

L'analyse du discours et l'intelligence artificielle pour réaliser une écriture inclusive : le projet E-MIMIC

Rachele Raus^{1,*}, Michela Tonti¹, Tania Cerquitelli², Luca Cagliero², Giuseppe Attanasio², Moreno La Quatra², et Salvatore Greco²

¹Département de Traduction et d'Interprétation, Université de Bologne, Italie

²Department of Control and Computer Engineering, École Polytechnique de Turin, Italie

Résumé. Cet article présente le projet E-MIMIC, une application qui vise à éliminer les préjugés et la non-inclusion dans les textes administratifs rédigés dans les pays européens, à commencer par ceux qui sont rédigés dans les langues romanes. Il présente une méthodologie conçue à partir de critères discursifs inspirés de l'analyse du discours française et utilisés pour étiqueter un corpus de documents institutionnels, qui sont utilisés pour l'apprentissage profond des réseaux neuronaux. Des architectures de modélisation profonde du langage sont exploitées pour identifier automatiquement les extraits de texte non inclusifs, suggérer des formes alternatives et produire des reformulations inclusives. Une évaluation préliminaire menée sur un ensemble de données de référence pour la langue italienne montre des résultats prometteurs, qui poussent à finaliser l'application et à la réaliser également pour d'autres langues, tel le français.

Abstract. Discourse Analysis and Artificial Intelligence to empower inclusive writing: the project E-MIMIC. This paper presents the E-MIMIC project, an application that aims to eliminate non-inclusive, prejudiced language forms in administrative texts written in European countries, starting with those written in Romance languages. It presents a methodology based on discourse criteria inspired by French discourse analysis and used to label a corpus of institutional documents, which are used for the deep learning of neural networks. Deep Language Modelling architectures are exploited to automatically identify non-inclusive text snippets, suggest alternative forms, and produce inclusive text rephrasing. A preliminary evaluation conducted on a benchmark dataset in Italian shows promising results and encourages us to finalise the application and to implement it also for other languages, such as French.

Introduction

La communication inclusive a connu récemment un essor important sur le plan international (Conseil de l'UE 2018), ce qui s'est également répercuté tant dans le monde

* rachele.raus@unibo.it

académique que dans celui des entreprises (Canfield *et alii* 2020). L'intérêt porté notamment à l'écriture inclusive, avant tout mais pas exclusivement à la féminisation, a déclenché des débats sur les plans nationaux¹ et a donné l'occasion de réfléchir sur l'importance des discours et des usages en relation à la création des relations sociales (Charaudeau 2021). Outre les problèmes des dissymétries lexicales (grammaticale et sémantique) entre les genres masculin et féminin qui ont été analysés et dénoncés dès les années 1970² par les premières recherches qui s'intéressaient au sexisme socio-discursif, les recherches actuelles incluent la dénonciation de stéréotypes liés à la circulation discursive des mots qui peuvent déboucher sur des préjugés et même justifier des formes de violence (Conseil de l'Europe 2014). Dans ce contexte, l'utilisation d'algorithmes d'intelligence artificielle³ dans l'industrie des langues risque d'empirer les choses, comme l'ont dénoncé plusieurs chercheur·e·s, en parlant de *discrimination algorithmique* (Bartoletti 2020) et de la manière dont l'utilisation de données préalables déjà problématiques, finit par accentuer les discriminations, notamment dans le traitement automatique des langues (voir, par exemple, Marzi 2021, Savoldi *et alii* 2021...). Les réseaux neuronaux, en effet, s'entraînent sur des corpus, et donc des discours, déjà biaisés en devenant des dispositifs capables d'alimenter ce qu'Eni Orlandi appelle la « mémoire métallique » (1996), c'est-à-dire la mémoire des « machines » qui contribue à naturaliser des discours et, dans ce cas, à enrainer des pratiques discriminatoires.

C'est justement pour intervenir sur ces questions que nous proposons un nouveau paradigme de travail entre les linguistes adoptant l'approche de l'analyse du discours dite à la française⁴ et le personnel expert en apprentissage profond de manière à proposer des outils linguistiques au service de l'écriture inclusive. Dans le cadre du projet *Empowering a Multilingual Inclusive Communication* (E-MIMIC) qui est financé par l'École Polytechnique de Turin et réalisé en collaboration avec l'Université de Bologne, nous sommes en train de réaliser une application à destination des administrations publiques pour le support à la rédaction inclusive. Ce projet est également soutenu par le Centre d'excellence Jean Monnet *Artificial Intelligence for European Integration* de l'Université de Turin⁵.

Cet article se donne l'objectif de présenter cette initiative afin de lancer un débat sur l'importance d'insérer des critères linguistiques et discursifs en amont de l'entraînement des algorithmes d'IA. Pour l'instant, l'application est en train d'être finalisée pour la langue italienne d'Italie⁶, nous comptons ensuite l'adapter pour la langue française de France.

À cette fin, nous allons tout d'abord présenter la notion d'écriture inclusive qui ne fait pas l'unanimité lors des pratiques de rédaction et qui dépend également des politiques linguistiques menées dans les pays concernés. Ensuite, nous allons encadrer l'ADF et notre approche par rapport à l'informatique, notamment à l'intelligence artificielle qui désormais est utilisée majoritairement dans l'industrie des langues, pour présenter les critères linguistico-discursifs préalablement choisis en relation à l'« apprentissage supervisé »⁷ des réseaux neuronaux. Enfin, nous présentons l'application et un test que nous avons mené sur celle-ci, tout en précisant que, comme pour l'instant nous n'avons finalisé que la partie d'étiquetage de l'application italienne, nous nous bornons à présenter dans cette étude les résultats de cette partie. Par ailleurs, il s'agit de résultats très prometteurs qui nous encouragent à finaliser l'application pour ensuite la localiser pour d'autres langues européennes à commencer par la langue française.

1 Quelle écriture inclusive ?

Les politiques linguistiques portant sur la non-discrimination dans le langage sont de longues dates en Europe et à l'international.

Pour ce qui est des langues romanes, notamment du français (et/ou des français des pays francophones) et de l'italien (d'Italie et de Suisse), les débats sur la question ont concerné avant tout la lutte contre le sexisme dans le langage (Houdebine-Gravaud 2001, Vecchiato 2004, Raus 2004), qui en France a débuté pendant les années 1970 et en Italie environ dix ans après. Il s'agissait d'intervenir sur les dissymétries lexicales, notamment d'accord grammatical et de stéréotypie sémantique, pour rendre l'écriture paritaire (voir Yaguello 1978[2006], Sabatini 1987). Ce type d'intervention a fini par poser une évidence qui a été longtemps débattue jusqu'à aujourd'hui, à savoir que l'usage des langues – le discours – est encore plus sexiste que la langue elle-même (Sabatini 1987, Charaudeau 2021).

L'introduction sur le plan international de la notion de « genre » à partir de 1995, lors de la 4^e conférence mondiale sur les femmes qui a été organisée à Pékin en 1995 par l'ONU, a produit des réflexions nouvelles sur la question. Les réticences françaises à l'utilisation du calque de l'anglais⁸ s'expliquent sans doute, et entre autres, par le fait que la catégorie nouvelle entraînait une approche différente des questions linguistiques, en privilégiant des formes d'écritures de « neutralisation » (UNESCO 1999), ce qui visait la déconstruction du binarisme prôné par la féminisation des approches non-sexistes et paritaires (Raus 2020). L'inclusion, bien que la question ne se posait pas encore en ces termes, ne concernait pas que les femmes mais également des identités non binaires, qui ne pouvaient pas se retrouver dans les dichotomies de la langue.

L'évolution du lexique à l'international (Raus, 2013 : 77-79) témoigne de la diffusion plus récente d'un langage sensible aux différences et qui essaie donc de marquer celles-ci, ce qui s'est concrètement réalisé par les formes de créativité syntaxique, grammaticale et autres qu'on a vues en France (de nouveaux pronoms⁹, l'utilisation de l'arobase ou de l'astérisque ou plus récemment du point médian), en Italie (l'utilisation des astérisques, du point médian, de la voyelle « u » à la place des désinences « -o/-a » ou, plus récemment, du schwa¹⁰) et ailleurs.

Ce n'est que dans les toutes dernières années qu'on a commencé à parler de « communication inclusive » et d'un langage inclusif. Si l'Union européenne entend par ce type de communication la volonté d'éliminer toute forme de discrimination à l'égard des genres, des ethnies, des religions, de l'âge, etc., il suffit d'aller voir dans Wikipedia pour se rendre compte du fait qu'en France la question est reformulée comme renvoyant à la seule égalité de genre, à savoir aux « moyens linguistiques visant à assurer une égalité de genre dans la langue » (Wikipedia, 2021¹¹) et d'aller faire un tour sur Internet pour se rendre compte qu'en Italie, non seulement le débat public mais même les autorités comme Treccani ou La Crusca¹² reviennent sur la question de manière similaire. Comme au tout début, les politiques linguistiques portant sur le genre finissaient par ne concerner que les femmes¹³, aujourd'hui le langage inclusif est restreint aux questions de genre, voire, parfois, à la seule féminisation linguistique.

Par ce dernier constat, il ne s'agit pas d'aplatir des questions qui sont de loin plus complexes et mériteraient d'être étudiées de plus près par rapport aux contextes socio-politiques des pays concernés, mais de souligner que le débat actuel sur l'écriture inclusive est un débat largement positionné sur le plan de l'interdiscours¹⁴ et a perdu de vue le trait linguistique fondamental qui est présent dans le syntagme, à savoir que l'inclusion devrait concerner tout le monde.

Sur le plan européen, parler plus largement de communication inclusive veut dire avant tout pratiquer la non-exclusion (Conseil de l'UE 2018), ce qui se localise dans les versions italienne et française du document de 2018 par la promotion de politiques linguistiques qui ne soient pas discriminatoires à l'égard des femmes (par exemple, en évitant des formules comme « les droits de l'Homme »), des identités non binaires (par exemple, par la promotion de noms collectifs, de mots épécènes, de mots génériques – comme individu, personne ou autre – ou des reformulations neutralisant l'appartenance sexuelle, etc.), mais également à l'égard des soi-disant minorités, comme les personnes handicapées, plus âgées,

etc. Cette approche est redevable d'un discours européen qui s'est de plus en plus focalisé sur la lutte contre les violences et les discriminations de manière large (voir notamment la Convention d'Istanbul du Conseil de l'Europe de 2014).

Parler alors d'écriture inclusive ne peut pas se faire sans soulever des débats sur ce que c'est une communication inclusive plus généralement ou sur ce qui fait qu'un langage (ou plutôt un discours) soit tel. En effet, l'existence des déjà-là encombrants sur ce sujet fait que désormais le syntagme même « écriture inclusive » s'actualise de manières multiformes et souvent contradictoires en France comme en Italie.

Pour ces raisons, il nous faut à notre tour nous positionner et prendre des décisions sur ce sujet, ce que nous faisons en considération des critères suivants :

1) comme notre intervention se limite à la reformulation des textes administratifs, il nous a fallu tenir compte des politiques linguistiques des institutions (les universités mais aussi les institutions qui, comme la DGLFLF en France ou l'Accademia della Crusca en Italie, jouent un rôle plus ou moins contraignant dans les politiques linguistiques du pays concerné) et également des politiques prônées par l'Union européenne (notamment, en tenant compte du texte de 2018 sur la communication inclusive) ;

2) nous avons également donné la possibilité aux personnes qui utilisent l'application de pouvoir choisir parmi plusieurs possibilités de reformulation de la tournure ou du mot non-inclusifs en privilégiant celle qui est perçue comme la plus appropriée par rapport au contexte et à la cohérence du document final. De cette manière, l'application signale à la personne qui s'en sert quel est le segment ou le mot inapproprié, en contribuant à la sensibilisation vis-à-vis de la discrimination et non-inclusion dans le langage, et suggère des reformulations possibles, la personne pouvant donc choisir parmi les solutions proposées.

Par rapport à ces dernières, par exemple, l'application peut proposer des reformulations binaires, pour rendre visible les femmes, ou/et non-binaires, comme on le verra plus en détail dans la section 2.1. Cela dit, la première des formes affichées dans la liste des reformulations possibles tient compte d'une part de la politique spécifique du pays concerné et de l'autre des statistiques d'usage des usagers, de manière à essayer de normaliser les usages au cas où on choisirait cette forme spécifique.

Avant de rentrer dans le détail dans l'application, il nous faut présenter les critères linguistiques et discursifs que nous avons privilégiés pour entraîner les algorithmes d'IA et en quoi l'analyse du discours dite à la française (dorénavant ADF) s'est révélée fondamentale à cette fin.

2 ADF et informatique : un mariage de longue date

L'ADF a toujours eu un penchant privilégié pour l'informatique. D'ailleurs, dès le début, l'ADF a servi pour traiter automatiquement les langues, et la revue *Mots*, ainsi que le laboratoire de lexicométrie de Saint-Cloud, dirigé pendant longtemps par Maurice Tournier, sont nés justement de l'intérêt porté à la lexicométrie et/ou aux formations discursives¹⁵ qui peuvent être traitées de manière informatique¹⁶. Cette tendance à lier informatique et ADF a survécu jusqu'à aujourd'hui, nous semble-t-il, de deux manières, directe et indirecte.

Pour ce qui est du premier cas, nous signalons, entre autres, les recherches que Damon Mayaffre¹⁷ a menées en lexicométrie et en logométrie, en privilégiant récemment l'utilisation des algorithmes d'IA (Mayaffre 2021) ; à l'égard du deuxième cas, nous renvoyons aux analystes qui ont voulu analyser les conditions de production des discours médiés par des dispositifs numériques plus généralement¹⁸ et qui ont donc enquêté sur la manière dont ces dispositifs contribuent à façonner le discours, par exemple, par la présence

de routines et de formations de plus en plus figées de la matérialité de la langue dans certains genres de discours (Née 2017), ou par la manière dont la reprise interdiscursive des mots et des discours par ces dispositifs permettent de parler de « mémoire métallique » (Orlandi 1996).

Ce sont justement ces dernières recherches qui nous ont aidé à mieux focaliser les critères qu'il fallait privilégier pour permettre aux algorithmes d'apprendre de manière optimale les critères d'écriture inclusive que nous avons présentés auparavant.

2.1 Des critères linguistiques et discursifs pour l'apprentissage de l'IA

Si nous parlons de critères linguistiques et discursifs en amont de l'apprentissage de l'IA, c'est pour souligner que nous ne pouvons plus opposer la langue au discours, puisque les réseaux d'apprentissage profond apprennent des formes récurrentes du dit (des formations discursives ou des routines, dirait l'analyste du discours) mais qu'elle arrive à « abstraire » les formes linguistiques derrière les manifestations discursives de la matérialité des discours qui forment le corpus d'où elle apprend. Il faut donc imaginer un lien de va-et-vient constant entre langue et discours, dans la mesure où l'intelligence mathématique du dispositif des réseaux artificiels (la binarité 0-1) tend à réduire les possibilités du discours à des formes linguistiques structurées préformatées. Dans les systèmes traditionnels de traduction automatique sur des bases statistiques, avant l'utilisation de l'IA, ce constat devient une évidence. Le dispositif, ou la « machine » comme on l'appelle plus généralement¹⁹, propose des segments pré-traduits en exploitant les mémoires de traduction qu'elle traite à partir de critères statistiques. Ce système produit l'aplatissement de la matérialité discursive et privilégie la reprise interlinguistique, à savoir une reprise interdiscursive décontextualisée (un *pattern*). C'est l'usager (traductrice ou traducteur) qui intervient sur le segment proposé par la machine pour l'éditer et assurer la bonne qualité de la traduction. Lors de l'apprentissage automatique, le dispositif neuronal décontextualise mais abstrait la norme linguistique pour (re)produire ensuite du nouveau discours. Dans ce cas aussi, l'entraînement reste fondé sur des critères statistiques. Par conséquent, il est avant tout fondamental d'intervenir en amont sur l'apprentissage pour éviter que le dispositif apprenne de manière erronée, par exemple, en abstrayant des règles grammaticales erronées. C'est, par exemple le cas de l'élimination de la forme féminine des noms et des adjectifs dans les langues romanes et cela pour deux raisons essentiellement :

1) le fait que souvent les algorithmes s'entraînent sur des corpus internationaux, notamment des organisations internationales (dorénavant O.I.), qui privilégient l'utilisation du masculin « neutre ». N'oublions pas, en effet, que les O.I. ont tendance à rédiger en priorité des documents législatifs, ou politiques mais finalisés à la rédaction de textes législatifs, où le masculin est normalement privilégié. C'est par exemple, le cas du traducteur automatique *DeepL*, qui utilise et exploite les corpus multilingues du concordancier *Linguee*²⁰.

Voici, en tableau 1, les résultats affichés le 27 novembre 2021 pour la traduction FR-IT et IT-FR d'une phrase simple :

Tableau 1 : Résultat de traduction (FR-IT et IT-FR) de l'adjectif « fière » dans *DeepL*²¹.

Traduction IT-FR	Traduction FR-IT
IT : Sono fiera di essere qui [= Je suis fière d'être ici]	FR : Je suis fier d'être ici <i>Traduction alternative</i> Je suis fière d'être ici

FR : Je suis fière d'être ici	IT : Sono orgoglioso di essere qui [= Je suis orgueilleux d'être ici] <i>Traduction alternative</i> Sono fiero di essere qui [=Je suis fier d'être ici] Sono orgogliosa di essere qui [=Je suis fière d'être ici]
-------------------------------	--

La forme féminine apparaît en dernier comme alternative. Nous faisons également remarquer que, jusqu'à il y a quelque temps, elle était absente.

2) Le dispositif utilise l'anglais comme langue pivot et, comme souvent l'anglais utilise des mots épiciènes pour les acteurs ou ne présente pas de formes binaires des adjectifs, la retraduction vers la langue d'arrivée finit par privilégier le masculin, faute de pouvoir attribuer un genre précis (Sono fiero [+FEM]...> *I'm proud* [Ø] > Je suis fier [+MAS]).

À ces dissymétries grammaticales s'accompagnent souvent des dissymétries sémantiques liées aux stéréotypes véhiculés par l'interdiscours caractérisant les corpus d'apprentissage des réseaux neuronaux, comme l'ont montré récemment plusieurs études (Marzi 2021, Savoldi *et alii* 2021, etc.).

Il est clair que l'un des problèmes caractérisant ces dispositifs est l'utilisation de la langue anglaise comme langue véhiculaire²². Cela dit, ce qui nous intéresse est d'introduire des critères linguistiques et discursifs qui permettent d'intervenir en amont de l'apprentissage pour éliminer ces dysfonctionnements.

Pour ce faire, nous avons tout d'abord créé un corpus de documents²³ produits par l'administration nationale du pays concerné, donc en italien dans le cas ci-traité, et les avons étiquetés selon des critères linguistico-discursifs. Le corpus choisi est homogène (discours administratif, plusieurs genres à son intérieur), ce qui n'a pourtant pas exclu la prise en compte des formes d'hétérogénéité constitutive et/ou montrée (Authier 1984), notamment les interférences diaphasiques (Maingueneau 1991 : 143), lors de l'étiquetage des autres discours qui traversent les documents choisis (voir *infra*). Ce corpus permet à la machine de se « spécialiser », en apprenant les routines et, plus généralement, les caractéristiques discursives définissant le type et les genres du discours administratif qui feront ensuite l'objet des reformulations inclusives. C'est justement la compétence de l'analyste du discours qui, dès ce stade, permet de déterminer les genres et les interférences diaphasiques dont on doit tenir compte ou les parties du discours qu'on doit étiqueter pour que la machine puisse ensuite détecter le problème de non-inclusion. Pour mieux comprendre ce point fondamental, nous donnons l'exemple d'un problème qui se vérifie souvent dans des logiciels de traduction automatique s'appuyant sur des réseaux d'apprentissage profond : si on fait traduire à des traducteurs automatiques comme *Google*, *DeepL*, *Reverso* la phrase italienne « *la capitana della squadra ha visto l'infermiera, che è uscita dall'ospedale* » (la capitaine de l'équipe a vu l'infirmière, qui est sortie de l'hôpital) en français, on obtient « le capitaine de l'équipe a vu l'infirmière, qui est sortie de l'hôpital ». Vice-versa, si on fait traduire « *la capitana della squadra ha visto l'infermiere, che è uscito dall'ospedale* » (la capitaine de l'équipe a vu l'infirmier, qui est sorti de l'hôpital), la phrase est traduite toujours par « le capitaine de l'équipe a vu l'infirmière, qui est sortie de l'hôpital »²⁴. Ces exemples, tirés de cas de traduction, permettent de revenir sur les discours monolingues de manière avisée par rapport aux questions que nous traitons ici. En français, l'omission de l'accord au féminin de l'article « le » pour capitaine dans le premier cas, alors que la forme féminine d'infirmière est traduite de manière correcte, ainsi que, à l'inverse, le non-accord du masculin pour l'infirmier du second exemple, montrent qu'il faut donner à la machine les moyens d'appréhender le discours. En effet, sur le plan

de la langue, le problème ne se poserait point car les formes respectivement au masculin au féminin de l'article et de ces noms de professions existent bel et bien. Il s'agit plutôt d'un problème interdiscursif, les mots au féminin circulant plus fréquemment dans certains discours que dans d'autres, tout comme il arrive pour leurs homologues masculins. Ce type de problème, il est vrai, pourrait être résolu par l'implémentation de dictionnaires et de sources lexicographiques existantes et qui accordent les noms au féminin, comme l'a fait le traducteur automatique *Systran*, mais cela ne résout pas le problème de la traduction d'infirmier par infirmière²⁵ et, plus généralement, résout seulement une partie des problèmes dont il est question ici. Par exemple, cela ne permet pas forcément de prendre en compte la question des stéréotypes qui peuvent se déclencher en discours. En effet, les dictionnaires peuvent ne pas suffire à trouver les mots biaisés, lorsque c'est l'utilisation discursive qui les rend tels, et il faut plutôt utiliser des systèmes d'exploration des données pour les détecter. C'est, par exemple, le cas du site d'une collectivité locale italienne qui est présent dans notre corpus, où l'invitation à se rendre à la journée porte ouverte d'une école primaire concerne « les mères, les pères, les grands-parents et toute la famille » de l'élève, frôlant ainsi le danger de ne pas inclure les enfants orphelins ou de véhiculer l'image stéréotypée d'un certain type de famille (« mères et pères ») à travers les implicites transmis en et par le discours.

Le fait d'utiliser l'approche discursive en amont permet de tenir compte plus généralement de toutes ces questions à la fois, en focalisant l'attention sur les conditions de production des discours biaisés, notamment l'interdiscours, et sur les parties du discours qui peuvent devenir problématiques, comme les formations discursives qui explicitent le sujet humain normalement au masculin, surtout si elles deviennent fréquentes d'un point de vue statistique, ce qui empêche la machine d'apprendre des usages inclusifs.

Après avoir souligné en quoi l'apport de l'analyste du discours devient fondamental dès le choix des corpus, il nous faut maintenant préciser que nous n'avons finalisé que la partie d'étiquetage de l'application et que donc c'est par rapport à cette partie que nous avons pu tester les critères linguistiques et discursifs que nous allons présenter tout de suite.

Ces critères ont tout d'abord été implémentés en permettant au dispositif de différencier les variantes diatopiques des langues concernées, par exemple, pour l'italien, l'italien d'Italie, de Suisse et de San Marino. Chaque mot circulant dans des interdiscours différents dans les différents pays, il nous a semblé bon de tenir compte du fait que non seulement les politiques linguistiques des pays concernés peuvent varier (par exemple, la tendance de la Suisse à privilégier des politiques « paritaires » dans le langage par la visibilité des formes féminines²⁶) mais aussi et surtout que les corpus administratifs à disposition dans ces deux pays seraient forcément différents en raison des productions discursives différentes (dispositifs de recommandation différents, discours procéduraux différents, interdiscours national différent...) et de l'interdiscours caractérisant les discours des administrations italiennes et suisses, par exemples.

Un deuxième critère a été celui de tenir compte du support web ou non web du document, en raison, là encore, de la nécessité de prendre en considération le type de dispositifs numérique, notamment du dispositif lié à la production des pages dans le réseau et qui conditionne la matérialité discursive (Paveau 2017). Un discours produit directement pour être publié dans une page Internet présente des conditions de production matérielles et discursives²⁷ spécifiques et différentes par rapport à un discours produit sur un support numérique traditionnel (par exemple, un fichier .docx)²⁸. Puis, il a fallu tenir compte du type et des genres du discours caractérisant les corpus utilisés. Bien que nous ayons privilégié le discours administratif, il a fallu considérer que ce dernier peut être traversé par d'autres discours faisant interférence par rapport au premier, par exemple, le discours de type législatif, qui est souvent utilisé à l'intérieur du discours de l'administration, avant tout lors des citations interdiscursives. Il est alors possible d'étiqueter la phrase, qui est l'unité

de base de l'application comme nous verrons ensuite, selon le type de discours, et non selon le type de texte plus généralement qui, lui, reste administratif ou législatif.

Nous avons également tenu compte des genres discursifs (par exemple, avis, code de conduite, rapport, règlement...), en raison du fait que chaque genre présente à son intérieur des routines discursives spécifiques.

Par rapport aux étiquettes à utiliser, nous avons tenu compte des éléments qui étaient inclusifs, – dans ce cas, la phrase a été étiquetée comme « inclusive » – ou pas. Dans ce dernier cas, la phrase est étiquetée comme non-inclusive et la personne qui prend en charge l'étiquetage du corpus a la possibilité d'insérer une ou plusieurs reformulations, en donnant la priorité aux formes recommandées par les choix politiques nationaux concernés.

Les étiquettes concernant les mots ou les segments non inclusifs dans la phrase sont les suivantes : « titre », « citation », « appellation qui peut être reformulée », « appellation qui ne peut pas être reformulée », « syntagme », « liste noire ». Ces étiquettes tiennent compte de plusieurs éléments liés à la fois à la langue et au discours.

Par titre, il s'agit des titres des documents qui sont signalés comme tels parce qu'ils peuvent respecter des contraintes liées à des critères rédactionnels spécifiques ou à la reprise interdiscursive.

L'étiquette « citation » concerne tous les segments venant d'autres sources qui sont repris de manière directe et qu'il faut donc laisser tels quels. Cette étiquette permettra au dispositif de ne pas corriger, par exemple, le syntagme venant du discours européen « *notte dei ricercatori* » (= nuits des chercheurs) s'il est présent dans la citation d'un document du Ministère italien de la recherche, mais il devra en indiquer la non-inclusion et suggérer la reformulation inclusive (par exemple, « *notte della ricerca* », c'est-à-dire la nuit de la recherche, ou d'autres reformulations possibles²⁹) quand il est utilisé dans une phrase du document qui n'est pas une citation.

Les appellations renvoient aux acteurs et aux parties du discours qui demandent une reformulation inclusive, en considérant le problème du masculin « neutre » qui normalement est utilisé dans le discours administratif italien pour renvoyer à toute personne. Il faut également que le dispositif arrive à bien décliner le féminin au cas où l'acteur spécifique serait une femme, ce qui nous a obligée à introduire parallèlement l'étiquette « appellation qui ne peut pas être reformulée » pour tenir compte des cas qui demandent au dispositif de ne rien toucher à des appellations correctes comme, par exemple, la phrase qui se compose de l'appellation suivie par le nom féminin (« *la ministra Maria Rossi* » = la ministre Maria Rossi). En effet, comme le dispositif apprend avant tout à partir d'un critère statistique, il pourrait privilégier des reformulations différentes, induites par l'apprentissage sur les autres données.

L'étiquette « syntagme » permet de signaler des expressions figées qui pourraient être plus ou moins inclusives. Par exemple, « *persona diversamente abile* » (= personne qui a des capacités différentes pour désigner les personnes handicapées) est un syntagme qui n'est ni inclusif ni clair et qui a été donc contesté³⁰ en raison du fait que n'importe quelle personne a des capacités spécifiques et différentes des autres. On peut donc le signaler comme non inclusif et proposer des reformulations possibles (par exemple, « *persona con disabilità* » comme recommandé par le Conseil de l'UE 2018).

Enfin, l'étiquette « liste noire » permet de tenir compte des stéréotypes que certains mots pourraient déclencher dans certains discours. Par exemple, l'expression « *i partner istituzionali dovrebbero...* » (= les partenaires institutionnels devraient..., nous signalons en gras la non-inclusion de l'article masculin italien « i » au pluriel) qui au féminin donnerait « *le partner istituzionali dovrebbero...* » mais qui déclencherait un stéréotype de genre en italien.

Lors des reformulations, nous avons donné la priorité aux reformulations inclusives par l'utilisation de mots épécènes, de pronoms indéfinis, de noms collectifs, de reformulations sans agent obtenues par métonymie³¹, de mots génériques (personne, individu...). Nous

avons également omis des segments non inclusifs dans les cas qui le permettaient, c'est-à-dire quand le sens de la phrase demeurait clair et compréhensible.

La forme binaire H/F a été insérée de manière résiduelle dans la langue italienne d'Italie, comme formulation non prioritaire, mais nous avons quand même permis que le dispositif apprenne l'accord au féminin dans le cas des appellations suivies par un nom féminin. Cela n'élimine pas la possibilité de laisser le segment au masculin au cas où la personne concernée privilégierait cette forme pour des raisons de dysphorie de genre, mais exclut la possibilité d'utiliser des formes d'écriture alternative de type expérimental. En tous les cas, ce système permet de signaler la présence d'un problème d'exclusion, ce qui plus généralement rend le dispositif très utile au développement d'une conscience métalinguistique à des fins de formation sur ces questions.

3 Présentation du projet E-MIMIC

L'application E-MIMIC se fonde sur l'idée qu'il est possible de traiter la reformulation des données non-inclusives du discours administratif italien comme s'il s'agissait de traduction intralinguistique. Il faut donc donner à la « machine », c'est-à-dire aux réseaux d'apprentissage profond, la possibilité d'apprendre à reconnaître le segment non inclusif de manière à en proposer des « traductions » inclusives. Ce système permettrait d'utiliser l'application non seulement à l'intérieur des administrations mais également à l'extérieur, par exemple pour améliorer et peaufiner de manière inclusive des documents traduits en italien et qui présentent des problèmes d'inclusion, comme il arrive pour la majorité des traducteurs automatiques actuels qui sont disponibles de manière gratuite sur Internet.

L'apprentissage profond d'E-MIMIC se fonde sur une architecture de transformateurs (Vaswani *et alii*), qui exploitent un mécanisme d'attention. Les transformateurs sont basés sur une architecture codeur-décodeur, qui produit une représentation vectorielle mathématique des séquences d'entrée, puis exploite la version codée pour générer de nouvelles séquences. Par l'utilisation d'une approche supervisée, cette architecture peut s'entraîner sur des données. Pour l'application E-MIMIC, nous avons spécialisé un modèle pré-entraîné afin de l'adapter à la réalisation d'une nouvelle tâche. Les connaissances générales acquises pendant la phase de préformation sont exploitées pendant l'étape de spécialisation en utilisant beaucoup moins de données étiquetées. Les réseaux de l'application se sont donc d'abord pré-entraînés de manière auto-supervisée sur l'italien à partir d'une grande collection générique de données non étiquetées. Pour ce faire, nous avons utilisé un modèle BERT (acronyme anglais de *Bidirectional Encoder Representations from Transformers* ; voir Devlin *et alii* 2019) pré-entraîné par une équipe de la Bavarian State Library³² à partir des données de Wikipedia 2019 et d'une collection de textes tirés du corpus OPUS (Tiedmann 2012)³³. Ces données sont divisées en phrases de maximum 512 mots par une librairie NLTK³⁴.

Une fois pré-entraînés sur ce corpus, les réseaux d'E-MIMIC se spécialisent de manière supervisée sur un ensemble plus restreint de documents administratifs étiquetés à partir des critères précédemment décrits. E-MIMIC utilise donc l'apprentissage par transfert : le modèle spécialisé est basé sur des modèles linguistiques déjà formés. Ces derniers permettent d'encoder les connaissances lexicales de la langue et sert donc de point de départ solide.

Tout d'abord, une étape de pré-entraînement nous permet de spécialiser le modèle linguistique préexistant à la langue administrative. Après cela, le modèle résultant sera capable de comprendre les associations de mots spécifiques au domaine et pourra donc être utilisé pour générer des phrases cohérentes. Les annotations supplémentaires recueillies pendant la phase d'étiquetage peuvent être utilisées pour la reconnaissance des entités. Ce faisant, le modèle apprend à contextualiser en fonction du type de mots traités. Les modèles

peuvent également être spécialisés par le biais de la classification du contenu, par exemple en prédisant si une phrase contient un contenu juridique, administratif, technique ou informatif. Comme pour la contextualisation basée sur les entités, les modèles seraient capables d'apprendre les facettes sémantiques adaptées à chaque domaine d'application.

Pour la réalisation de l'application finale, nous envisageons une dernière phase de spécialisation qui permettra à la machine de générer les segments inclusifs. Dans cette phase, le modèle apprendra les caractéristiques syntaxiques et la sémantique du langage inclusif, sur la base des annotations recueillies. Pour l'instant, cette phase n'a pas été encore réalisée.

Le corpus des documents administratifs que nous avons étiquetés se constitue de 68 textes pour un nombre total de 3661 phrases. Ces documents ont été téléchargés à partir des sites de 12 Universités italiennes³⁵ et de 5 collectivités locales³⁶. Ces documents, qui peuvent être de plusieurs genres discursifs³⁷, ont été divisés en phrases par un logiciel³⁸ qui a été réalisé exprès en Python avant de les étiqueter.

En effet, nous avons choisi comme unité de base la phrase, en tant que segment terminant par un point final. À ce sujet, bien que d'autres choix soient également possibles (par exemple, segmentation en paragraphes ou en mots), nous avons tenu compte des critères suivants :

- 1) ce n'est pas forcément le mot en soi qui peut être non-inclusif. Comme nous l'avons dit, c'est souvent l'usage de la langue en contexte qui produit la non-inclusion, bref le discours ;
- 2) les phrases représentent l'unité de base la plus courte qui contient des segments non inclusifs ;
- 3) les phrases peuvent aisément être regroupées en paragraphes, ce qui permet ensuite de prendre éventuellement en compte un contexte plus large (Yasunaga *et alii* 2019).

4 Résultats des premiers tests

Pour tester l'apprentissage des algorithmes, nous avons utilisé comme ensemble de données de référence des phrases générées synthétiquement en langue italienne, en utilisant une procédure de remplissage de modèles. Plus précisément, nous avons formulé un ensemble de modèles consistant en une phrase avec une partie vide à remplir qui indique si la phrase est inclusive ou pas. Le test a été effectué en soumettant 112 phrases à la machine à partir de 4 modèles différents.

Dans le tableau 2, nous montrons un exemple de ces modèles. Nous avons collecté un corpus parallèle de graines utilisées pour remplir les parties en blanc des modèles. Chaque graine comprend deux phrases, l'une qui répond aux critères linguistiques de l'écriture inclusive (I) et l'autre qui n'y répond pas (NI). Par conséquent, le processus de remplissage d'un modèle avec une graine aléatoire génère deux phrases, inclusive et non inclusive.

Tableau 2 : Exemple de modèle pour tester l'application E-MIMIC.

Modèle
Occorre richiedere la firma [blank] [= Il faut demander la signature (blanc)]
Exemples de synthèse
NI = Occorre richiedere la firma agli interessati [= Il faut demander la signature aux intéressés]
I = Occorre richiedere la firma alle persone interessate

[= Il faut demander la signature aux personnes intéressées]

Lors du test, le modèle spécialisé a atteint une exactitude³⁹ de 95,53% pour la reconnaissance des phrases inclusives ou non inclusives sur les données testées. Des extraits de ces résultats sont transcrits dans le tableau 3.

Tableau 3 : Des exemples des résultats des prévisions d'E-MIMIC.

Phrase	Inclusif (I)/Non inclusif (NI)	Prévision d'E-MIMIC
I candidati sono invitati a consultare il bando prima di presentare domanda [= les candidats sont invités à lire l'avis avant d'envoyer leur candidature]	NI	NI
Lo studente è invitato a prender visione della sede dei colloqui [= L'étudiant est invité à vérifier le lieu de l'entretien]	NI	NI
Sono stati analizzati i riscontri dei docenti [= Le retour d'information des enseignants a été analysé]	NI	NI
Si invita a consultare il bando prima di presentare domanda [= Il est recommandé de lire l'avis avant de candidater]	I	I
Si invita a prendere visione della sede dei colloqui [= Il est recommandé de vérifier le lieu de l'entretien]	I	I
Sono stati analizzati i riscontri del personale docente [= Le retour d'information du personnel enseignant a été analysé]	I	I

Rappelons que les données de test sont des phrases inclusives et non inclusives que le modèle n'a jamais vues lors de son apprentissage. Par conséquent, bien que le test ne permette pas encore de valider un par un les différents critères discursifs prévus préalablement, le résultat obtenu démontre que la machine a pu généraliser et apprendre la structure du langage inclusif.

Conclusions

Le test présenté montre que l'utilisation de critères discursifs en amont de l'apprentissage profond de la machine permet des performances prometteuses et démontre qu'un autre paradigme est possible quand on s'intéresse à l'utilisation de l'intelligence artificielle dans l'industrie des langues. Ce paradigme se veut novateur de deux manières différentes : d'une part, il permet de modéliser des architectures d'apprentissage profond à partir de critères qui ne tiennent pas compte de l'anglais mais des langues romanes, ce qui est en contenance avec l'utilisation majoritaire de l'anglais dans les technologies linguistiques

(Vetere en préparation) et qui donne la possibilité de préserver le multilinguisme ; de l'autre, il renverse la tendance actuelle à considérer l'intervention de l'expert-e en linguistique en aval de l'apprentissage profond en tant que personne qui révise le texte produit par la machine. Au contraire, l'analyste du discours doit travailler à côté de la personne experte en informatique pour choisir les critères d'étiquetage dans la langue concernée, ce qui permet à la machine d'apprendre de manière supervisée en évitant les tournures ou les mots utilisés de manière non inclusive et/ou qui peuvent déclencher des préjugés en discours. Cela nous semble d'autant plus important que ces critères ne sont pas purement linguistiques mais justement discursifs, tenant donc compte des usages non-inclusifs à éviter. Cette piste de travail répond non seulement à la demande que les personnes expertes en apprentissage profond posent relativement à la nécessité de bien choisir les données spécialisées pour l'apprentissage de la machine afin d'en améliorer les performances (voir, entre autres, Langlais en préparation), mais elle répond également au besoin de superviser l'apprentissage de la machine par des critères qui puissent l'optimiser de manière à obtenir des résultats de qualité.

Enfin, les résultats obtenus nous encouragent à mener à bien la finalisation de l'application pour la localiser ensuite pour d'autres langues, — à commencer par la langue française de France —, ce qui donnera l'opportunité de montrer comment l'ADF peut et pourra jouer un rôle fondamental à l'avenir, en contribuant à l'élaboration d'une intelligence artificielle respectueuse des valeurs humaines.

Références bibliographiques

- Abbou, J., Arnold, A., Candea, M., Marignier, N. (2018). Qui a peur de l'écriture inclusive ? Entre délire eschatologique et peur d'émasculatation. *Semen*, 44. <https://journals.openedition.org/semen/10800>
- Authier, J. (1984). Hétérogénéité(s) énonciative(s). *Langages*, 73, 98-111.
- Bartoletti, I. (2021). *An Artificial Revolution. On Power, Politics and AI*. Édimbourg : Indigo.
- Charaudeau, P. (2021). *La langue n'est pas sexiste. D'une intelligence du discours de féminisation*. Lormont : Le bord de l'eau.
- Charaudeau, P., Maingueneau, D. (2002). *Dictionnaire d'analyse du discours*. Paris : Éditions du Seuil.
- Conseil de l'Europe (2014). *Convention d'Istanbul*. <https://rm.coe.int/1680084840>
- Conseil de l'Union européenne (2018). *Communication inclusive*. Union européenne. https://www.consilium.europa.eu/media/35450/fr_brochure-inclusive-communication-in-the-gsc.pdf
- Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Vol.1 (Long and Short Papers), 4171–4186.
- Greco, L. (2013). Langage et pratiques « transgenres ». *Bulletin de l'Observatoire des pratiques linguistiques « Langue et cité »*, 24. https://www.culture.gouv.fr/content/download/93563/file/lc_24_feminin-masculin_def.pdf?inLanguage=fr-FR
- Houdebine-Gravaud, A.-M. (2001). La féminisation des nomes de métiers en français. In Stistrup, M. (dir) *Nature, langue, discours*. Lyon : PUL, 95-107.

- Houdebine-Gravaud, A.-M. (2003). Trente ans de recherche sur la différence sexuelle ou Le langage des femmes et la sexuation dans la langue, le discours, les images. *Langage et Société*, 106, 33-61.
- Langlais, P. (en préparation). « A Journey in Neural Machine Translation ». *De Europa, Special Issue 2022*.
- Le Cun, Y. (2019). *Quand la machine apprend*. Paris : Odile Jacob.
- Maingueneau, D. (1991). *L'analyse du discours*. Paris : Hachette.
- Maingueneau, D. (2011). Pertinence de la notion de formation discursive en analyse du discours. *Langage & Société*, 135. <https://www.cairn.info/revue-langage-et-societe-2011-1-page-87.htm>
- Marzi, E. (2021). La traduction automatique neuronale et les biais de genre : le cas des noms de métiers entre l'italien et le français. *Synergies Italie*, 17, 19-36. <http://gerflint.fr/Base/Italie17/marzi.pdf>
- Mayaffre, D. (2021). *Macron ou le mystère du verbe. Ses discours décryptés par la machine*. Paris : Éditions de l'aube.
- Mazière, F. (2015[2005]). *L'analyse du discours : histoire et pratiques*. Paris : PUF.
- Moirand, S. (2020). Retour sur l'analyse du discours française. *Pratiques linguistique, littérature, didactique*, 185-186. <https://journals.openedition.org/pratiques/8721>
- Née, E. (éd.) (2017). *Méthodes et outils informatiques pour l'analyse des discours*. Rennes : PUR.
- Orlandi, E. (1996). *Interpretação: Autoria, leitura e efeitos do trabalho simbólico*. Petrópolis : Vozes.
- Paveau, M.-A. (2017). *L'analyse du discours numérique. Dictionnaire des formes et des pratiques*. Paris : Hermann.
- Raus, R. (2004). La linguistique française et les études de genre. In Raus, R. (éd) *Linguaggi e discriminazioni*. <http://www.cirsde.unito.it>
- Raus, R. (2013). *La terminologie multilingue. La traduction des termes de l'égalité H/F dans le discours international*. Bruxelles : De Boeck.
- Raus, R. (2020). La traduction française de *gender* et d'*intersectionality* dans les organisations internationales entre mémoires 'féministes' et déconstruction des différences. *De Genere*, 5, <https://www.degenere-journal.it/index.php/degenere/issue/view/6>
- Raus, R. (en préparation). La traduction des discours européens sur l'intelligence artificielle entre effets de sens et 'capitalisme de surveillance'. *Mots. Les langages du politiques*, 128.
- Sabatini, A. (1987). *Il sessismo nella lingua italiana*. Presidenza del Consiglio dei Ministri. http://www.funzionepubblica.gov.it/sites/funzionepubblica.gov.it/files/documenti/Normativa%20e%20Documentazione/ Dossier%20Pari%20opportunità/linguaggio_non_sessista.pdf
- Savoldi B., Gaido M., Bentivoglio L., Negri M., Turchi M. (2021), Gender Bias in Machine Translation, 9, 845-874. https://direct.mit.edu/tacl/article/doi/10.1162/tacl_a_00401/106991/Gender-Bias-in-Machine-Translation
- Tiedemann, J. (2012). Parallel data, tools and interfaces in OPUS. In *Lrec, 2012*, 2214-2218.
- UNESCO (1999). *Guidelines on Gender-Neutral Language*. https://eige.europa.eu/sites/default/files/unesco_guidelines_gender-neutral_language_0.pdf
- Vaswani, A., Shazeer, N., Parmar, N., J. Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., Polosukhin, I. (2017). Attention is all you need. *CoRR*, vol. abs/1706.03762.
- Vecchiato, S. (2004). Le sexisme dans le langage. Notes sur l'italien et le français. In Raus, R. (éd) *Linguaggi e discriminazioni*. <http://www.cirsde.unito.it>

Vetere, G. (en préparation). Elaborazione automatica dei linguaggi diversi dall'inglese: introduzione, stato dell'arte e prospettive. *De Europa, Special Issue 2022*.

Yaguello, M. (1978[2006]). *Les mots et les femmes. Essai d'approche sociolinguistique de la condition féminine*. Paris : Payot.

Yasunaga, M., Kasai, J., Zhang, R., Fabbri, A., Li, I., Friedman, D., Radev, D. (2019). ScisummNet: A large annotated corpus and content-impact models for scientific paper summarization with citation networks, in *Proceedings of AAAI 2019*.

Notes

¹ Voir, en France, les réactions à la tribune parue le 18 septembre 2020 dans le magazine *Marianne*, <https://www.marianne.net/agora/tribunes-libres/une-ecriture-excluant-qui-s-impose-par-la-propagande-32-linguistes-listent-les>

² Voir, en France, Marine Yaguello (1978[2006]) et Houdebine-Gravaud (2003) ou les études d'Alma Sabatini (1987) en Italie.

³ Nous tenons à préciser que la métaphore biologique qui a permis la parution du syntagme « intelligence artificielle » est désormais un véritable « mythe » naturalisé (voir Raus en préparation) : il suffit de voir des expressions utilisées couramment en informatique et en ingénierie comme « l'algorithme apprend », « apprentissage auto-supervisé », etc. pour s'en apercevoir.

⁴ Pour une présentation de la discipline, nous renvoyons à Mazière (2015) et Moirand (2020).

⁵ Cf. <http://www.jmcoe.unito.it/home>

⁶ Nous avons tenu également compte de la variation diatopique des langues et des politiques linguistiques différentes existant entre pays utilisant des variantes de la même langue (par exemple, l'italien d'Italie et de Suisse, le français de France, de Belgique ou du Luxembourg...).

⁷ Selon la banque terminologique de l'Union européenne IATE, l'apprentissage supervisé se définit de la manière suivante : « dans le contexte de l'intelligence artificielle et de l'apprentissage automatique, système qui fournit à la fois les données en entrée et les données attendues en sortie, les données en entrée et en sortie étant étiquetées en vue de leur classification, afin d'établir une base d'apprentissage pour le traitement ultérieur des données » [fiche 3676889 consultée le 16 décembre 2021]. Cf. <https://iate.europa.eu/home>

⁸ Voir, entre autres, voir la fiche de l'alors Commission générale de terminologie et de néologie (COGETER) de juillet 2005 (<http://www.culture.fr/Ressources/FranceTerme/Recommandations-d-usage/GENDER>)

⁹ Voir, entre autres, Greco (2013) et Abbou, Arnold, Candeia, Marignier (2018).

¹⁰ Voir, entre autres, <https://italianoinclusivo.it> et <https://thesubmarine.it/2020/08/03/schwa-linguaggio-inclusivo-vera-gheno> Les réactions au schwa ne se sont pas faites attendre : <https://accademiadellacrusca.it/it/consulenza/un-asterisco-sul-genere/4018>

¹¹ Cf. https://fr.wikipedia.org/wiki/Langage_inclusif_en_français (consulté le 16 décembre 2021).

¹² Voir <https://www.treccani.it/magazine/chiasmo/extra/linguaggioinclusivo.html> et <https://accademiadellacrusca.it/it/consulenza/un-asterisco-sul-genere/4018> (consultés le 29 novembre 2021)

¹³ Voir, entre autres, la traduction française de la notion de genre dans ONU 1995 (Raus 2020).

¹⁴ Pour le débat sur l'écriture inclusive et les questions idéologiques en France, signalons l'article d'Abbou, Arnold, Candeia, Maregnier (2018).

¹⁵ Nous utilisons cette notion, tout en sachant qu'elle est débattue en ADF (voir Mainueneau 2011)

¹⁶ Voir, entre autres, le n°37 de la revue *Langages* de 1975 au lien https://www.persee.fr/issue/lgge_0458-726x_1975_num_9_37 (consulté le 27 novembre 2021).

¹⁷ Voir, entre autres, le site <http://mesure-du-discours.unice.fr/>

¹⁸ Voir le dictionnaire de Marie-Anne Paveau (2017).

¹⁹ Voir par exemple, le livre de vulgarisation de l'un des experts le plus célèbre de ce domaine, Yan Le Cun (2019). Dans cet article, nous utiliserons donc l'expression « dispositif neuronal » pour renvoyer justement aux réseaux des neurones d'apprentissage profond qui permettent à la « machine » de performer.

²⁰ Ce type d'erreur est actuellement en cours de résolution, au sens où on en a pris conscience et on est en train d'intervenir pour le résoudre.

²¹ Les exemples tirés de *Deepl* ont été consultés pendant la rédaction de l'article et vérifiés pour la dernière fois le 4 janvier 2022.

²² Nous n'avons pas la possibilité de décrire en détail les dérives que l'utilisation de l'anglais entraînerait. Nous nous limitons à suggérer la lecture d'un article de Guido Vetere (en préparation) sur ce sujet.

²³ Le corpus est décrit en détail dans la section 3.

²⁴ L'exemple tiré de *Google, Deepl* et *Reverso* ont été consultés pendant la rédaction de l'article et vérifiés pour la dernière fois le 2 mars 2022.

²⁵ En effet, la phrase « *la capitana della squadra ha visto l'infermiera, che è uscita dall'ospedale* » (la capitaine de l'équipe a vu l'infirmière, qui est sortie de l'hôpital) est traduite par le traducteur *Systran* « la capitaine de l'équipe a vu l'infirmière, qui est sortie de l'hôpital ».

²⁶ Cf. la *Guida al pari trattamento linguistico di donna e uomo nei testi ufficiali della Confederazione*.

https://www.bk.admin.ch/dam/bk/it/dokumente/sprachdienste/Sprachdienst_it/02/objekt_40366.pdf.dowload.pdf/guida_al_pari_trattamentolinguisticodidonnaeuomo.pdf

²⁷ Sur ces deux types de conditions de production, voir Charaudeau, Maingueneau (2002 : 119). Les conditions matérielles sont celles que Charaudeau considère comme étant des conditions communicationnelles, à savoir « l'ensemble des données non-linguistiques qui président à l'acte d'énonciation ».

²⁸ Nous avons tenu compte de l'original du document : par exemple, un document rédigé par le logiciel word, qui ensuite a été publié en ligne pour des questions de transparence administrative, a été quand même considéré comme document « non web ».

²⁹ Signalons que nous n'avons pas pu suggérer des formulations avec le point médian ou d'autres choix expérimentaux pour plusieurs raisons, la première étant que ces reformulations excluent d'autres catégories comme, par exemple, les personnes malvoyantes. En outre, ces reformulations, qui sont éventuellement possibles pour la seule communication interne pour ne pas créer plus généralement des problèmes de lecture qui excluraient la majorité des personnes (plusieurs de ces formes ne sont reconnaissables que par un public cultivé, alors que l'administration doit pouvoir s'adresser à tout le monde, voir à ce sujet la question débattue de la clarté du langage de l'administration), ne sont pas acceptées par les institutions qui font autorité sur le plan des politiques linguistiques.

³⁰

Cf.

<https://www.unipd.it/sites/unipd.it/files/2018/Le%20parole%20delle%20disabilita%20e%20inclusion e.pdf>

³¹ Par exemple, « *ricerca* » (=recherche) à la place de « *ricercatore* » (chercheur).

³² Voir <https://github.com/dbmdz/berts>

³³ Le corpus a une dimension de 13 GB et se compose d'environ 2 milliards de mots.

³⁴ Cf. <https://www.nltk.org/>

³⁵ Il s'agit des Universités de Bari, Bologne, Turin, Catane, Enna, Modène et Reggio d'Émilie, Naples « *L'Orientale* », Pise, Rome « *La Sapienza* », Turin, Urbino et de l'École polytechnique de Milan.

³⁶ Il s'agit des « villes métropolitaines » de Bologne, de Rome, de Turin, de Venise, la mairie d'Ancône. Nous tenons à préciser que la ville métropolitaine de Turin a adhéré au projet E-MIMIC et a mis à disposition ses propres documents internes.

³⁷ Pour l'instant, nous avons chargé 56 appels à concours, 4 rapports et 8 procès-verbaux.

³⁸ Ce logiciel permet de traiter des fichiers avec des extensions différentes (.docx, .odt, .pdf...).

³⁹ L'exactitude renvoie au nombre de documents étiquetés par rapport au nombre total des documents analysés.