀĤ me  soucier d'elle ; je vous  donne un dĀŦlai de trois  jours pour chercher une  ombre ; dans trois jours,  prĀŦsentez-vous ĀĤ moi avec  une ombre qui vous aille  bien, et vous serez le  bienvenu ; mais le  quatriĀĤme jour âĀŦ ĀŦcou-  tez bien ce que je vous dis  âĀŦ ma fille sera la femme  d'un autre. ĀŦ \n Je  voulus encore adresser un  mot ĀĤ Mina ; mais, redou-\  nblant de sanglots, elle se  serra plus fort contre sa  mĀĤre et celle-ci, sans mot  dire, me fit signe de m'  ĀŦloigner."</p>	<p>1 "Il prit \naussitĀŦt la parole  : \n ĀĤ Je m'ĀŦtais  annoncĀŦ pour aujourd'hui,  vous n'avez pas eu \nl a p  a t i e n c e d ' a t t e  n d r e l e m o m e n t  f i x ĀŦ . M a i s t o u  t p e u t e n c o r e s  'arranger : acceptez mon  conseil, votre ombre est  encore ĀĤ vo-tre  disposition, et aussitĀŦt  vous revenez sur vos pas.  Vous serez \nle bienvenu  dans le jardin de l'  inspecteur des forĀŦts, et  tout cela n'aura ĀŦtĀŦ qu'  une plaisanterie ; quant ĀĤ  Rascal, qui vous a trahi e  t d e m a n d e l a m a i  n d e v o t r e f i a n c ĀŦ  e , j e m e c h a r g e  d e l u i e t l a c h o s e  est pratiquement rĀŦglĀŦe.  ĀŦ \n J'avais la sensation  de vivre un rĀŦve : ĀĤ  Vous vous ĀŦtiez annoncĀŦ  pour aujourd'hui ? ĀŦ "</p>
--	---

traitement(exemple : le LLM doit-être capable de reconnaître un dialogue, alors que cette information est superflu, voir gênante dans notre analyse.

4. Nous voulons simplifier le plus possible, la donnée que le LLM doit prédire plutôt que de complexifier son modèle d'apprentissage.

5. "ab".replace("abd", "cccc") == "ccccd" est une proposition vraie.



## définitions

### Encodage

Soit  $\Sigma_1$  et  $\Sigma_2$  2 alphabets.

On appelle *ENCODAGE* toute application :

$$e : \Sigma_1^* \rightarrow \Sigma_2^*$$

### exemple

Dans cet exemple, il s'agit juste de passer d'un langage binaire a un autre  $a \mapsto 0$  et  $b \mapsto 1$ .

$\Sigma_1 = \{a, b\}$	$\Sigma_2 = \{0, 1\}$
abbbabbabbabbabb	0111011011011011
aaaabbbaabbbaaaa	0000110001100000
aaaaabbbbbaabbabb	0000011110011101
bbababbbaabaabbbb	1101011000100111
aaabbaab	00011001

### Encodeur

Soit  $\Sigma_1$  et  $\Sigma_2$  2 alphabets.

Soit  $n \in \mathbb{N}$ .

Soit  $((u_1, u_2, \dots, u_n), (v_1, v_2, \dots, v_n)) \in (\Sigma_1^* \times \Sigma_2^*)$  ON NOTE :

$$\sigma = \begin{pmatrix} u_1 & u_2 & \cdots & u_n \\ v_1 & v_2 & \cdots & v_n \end{pmatrix}.$$

On appelle *ENCODEUR* tout  $e : \Sigma_1^* \rightarrow \Sigma_2^*$  :

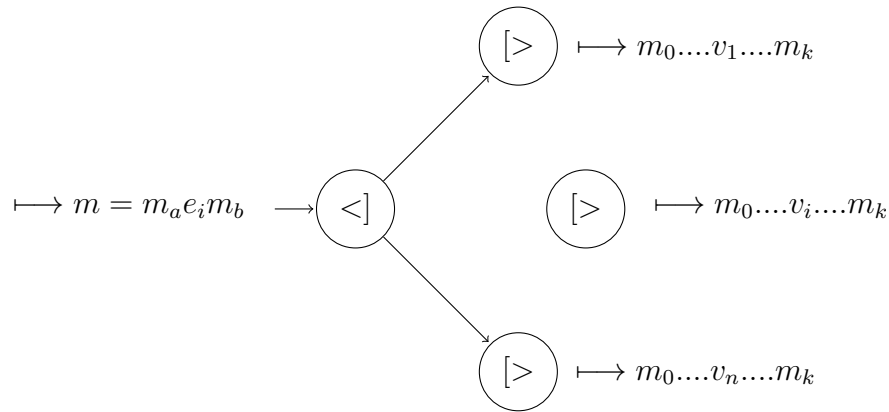
$$e : \begin{cases} \Sigma_1^* \rightarrow \Sigma_2^* \\ m \mapsto \begin{cases} e(m_a).v_i.e(m_b) & \text{si } m = m_a.u_i.m_b \\ m & \text{sinon} \end{cases} \end{cases}.$$

$e$  est dit *CORRECTEUR* si  $\Sigma_1 = \Sigma_2$ .

## clarifications

Les encodeurs ne sont en outre que des encodages spécifiques, qui changent certains patterns au profit de d'autres. Ce sont des applications sur notre magma, mais ceux-ci ne se comportent pas comme des morphismes de groupes (si l'analogie nous est permise). Soit  $\forall f$  encodeur, il est FAUX de dire que :

$$\forall (m_a, m_b) \in (\Sigma_1^*)^2 f(m_a.m_b) = f(m_a).f(m_b)$$



## exemple

```

» > coder = EncodeAutomate("a" : "c" , "aa" : "b")
» > coder.sequence("aaaaaaaaabaaaaaba")
» > 'bbbbbcbbbbba'

```

### Activateur

Soit  $\Sigma$  un alphabet.  
Soit  $S_L \subseteq \Sigma^*$ .

On appelle *ACTIVATEUR*, tout encodeur  $\varphi$  tel que :

$$\forall m = m_1.m_2....m_n \in \Sigma^* \quad \varphi : \begin{cases} \Sigma^* \rightarrow \{0, 1\} \\ m \mapsto \begin{cases} 1 \text{ si } \exists i \in [1; n] \text{ tel que } m_i \in S_L \\ 0 \text{ sinon} \end{cases} \end{cases}$$

.

### clarifications

L'activateur est simplement une fonction booléenne sur notre magma. Ils vont nous être utiles, afin de mener des études statistiques sur de larges portions de textes, et d'ainsi nous ramener à des résultats numériques.

## Automates

Soit  $n \in \mathbb{N}$  on pose  $E = \llbracket 1; n \rrbracket$ , un ensemble d'états fini.

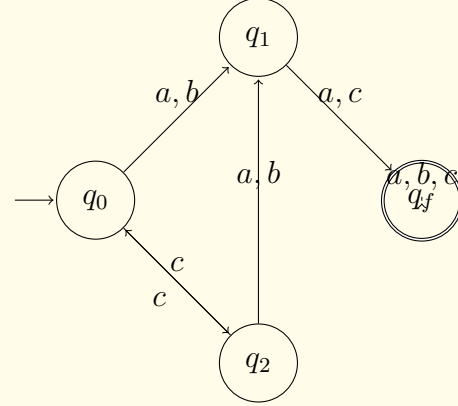
Soit  $(q_0, q_f) \in E^2$ , respectivement son état initial et acceptant.

Soit  $\Sigma$  un alphabet.

Soit  $f$  une fonction dite de *transition*, tel que :

$$f : \begin{cases} E \times \Sigma \rightarrow P(E) \\ (x, a) \mapsto f(x, a) \subseteq E \end{cases}$$

Alors  $(E, \Sigma, f)$  est un *AUTOMATE*.



## Automate Structurel

Soit  $\Sigma$  un alphabet.

Soit  $L_\Sigma$  un langage sur  $\Sigma$ , à base disjointe  $B = (E_1, E_2, \dots, E_n)$  avec  $n \in \mathbb{N}$ .

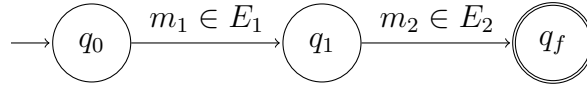
Soit  $S = \llbracket 1, n \rrbracket$ , un ensemble de sommets.

Soit  $f : (S, \Sigma^*) \rightarrow S$ , une fonction de transition sur  $(S, \Sigma)$ .

Soit  $(q_0, q_f) \in S^2$ , respectivement son état initial et acceptant.

$$f' : \begin{cases} S \times \Sigma^* \rightarrow P(E) \\ (s, m) \mapsto \begin{cases} f(s, m) \text{ si } m \in E_{f(s, m)} \\ \emptyset \text{ sinon} \end{cases} \end{cases}$$

$(S, \Sigma^*, f')$  est donc un automate structurel sur  $B$ .



# Théorie et Modélisation

Il est ardu de définir la langue française, comme étant l'un de nos objets précédemment défini. Ceci pour la simple et bonne raison, qu'il est difficile d'en déterminer les frontières. Une phrase est-elle française pour la simple raison qu'elle est syntaxiquement correcte<sup>1</sup>, ou faut-il nécessairement que sa sémantique soit bonne<sup>2</sup> ?

En d'autres termes, suffit-il de former une phrase à partir d'un n-uplet quelconque de mots de la langue française, pour former une phrase correcte (indépendamment de son sens) ? Si oui, notre langue française peut être considérée comme une partie stable d'un autre langage plus large, en cas contraire il faudra en énumérer les structures<sup>3</sup> ?

Cela pose un problème autrement plus fondamental. Nous serons amenés à travailler avec des textes de fantaisie, dont les auteurs sont particulièrement friands de noms inexistantes. Considérant la phrase : "La compagnie prenait alors la route de Kathovar." , il nous faut faire un choix.

Il nous est alors soit interdit de considérer cette phrase comme étant française, ou dans un second cas, il nous faut faire le constat qu'une phrase française peut-être constituée d'éléments qui ne le sont pas nécessairement. Cela tient au fait que la langue française, est une langue naturelle et non une langue construite.

## Mots et Phrases

Nous voulons entraîner un modèle à produire des phrases françaises. Si comme le veut la grammaire française : "une phrase commence par une majuscule et finie par un point", il nous faut alors nuancer son applicabilité à notre base de données.

En effet, si une phrase est supposée commencer par une majuscule et que sa conclusion est suivie d'un point, toutes majuscules n'impliquent pas

- 
1. "Pierre Chat la lancer a-t-il ?" est-elle une phrase de la langue Française ?
  2. "Le chat a-t-il lancer la pierre ?" est-elle une phrase française pour la simple et bonne raison qu'elle a du sens ?
  3. Pour espérer approximer la langue française par un langage à bases disjointes ?

forcément la formation d'une phrase, et tous les points n'en concluent pas nécessairement une !

Or extraire les phrases puis les mots d'un texte, est une condition nécessaire à l'étude préalable de nos textes.

Par ailleurs, extraire les phrases d'un texte d'origine, se révèle être un exercice bien moins intuitif qu'il n'y paraît ! "Nous pouvons citer Gorge R.R Martin et Philipe K. Dick !", cette seule phrase suffit à expliquer l'inefficacité d'une méthode naïve d'extraction de phrase.

Extraire les mots d'une phrase, met également la méthode naïve (la simple utilisation de la fonction `split` sur notre texte) en échec (en raison des espaces insécables (certains mots contiennent un ou plusieurs espaces)).

## Parlons sémantique

Dire qu'une phrase est bien formée ou non, est non seulement l'objet de notre étude, mais également une question délicate ! Nous ne pouvons le faire que par une succession d'exemples, puis parier sur le fait que cette énumération soit suffisamment exhaustive (bien que cela, soit peu probable). C'est pour cela, que nous préférons, approximer la langue française (qui est une langue naturelle), par un langage formel dont nous aurons déterminer au préalable la base.

Il nous faut composer avec l'hypothèse, que la langue française est composée de structures intelligibles, que nous pouvons énumérées, à l'aide d'études statistiques sur des phrases bien formées. Pour ainsi construire les bases disjointes de notre langage formel. Pour ce faire, nous comptons séparer les mots de la langue française, selon leur fonctions (fonctions sémantiques), puis énumérer leurs différents enchainements, dans les phrases à analyser.

A priori, trouver des textes en bon français, n'est pas chose bien compliquée ; de façon intuitive, il conviendra de collecter nos exemples à partir de textes d'auteurs et/ou d'académistes. Notre capacité à en extraire correctement les phrases puis les mots, selon leurs fonctions, déterminera la qualité de notre analyse.

## fonctions des mots

Nous possédons une liste (conséquente, mais non exhaustive) de mots de la langue française (en fichier `texte`). Nous remarquons également que le

dictionnaire français en ligne : "Larousse"<sup>1</sup>, nous fournit dans ses descriptions : la fonction (ou les fonctions) de chaque mots. En s'inspirant d'une API, qui lui est dédiée<sup>2</sup>, nous programmons une fonction, capable de retourner la fonction d'un mot donne, si tenter que celui-ci ait été décrit dans le Larousse.

## implémentation

```
1 def fundamental(word : str):
2     """obtenir la fonction d'un mot"""
3     url = "https://www.larousse.fr/dictionnaires/francais/" +
4         word.lower()
5     soup = BeautifulSoup(requests.get(url=url).text, 'html.
6         parser')
7     if soup.select('p') != []:
8         return unicodedata.normalize("NFKD", soup.select('p')
9         [0].text)
10    else:
11        return None
```

Listing 1 – Python classification Implementation

Nous construisons et sauvegardons, en un premier temps, un dictionnaire (python), contenant les mots de notre liste et leurs fonctions. Pour pallier la lenteur problématique de notre algorithme de web scrapping, nous multi-threadrons ce handicap pour en faire un problème mineur.

Sur les 22735 mots de notre liste, nous comptons 22661 fonctions *trouvees*.

---

1. <https://www.larousse.fr/dictionnaires/francais>  
2. [https://github.com/quentin-dev/larousse\\_api](https://github.com/quentin-dev/larousse_api)

# Parlons Aquisition de Donnees

## Préambule

Pour des raisons de temps d'exécution, notre analyse sémantique (qui se fera en plusieurs temps), sera comprise dans notre grand model d'acquisition de données, chargé de collecter des livres sur internet, puis d'en extraire les textes. Celui-ci traitera le texte une première fois, l'analysera, puis le traitera une seconde fois. Pour ce faire, nous aurons besoin de nos structures précédemment définies d'activateurs , d'automates et d'encodeurs.

Nous avons peu de moyens pour influencer sur la qualité du model d'entraînement, c'est pour cela que nous baserons pratiquement l'entièreté de notre stratégie sur l'amélioration de la qualité de notre donnée d'entraînement. Ainsi, notre étude, se portera principalement sur le processus d'acquisition de données, de leur traitement et de leur analyse.

Traiter de grandes quantités de données, est en général un exercice préparatoire a une course contre la montre. Il nous faut donc, veiller à conserver une vigilance absolue au niveau de la complexité de nos algorithmes. On la vue précédemment, faire appel à de la programmation asynchrone, peut nous permettre dans une large mesure, de pallier ce problème. Cependant démultiplier le nombre de thread qui travaillent en même temps, démultiplie également le risque que certains d'entre eux échouent, dans leur requêtes(car l'on envoie et l'on reçoit des paquets d'informations). Il nous faut donc conserver un certain équilibre.

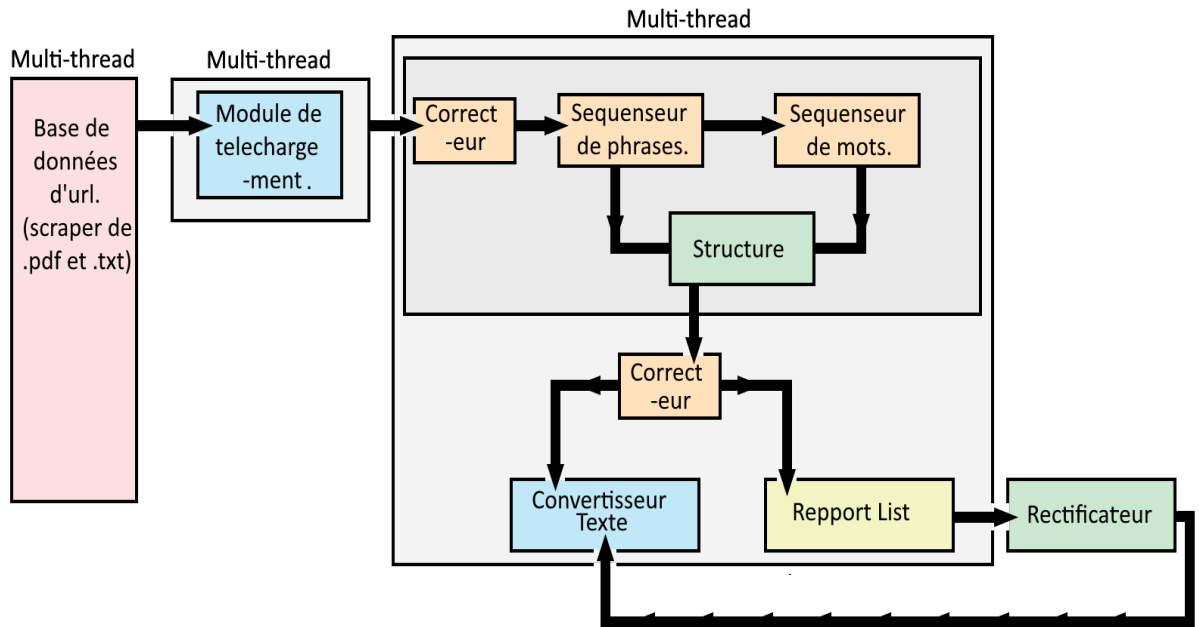
Pour ce faire, nous utiliserons fréquemment notre fonction : "src manager", que voici :

```
1 def src_manager(file_procc : [str], n : int):
2     """segmente notre liste en plusieurs listes plus petites
3     """
4     if len(file_procc) <= n:
5         return [file_procc]
6     else:
7         return [file_procc[:n]] + src_manager(file_procc[n:],
8         n)
```

Listing 2 – Python ramification Implementation



## Architècture de notre Modèle de Web Scraping



Celle-ci possède 4 phases :

1. Nos robots web scrolleurs parcourent nos pages pour en répertorier les url.
2. Nous téléchargeons lesdits documents (pdf ou texte).
3. Après analyse, le modèle décide s'il doit l'enregistrer dans nos données d'entraînement (s'il a réussi à cartographier les structures sémantiques), ou s'il l'envoie au rectificateur.
4. le rectificateur enregistre dans un fichier les phrases ou les mots que le modèle n'arrive pas à cataloguer.

## Web Scrollers

Après délibérations, nous avons conservé 4 sources de données majeures, que nous serons amenés à expliciter par la suite. Pour chacun de ces sites, nous avons analysé leur architecture, puis à l'aide d'automates, avons classé leurs liens pertinents, en les répertoriant dans des listes. Ces listes, nous servirons par la suite, afin de télécharger les ouvrages gratuits, proposés par ces sites.

Pour extraire les phrases de notre texte, nous le parcourons tout entier, en instanciant notre tête à chaque majuscule précédée d'un point (l'initialisant à 0), puis en instanciant notre queue à chaque point précédé d'une majuscule. En excluant certains patterns récurrents, tel que : "M. ", puis en ajoutant à une liste chacune de ces séquences, nous ramifions notre texte par phrases.

Nous partons du principe, que nous avons une phrase et que nous voulons en extraire les mots. Nous initialisons alors au préalable une liste des séquences de notre phrase, en utilisant la fonction native de python : "split()", sur notre texte. Nous avons ainsi découpé notre phrase, à chaque espace. Puis nous nous munissons d'un dictionnaire référençant tous les mots (de notre liste de mots français), contenant des espaces non sécables. Nous additionnons les séquences ainsi obtenues, tant que cette somme appartient à notre dictionnaire, puis, appliquons le même procédé de façon récursive : sur le reste de notre liste de séquences.

*EXEMPLE* : "Je m'étais assoupi parce que la journée arrivait à son terme." Nous avons notre liste de séquences suivante :

["Je", "m'étais", "assoupi", "parce", "que", "la", "journée", "arrivait", "à", "son", "terme"]

Nous obtenons après analyse :

["Je", "m'étais", "assoupi", "parce que", "la", "journée", "arrivait", "à", "son", "terme"]

## Rôle de nos automates

Nos automates, ont un rôle double : Après analyse de nos données, nous obtenons un dictionnaire contenant : les structures sémantiques de notre texte, et de leur occurrence dans notre donnée d'origines. Nous initialisons alors, une liste d'automates structurels à partir de ce dictionnaire. De façons périodiques, notre modèle LLM, sauvegardera ses paramètres lors de son entraînement, puis une étude statistique sera menée à partir de notre liste d'automates structurels, qui mesurons l'occurrence de ces structures sémantiques, dans ses productions. En outre à l'aide de ces fréquences, nous pourrions nous ramener à des notions de distance sur un espace vectoriel à dimension finie  $R^n$  (en mesurant la distances entres les fréquences obtenues sur nos données d'entraînement et celles obtenues à partir de notre model), pour mesurer sa capacité, à produire des phrases qui ont du sens (correctes sémantiquement).

Les automates ont un second rôle, moins abstrait : trier les url obtenues par nos web scrollers. Prenons notre cas de base : Nous nous trouvons sur une page web référençant un e-book gratuit, en format PDF. Celle-ci contient, plusieurs liens plus ou moins pertinents : -Des liens extérieurs menants à d'autres sites (comme des pages Facebook ou X (anciennement twitter) ) -Des liens menants à la page d'accueil du site, ou à ses diverses catégories. -Et bien sur des liens contenant des informations sur le dit livre, dont le plus important se trouve être le lien d'accès à sa version PDF. Faire la différence entre ces liens, est aisé pour un opérateur humain et l'est beaucoup moins, pour un robot. Cependant, analyser la structure de ce site, nous permet de repérer des patterns récurrents dans ces liens.

Par exemple les liens donnent accès à ses livres numériques, commencent tous par un "https", suivit du nom de domaine du site, puis par un "ebook" et fini par un ".pdf". Nous faisons alors appel à notre constructeur d'automate, à l'aide de ces données.

Il nous vient immédiatement à l'esprit la norme 2 sur  $R^n$ , cependant (et malgres une equivalence des normes), rien ne nous dit que celle-ci soit la plus precise, pour mesurer nos ecarts semantiques.

```
1 def Norme_2_Distance(vect_1 : (int) , vect_2 : (int)):  
2     assert (len(vect_1) != len(vect_2))  
3     return sum([(vect_1[i] - vect_2[i])**2 for i in range(len  
    (vect_1))])**(1/2)
```

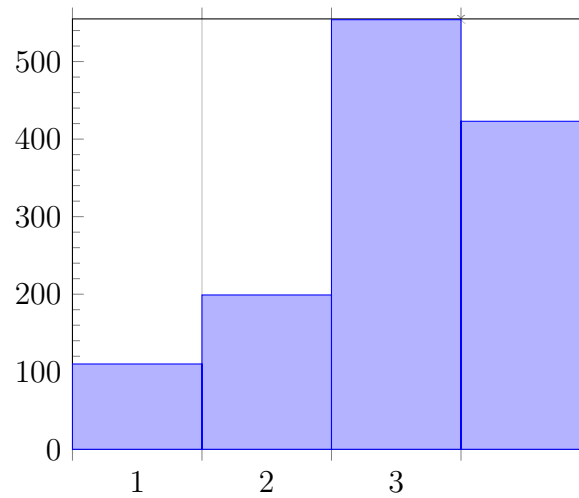
Listing 3 – Python Norme 2 Implementation

## Jeu de données

Regardons de plus près notre donnée.

Nous avons ci-dessous, la répartition de nos données, scrappées par nos algorithme. Celles-ci ont été récoltées sur 4 sites distincts.

1. <http://www.livrespourtous.com>, qui représente 8.5% de nos données d'entraînement.
2. <https://www.vousnousils.fr>, qui représente 15.4% de nos données d'entraînement.
3. <https://www.edition999.info>, qui représente 43% de nos données d'entraînement.
4. <https://www.gutenberg.org>, qui représente 32.8% de nos données d'entraînement.



Il est a noté, que seul Gutenberg, présente des livres numériques en format texte(.txt) standard. On a ainsi 68% de nos données, extraite de pdf. Concernant la part d'œuvres de fantaisies (ou fantastiques) dans nos données d'entraînement, celle-ci devrait excéder les 50%, sachant qu'elle compose en totalité les œuvres d'édition 999 et une bonne moitié de celle de livres-pourtous. Il est cependant impossible d'en connaître la composition, pour Gutenberg.

Plus précisément nous avons utiliser nos web scrollers sur les liens ci-dessous :

1. [http://www.livrespourtous.com/e-books/list/onecat/Livres-electroniques+Contes-et-nouvelles+Fantastique-et-Science-Fiction/0/all\\_items.html](http://www.livrespourtous.com/e-books/list/onecat/Livres-electroniques+Contes-et-nouvelles+Fantastique-et-Science-Fiction/0/all_items.html).
2. <https://www.edition999.info/-Fantastique-.html>.
3. [http://www.livrespourtous.com/e-books/list/onecat/Livres-electroniques+Romans+Fantastique-et-SF/0/all\\_items.html](http://www.livrespourtous.com/e-books/list/onecat/Livres-electroniques+Romans+Fantastique-et-SF/0/all_items.html).
4. [http://www.livrespourtous.com/e-books/list/onecat/Livres-electroniques+Lettres-et-memoires/0/all\\_items.html](http://www.livrespourtous.com/e-books/list/onecat/Livres-electroniques+Lettres-et-memoires/0/all_items.html).
5. [http://www.livrespourtous.com/e-books/list/onecat/Livres-electroniques+Documents-et-essais/0/all\\_items.html](http://www.livrespourtous.com/e-books/list/onecat/Livres-electroniques+Documents-et-essais/0/all_items.html).
6. [http://www.livrespourtous.com/e-books/list/onecat/Livres-electroniques+Romans+Historique/0/all\\_items.html](http://www.livrespourtous.com/e-books/list/onecat/Livres-electroniques+Romans+Historique/0/all_items.html).

## Du code et encore du code

### séquenceurs

```
1  """
2  All Sequensers Models
3  """
4  #Imports Pannel requierement
5  import matplotlib.pyplot as plt
6  import threading
7  import os
8  import io
9  import pickle
10 import csv
11 import requests
12 import unicodedata
13 from bs4 import BeautifulSoup
14 accept_dict = {".": None, "!" : None, "?" : None}
15 class EncodeAutomate():
16     def __init__(self, activation_function : dict):
17         self.af = activation_function
18         self.activation = {}
19         def insere(sequence : str , sequence_list : [str]):
20             if sequence_list == []:
21                 return [sequence]
22             for inc in range(len(sequence_list)):
23                 if (len(sequence) >= len(sequence_list[inc]))
24 :
25                 return sequence_list[:inc] + [sequence] +
26                 sequence_list[inc:]
27         for sequence in self.af:
28             if sequence[0] in self.activation:
29                 self.activation[sequence[0]] = insere(
30 sequence, self.activation[sequence[0]])
31             else:
32                 self.activation[sequence[0]] = [sequence]
33         def sequence(self, sentence : str):
34             def activation(work : str , sentence : str):
35                 assert len(work) == len(sentence)
36                 if len(work) == 0:
37                     return True
38             else:
39                 if (work[0] == sentence[0]) or (sentence[0]
40 == "ÂÍ"):
41                     return activation(work[1:], sentence[1:])
```

```

38         else:
39             return False
40         my_new_sentence = ""
41         increase = 0
42         while increase < len(sentence):
43             if sentence[increase] in self.activation:
44                 NotEqual = True
45                 sec_increase = 0
46                 while NotEqual and (sec_increase < len(self.
activation[sentence[increase]])):
47                     sequence = self.activation[sentence[
increase]][sec_increase]
48                     if increase + len(sequence) <= len(
sentence):
49                         if activation(sentence[increase:
increase + len(sequence)],sequence):
50                             my_new_sentence = my_new_sentence
+ self.af[sequence]
51                             increase = increase + len(
sequence)
52                             NotEqual = False
53                         else:
54                             sec_increase += 1
55                         else:
56                             sec_increase += 1
57                     if NotEqual:
58                         my_new_sentence = my_new_sentence +
sentence[increase]
59                         increase += 1
60                     else:
61                         my_new_sentence = my_new_sentence + sentence[
increase]
62                         increase += 1
63                 return my_new_sentence
64 def shearch_patterns(txt : str, i : int, pattern : str):
65     return (((i + len(pattern)) <= len(txt)) and (txt[i:i +
len(pattern)] == pattern))
66 anti_point_encoder = EncodeAutomate({"." : "", "!" : "", "?"
: ""})
67 coder = EncodeAutomate({"\n" : " ", "--" : "", "_" : "", "Âñ" :
" ", "Âž" : " ", ", " : " ", ";" : " ", ":" : " ", "n'" : "ne
pas", "qu'" : "que ", "Qu'" : "Que ", "quâĂŽ" : "que ", "QuâĂŽ
" : "Que ", "nâĂŽ" : "ne pas"})
68 def get_sentence_list(my_txt : str):
69     if len(my_txt) > 1:

```

```

70     initialize_list = []
71     head_play = 0
72     in_stand = True
73     for i in range(1,len(my_txt)):
74         if (my_txt[i] != (my_txt[i].lower())) and (not
shearch_patterns(my_txt,i,"M. ")) and (not in_stand):
75             head_play = i
76             in_stand = True
77             elif (my_txt[i] in accept_dict) and (my_txt[i-1]
!= "M") and in_stand:
78                 initialize_list.append(coder.sequence(my_txt[
head_play:i] + my_txt[i]))
79                 in_stand = False
80                 elif i == (len(my_txt) - 1):
81                     initialize_list.append(coder.sequence(my_txt[
head_play:i]))
82             return initialize_list
83     else:
84         return []
85 def get_space_words():
86     fileIn = open("data collections\liste_francais.txt","r")
87     content = [(line[:len(line) - 1]) for line in fileIn.
readlines() if (" " in line[:len(line) - 1])]
88     fileIn.close()
89     return dict(zip(content,[None for a in content]))
90 def sub_analyse(word : str, work_list : [str] , word_dict :
dict):
91     word = anti_point_encoder.sequence(word)
92     if work_list == []:
93         return (word,[])
94     elif (word + " " + work_list[0]) not in word_dict:
95         return (word,work_list)
96     else:
97         return sub_analyse(word + " " + work_list[0],
work_list[1:],word_dict)
98 def init_analyse(work : [str], non_list : [str],
words_with_space : dict):
99     if work == []:
100         return non_list
101     else:
102         results = sub_analyse(work[0],work[1:],
words_with_space)
103         return init_analyse(results[1],non_list + [results
[0]],words_with_space)
104 def get_words_list_of_a_sentence(sentence : str,

```



```

105     words_with_space : dict):
106     my_stock = sentence.split()
107     return init_analyse(my_stock, [], words_with_space)
pass

```

Listing 4 – Python Code des Sequenceurs

## activateurs

```

1  """
2  All Activators Models
3  """
4  class ActivatorAutomate():
5      def __init__(self):
6          pass
7      def recognize(self, sequence : str):
8          return True
9  class ActivatorName(ActivatorAutomate):
10     def __init__(self):
11         ActivatorAutomate.__init__(self)
12         pass
13     def recognize(self, word : str):
14         return ((word[0] != (word[0].lower())) or word.
isdigit() or ((len(word)>2) and (word[0]+word[1] in ["d'",
"l'"])) and (word[2] != word[2].lower())) or ((len(word) >
1) and word[0].isdigit and (word[1] != (word[1].lower()))))
15 class TrivialAutomate(ActivatorAutomate):
16     def __init__(self, transitions : [dict], final_state :
int, initial_state=0):
17         ActivatorAutomate.__init__(self)
18         self.transitions = transitions
19         self.final_state = final_state
20         self.initial_state = initial_state
21         pass
22     def recognize(self, word : str, state=0):
23         if word == "":
24             return (state== self.final_state)
25         elif state > (len(self.transitions) - 1):
26             return False
27         else:
28             if word[0] in self.transitions[state]:
29                 return self.recognize(word[1:], self.
transitions[state][word[0]])

```

```

30         else:
31             return False
32 class LoopedAutomate(ActivatorAutomate):
33     def __init__(self, transitions : [dict], final_state :
34     int, initial_state=0):
35         ActivatorAutomate.__init__(self)
36         self.transitions = transitions
37         self.final_state = final_state
38         self.initial_state = initial_state
39     pass
40     def recognize(self, word : str, state=0):
41         if word == "":
42             return (state== self.final_state)
43         elif state > (len(self.transitions) - 1):
44             return False
45         else:
46             if word[0] in self.transitions[state]:
47                 return self.recognize(word[1:], self.
48                 transitions[state][word[0]])
49             else:
50                 return self.recognize(word[1:], state)
51 class AsyncAutomate(ActivatorAutomate):
52     def __init__(self, transitions : [[bool,int,dict]],
53     final_state : int, initial_state=0):
54         ActivatorAutomate.__init__(self)
55         self.transitions = transitions
56         self.final_state = final_state
57         self.initial_state = initial_state
58     pass
59     def recognize(self, word : str, state=0):
60         if word == "":
61             return (state== self.final_state)
62         elif state > (len(self.transitions) - 1):
63             return False
64         else:
65             if word[0] in self.transitions[state][2]:
66                 return self.recognize(word[1:], self.
67                 transitions[state][2][word[0]])
68             else:
69                 if not self.transitions[state][0]:
70                     return self.recognize(word[1:], state)
71                 else:
72                     return self.recognize(word[1:], self.
73                     transitions[state][1])
74 class AsyncAutomateMono(ActivatorAutomate):

```

```

70     def __init__(self, transitions : [[bool,int,dict]],
final_state : int, initial_state=0):
71         ActivatorAutomate.__init__(self)
72         self.transitions = transitions
73         self.final_state = final_state
74         self.initial_state = initial_state
75         pass
76     def recognize(self, word : str, state=0):
77         if word == "":
78             return (state== self.final_state)
79         elif state > (len(self.transitions) - 1):
80             return False
81         else:
82             if word[0] in self.transitions[state][2]:
83                 return self.recognize(word[1:], self.
transitions[state][2][word[0]])
84             else:
85                 if not self.transitions[state][0]:
86                     return self.recognize(word[1:], state)
87                 else:
88                     return self.recognize(word, self.
transitions[state][1])
89 def constructEndAutomata(construct_lists : [[str]]):
90     basic_graph_construct = [[False,0,{}]]
91     tails_list = [0]
92     for construct_list in construct_lists:
93         sub_tail_list = []
94         for word in construct_list:
95             state = tails_list[-1]
96             indice = len(word) - 1
97             for i in range(indice):
98                 if word[i] in basic_graph_construct[state
][2]:
99                     state = basic_graph_construct[state][2][
word[i]]
100                 else:
101                     new_node = len(basic_graph_construct)
102                     basic_graph_construct.append([True,
tails_list[-1],{}])
103                     basic_graph_construct[state][2][word[i]]
= new_node
104                     state = new_node
105                     sub_tail_list.append((state,word[indice]))
106                     new_big_node = len(basic_graph_construct)
107                     basic_graph_construct.append([False,new_big_node,{}])

```

```

108         tails_list.append(new_big_node)
109         for tail in sub_tail_list:
110             basic_graph_construct[tail[0]][2][tail[1]] =
new_big_node
111         basic_graph_construct[tails_list[-1]][0] = True
112         basic_graph_construct[tails_list[-1]][1] = tails_list[-2]
113         return AsyncAutomate(basic_graph_construct, tails_list
[-1])
114 def constructEndAutomataNotEq(construct_lists : [[str]]):
115     basic_graph_construct = [[False, 0, {}]]
116     tails_list = [0]
117     for construct_list in construct_lists:
118         sub_tail_list = []
119         for word in construct_list:
120             state = tails_list[-1]
121             indice = len(word) - 1
122             for i in range(indice):
123                 if word[i] in basic_graph_construct[state
][2]:
124                     state = basic_graph_construct[state][2][
word[i]]
125                 else:
126                     new_node = len(basic_graph_construct)
127                     basic_graph_construct.append([True,
tails_list[-1], {}])
128                     basic_graph_construct[state][2][word[i]]
= new_node
129                     state = new_node
130                     sub_tail_list.append((state, word[indice]))
131                     new_big_node = len(basic_graph_construct)
132                     basic_graph_construct.append([False, new_big_node, {}])
133                     tails_list.append(new_big_node)
134                     for tail in sub_tail_list:
135                         basic_graph_construct[tail[0]][2][tail[1]] =
new_big_node
136         return AsyncAutomate(basic_graph_construct, tails_list
[-1])
137 class DictAutomate():
138     def __init__(self, transitions : [dict], final_state :
int, initial_state=0):
139         self.transitions = transitions
140         self.final_state = final_state
141         self.initial_state = initial_state
142     def ddelete_occurence(my_list : list):
143         new_list = []

```

```

144         for k in my_list:
145             if k not in new_list:
146                 new_list.append(k)
147         return new_list
148     self.links = [delete_occurrence([k for k in tr.values
149 (()) for tr in transitions]
150 pass
151 def recognize(self, sentence : [str], state=0):
152     if sentence == []:
153         return (state== self.final_state)
154     else:
155         if sentence[0] in self.transitions[state]:
156             return self.recognize(sentence[1:], self.
157 transitions[state][sentence[0]])
158         else:
159             return False
160 def recognize_sentence(self, sentence : str):
161     return self.recognize(sentence.split(" "), self.
162 initial_state)
163 def constructEndAutomataMoloList(construct_lists : [str]):
164     basic_graph_construct = [[False,0,{}]]
165     tails_list = [0]
166     for construct_list in construct_lists:
167         sub_tail_list = []
168         for word in construct_list:
169             state = tails_list[-1]
170             indice = len(word) - 1
171             for i in range(indice):
172                 if word[i] in basic_graph_construct[state
173 ][2]:
174                     state = basic_graph_construct[state][2][
175 word[i]]
176                 else:
177                     new_node = len(basic_graph_construct)
178                     basic_graph_construct.append([True,
179 tails_list[-1],{}])
180                     basic_graph_construct[state][2][word[i]]
181 = new_node
182                     state = new_node
183                     sub_tail_list.append((state,word[indice]))
184     new_big_node = len(basic_graph_construct)
185     basic_graph_construct.append([False,new_big_node,{}])
186     tails_list.append(new_big_node)
187     for tail in sub_tail_list:
188         basic_graph_construct[tail[0]][2][tail[1]] =

```

```

new_big_node
182     basic_graph_construct[tails_list[-1]][0] = True
183     basic_graph_construct[tails_list[-1]][1] = tails_list[-2]
184     return AsyncAutomateMono(basic_graph_construct,
tails_list[-1])

```

Listing 5 – Python Code des Activateurs

## web scrapeur

```

1 import all_sequensers_models as asqs
2 import requests
3 import unicodedata
4 from bs4 import BeautifulSoup
5 import os
6 import io
7 import pickle
8 import csv
9 import threading
10 text_cleaner_encoder = asqs.EncodeAutomate({"Ã" : "e", "Ã" :
    "e", "Ã" : "e", "Ã" : "e",
11                                           "-\n" : "", "- "
    : "\n - ", "aÃ Ã aÃ " : "",
12                                           "aÃ ÃÃÃ aÃ " :
    "" , "aÃ ÃÃÃ aÃ " : "",
13                                           " aÃ " : "\n aÃ
    ", "-\n" : "", "Ã" : "a",
14                                           "Ã" : "a", "Ã" :
    "a", "Ã" : "o", "Ã" : "o",
15                                           "Ã" : "i", "Ã" :
    "i", "Ã" : "u", "Ã" : "u", " " : ""})
16 def save(Dictionnaire : dict ,path="data collections\
database_dictionnaire.obj"):
17     fileIn = open(path,"wb")
18     pickle.dump(Dictionnaire,fileIn)
19     fileIn.close()
20 def load(path="data collections\database_dictionnaire.obj"):
21     fileout = open(path,"rb")
22     Load = pickle.load(fileout)
23     fileout.close()
24     return Load
25 def MultiThreadFragmentation():
26     Words_Functions = {}

```

```

27     Interrupt_Dict = {'conjonction de coordination' : None, '
conjonction' : None, "adverbe" : None, 'proverbe': None, '
interjection' : None, "adjectif" : None,
28         "article" : None, "preÃAposition" : None,
'preÃAfixe' : None, "nom" : None, "verbe" : None, "Participe"
: None, "pronom" : None, "locution" : None}
29     register_list = ['conjonction de coordination', '
conjonction', "adverbe", 'proverbe', 'interjection', "adjectif
", "article", "preÃAposition", 'preÃAfixe', "nom", "verbe", "
Participe", "pronom", "locution"]
30     def initializeRegister(word_path : str , word : str):
31         if word_path != None:
32             if word_path in Words_Functions:
33                 Words_Functions[word_path][word] = None
34             else:
35                 Words_Functions[word_path] = {word : None}
36         else:
37             initializeRegister("nom", word)
38     def three_activation(word_path : str, word : str,
pattern_list : [str]):
39         if word_path in Interrupt_Dict:
40             initializeRegister(word_path, word)
41         else:
42             if word_path != None:
43                 for pattern in pattern_list:
44                     if pattern in word_path:
45                         initializeRegister(pattern, word)
46                         return None
47             initializeRegister("nom", word)
48         pass
49     def fundamental(word : str):
50         """obtenir la fonction d'un mot"""
51         url = "https://www.larousse.fr/dictionnaires/francais
/" + word.lower()
52         soup = BeautifulSoup(requests.get(url=url).text, '
html.parser')
53         if soup.select('p') != []:
54             return unicodedata.normalize("NFKD", soup.select(
'p')[0].text)
55         else:
56             return None
57     def encapsule(n : int , file_procc : [str]):
58         """craie des ramifications d'une meme liste par
paquet de n valeurs"""
59         if len(file_procc) <= n:

```

```

60         return [file_procc]
61     else:
62         return [file_procc[:n]] + encapsule(n,file_procc[
n:])
63     class MyThread(threading.Thread):
64         def __init__(self, words : [str]):
65             threading.Thread.__init__(self)
66             self.words = words
67         def run(self):
68             for word in self.words:
69                 path = fundamental(word)
70                 three_activation(path,word,register_list)
71             pass
72     class MyThreadRelation(threading.Thread):
73         def __init__(self, words : [str]):
74             threading.Thread.__init__(self)
75             self.words = words
76         def run(self):
77             i = 0
78             for word in self.words:
79                 print(str(i) + "%")
80                 path = fundamental(word)
81                 i = i + 1
82                 three_activation(path,text_cleaner_encoder.
sequence(word),register_list)
83             pass
84     fileIn = open("data collections\liste_francais.txt","r")
85     content = [line[:len(line) - 1].lower() for line in
fileIn.readlines()]
86     fileIn.close()
87     partition = encapsule(100,content)
88     Multi_Thread = [MyThreadRelation(partition[0])] + [
MyThread(partition[j]) for j in range(1,len(partition))]
89     for th in Multi_Thread:
90         th.start()
91     for th in Multi_Thread:
92         th.join()
93     save(Words_Functions)
94 def MultiThreadFragmentationCorrecteur(content : list):
95     Words_Functions = load()
96     Interrupt_Dict = {'conjonction de coordination' : None,'
conjonction' : None,"adverbe" : None,'proverbe': None,'
interjection' : None,"adjectif" : None,
97                     "article" : None,"preposition" : None,
'prefixe' : None,"nom" : None,"verbe" : None,"Participe"

```



```

: None, "pronom" : None, "locution" : None}
98 register_list = ['conjonction de coordination', '
conjonction', "adverbe", 'proverbe', 'interjection', "adjectif
", "article", "prelAposition", 'prelAfixe', "nom", "verbe", "
Participe", "pronom", "locution"]
99 def initializeRegister(word_path : str , word : str):
100     if word_path != None:
101         if word_path in Words_Functions:
102             Words_Functions[word_path][word] = None
103         else:
104             Words_Functions[word_path] = {word : None}
105     else:
106         initializeRegister("nom", word)
107 def three_activation(word_path : str, word : str,
pattern_list : [str]):
108     if word_path in Interrupt_Dict:
109         initializeRegister(word_path, word)
110     else:
111         if word_path != None:
112             for pattern in pattern_list:
113                 if pattern in word_path:
114                     initializeRegister(pattern, word)
115                 return None
116             initializeRegister("nom", word)
117         pass
118 def fundamental(word : str):
119     """obtenir la fonction d'un mot"""
120     url = "https://www.larousse.fr/dictionnaires/francais
/" + word.lower()
121     soup = BeautifulSoup(requests.get(url=url).text, '
html.parser')
122     if soup.select('p') != []:
123         return unicodedata.normalize("NFKD", soup.select(
'p')[0].text)
124     else:
125         return None
126 def encapsule(n : int , file_procc : [str]):
127     """craie des ramifications d'une meme liste par
paquet de n valeurs"""
128     if len(file_procc) <= n:
129         return [file_procc]
130     else:
131         return [file_procc[:n]] + encapsule(n, file_procc[
n:])
132 class MyThread(threading.Thread):

```

```

133         def __init__(self, words : [str]):
134             threading.Thread.__init__(self)
135             self.words = words
136         def run(self):
137             for word in self.words:
138                 path = fundamental(word)
139                 three_activation(path,word,register_list)
140             pass
141     class MyThreadRelation(threading.Thread):
142         def __init__(self, words : [str]):
143             threading.Thread.__init__(self)
144             self.words = words
145         def run(self):
146             i = 0
147             for word in self.words:
148                 print(str(i) + "%")
149                 path = fundamental(word)
150                 i = i + 1
151                 three_activation(path,text_cleaner_encoder.
sequence(word),register_list)
152             pass
153     partition = encapsule(100,content)
154     Multi_Thread = [MyThreadRelation(partition[0])] + [
MyThread(partition[j]) for j in range(1,len(partition))]
155     for th in Multi_Thread:
156         th.start()
157     for th in Multi_Thread:
158         th.join()
159     save(Words_Functions)

```

Listing 6 – Python Code des Web Scrappers

## main

```

1  """
2  Big Scrapper Model
3  """
4  #Imports Pannel requierement
5  from bs4 import BeautifulSoup
6  import requests
7  import random
8  import os
9  import io
10 import pickle

```

```

11 import csv
12 import threading
13 import unicodedata
14 from PyPDF2 import PdfReader
15 import all_activators_models as actv
16 import all_sequensers_models as asqs
17 import dictionnary_scrapper as dscp
18 #globals Variables
19 text_cleaner_encoder = asqs.EncodeAutomate({"Ã" : "e", "Ã" :
    "e", "Ã" : "e", "Ã" : "e",
20
    " -\n" : "", " - "
    : "\n - ", "aÃ Ã aÃ " : "",
    "aÃ ÃÃÃ aÃ " :
21
    "" , "aÃ ÃÃÃ aÃ " : "",
    " aÃ " : "\n aÃ
22
    ", " -\n" : "", " " : "", "Ã" : "a",
    "Ã" : "a", "Ã" :
23
    "a", "Ã" : "o", "Ã" : "o",
    "Ã" : "i", "Ã" :
24
    "i", "Ã" : "u", "Ã" : "u", "Ã" : "c"})
25 structurize_cleaner = asqs.EncodeAutomate({" -\n" : "", " - " :
    "\n - ", "aÃ Ã aÃ " : "",
    "aÃ ÃÃÃ aÃ " :
26
    "" , "aÃ ÃÃÃ aÃ " : "",
    " aÃ " : "\n aÃ
27
    ", " -\n" : ""})
28 simplification_cleaner = asqs.EncodeAutomate({"Ã" : "e", "Ã" :
    "e", "Ã" : "e", "Ã" : "e", "Ã" : "a",
    "Ã" : "a", "Ã" :
29
    "a", "Ã" : "o", "Ã" : "o",
    "Ã" : "i", "Ã" :
30
    "i", "Ã" : "u", "Ã" : "u", "Ã" : "c"})
31 analyse_coder = asqs.EncodeAutomate({"\n" : " ", "--" : "", "_":
    " ", "Ã" : " ", "Ã" : " ", " " : " ", " " : " ", " " : " ",
    "aÃ" : ""})
32 txt_recorder_coder = asqs.EncodeAutomate({" " : "\t", "Ã" : "
    e", "Ã" : "e", "Ã" : "e", "Ã" : "e", "Ã" : "a",
    "Ã" : "a", "Ã" :
33
    "a", "Ã" : "o", "Ã" : "o",
    "Ã" : "i", "Ã" :
34
    "i", "Ã" : "u", "Ã" : "u", "Ã" : "c"})
35 verb_automata = actv.constructEndAutomataMoloList([["es", "ons
    ", "ez", "ent", "ont", "ait", "ais", "ions", "iez", "aient", "erai"
    , "eras", "era", "erons", "erez", "eront", "ant",
36
    "asse", "asses", "

```

```

37     Ãt", "assions", "assiez", "assent", "rer", "a", "ais", "etait", "
    est", "as", "ÃtÃ", "ete", "es", "Ãzmes", "fumes",
    "Ãtes", "etes", "
38     sommes", "sois", "soit", "aient"]])
39 is_a_word_activator = actv.ActivatorName()
40 words_with_spaces_dict = asqs.get_space_words()
41 abbreviation_dict = {"l'" : None, "s'" : None, "L'" : None, "S'"
    : None, "c'" : None, "C'" : None, "j'" : None, "J'" : None
    , "t'" : None, "T'" : None, "d'" : None, "D'" : None, "m'" :
    None, "M'" : None,
42     "laÃ" : None, "saÃ" : None, "LaÃ" : None
    , "SaÃ" : None, "caÃ" : None, "CaÃ" : None, "jaÃ" :
    None, "JaÃ" : None, "taÃ" : None, "TaÃ" : None, "daÃ" :
    None, "DaÃ" : None, "maÃ" : None, "MaÃ" : None}
43
44
45
46
47
48
49 #obj stockage
50 def save_dict(Dictionnaire : dict ,path="data collections\
    database_dictionary.obj"):
51     fileIn = open(path,"wb")
52     pickle.dump(Dictionnaire,fileIn)
53     fileIn.close()
54 def load_dict(path="data collections\database_dictionary.obj
    "):
55     fileout = open(path,"rb")
56     Load = pickle.load(fileout)
57     fileout.close()
58     return Load
59
60
61
62
63
64
65 #rsc management for multitask
66 def src_manager(file_procc : [str], n : int):
67     """segmente notre liste en plusieurs listes plus petites
    """
68     if len(file_procc) <= n:
69         return [file_procc]

```

```

70     else:
71         return [file_procc[:n]] + src_manager(file_procc[n:],
n)
72 def get_all_links(url : str):
73     my_request = requests.get(url)
74     bsoup = BeautifulSoup(my_request.text,"html.parser")
75     all_links = [l.get("href") for l in bsoup.findAll("a")]
76     return all_links
77 def get_all_links_filter(url : str, filter_automata :
callable):
78     return [link for link in get_all_links(url) if (link!=
None) and filter_automata(link)]
79 def get_all_links_anti_filter(url : str, filter_automata :
callable):
80     return [link for link in get_all_links(url) if (link!=
None) and not filter_automata(link)]
81 def turn_a_pdf_innto_a_txt(path : str):
82     pass
83 fast_dict = load_dict("data collections\database_fast_dict.
obj")
84
85
86
87
88
89
90
91
92
93
94 #web scroller collect urls
95 def collection_0002():
96     print("aquisition url, phase 02")
97     filter_automata_01 = actv.constructEndAutomata([[ "
vousnousils.fr"],[".pdf"]])
98     links_list = get_all_links_filter("https://www.
vousnousils.fr/ebooks-gratuits", filter_automata_01.
recognize)
99     return links_list
100 def collection_0004():
101     print("aquisition url, phase 04")
102     mapping = {}
103     filter_automata_01 = actv.constructEndAutomataNotEq([[ "
ebooks/"]])
104     filter_automata_02 = actv.constructEndAutomata([[ "ebooks/

```

```

    ],["txt.utf-8"]])
105     links_list = src_manager(["https://www.gutenberg.org" + a
        for a in get_all_links_filter("https://www.gutenberg.org/
browse/languages/fr#a50488", filter_automata_01.recognize)
    ], 20)
106     class TextScraper(threading.Thread):
107         def __init__(self, work : [str]):
108             threading.Thread.__init__(self)
109             self.work = work
110             self.filter = filter_automata_02.recognize
111         def run(self):
112             for task in self.work:
113                 for url in get_all_links_filter(task,self.
filter):
114                     if url not in mapping:
115                         mapping["https://www.gutenberg.org"+
url] = None
116             bot_list = [TextScraper(task_list) for task_list in
links_list]
117             for bot in bot_list:
118                 bot.start()
119             for bot in bot_list:
120                 bot.join()
121             return list(mapping.keys())
122 def collection_0003():
123     print("aquisition url, phase 03")
124     basic_mapping = ["https://www.edition999.info/-
Fantastique-.html?debut_articles=" +str(k*20)+"#
pagination_articles" for k in range(1,16)] + ["https://www
.edition999.info/-Fantastique-.html",
125         "https://www.edition999.info/-Heroic-
Fantasy-.html?debut_articles=10#pagination_articles",
126         "https://www.edition999.info/-Heroic-
Fantasy-.html?debut_articles=20#pagination_articles"]
127     mapping = {}
128     filter_automata_01 = actv.constructEndAutomataNotEq([[ '
https', 'http', "pagination_articles", '-Fantastique-.html',
129         'Avant-premiere-.html', '-Science-Fiction-Anticipation-.
html',
130         'Heroic-Fantasy-.html', '-Litterature-.html', '-Policier-et
-suspense-.html',
131         'Poesie-.html', '-Classique-.html',

```

```

132         'Autobiographie-Temoignage-Autofiction-.html', '-Jeunesse-.html',
133         'Litterature-Erotique-.html', '-Biographie-.html', '-Histoire-courte-.html',
134         'Langues-etrangeres-.html', '-Saga-.html',
135         'Essai-Politique-Scolaire-Education-.html', '-Philosophie-et-spiritualite-.html',
136         'Adolescents-et-Jeunes-Adultes-.html', '#']]
137     filter_automata_02 = actv.constructEndAutomataNotEq(["edition999.info"], ["spip.php?action=telecharger&arg"])
138     links_list = []
139     for main_link in basic_mapping:
140         links_list = links_list + ["https://www.edition999.info/" + link for link in get_all_links_anti_filter(
141             main_link, filter_automata_01.recognize)]
142         links_list = src_manager(links_list, 20)
143     class TextScrapper_03_01(threading.Thread):
144         def __init__(self, work : [str]):
145             threading.Thread.__init__(self)
146             self.work = work
147             self.filter = filter_automata_02.recognize
148         def run(self):
149             for task in self.work:
150                 for url in get_all_links_filter(task, self.filter):
151                     if url not in mapping:
152                         mapping[url] = None
153     class TextScrapper_03_02(threading.Thread):
154         def __init__(self, work : [str]):
155             threading.Thread.__init__(self)
156             self.work = work
157             self.filter = filter_automata_02.recognize
158         def run(self):
159             for k in range(len(self.work)):
160                 for url in get_all_links_filter(self.work[k], self.filter):
161                     if url not in mapping:
162                         mapping[url] = None
163                     print("etape : " + str(k))
164     bot_list = [TextScrapper_03_02(links_list[0])] + [

```

```

TextScrapper_03_01(links_list[k]) for k in range(1,len(
links_list))]]
164     for bot in bot_list:
165         bot.start()
166     for bot in bot_list:
167         bot.join()
168     return list(mapping.keys())
169 def collection_0001():
170     print("aquisition url, phase 01")
171     basic_mapping = []
172     scroll_mapping = ["http://www.livrespourtous.com/e-books/
list/onecat/Livres-electroniques+Contes-et-nouvelles+
Fantastique-et-Science-Fiction/"+str(k)+"/all_items.html"
for k in range(5)] + ["http://www.livrespourtous.com/e-
books/list/onecat/Livres-electroniques+Romans+Fantastique-
et-SF/"+str(k)+"/all_items.html" for k in range(11)]+["
http://www.livrespourtous.com/e-books/list/onecat/Livres-
electroniques+Documents-et-essais/"+str(k)+"/all_items.
html" for k in range(6)]+["http://www.livrespourtous.com/e-
books/list/onecat/Livres-electroniques+Lettres-et-
memoires/"+str(k)+"/all_items.html" for k in range(5)]+["
http://www.livrespourtous.com/e-books/list/onecat/Livres-
electroniques+Romans+Historique/"+str(k)+"/all_items.html"
for k in range(17)]
173     mapping = {}
174     filter_automata_00 = actv.constructEndAutomata(["e-books
/detail"],["all_items.html"])
175     filter_automata_02 = actv.constructEndAutomata([".pdf"
]])
176     class TextScrapper_01_01(threading.Thread):
177         def __init__(self, work : [str]):
178             threading.Thread.__init__(self)
179             self.work = work
180             self.filter = filter_automata_02.recognize
181         def run(self):
182             for task in self.work:
183                 for url in get_all_links_filter(task,self.
filter):
184                     if url not in mapping:
185                         mapping[url] = None
186     links_list = []
187     for link in scroll_mapping:
188         links_list = links_list + get_all_links_filter(link,
filter_automata_00.recognize)
189     print("etape 1")

```



```

190     bot_list = [TextScrapper_01_01(work) for work in
src_manager(links_list,2)]
191     for bot in bot_list:
192         bot.start()
193     for bot in bot_list:
194         bot.join()
195     return list(mapping.keys())
196
197
198
199
200
201
202
203
204 #data manager
205 def get_the_file_pdf(url : str, name : str):
206     request = requests.get(url)
207     try:
208         with open("back room\pdf\ " +str(name)+".pdf","wb")
as pdf:
209             for chunk in request.iter_content(chunk_size
=1024):
210                 if chunk:
211                     pdf.write(chunk)
212     except:
213         pass
214 def get_the_file_txt(url : str):
215     request = requests.get(url)
216     openFile = open(r"back room\gutenberg_unprocess_data.txt"
,"a+",encoding="utf-8")
217     if request.text != None:
218         openFile.write(request.text)
219     openFile.close()
220 def get_a_generate_key():
221     txt = ""
222     for k in range(14):
223         txt = txt + str(random.randint(0,9))
224     return txt
225
226
227
228
229
230

```

```

231
232
233
234
235
236
237
238
239
240 #Big Scrapper
241 def massive_scrapper():
242     l1_links = collection_0001()
243     l2_links = collection_0002()
244     l3_links = collection_0003()
245     the_txt_url_list = src_manager(collection_0004(),100)
246     the_pdf_url_list = src_manager(l1_links + l2_links +
247                                   l3_links,100)
248     class DownloadThread(threading.Thread):
249         def __init__(self, work : [str],file_type : str):
250             threading.Thread.__init__(self)
251             self.work = work
252             self.file_type = file_type
253         def run(self):
254             if self.file_type == "pdf":
255                 for task in self.work:
256                     get_the_file_pdf(task,get_a_generate_key
257                                     ())
258             else:
259                 for task in self.work:
260                     get_the_file_txt(task)
261     pdf_bots_list = [DownloadThread(work,"pdf") for work in
262                     the_pdf_url_list]
263     for bot in pdf_bots_list:
264         bot.start()
265     for bot in pdf_bots_list:
266         bot.join()
267     print("section pdf finie")
268     txt_bots_list = [DownloadThread(work,"texte") for work in
269                     the_txt_url_list]
270     for bot in txt_bots_list:
271         bot.start()
272     for bot in txt_bots_list:
273         bot.join()
274     print("section texte finie")

```

```

272
273
274
275
276
277 #analyseur
278 def pre_path_analyse():
279     words_functions = load_dict()
280     new_dict = {}
281     for ver in words_functions.keys():
282         for word in (words_functions[ver].keys()):
283             new_dict[word] = ver
284     save_dict(new_dict, path="data collections\
database_fast_dict.obj")
285 def path_analyse(word : str, fast_words_dict : dict):
286     new_word = simplification_cleaner.sequence(word)
287     if new_word in fast_words_dict:
288         return fast_words_dict[new_word]
289     elif (len(new_word) > 2) and new_word[:2] in
abreviation_dict:
290         cut_word = new_word[2:]
291         return path_analyse(cut_word, fast_words_dict)
292     else:
293         wo , ve = is_a_word_activator.recognize(new_word) ,
verb_automata.recognize(new_word)
294         if (wo and ve) or ve or (new_word == "a") or ((len(
new_word) > 2) and (new_word[:1] == "j'))):
295             return "verbe"
296         elif wo:
297             return "nom"
298         else:
299             openFile = open(r"back room\uknow_words.txt", "a+"
)
300             openFile.write(word+"\n")
301             openFile.close()
302             return "unknow"
303     pass
304 def get_list_recompiler(my_list : [str]):
305     new_list = []
306     for path in my_list:
307         if new_list == []:
308             new_list.append(path)
309         elif (new_list[-1] == "nom") and (path == "nom"):
310             pass
311         elif (new_list[-1] == "unknow") and (path == "unknow"

```

```

):
312         pass
313     else:
314         new_list.append(path)
315     return new_list
316 def sentence_analyser(sentence : str, fast_words_dict : dict,
317     structures_referencer : dict):
318     all_words = asqs.get_words_list_of_a_sentence(sentence,
319     words_with_spaces_dict)
320     recompile_tuple = tuple(get_list_recompiler([path_analyse
321     (word,fast_words_dict) for word in all_words]))
322     if "unknown" in recompile_tuple:
323         openFile = open(r"back room\unreferenced_data.txt","a
324         +")
325         openFile.write(sentence+" ||\n")
326         openFile.close()
327         pass
328     else:
329         if recompile_tuple in structures_referencer:
330             structures_referencer[recompile_tuple] =
331             structures_referencer[recompile_tuple] + 1
332         else:
333             structures_referencer[recompile_tuple] = 1
334             openFile = open(r"back room\training_data.txt","a+")
335             openFile.write(simplification_cleaner.sequence(
336             sentence)+" ||\n")
337             openFile.close()
338         pass
339 def txt_sentence_analyser(sentence : str, fast_words_dict :
340     dict, structures_referencer : dict):
341     all_words = asqs.get_words_list_of_a_sentence(sentence,
342     words_with_spaces_dict)
343     recompile_tuple = tuple(get_list_recompiler([path_analyse
344     (word,fast_words_dict) for word in all_words]))
345     if "unknown" in recompile_tuple:
346         openFile = open(r"back room\unreferenced_data.txt","a
347         +")
348         openFile.write(sentence+" ||\n")
349         openFile.close()
350         pass
351     else:
352         if recompile_tuple in structures_referencer:
353             structures_referencer[recompile_tuple] =
354             structures_referencer[recompile_tuple] + 1
355         else:

```

```

345         structures_referencer[recompile_tuple] = 1
346         openFile = open(r"back room\training_data.txt","a+")
347         openFile.write(txt_recorder_coder.sequence(sentence)+
" ||\n")
348         openFile.close()
349         pass
350 def all_pdf_analyse():
351     pdf_references = os.listdir(r"back room\pdf")
352     structure_dict = load_dict("data collections\
database_big_vector.obj")
353     text = ""
354     for reference in pdf_references:
355         try:
356             reader = PdfReader("back room\pdf\ " + reference
[1:])
357             number_of_pages = len(reader.pages)
358             for i in range(number_of_pages):
359                 page = reader.pages[i]
360                 text = text + structurize_cleaner.sequence(
page.extract_text())
361         except:
362             pass
363     class Structurer(threading.Thread):
364         def __init__(self, work : [str]):
365             threading.Thread.__init__(self)
366             self.work = work
367         def run(self):
368             for task in self.work:
369                 sentence_analyser(task,fast_dict,
structure_dict)
370             pass
371     print(text)
372     all_sentences = asqs.get_sentence_list(text)
373     task_list = src_manager(all_sentences,200)
374     bot_list = [Structurer(work) for work in src_manager(
all_sentences,200)]
375     for bot in bot_list:
376         bot.start()
377     for bot in bot_list:
378         bot.join()
379     print("transfert fini")
380     save_dict(structure_dict,"data collections\
database_big_vector.obj")
381     pass
382 def all_txt_analyse():

```

```

383     structure_dict = load_dict("data collections\
database_big_vector.obj")
384     fileIn = open(r"back room\gutenberg_unprocess_data.txt",
r",encoding="utf-8")
385     content = fileIn.readlines()
386     fileIn.close()
387     def split(a, n):
388         k, m = divmod(len(a), n)
389         return (a[i*k+min(i, m):(i+1)*k+min(i+1, m)] for i in
range(n))
390     text = ""
391     for line in content:
392         text = text + structurize_cleaner.sequence(line)
393     class Structurer(threading.Thread):
394         def __init__(self, work : [str]):
395             threading.Thread.__init__(self)
396             self.work = work
397         def run(self):
398             for task in self.work:
399                 txt_sentence_analyser(task,fast_dict,
structure_dict)
400                 pass
401     all_sentences = asqs.get_sentence_list(text)
402     bot_list = [Structurer(work) for work in split(
all_sentences,200)]
403     for bot in bot_list:
404         bot.start()
405     for bot in bot_list:
406         bot.join()
407     print("transfert fini")
408     save_dict(structure_dict,"data collections\
database_big_vector.obj")
409     pass
410 def rectification_01():
411     anti_n = asqs.EncodeAutomate({"\n" : ""})
412     fileIn = open(r"back room\uknow_words.txt", "r")
413     content = fileIn.readlines()
414     fileIn.close()
415     words_dict = {}
416     for word in content:
417         new_word = anti_n.sequence(word)
418         if new_word not in words_dict:
419             words_dict[new_word] = None
420     words_list = list(words_dict.keys())
421     dscp.MultiThreadFragmentationCorrecteur(words_list)

```

```

422 def recompile():
423     my_dict = load_dict()
424     new_dict = {}
425     for cat in my_dict:
426         for word in my_dict[cat]:
427             new_dict[word] = cat
428     fast_dict = save_dict(new_dict, "data collections\
database_fast_dict.obj")
429 data_decoder = asqs.EncodeAutomate({" ||\n" : "", "\t" : " "
}).sequence
430 def get_sentence_of_data(path=r"back room\training_data.txt")
:
431     with open(path, "r") as f:
432         text = f.read()
433     return text.split(" ||\n")
434 decode_sentence = asqs.EncodeAutomate({"\t" : " "}).sequence
435 def rectification_02():
436     structure_dict = {}
437     class Structurer(threading.Thread):
438         def __init__(self, work : [str]):
439             threading.Thread.__init__(self)
440             self.work = work
441         def run(self):
442             for task in self.work:
443                 sentence_analyser(task, fast_dict,
structure_dict)
444             pass
445     all_sentences = [data_decoder(a) for a in
get_sentence_of_data(r"back room\unreferenced_data.txt")]
446     task_list = src_manager(all_sentences, 200)
447     bot_list = [Structurer(work) for work in src_manager(
all_sentences, 200)]
448     for bot in bot_list:
449         bot.start()
450     for bot in bot_list:
451         bot.join()
452     print("transfert fini")
453     return structure_dict
454 def rectification():
455     rectification_01()
456     recompile()
457     rectification_02()
458 def massive_analyse():
459     all_pdf_analyse()
460     all_txt_analyse()

```

```

461     rectification()
462
463
464 def main_scraper():
465     massive_scrapper()
466     massive_analyse()
467     pass

```

Listing 7 – Python Code du collecteur de donnees

## Biaisons les données

# Modèle

```

1 import random
2 import pickle
3 import heapq
4 import numpy as np
5 import pandas as pd
6 import matplotlib.pyplot as plt
7 from nltk.tokenize import RegexpTokenizer
8 import tensorflow as tf
9 from tensorflow.keras.models import Sequential, load_model
10 from tensorflow.keras.layers import LSTM, Dense, Activation
11 from tensorflow.keras.optimizers import RMSprop
12 text = open("training_data.txt", 'r').read()
13 partial_text = text[:1000000]
14 tokenizer = RegexpTokenizer(r"\w+")
15 tokens = tokenizer.tokenize(partial_text)
16 unique_tokens = np.unique(tokens)
17 unique_token_index = {token: index for index, token in
18     enumerate(unique_tokens)}
19 n_words = 10
20 input_words = []
21 next_word = []
22 for i in range(len(tokens) - n_words):
23     input_words.append(tokens[i:i + n_words])
24     next_word.append(tokens[i + n_words])
25 X = np.zeros((len(input_words), n_words, len(unique_tokens)),
26     dtype=bool)

```



```

25 y = np.zeros((len(next_word), len(unique_tokens)), dtype=bool
    )
26 for i, words in enumerate(input_words):
27     for j, word in enumerate(words):
28         X[i, j, unique_token_index[word]] = 1
29     y[i, unique_token_index[next_word[i]]] = 1
30 model = Sequential()
31 model.add(LSTM(128, input_shape=(n_words, len(unique_tokens))
    , return_sequences=True))
32 model.add(LSTM(128))
33 model.add(Dense(len(unique_tokens)))
34 model.add(Activation("softmax"))
35 optimizer = RMSprop(learning_rate=0.01)
36 model.compile(loss="categorical_crossentropy", optimizer=
    optimizer, metrics=["accuracy"])
37 history = model.fit(X, y, batch_size=256, epochs=10, shuffle=
    True).history
38 model.save("textgenerator1.h5")
39 model = tf.keras.models.load_model("textgenerator1.h5")
40 history = model.fit(X, y, batch_size=256, epochs=10, shuffle=
    True).history
41 model.save("textgenerator2.h5")
42 model = tf.keras.models.load_model("textgenerator2.h5")
43 history = model.fit(X, y, batch_size=256, epochs=40, shuffle=
    True).history
44 model.save("textgenerator3.h5")
45 model = tf.keras.models.load_model("textgenerator3.h5")
46 history = model.fit(X, y, batch_size=256, epochs=40, shuffle=
    True).history
47 model.save("textgenerator4.h5")
48 def predict_next_word(input_text, n_best):
49     input_text = input_text.lower()
50     X = np.zeros((1, n_words, len(unique_tokens)))
51     for i, word in enumerate(input_text.split()):
52         X[0, i, unique_token_index[word]] = 1
53
54     predictions = model.predict(X)[0]
55     return np.argsort(predictions, -n_best)[-n_best:]
56 def generate_text(input_text, n_words, creativity=3):
57     word_sequence = input_text.split()
58     current = 0
59     for _ in range(n_words):
60         sub_sequence = " ".join(tokenizer.tokenize(" ".join(
            word_sequence).lower()))[current:current+n_words]
61         try:

```

```

62         choice = unique_tokens[random.choice(
        predict_next_word(sub_sequence, creativity))]
63     except:
64         choice = random.choice(unique_tokens)
65         word_sequence.append(choice)
66         current += 1
67     return " ".join(word_sequence)

```

Listing 8 – Réseau Recurrent

## Moteur de Jeu

Notre environnement de jeu, n'étant pas bien complexe, notre moteur de jeu ne le sera pas non plus. Pour ce faire, listons nos besoins :

1. Afficher des images en arrière-plan.
2. Afficher des images de personnages.
3. Afficher une boîte de dialogue.
4. Possiblement post-traiter notre rendu final.
5. Une fenêtre de saisie de texte.



## moteur de jeu

```
1  """
2  pygame model
3  """
4  from pygame.locals import *
5  import pygame
6  from PIL import Image
7  import tensorflow as tf
8  import random
9  import numpy as np
10 import os
11 import time
12
13
14 from tensorflow.keras.models import Sequential
15 from tensorflow.keras.optimizers import RMSprop
16 from tensorflow.keras.layers import Activation, Dense, LSTM
17 from keras.layers import Dropout
18 import all_sequensers_models as asqs
19
20
```

```

21 decode_sentence = asqs.EncodeAutomate({"\t" : " ", "\n" : " "
    }).sequence
22 encode_sentence = asqs.EncodeAutomate({" " : "\t"}).sequence
23
24
25 text = open("test.txt", 'rb').read().decode(encoding='utf-8')
    .lower()
26 #text = text[]
27
28
29 characters = sorted(set(text))
30
31
32 char_to_index = dict((c,i) for i, c in enumerate(characters))
33 index_to_char = dict((i,c) for i, c in enumerate(characters))
34
35
36 SEQ_LENGTH = 40
37 STEP_SIZE = 3
38 model = tf.keras.models.load_model("textgenerator.model")
39
40
41 def sample(preds, temperature=1.0):
42     preds = np.asarray(preds).astype('float64')
43     preds = np.log(preds) / temperature
44     exp_preds = np.exp(preds)
45     preds = exp_preds / np.sum(exp_preds)
46     probas = np.random.multinomial(1, preds, 1)
47     return np.argmax(probas)
48 def generate_text(length, temperature):
49     start_index = random.randint(0, len(text) - SEQ_LENGTH -
50     1)
51     generated = ""
52     sentence = text[start_index: start_index + SEQ_LENGTH]
53     generated += sentence
54     for i in range(length):
55         x_predictions = np.zeros((1, SEQ_LENGTH, len(
56         characters)))
57         for t, char in enumerate(sentence):
58             x_predictions[0, t, char_to_index[char]] = 1
59
60         predictions = model.predict(x_predictions, verbose=0)
61         next_index = sample(predictions,
62                             temperature)

```

```

61         next_character = index_to_char[next_index]
62
63         generated += next_character
64         sentence = sentence[1:] + next_character
65     return generated
66
67
68 def generate_text_aux(length, temperature, my_sent):
69     start_index = random.randint(0, len(text) - SEQ_LENGTH -
70     1)
71     generated = my_sent
72     sentence = text[start_index: start_index + SEQ_LENGTH]
73     generated += sentence
74     for i in range(length):
75         x_predictions = np.zeros((1, SEQ_LENGTH, len(
76         characters)))
77         for t, char in enumerate(sentence):
78             x_predictions[0, t, char_to_index[char]] = 1
79
80         predictions = model.predict(x_predictions, verbose=0)
81         next_index = sample(predictions,
82                             temperature)
83         next_character = index_to_char[next_index]
84
85         generated += next_character
86         sentence = sentence[1:] + next_character
87     return generated
88
89 def dialog_manager(historic : list):
90     old_text = encode_sentence("\n".join(historic))
91     new_text = (generate_text_aux(random.randint(0,200),0.8,
92     old_text))[len(old_text):]
93     return new_text
94
95
96 def src_manager(file_procc : [str], n : int):
97     """segmente notre liste en plusieurs listes plus petites
98     """
99     if len(file_procc) <= n:
100         return [file_procc]
101     else:
102         return [file_procc[:n]] + src_manager(file_procc[n:],
103         n)

```

```

100 def load_animation_images_sky():
101     images = []
102     path = "src\\sprite01\\my_type01_sprite"
103     for num in range(1,17):
104         image_path = path + str(num) + ".png"
105         img = pygame.image.load(image_path)
106         images.append(img)
107     return images
108 animations = {"sky" : load_animation_images_sky()}
109 class AnimationSky(pygame.sprite.Sprite):
110     def __init__(self, rect_x, rect_y):
111         super().__init__()
112         self.image = pygame.image.load("src\\sprite01\\
my_type01_sprite1.png")
113         self.rect = self.image.get_rect()
114         self.rect.x = rect_x
115         self.rect.y = rect_y
116         self.current_image = 1
117         self.animation = True
118         self.images = animations["sky"]
119     def animate(self, loop=True):
120         if self.animation:
121             self.current_image += 1
122             if self.current_image >= len(self.images):
123                 self.current_image = 1
124                 if not loop:
125                     self.animation = False
126             self.image = self.images[self.current_image]
127
128
129 class GraphicModel():
130     def __init__(self ,title : str ,bg_color="dark gray",fuls
= False ,size=(800,500), tick=6):
131         #pre initialize
132         pygame.init()
133         if fuls:
134             self.windows = pygame.display.set_mode(size,
pygame.FULLSCREEN)
135         else:
136             self.windows = pygame.display.set_mode(size)
137
138         #Initialize
139         self.meta = {"title" : title , "bg-color" : bg_color
, "tick" : tick}

```

```

141         pygame.display.set_caption(title)
142         self.timer = pygame.time.Clock()
143         self.font = pygame.font.Font("police.ttf",24)
144         self.snip = self.font.render("",True,"white")
145         self.run = True
146         self.last_messages = [generate_text(100,0.8)]
147         self.dialog_mode = False
148         self.images = {"bg" : pygame.image.load("src\img0004.
png"),
149                        "fg" : pygame.image.load("src\img0007.
png"),
150                        "db" : pygame.image.load("src\img0003.
png"),
151                        "char1" : pygame.image.load("src\
img0001.png"),
152                        "char2" : pygame.image.load("src\
img0006.png")}
153         self.counter = 0
154         self.done = False
155         self.speed = 1
156         self.user_text = ""
157         self.user_dialog = False
158         self.char = 1
159         self.sky_sprite = pygame.sprite.Group()
160         self.head = 0
161         self.display_dialog = src_manager(self.last_messages
[-1],54)
162         self.sky_sprite.add(AnimationSky(0,0))
163         def add_a_dialog(self, new_txt : str):
164             self.last_messages.append(new_txt)
165         def open_dialog_box(self):
166             if self.dialog_mode:
167                 self.windows.blit(self.images["db"],(0,400))
168                 self.windows.blit(self.images["char"+str(self.
char)],(0,272))
169                 message = self.display_dialog[self.head]
170                 if self.counter < self.speed * len(message):
171                     self.counter += 1
172                 elif self.counter >= self.speed*len(message):
173                     self.done = True
174                     self.snip = self.font.render(message[0:self.
counter//self.speed],True, "white")
175                     self.windows.blit(self.snip,(30,440))
176                     pass
177                 else:

```

```

178         self.counter = 0
179         self.done = False
180     pass
181     def input_text(self, event):
182         if event.type == pygame.KEYDOWN and self.user_dialog:
183             if event.key == pygame.K_BACKSPACE:
184                 self.user_text = self.user_text[:-1]
185             elif event.key == pygame.K_RETURN:
186                 self.last_messages.append(self.user_text)
187                 self.last_messages.append(dialog_manager(self
188 .last_messages))
189                 self.user_text = ""
190                 pygame.time.wait(1000)
191                 self.display_dialog = src_manager(self.
192 last_messages[-1], 54)
193                 self.user_dialog = False
194                 self.dialog_mode = True
195             else:
196                 self.user_text = self.user_text + event.
197 unicode
198     def self_open_dialog_box(self):
199         if self.user_dialog:
200             self.windows.blit(self.images["db"], (0, 400))
201             self.windows.blit(self.images["char"+str(self.
202 char)], (0, 272))
203             self.snip = self.font.render(decode_sentence(self
204 .user_text), True, "black")
205             self.windows.blit(self.snip, (30, 440))
206     pass
207     def display_pre_initialize(self):
208         rect = self.images["bg"].get_rect()
209         self.windows.blit(self.images["bg"], (0, 0))
210         for char in self.sky_sprite:
211             if char.animation:
212                 char.animate()
213             else:
214                 self.sky_sprite.remove(char)
215         self.sky_sprite.draw(self.windows)
216     pass
217     def display_initialize(self):
218         self.windows.blit(self.images["fg"], (0, 0))
219         self.open_dialog_box()
220         self.self_open_dialog_box()
221     pass
222     def display_post_initialize(self):

```



```

218         pygame.image.save(self.windows, "p_t.png")
219     pass
220     def game_loop(self):
221         while self.run:
222             self.timer.tick(self.meta["tick"])
223             self.display_pre_initialize()
224             self.display_initialize()
225             self.display_post_initialize()
226             pygame.display.flip()
227             pygame.display.update()
228             for event in pygame.event.get():
229                 if event.type == pygame.QUIT:
230                     self.run = False
231                 if event.type == pygame.KEYDOWN:
232                     if event.key == pygame.K_DOWN:
233                         if self.dialog_mode == False:
234                             self.dialog_mode = True
235                             self.user_dialog = False
236                     else:
237                         self.dialog_mode = False
238                         self.user_dialog = False
239                 if event.key == pygame.K_UP:
240                     if self.user_dialog == False:
241                         self.user_dialog = True
242                         self.dialog_mode = False
243                     else:
244                         self.user_dialog = False
245                         self.dialog_mode = False
246                 if event.key == pygame.K_a:
247                     if self.head == len(self.
display_dialog) - 1:
248                         self.head = 0
249                     else:
250                         self.head = self.head + 1
251                         self.counter = 0
252                         self.done = False
253                 if (event.key == pygame.K_LEFT) or (event
.key == pygame.K_RIGHT):
254                     if self.char == 1:
255                         self.char = 2
256                     else:
257                         self.char = 1
258                 if self.user_dialog or self.
dialog_mode:
259                     self.user_dialog = True

```

```

260             self.dialog_mode = False
261             if event.key == pygame.K_TAB:
262                 pygame.image.save(self.windows, "
screenshot.png")
263                 self.input_text(event)
264                 pygame.quit()
265                 pass
266
267
268 my_game = GraphicModel("tipe",)
269 my_game.game_loop()

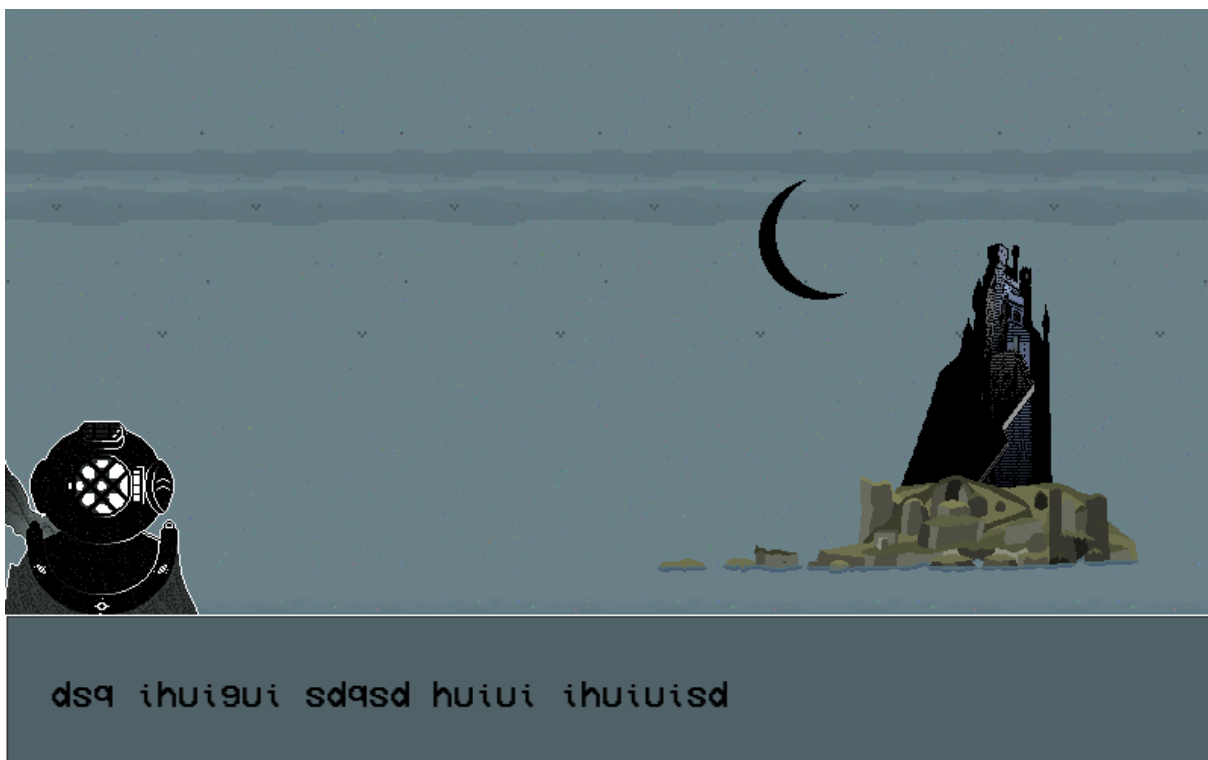
```

Listing 9 – Python Moteur de Jeu Implementation

Pour commencer affichons notre arriere plan.



Implementons un systeme d'affichage de texte (pour les dialogues).



Implementons un systeme de saisie de texte.



# Générateur de texte

## Chaine de markov

### implémentation

```
1 import markovify
2 import json
3 import os
4 fileIn = open(r"training_data.txt", "r")
5 content = fileIn.readlines()
6 fileIn.close()
7 training_data = "./n".join([a[:-4] for a in content])
8 text_model = markovify.Text(training_data)
9 print(text_model.make_sentence(tries=1000))
```

Listing 10 – Chaine de Markov Implementation

performances

## Réseau LSTM

implémentation

```
1 import tensorflow as tf
2 import random
3 import numpy as np
4 import os
5 import time
6
7
8 from tensorflow.keras.models import Sequential
9 from tensorflow.keras.optimizers import RMSprop
10 from tensorflow.keras.layers import Activation, Dense, LSTM
11 from keras.layers import Dropout
12
13
14
15 strategy = tf.distribute.MirroredStrategy()
16 print('Number of devices: {}'.format(strategy.
    num_replicas_in_sync))
17 text = (open("training_data.txt").read())
18 text = text[:1000000]
19 characters = sorted(list(set(text)))
20
21
22 char_to_index = dict((c,i) for i, c in enumerate(characters))
23 index_to_char = dict((i,c) for i, c in enumerate(characters))
24
25
26 SEQ_LENGTH = 40
27 STEP_SIZE = 3
28 sentences = []
29 next_characters = []
30 for i in range(0, len(text) - SEQ_LENGTH , STEP_SIZE):
31     sentences.append(text[i : i + SEQ_LENGTH])
32     next_characters.append(text[i + SEQ_LENGTH])
33 x = np.zeros((len(sentences), SEQ_LENGTH,
34               len(characters)), dtype=bool)
35 y = np.zeros((len(sentences),
36               len(characters)), dtype=bool)
37 for i, satz in enumerate(sentences):
```

```

38     for t, char in enumerate(satz):
39         x[i, t, char_to_index[char]] = 1
40         y[i, char_to_index[next_characters[i]]] = 1
41 model = Sequential()
42 model.add(LSTM(128, input_shape=(SEQ_LENGTH, len(characters))
43             ))
44 model.add(Dropout(0.2))
45 model.add(Activation("softmax"))
46 model.add(Dropout(0.2))
47 model.add(Dense(len(characters)))
48 model.compile(loss= "categorical_crossentropy" , optimizer=
49               RMSprop(learning_rate=0.01))
50 model.fit(x,y,batch_size=256,epochs=4)
51 model.save("textgenerator1.model")
52 model = tf.keras.models.load_model("textgenerator1.model")
53 model.fit(x,y,batch_size=256,epochs=10)
54 model.save("textgenerator1.model")
55 model = tf.keras.models.load_model("textgenerator1.model")
56 model.fit(x,y,batch_size=256,epochs=40)
57 model.save("textgenerator2.model")
58 model = tf.keras.models.load_model("textgenerator2.model")
59 model.fit(x,y,batch_size=256,epochs=20)
60 model.save("textgenerator3.model")
61 model = tf.keras.models.load_model("textgenerator0.model")
62 def sample(preds, temperature=1.0):
63     preds = np.asarray(preds).astype('float64')
64     preds = np.log(preds) / temperature
65     exp_preds = np.exp(preds)
66     preds = exp_preds / np.sum(exp_preds)
67     probas = np.random.multinomial(1, preds, 1)
68     return np.argmax(probas)
69 def generate_text(length, temperature):
70     start_index = random.randint(0, len(text) - SEQ_LENGTH -
71     1)
72     generated = ""
73     sentence = text[start_index: start_index + SEQ_LENGTH]
74     generated += sentence
75     for i in range(length):
76         x_predictions = np.zeros((1, SEQ_LENGTH, len(
77         characters)))
78         for t, char in enumerate(sentence):
79             x_predictions[0, t, char_to_index[char]] = 1
80
81         predictions = model.predict(x_predictions, verbose=0)
82         [0]

```

```

78         next_index = sample(predictions,
79                               temperature)
80         next_character = index_to_char[next_index]
81
82         generated += next_character
83         sentence = sentence[1:] + next_character
84     return generated
85 def generate_text_aux(length, temperature, my_sent):
86     start_index = random.randint(0, len(text) - SEQ_LENGTH -
87     1)
88     generated = my_sent
89     sentence = text[start_index: start_index + SEQ_LENGTH]
90     generated += sentence
91     for i in range(length):
92         x_predictions = np.zeros((1, SEQ_LENGTH, len(
93         characters)))
94         for t, char in enumerate(sentence):
95             x_predictions[0, t, char_to_index[char]] = 1
96
97         predictions = model.predict(x_predictions, verbose=0)
98         [0]
99         next_index = sample(predictions,
100                               temperature)
101         next_character = index_to_char[next_index]
102
103         generated += next_character
104         sentence = sentence[1:] + next_character
105     return generated

```

Listing 11 – RÅlseau RÅlccurent Implementation

performances

## Conclusions

## Bibliographie Commentée

Depuis la parution de Dungeon et dragon (DND) en 1974, nombreuses ont été les tentatives de portages sur ordinateur, du célèbre jeu de rôle sur table.

Ces bien heureuses tentatives, de par leur contraintes techniques, ont donné naissance aux mécaniques de RPG, très largement rependues dans l'industrie du jeu vidéo. Cependant, après tant d'années, il nous faut faire l'amère constat, de l'échec partiel de ces portages, car les grandes qualités du jeu de rôle, sont à l'origine de ces semi-réussites. Pour cause, le « CRPG » ( pour « computer role playing game » ), perd en flexibilité, dans ses tentatives de réduire le jeu de rôle à son aspect statistique.

C'est là : une contrainte que le jeu de rôle, partage avec le traitement naturel du langage. Car les faiblesses du jeu de rôle, sont également ses plus grands atouts :

- Son côté humain, collaboratif et communautaire, rend son organisation pénible et bien des fois hasardeuse.

- Les participants n'ont bien souvent, pas les mêmes disponibilités, voir volontés de s'y investir.

- Le jeu de certains joueurs, restreint bien souvent celui des autres.

- Ou encore des erreurs humaines, peuvent être commises lors de l'application des règles du jeu.

C'est pour pallier ces problèmes, tout en conservant leur côté avantageux, qu'il nous est venu à l'esprit de coupler un système de jeu (aussi rudimentaire soit-il) à un « LLM ».

Le soudain essor des techniques de « deep learning », est aller de pair avec un développement accru, du traitement du langage naturel (NLP).<sup>5</sup> Longtemps considéré comme une des limites, du « machine learning » ; en raison des difficultés à vectoriser : mots et autres morceaux de phrases (<sup>4</sup> la



tokenisation), le NLP a néanmoins connu de récentes avancées majeurs dans le domaine de la complétion de texte.<sup>4</sup>

Générer du texte en appliquant de la statistique à la linguistique, est une méthode conçue en 1948 par Claude Shannon.<sup>6</sup> Bien que, le model (basé sur des principes de probabilités Markoviennes<sup>6</sup>), produise la plupart du temps des phrases correcte au niveau de la syntaxe, cela se fait au détriment de la sémantique et du sens implicite de celles-ci. En 2017 est publier un article scientifique, nommé : « Attention Is All You Need »<sup>2</sup>, qui a l'aide de « transformers », parvient à générer des phrases dont la sémantique est correcte, tout en comprenant le contexte des phrases précédentes (en simulant le mécanisme d'oubli du cerveau humain). C'est sur cette base nouvelle, que seront battis les récents LLM tel que : « GPT-2, GPT-3, GPT-4, Bloom, Bert, e.t.c. » « NovelAI » et « AiDungeon », pour citer les exemples les plus connus d'utilisations de LLM dans le jeu de rôle, sont basés sur ces derniers.

---

1. [https://proceedings.neurips.cc/paper\\_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf)

2. <https://datascientest.com/introduction-au-nlp-natural-language-processing>

3. <https://larevueia.fr/introduction-aux-reseaux-de-neurones-transformers/>

4. <https://larevueia.fr/introduction-au-nlp-avec-python-les-ia-prennent-la-parole/>

5. [https://fr.wikipedia.org/wiki/Traitement\\_automatique\\_du\\_langage\\_naturel](https://fr.wikipedia.org/wiki/Traitement_automatique_du_langage_naturel)

6. <https://www.wikidata.fr-fr.nina.az/...>