

[SLIDE 1/36]

Bonjour à toutes et à tous, aujourd'hui nous allons parler d'utilisation pédagogiques des IA génératives. La présentation que je vais vous faire aujourd'hui a été préparée par moi-même (Guillaume MULLER), Karine RICHOU et Julien MORICE.

[SLIDE 3/36]

Tout d'abord, commençons par de rapides présentations. Je m'appelle Guillaume MULLER, je suis Maître de Conférences à l'École des Mines de Saint-Étienne. Mes spécialités, autant en recherche qu'en enseignement, sont l'Intelligence Artificielle et la Cyber-Sécurité.

Karine RICHOU est ingénieure pédagogique à l'École des Mines de Saint-Étienne.

Julien MORICE est ingénieur pédagogique à l'IMTBS (Institut Mines-Telecom Business School) et en charge du projet Practice.

[SLIDE 3/36]

Puisqu'on ne se connaît pas je vous propose de faire un petit questionnaire pour identifier ce que vous savez de l'IA, en général, et des IA Génératives, en particulier.

Vous pouvez vous rendre sur le site WooClap. Le nom du quiz est QUIZAI, tout attaché et tout en majuscules, ou vous pouvez scanner le QR code qui s'affiche actuellement.

[SLIDE 4/36]

[Résumé du Quiz:

+ Pour quelles tâches l'IA est-elle aujourd'hui déjà utilisée dans nos quotidiens ?

- La reconnaissance des visages sur les photos
- La conduite des enfants à l'école en voiture autonome.
- L'analyse des émotions selon les expressions chez le médecin.
- Le choix d'itinéraire est optimisé selon le trafic en temps réel.
- La recommandation de Chanson par rapport à nos écoutes précédentes.

+ Quels outils d'IA connaissez vous ?

+ J'utilise une IA...

- Tous les jours
- 1 fois par mois
- 1 fois par semaine
- Jamais

+ Pour...

+ Quelles sont les images générées par IA ?

[3 exemples d'images générées par des IA, notamment celle du pape en doudoune "balenciaga"]

]

[SLIDE 6/36]

Maintenant que je sais un peu mieux ce que vous savez de l'IA Générative (apparemment une grande majorité d'entre vous, a déjà entendu parler de pas mal d'outils), la première chose qu'on va faire, c'est de voir un petit panel de tout ce qui existe pour mettre tout le monde au même niveau.

Premièrement ce qu'il faut savoir c'est que des outils s'appuyant sur l'intelligence artificielle, on en utilise depuis très longtemps !

Par exemple, pour les mathématiciens parmi vous, il y en a sûrement qui ont utilisé Wolfram Alpha. Cet outil existe depuis des années et s'appuie sur des systèmes d'IA. Même si ce ne sont pas des systèmes d'IA générative, comme les outils d'aujourd'hui.

Dans les outils d'aujourd'hui il y a plusieurs catégories d'outils :

- On a d'abord les outils qui manipulent du *texte*. Ici aussi, il y en a au moins un qui existe depuis 'très' longtemps (2017) que beaucoup d'entre vous doivent utiliser, sans savoir qu'il s'appuie de l'IA. C'est DeepL, qui permet de faire de la traduction. Vous avez évidemment entendu parler de ChatGPT, que Microsoft utilise dans son outil qui s'appelle BingChat (renommé Copilot depuis fin 2023). Google a un outil qui s'appelle Bard, Facebook a un outil qui s'appelle LLama/Alpaca et puis il y a d'autres sites comme Chatsonic ou Perplexity. Comme c'est écrit en dessous, tous ces outils, manipulent du texte. Ils peuvent faire plein de choses : générer, de quasiment rien, du texte ; résumer du texte ; traduire du texte ; reformuler du texte (par exemple dans un certain style qu'on lui demande) ; et puis éventuellement poser/répondre à des questions.

- Il existe des outils qui peuvent générer du *code* informatique. Le code étant une forme de texte très structurée, ça marche plutôt bien. ChatGPT sait générer du code, mais il y a des outils spécialisés, comme Copilot, qui est produit par Microsoft, et un autre qui s'appelle TabNine. Ils sont capables de générer du code dans tout un tas de langages. On leur demande par exemple, d'écrire un programme qui fait si ou ça, et il va écrire le code qui correspond à ça dans le langage qu'on lui a demandé.

[SLIDE 7/36]

- Ensuite, il y a des outils qui permettent de manipuler des *images*. Parmi les plus connus, donc il y a Dall-E (qui est fait par OpenAI, la même société que celle qui fait ChatGPT). Un autre très connu qui s'appelle Stable Diffusion (qui est accessible librement), un site qui s'appelle Mid-Journey (fait par des anciens d'OpenAI) et comme je disais, il y en a toujours plein, donc on a un autre que peut-être moins connus, qui s'appelle InstantART.

L'utilisation est assez similaire : on écrit un bout de texte (on dit par exemple : "génère-moi une image de chaise en forme d'avocat" et l'outil va générer une image qui représente une chaise en forme d'avocat). Ça marche assez bien, c'est plutôt bluffant.

- On a aussi tout ce qui permet de générer du *son*. Là j'ai mis deux outils qui ont deux utilisations très différentes :

+ MuBERTE qui est un outil qui permet de /générer de la musique/. Je peux vous faire écouter à un petit extrait de musique générée par MuBERT, si vous voulez, en cliquant sur le lien. MusicLM proposé par Google (sous code libre) fait pareil.

+ Il y a une autre catégorie d'outils, comme Vall-E (clin d'œil à Dall-E pour les images) qui sont les outils qui permettent de /générer de la voix/. Avec cet outil, on peut enregistrer trois secondes de sa propre voix et après taper un texte sur le clavier de l'ordinateur et l'ordinateur va le lire avec notre propre voix.

- De la même façon, les via génératives sont maintenant capables de générer des *vidéos*. Il y a un site comme D-ID qui permet d'animer un avatar en fonction de la texte qu'on écrit au clavier. On va donc pouvoir faire dire à l'avatar (potentiellement une photo de vous qui va être animée automatiquement) ce que vous avez tapé.

Il a même déjà des sites qui vont encore plus loin que ça et qui permettent, de générer des MOOCs (des cours en ligne) : Vous préparer la structure du cours,

vous écrivez/faites générer les textes qu'il faudrait dire dans les cours, et un avatar (qui est potentiellement votre photo) va dire (avec potentiellement votre voix) tout le contenu du cours.

La plupart des outils dont je vous parlais font une seule chose : soit générer des images, soit générer du son, soit générer la vidéo. Mais en fait tous ces outils commencent à être intégrés les uns avec les autres, dans des solutions plus complètes. Je vous parlais de celles pour générer des MOOCs. Mais on commence à voir de plus en plus de solutions très complètes.

Par exemple, les différents outils qui vous permettent de faire des visio-conférences, par exemple Meet de Google ou Teams de Microsoft, commencent à intégrer des outils d'IA. D'une part, pour la partie texte, par exemple, pour générer les transcripts des réunions automatiquement, ou pour améliorer l'aspect de la "chambre virtuelle" dans laquelle vous êtes. Par exemple, on peut demander que la "chambre virtuelle" contienne 5 sièges, car on est 5 personnes.

Comme je disais, Microsoft a très gros contrat avec OpenAI (la société qui fait ChatGPT et Dall-E), et ils sont en train d'intégrer l'outil, dans leur système d'exploitation (Windows). Ils l'intègrent à plein d'endroits. Il faut noter que tous les outils Microsoft qui utilisent ChatGPT ont été renommés Copilot récemment. L'assistant qu'on avait dans Windows jusqu'à maintenant, qui s'appelaît Cortana, va s'appeler Copilot, de même que le moteur de recherche Bing ou que le générateur de code de GitHub.

[SLIDE 8/36]

Comme vous pouvez le voir, il y a des tas de services qui existent ou se construisent autour des IA Génératives, notamment à destination des enseignants.

Le paysage change sans arrêt. Je viens de vous donner une image de ce qui existe aujourd'hui, en janvier 2024.

Cette liste d'outils, elle est impossible à maintenir pour moi, parce qu'il y en a des centaines qui sortent tous les jours.

Pour pouvoir suivre ça, je vous propose 2 sites ici :

- Il y en a un qui s'appelle "There's an AI for that", avec un moteur de recherche où vous mettez des mots clés et vous récupérez une liste d'outils d'IA qui le font ce que vous avez demandé. Par exemple, j'ai écouté un podcast

récemment où une dame parlait d'une IA qui permet de designer des vêtements :
vous décrivez un peu ce que vous voulez, ça vous propose des patrons de vêtements, à partir du patron de vêtements, vous écrivez encore du texte pour configurer en disant par exemple que vous voulez que ce soit plus serré à la taille, de telle couleur, ou des choses comme ça, ou avoir un col de telle ou telle sorte. Le site adapte le patron automatiquement. À la fin, soit vous pouvez, télécharger le patron et construire vous-même l'habit, soit vous pouvez carrément demander de l'acheter et le site, qui vous le produit et vous l'envoie directement chez vous.

- Un autre site du même style, s'appelle "Future Tools".

[SLIDE 10/36]

Voilà, alors maintenant, j'ai fait un peu un panorama de tout ce qui existe. On va rentrer un tout petit peu dans les détails, notamment de l'outil ChatGPT (outil de génération de texte) pour voir un peu comment ça marche à l'intérieur, et notamment pour discuter un peu des avantages et des limites (intrinsèques) que ces outils peuvent avoir.

L'outil qu'il y a derrière ChatGPT, c'est ce qu'on appelle un "modèle de langage" et ce dernier s'appelle "GPT", pour "Generative Pre-trained Transformer".

Ce système s'appuie sur un mécanisme qu'on appelle l'attention et l'idée c'est d'arriver à lui faire prédire le prochain élément d'une séquence, par exemple le prochain mot d'une phrase.

Je vous ai mis un petit exemple ici : on a une phrase "The FBI is chasing a criminal on the run". L'idée c'est qu'on a des tas de phrases comme ça, mais je vous montre l'exemple pour une seule phrase.

L'algorithme commence par le premier mot (celui qui est en rouge).

Pour l'instant le modèle ne voit pas la suite. Il voit le mot "The", et puis c'est tout. À force de voir des mot comme cela en début de phrase, il apprend que c'est un mot qui peut apparaître en début de phrase.

Ensuite on passe au mot suivant "FBI". Là le modèle, on voit (en bleu), qu'il

porte son *attention* sur le mot "The", parce qu'il estime que c'est important de mettre en combinaison "FBI" et "The". Le modèle apprend le concept d'article défini(?).

En anglais, ça ne se voit pas trop, mais en français par exemple on pourrait avoir "le chien" ou "la chienne". À force de voir "le chien", toujours en association et jamais "la chien", ou de voir toujours en association "la chienne" et jamais "le chienne", le modèle va apprendre, en gros, ce que c'est que le féminin, le masculin, le pluriel, etc. de

De la même façon, donc si on continue avec le troisième mot "is" c'est le verbe, et donc là il va porter son attention sur FBI, qui est le sujet. Là c'est un singulier, mais on pourra voir un pluriel, donc si on avait quelque chose au pluriel, on n'aurait pas "is" mais "are", et à force de voir des phrases ou à chaque fois qu'il y a un singulier, il y a "is" et à chaque fois qu'il y a un pluriel il y a "are", le modèle va créer des associations.

Et l'ensemble de toutes ces associations (statistiques) : masculin/féminin/pluriel/conjugaisons, etc. le système va créer un *modèle de langage*. C'est comme cela qu'on appelle GPT, le modèle sous-jacent à ChatGPT.

Je passe les détails pour la suite du processus, mais quand on arrive à la fin de la phrase, que quand le modèle arrive sur le point final, pour lui ce qui est important, ce qu'il a retenu, c'est qu'il y a un "criminal" qui est "on the run" et qui est "poursuivi" par le "FBI", et donc le système a, entre guillemets, compris le sens de la phrase, en associant en fait, quels étaient les mots de cette phrase, qui étaient importants pour chacun des autres mots.

[

Là je vous ai expliqué comment on *entraîne* le modèle en lui donnant à "manger" des centaines de millions de phrases

ChatGT, il a été entraîné sur quasiment tous les documents que l'on peut trouver sur Internet : tous les sites web, les messages sur les réseaux sociaux, quasiment tous les livres qui ont déjà été publiés (librement ou pas!), tous les messages Twitter, etc. La liste officielle n'est pas publique, mais c'est ce qu'on peut deviner en fonction des réponses du modèle.

Maintenant pour *utiliser* le modèle, c'est-à-dire générer des phrases, c'est assez simple, on peut soit lui donner un début de phrase et automatiquement, mot par mot, il va compléter la phrase en collant le mot le plus probable derrière ou alors on ne lui donne rien du tout, il va trouver un premier mot le plus probable et puis enchaîner, comme ça, de fil en aiguille.

]

Il se trouve que les gens qui ont créé ce modèle, se sont rendus compte, qu'en fait, c'est assez rigolo parce que si on donne une question à un tel modèle, la plus forte probabilité des mots qui vont suivre derrière une question, ce sont les mots qui correspondent à la réponse. Donc du coup, un effet de bord de ces modèles c'est qu'ils sont en partie capable de répondre à des questions.

Et donc les gens qui ont créé ces modèles se sont dit : "en fait, est-ce qu'on pourrait pas faire en sorte d'entraîner encore plus le modèle pour vraiment répondre à des questions ?

Et c'est comme cela qu'on va passer de *GPT, un modèle de langage*, donc qui sait comment est structurée une phrase, à un *système de Chat*, capable de répondre à des questions et de mener de véritables conversations.

[SLIDE 11/36]

Pour expliquer rapidement le passage de l'un à l'autre : , on a le fameux GPT qui est en haut, qui est un modèle de langage. Il a été ensuite entraîné à répondre à des questions. Le modèle résultant s'appelle "InstructGPT" et il a été entraîné selon une méthode qu'on nomme "RLHF", "Reinforcement Learning for Human With Human Feedback", donc l'idée, c'est de mettre le modèle dans une boucle avec un être humain : l'être humain va poser une question au chatbot, le chatbot va répondre et l'être humain va corriger ou proposer des réponses alternatives au le système. Qui va être entraîné à générer plutôt ce type de réponses-là que ces réponses initiales.

Vous pouvez voir les effet de ce processus typiquement quand vous posez des questions à ChatGPT, et qu'il va avoir des réponses assez structurées ou

stéréotypées, comme "je vous remercie beaucoup d'avoir posé cette question...", ou "je suis vraiment désolé de m'être trompé...", etc. C'est dans cette phase qu'on lui apprend à faire une réponse structurée, à être poli, etc.

Comme ce système a été entraîné sur tout un tas de données qui existent sur internet, qui vous le savez, sont particulièrement biaisée, par exemple racistes, ce InstructGPT est potentiellement biaisé/raciste.

L'entreprise OpenAI qui a créé ce système a donc rajouté une 3ème couche de sécurité pour éviter de répondre à certaines questions ou pour les reformuler "correctement".

Le résultat de ces 3 phases c'est le chatbot qu'on appelle ChatGPT : un modèle de langage "GPT", qui a été entraîné par feedback humain à répondre à des questions "InstructGPT", donc puis lissé par une couche de sécurité.

[SLIDE 12/36]

Vous l'avez bien compris, ce système, en fait, il n'est pas "intelligent" (au sens commun) dans le sens où il ne "comprend" pas les questions qu'on lui pose et ne "raisonne" aucunement pour répondre. Il génère juste les mots les plus probables derrière une question. Et les probabilités viennent de tout un ensemble de documents (corpus) qu'on lui a fourni dans sa phase d'entraînement. Ces documents étant potentiellement faux et/ou biaisés, le système résultant l'est tout autant.

Ce système génère des textes très bien faits, donc les gens se sont vite extasiés des résultats. Mais très rapidement d'autres personnes ont commencé à relativiser et à le critiquer.

Notamment, 2 chercheuses américaines l'ont surnommé "stochastic parot", i.e. un "perroquet probabiliste", puisqu'il ne répète que ce qu'il a déjà lu/vu, mais sous une forme paraphrasée, grâce aux probabilités.

Un youtubeur que j'aime bien, qui s'appelle "science étonnante", parle pour sa part d'un "baratineur". J'aime beaucoup cette définition, je trouve qu'elle image très très bien comment fonctionne ChatGPT et quelles sont ses limites.

Je vais vous illustrer cela avec l'exemple qu'il donne dans sa vidéo.

Mettons qu'on pose une question simple à ChatGPT, par exemple : "qu'est-ce qui a

le meilleur impact pour résoudre le changement climatique : manger moins de viande ou manger local ?".

Si on lui demande une centaine de fois de répondre, 70% des fois, il va vous dire que c'est de manger moins de viande, et dans 30% des cas, il va vous dire que c'est de manger local. Il se trouve que c'est le consensus scientifique : c'est effectivement manger moins de viande, qui permettrait le mieux d'aider à résoudre le changement climatique, notamment parce que les animaux eux-mêmes consomment génèrent du CO2 et que le transport génère du CO2. Donc la majorité du temps ChatGPT vous donnera la bonne réponse.

Pour comprendre sa réponse, c'est simple, il suffit dans s'imaginer l'ensemble des documents sur lesquels ChatGPT a été entraîné. Parmi cet ensemble, ne considérer que ceux qui parle du dérèglement climatique et de manger moins de viande et de manger local. Et bien, dans 70% de ces documents, il est dit que c'est manger moins de viande qui aurait le meilleur impact et 30% qui disent que c'est manger local.

Du coup, ChatGPT, quand il doit choisir les mot qu'il enchaîne pour générer ses phrases, il va s'appuyer sur ces probabilités et générer dans 70% des cas des phrases qui donnent la bonne réponse.

Maintenant, si je rajoute du contexte à ma question. Si par exemple je dis : "je suis un éleveur d'animaux, quelle est la chose qu'il a le plus d'impact pour résoudre le changement climatique, est-ce que c'est manger moins de viande ou manger local ?

À votre avis, qu'est-ce qu'il répond ?

Bien oui, vous l'avez deviné, il va répondre qu'il vaut mieux manger local. On voit que là, les probabilités sont pas hyper séparées, on est presque à du 50/50, mais il va quand même globalement favoriser la réponse qui est la mauvaise.

C'est pas qu'il a compris que j'avais dit que j'étais un éleveur d'animaux, et que du coup, j'ai une préférence pour le fait qu'il faudrait quand même continuer à manger de la viande, parce que ça arrange mon business. Ce n'est pas

du tout ça. Pour reprendre l'image de tout à l'heure, si on reconstituait encore une fois l'ensemble de tous les documents sur lesquels ChatGPT a été entraîné et que, cette fois-ci, parmi tous ces documents qui parlent de la question climatique et de manger moins de viande ou manger local, on ne retenait que ceux qui parlent aussi du fait d'être un éleveur d'animaux, et bien, les probabilités ne sont pas les mêmes : les documents penchent cette fois-ci pour l'autre réponse.

Le fait de rajouter un contexte, ça réduit le corpus des documents que l'algorithme va considérer pour donner sa réponse.

Évidemment, on peut continuer à rigoler comme ça : on peut dire par exemple, "je suis un fermier en biodynamie, ...", bah là, forcément, ChatGPT, il va passer à 92% pour dire qu'il vaut mieux manger local.

Si on fait complètement l'inverse, en disant "je suis antispéciste ...", là bien sûr le but, c'est de manger beaucoup moins de viande, donc les proportions des proportions monteront carrément à 94% pour dire qu'il faut manger moins de viande.

Pour résumer un peu cet exemple rigolo, pourquoi est-ce que Science Étonnante qualifie ChatGPT de "baratineur" ?

Parce que, en fait, comme je vous disais, ChatGPT ne comprend pas la question ni sa réponse, donc il ne cherche pas à vous donner la vérité, puisqu'il ne la connaît pas. Il ne cherche pas non plus à vous dire le contraire, c'est-à-dire à vous mentir.

Il cherche juste à vous donner la réponse qui est la plus probable ... dans un certain contexte ... et ce contexte, c'est vous qui donnez ... donc en fait, il cherche à vous donner la réponse qui vous satisfait le plus ... et c'est exactement la définition d'un baratineur : quelqu'un qui va toujours dans votre sens, qui vous brosse dans le sens du poil.

[SLIDE 14/36]

Pour revenir à ChatGPT et son utilisation.

Il y a deux utilisations de ChatGPT :

- La première, celle que j'évoquais jusqu'à maintenant et qu'on voit à gauche, c'est d'aller sur le site de OpenAI et d'utiliser l'interface Web de ChatGPT.

- La deuxième c'est d'utiliser ce qu'on appelle l'API. En fait c'est d'intégrer ChatGPT dans une autre application déjà existante. Par exemple, je vous disais que Microsoft a intégré ChatGPT dans le moteur de recherche Bing. Si vous utilisez le moteur de recherche Bing, Microsoft effectuera très certainement un appel à ChatGPT sans que vous vous en rendiez compte.

De la même façon, si vous utilisez le site Whimsical, qui sert à générer des mindmaps, il y a un petit bouton avec une baguette magique qui est apparu. Si vous cliquez sur un des nœuds de votre carte mentale et que vous cliquez ensuite sur ce petit bouton, il va y avoir un appel à ChatGPT pour générer les descendants possibles de votre nœud grâce à ChatGPT, qui seront directement affichés par Whisical dans son interface. Et vous ne verrez pas qu'il y a eu un appel à ChatGPT, même si c'est le cas.

ChatGPT a aussi été intégré dans Office365 de Microsoft et on peut, par exemple, demander à ChatGPT de synthétiser la colonne numéro 3 ou de faire générer le graphique qui va bien, ou à la place d'écrire les formules dans Excel nous-mêmes.

Au passage, il est à noter que, dans les dernières versions de ChatGPT (4, 4

Turbo), a été intégré un système de plug-in qui permet à ChatGPT de se connecter à différents sites ou services sur internet. Donc d'un point de vue

fonctionnalité vous pouvez avoir un peu la même chose avec les deux systèmes.

Par exemple, depuis l'interface de CatGPT vous pouvez lui demander de se

connecter à whimsical et de générer une mindmap directement dans l'interface de ce site.

Tout ça pour dire que vous allez bientôt faire du ChatGPT comme Mr Jourdain faisait de la prose, c'est-à-dire sans vous en prendre compte.

Il y a quand même une différence majeure entre ces deux approches, notamment en ce qui concerne les références.

Vous avez déjà dû entendre parler du fait que ChatGPT est très mauvais pour donner des références.

Notamment, si on essaye de lui faire générer un article scientifique, il va générer des super belles références, parce que dans son modèle de langage, il a bien compris qu'une référence, c'est généralement : ouvrez crochet, un nombre, fermé crochet, une liste de noms d'auteurs, un titre, un nom de conférence une année, etc. Il a bien la compris la structure générale d'une référence, mais il va générer des faux noms, des faux titres, fausses dates, etc. On pouvait donc très rapidement voir qu'un texte était généré, en tout cas un texte scientifique avec des références, car ces références étaient à 99% complètement fausses, n'existait pas. Malheureusement, cela n'est pas/plus tout à fait vrai :

- L'intégration de ChatGPT avec le moteur de recherche Bing de Microsoft se fait comme suit : quand vous posez une question, il reformule la question sous la forme d'une ou plusieurs requête(s), il l'envoie dans le moteur de recherches, il récupère un certain nombre de réponses (mettons les 5 premiers sites) et vous fait une synthèse de ces sites en guise de réponse. Du coup, Bing+ChatGPT est capable de vous donner les références des pages web où il a trouvé ses réponses, et celles-ci existent et sont correctes. Même si les sites qu'il utilise comme sources ne racontent que des conneries, Bing+ChatGPT est capable de vous donner ses sources.

- Jusqu'à maintenant le site ChatGPT ne pouvait pas faire ça. C'était la grosse différence entre les 2 approches. Il ne pouvait vous donner des réponses que par rapport à l'ensemble des documents sur lesquels il avait été entraîné. Par exemple, pour Chat GPT 3.5, qui n'avait été entraîné sur des documents jusqu'à 2021, il était incapable de vous donner des réponses sur des faits qui s'étaient produits après 2021.

Avec le nouveau système de plug-in, cela n'est plus vrai, puisque les nouveaux ChatGPT (4, 4 Turbo - payants) peuvent se connecter à Internet pour aller chercher ses références.

Pour résumer, les fausses références ne sont plus un moyen fiable de détecter un

texte généré ☺

[SLIDE 15/36]

Dans le présent transparent, j'ai essayé de mettre à jour un peu tout ce qui se passe avec la toute dernière version de ChatGPT qui s'appelle ChatGPT 4 Turbo.

Elle date de quelques semaines maintenant, les nouvelles sont quand même encore assez fraîches.

ChatGPT 4 Turbo a été entraîné sur des *données plus récentes, allant jusqu'à 2023*. Donc, de base, sans aucune connexion internet, il est capable de répondre à à peu près toutes les questions qui sont des événements qui seraient arrivés jusqu'en 2023.

OpenAI a augmenté le contexte jusqu'à *128 000 tokens*, (je rentre pas dans les détails mais un token c'est quasiment équivalent à un mot). Ça veut dire qu'on peut avoir une conversation d'avec ChatGPT qui contient à peu près 128 000 mots, et juste pour vous donner une idée, c'est de l'ordre de plusieurs livres !! Vous pouvez désormais avoir une conversation avec ChatGPT qui a la taille de plusieurs livres et ChatGPT se rappellera de tout ce que vous lui avez dit et de ce qu'il vous a répondu sur cette période !

Ensuite, une autre chose qu'ils ont rajoutée, c'est la *parole*. Et ils l'ont rajoutée dans les deux sens : à la fois avec la reconnaissance de la parole (STT), donc vous pouvez, au lieu de taper votre requête à ChatGPT, la lui dicter. Et, en réponse, ChatGPT, au lieu de vous afficher un texte à l'écran, peut vous la lire. Vous pouvez avoir une conversation complète avec ChatGPT en audio !!

Une autre chose très intéressante, surtout pour nous, enseignants, qui a été introduite, c'est la *reproductibilité des résultats*. Vous pouvez demander à ChatGPT, quand vous lui posez une question et qu'il répond, de vous donner un nombre unique (une référence) pour identifier cette conversation. Si, par exemple, un élève a une conversation avec ChatGPT et me donne ce numéro, je peux moi-même poser la même question à ChatGPT, et en lui précisant ce nombre, récupérer exactement la même réponse que l'élève, alors que, jusqu'à maintenant,

ChatGPT répondait des choses différentes à chaque fois, du fait de son caractère probabiliste.

Enfin, OpenAI a créé un nouveau truc qui s'appelle "GPTs". L'idée, c'est de vous permettre de spécialiser ChatGPT pour une fonctionnalité en particulier. Par exemple, je peux créer un ChatGPT, qui m'aide à designer des habits ou un ChatGPT spécialisé pour la cuisine, etc.

Ce qui est très rigolo, c'est que pour faire ce ChatGPT spécialisé je n'ai pas besoin d'écrire de code et un énorme prompt/requête. Je discute juste avec le ChatGPT "général" et il va me construire lui-même cette énorme requête.

Ensuite je vous ai mis des petits exemples rigolo de choses qu'on peut faire avec ChatGPT. Par exemple en haut à droite, vous pouvez prendre une photo de votre frigo, et ChatGPT va vous donner une recette que vous pouvez faire. Il est capable d'analyser dans l'image quels sont les différents ingrédients qui sont présents, et d'aller chercher la recette qui correspond le plus à cette liste.

Durant le grand show de présentation de GPT4 Turbo, ils montrent un truc absolument bluffant : ils donnent une vidéo sur YouTube, et ChatGPT décrit la vidéo en disant "Ah bah on voit une voiture qui circule dans un désert, et qui monte une dune, et la voiture elle est très grosse, etc, et à la fin il dit : c'est très probablement une pub pour BMW" !!!

Donc assez impressionnant, il arrive non seulement à décrire tout ce qu'il voit dans la vidéo, et en plus de synthétiser tout cela en "devinant" que c'est une publicité.

Vous avez certainement déjà entendu parler aussi que beaucoup de gens ont essayé de faire passer des examens à ChatGPT. Comme les examens aux États-Unis sont souvent des QCM, donc entre guillemets "faciles" quand on sait retenir par cœur - spécialité de ChatGPT -, donc ChatGPT en a validé pas mal, notamment celui pour être avocat dans je-ne-sais-plus-quel État des États-Unis.

[Mais plus proche de nous, on a par exemple Maxime Le François qui a donné le PDF d'un examen qu'il donne aux étudiants de Master 2 CPS2 à ChatGPT et ce

dernier a eu une meilleure note que le meilleur des étudiants de Master ! Malgré des questions non triviales, autrement plus difficile que du par cœur.]

Comme je vous le disais, on peut aussi faire générer du code. Par exemple avec une (longue) discussion avec ChatGPT, on peut faire générer le code complet d'un jeu.

[SLIDE 16/36]

Comme je vous le disais aussi, tous ces outils commencent à être intégrés, notamment Microsoft, via son énorme partenariat avec OpenAI essayent d'intégrer ça dans tous leurs outils.

Je vous ai parlé de Bing Search, l'intégration de ChatGPT au moteur de recherche de Microsoft.

Je vous ai parlé de Office365, donc l'intégration de ChatGPT à tous les outils d'édition : textes, tableur, présentation, etc.

Je vous ai parlé de Teams, l'outil pour faire des visio-conférences.

Je vous ai parlé de le système d'exploitation Windows, par exemple pour chercher des fichiers, par le contenu, en dialogant avec le système, plutôt que d'aller chercher les fichiers en ouvrant des dossiers ou en cherchant par noms.

Ils intègrent aussi ChatGPT à leur outil pour développer du code, VSCode, sous le nom de Copilot.

Donc tout ça, ça va bientôt être intégré encore plus, de manière "transparente" pour l'utilisateur. Ça sera directement dans Windows, et ça sera partout, sans qu'on s'en rende compte.

[Ça va être difficile d'échapper à ces IA Génératives !]

[SLIDE 18/36]

Comme on l'a évoqué tout à l'heure, si on veut des bonnes réponses avec ChatGPT, il faut lui donner le bon contexte. Est donc en train d'émerger toute une science qu'on appelle le "Prompt Engineering" sur la façon de bien poser les questions à ChatGPT pour en obtenir le meilleur.

Dans les prochains transparents, j'essaie de vous résumer les petites astuces/techniques que j'ai trouvées dans la littérature.

La première chose, c'est qu'il faut effectivement donner un bon contexte.

Souvent, un des trucs qui marche bien, c'est de donner un rôle, un "persona", à ChatGPT. Par exemple, lui dire : "tu es enseignant, tu écris un examen pour un cours de blablabla, pour des élèves de niveau blablabla, peut-tu me donner trois questions permettant d'évaluer leur compréhension de XXX".

Lui donner un rôle, ça permet très fortement contextualiser les réponses qu'il va donner.

Je vous ai mis à droite plusieurs ressources, sur comment on écrit des Prompts, de façon générale, mais aussi, plus spécifiquement, pour les enseignants. Notamment deux ebooks ici, que j'ai fait acheter la bibliothèque, donc que vous devriez bientôt pouvoir trouver la bibliothèque de l'école.

[SLIDE 19/36]

Pour résumer, il y a 2 petites images que je vais vous montrer, que j'ai trouvé sur le compte Twitter d'une personne qui s'appelle Matthieu, qui a aussi une chaîne YouTube qui s'appelle "OutilsIA", qui est très intéressante.

En gros, il a développé une méthode qui l'appelle "ACTIF", qui est assez très performante, qui consiste à décomposer le prompt en :

- une Action, puisqu'au final on demande toujours à ChGPT de faire quelque chose
 - (ici en jaune "répond aux 5 questions les plus fréquentes des utilisateurs"
- un Contexte
 - (ici en vert que l'entreprise est active dans le domaine XXX)
- une tonalité, car GPT a tellement bien compris la structure du langage qu'il peut générer du texte dans un certain style, par exemple sous la forme d'un poème, ou avec le style de Shakespeare, un email très formel pour le directeur de l'école, ou au contraire un truc de vulgarisation pour un enfant de trois ans.
- une identité, comme je vous disais, un persona(ge),
 - (ici en orange "agit comme un consultant en stratégie e-commerce"),
- finalement, un format. ChatGPT peut mettre un peu en forme le texte qu'il génère, par exemple, on peut lui demander de produire une liste, un tableau, etc.

La deuxième petite infographie, c'est pour vous donner des exemples de cette méthode. Par exemple, on peut lui dire, "agit comme un consultant", "tu dois

écrire un email de vente" "dans un style créatif", "sous la forme de slides".

D'autres astuces que j'ai trouvées sur internet, disent que souvent avoir un long Prompt comme ça, ça marche vous aurez des très bonnes réponses, mais c'est encore mieux si vous utilisez les facultés de conversation avec ChatGPT : commencer simplement, avec un petit prompt, et corriger ce prompt, par exemple en rajoutant un personnage ou en rajoutant du contexte, au fur et à mesure, en fonction des réponses que vous donne ChatGPT.

[SLIDE 20/36]

Pour finir sur l'histoire des Prompts.

Il y a des sites communautaires (gratuits ou payants), où les gens partagent des prompts qui fonctionnent bien pour certaines tâches.

Par exemple, j'en ai mis PromptVine, FlowGPT, etc. Ça peut être intéressant quand vous utilisez ChatGPT régulièrement pour les mêmes tâches.

Le dernier lien c'est une liste de sites de ce type.

La grande question qui se pose maintenant c'est est-ce que le Prompt Engineering va être la nouvelle science, est-ce que tout le monde va devoir savoir faire du Prompt Engineering ou pas.

D'une part, il existe des plugins soit de ChatGPT soit de votre navigateur, qui permettent d'insérer des prompts tout faits directement dans ChatGPT. Je vois que vous pouvez aussi avoir votre propre base de données de prompts directement dans les versions payantes de ChatGPT.

D'autre part, tous ces outils d'IA Générative s'améliorent, et comme je vous le disais pour les GPTs, ChatGPT commence lui-même capable de vous aider à écrire vos propres prompts (une sorte d'inception ou une IA nous aide à créer une IA !).

Finalement, se pose même la question aujourd'hui de savoir si on aura vraiment besoin de maîtriser ce Prompt Engineering ou pas.

[SLIDE 22/36]

Je vous ai beaucoup parlé de théorie. Maintenant, passons à la pratique.

Dans cette partie, on vous propose un petit exercice pour essayer

d'apprendre à maîtriser les outils d'IA Gen.

Comme on est plutôt entre enseignants, le petit exercice qu'on vous propose c'est d'essayer de concevoir un cours avec ChatGPT, soit pour la préparation du cours, soit pendant ou après le cours.

On va le faire en plus sous une forme pédagogique un peu "innovante" : les 1-2-tous. Je sais pas si vous connaissez le système 1-2-tous ? L'idée c'est que je vous pose une question/un problème, vous réfléchissez tout seul en votre coin, ensuite vous partagez avec votre voisin direct et puis finalement on partage tous ensemble nos feedbacks.

On peut faire l'exercice sous cette forme là pour que ce soit un petit peu rigolo.

Votre challenge ça va être donc de trouver des utilisations intéressantes des différents outils d'IA Gen (texte, images, vidéo, son) et d'écrire les prompts qui vont bien.

Pour la partie texte/images, je vous ai créé des comptes sur BingChat/Copilot (ChatGPT gratuit de Microsoft). Vous avez en bas du transparent les informations pour pouvoir vous connecter. Attention à ne pas réutiliser ces compte en dehors d'ici, car ils sont connus de beaucoup de gens et tous vos prompts sont stockés par Microsoft, donc tout le monde va pouvoir les voir.

Si jamais il y en a parmi vous qui bloquent, parce qu'il savent pas trop par où commencer, notamment s'ils ne maîtrisent pas trop les techniques de création d'un cours (c'est malheureusement mais en France nous ne sommes pas formés à l'enseignement, alors que c'est la moitié de notre travail !)

Pour ceux qui ne le sauraient pas, déjà il y a 5 composantes pour un cours :

- Définir/Partager un *thème* avec les apprenants (pour nous, en général, c'est l'École/la maquette qui le donne)
- Ensuite, on définit des *objectifs pédagogiques*, c'est une liste de choses que le prof voudrait que les élèves aient retenu à la fin du cours
- Après, il faut définir des *contenus théoriques* et *pratiques*

- Et enfin les mode d'évaluation. C'est là qu'apparaît le fameux
*"alignement
pédagogique"* dont l'idée est que le /mode d'évaluation/ doit
correspondre (/être aligné/) à ce qu'on cherche à faire retenir aux
élèves (/objectifs pédagogiques/).

Si par exemple, on est dans un cours de médecine et qu'on veut que les
élèves
retiennent tous les os de la main, c'est du par cœur, et donc, on peut
faire
un QCM pour l'évaluer, parce que le QCM, c'est très bien pour évaluer
l'apprentissage par cœur.

En revanche, si l'on est dans un cours de philosophie et que ce qu'on
veut,
c'est évaluer la capacité d'un élève à analyser un texte, par exemple,
le
faire en dissertation, est plus adapté qu'en QCM.

Pour la définition des objectifs pédagogiques, il y a six ou sept
éléments
différents qu'on peut essayer d'évaluer chez les élèves. Je vous ai mis
la
taxonomie de bloom qui est très connue.

Pour ce qui est la *définition des contenus pratiques et théoriques*, il
y a
dans le livret du participant de MÉDIANE 2023, toute une liste de façons
de
dérouter une session de "Cours" ou "TP" (par exemple, c'est là que j'ai
trouvé
la méthode 1-2-tous).

Pour les évaluations, pareil, il existe tout un tas de façons d'évaluer
les
élèves. Encore dans le livret du participant de MÉDIANE 2023, il y a
aussi toute
une section sur ces sujets-là.

[SLIDE 24/36]

Maintenant que vous avez bien réfléchi au sujet des utilisations
pédagogiques de
ChatGPT, je vais vous donner quelques exemples qui devraient résonner
avec ce
que vous avez trouvé.

Parmi les trucs de base :

- La première chose qui peut venir à l'esprit, si c'est un cours que vous
n'avez
jamais donné, c'est tout simplement de générer les idées de sujets à
aborder
ou un plan. On peut même demander un timing à ChatGPT !
- De plus en plus, les cours sont déposés sur un ENT/Moodle/eCampus, donc
on
peut vouloir générer une image pour illustrer le cours, une petite
icône, pour
que les élèves identifient très rapidement le cours.

- On génère une petite vidéo pour teaser le cours aux élèves, pour qu'ils aient envie de venir (ex. pour les Majeures/ToolBoxes/Défis qu'on présente à la journée des centres),
- Une fois qu'on a fait tous les contenus, on a tout un tas de documents, mais le syllabus du cours n'est plus forcément à jour. On peut utiliser ChatGPT pour créer un résumé du contenu sous la forme d'un syllabus et ainsi avoir un syllabus toujours à jour.

Pour les contenus & évaluation:

- On peut faire générer des quiz à ChatGPT. Il y a des sites entiers dédiés à ça (WooClap?). Avec un peu de pratique, on peut aussi demander à ChatGPT de générer une sortie en GIFT, qui est le format pour importer/exporter des quiz dans Moodle. Plus on a de questions plus on peut changer fréquemment pour éviter la triche. Dans Moodle, quand on a un pool de questions suffisant, on peut même lui demander de générer un ensemble de questions spécifique à chaque élève.

Ces quiz pourront soit sous forme pour des évaluations sommatives (pour donner une note) ou pour des évaluations formatives (juste pour aider élèves à savoir où ils en sont de leur compréhension)

- On peut préparer une session de cours entière : le plan du cours, et puis après, une fois qu'il a fait le plan, on peut lui demander pour chaque session de donner un peu plus de détails, et comme je vous le disais, on a maintenant des outils qui peuvent même générer les slides directement pour nous, une fois qu'on a validé le contenu avec ChatGPT.

- J'ai vu une utilisation très intéressante aussi, qui consiste à spécialiser le matériel. Avec 2 applications : 1) je peut suggérer le sujet général du TP (ex : un concept de base de données), et je peux demander à ChatGPT, de me générer le UseCase (thème, données, etc.) pour que les élèves aient tous des sujets différents et ne puissent pas tricher. 2) L'autre utilisation c'est pour la "pédagogie différenciée". ici on en n'a pas beaucoup, mais quand

j'étais à TSE, j'avais pas mal d'étudiants qui avaient des tiers temps pour cause de dyslexie, dyspraxie, etc. On peut demander à ChatGPT de réécrire notre matériel dans un style plus léger ou plus simple pour que ces élèves puissent mieux les appréhender.

- Enfin, on peut demander à ChatGPT de générer une grille d'évaluation. Je l'ai fait récemment car des élèves m'ont demandé la grille d'évaluation d'un projet, et j'avais écrit bêtement le sujet dans penser à l'évaluation. Du coup j'ai pu répondre au mail des élèves à 4 secondes, en demandant son avis à ChatGPT et en retouchant selon mes goûts. Ça m'a gagné pas mal de temps.

Pour ce qui est des activités elles-mêmes : Dans une autre (excellente) vidéo (Youtube) Sciences Étonnante. Il aborde le sujet des techniques qui marchent pour réviser ses cours pour les élèves. Je vous invite à le regarder, parce que c'est très intéressant, on apprend notamment que relire son cours, ça sert à rien, que fluoter ses notes, ça sert à rien, c'est scientifiquement prouvé.

En revanche, il y a des techniques qui marchent.

- Par exemple, le système de flashcard : se poser des questions régulièrement sur le cours, sans le relire. L'idée étant d'entraîner son cerveau à transférer des éléments de la mémoire à long terme vers la mémoire à court terme, plutôt que de travailler 100% dans le court-terme (relecture).

Une flashcard ça peut être une idée/concept que l'on doit décrire ou une question à laquelle on doit répondre. L'idée c'est on a un deck de telles cartes et on retourne les cartes les unes après les autres, en se posant les questions à soi-même. Celles auxquelles on sait répondre, on peut les sortir du deck. Celles auxquelles on ne sait pas répondre, on les remet dans le deck pour se les reposer de plus en plus souvent.

On peut demander à ChatGPT de générer le deck de cartes, soit simplement à partir du thème du cours, soit carrément en lui donnant le matériel du cours !

On peut aussi lui demander de jouer le rôle du deck de carte de de poser les questions les unes après les autres et même de commenter/corriger nos réponses

!

- On peut aussi demander de l'aide à ChatGPT pour concevoir Serious Game (thème, règles, cartes, etc.). On peut aussi lui donner les règles d'un jeu et lui demander de jouer certaines personnes, et nous, on jouera une autre personne.

- Comme ChatGPT répond toujours à ce qu'on lui demande, on peut le forcer à dire des grosses sottises. Ça peut être rigolo, d'utiliser ChatGPT pour développer le sens critique des élèves. Par exemple, Julien MORICE de l'IMT-BS a fait écrire un article de blog avec photos et tous sur l'aqua-poney. Et il l'a utilisé pour demander aux élèves de commenter, rechercher des sources, de croiser les faits, etc.

[SLIDE 25/36]

Dans le présent transparent, je vous ai mis encore deux ressources que j'ai trouvées sur les utilisations pédagogiques de ChatGPT. L

La première est gratuite, "AI de Classroom", où il y a cette jolie infographie, qui liste 20 exemples d'utilisation de ChatGPT de façon pédagogique.

Par exemple, pour les profs de langues (qui se posaient beaucoup de questions à l'arrivée de ChatGPT), si ils ont fait une session sur "comment passer un entretien d'embauche en anglais", il peuvent demander aux élèves, une fois rentrés chez eux, d'interagir avec ChatGPT pour simuler un entretien d'embauche en anglais et de leur fournir la trace de la discussion comme soumission. Dans un tel cas, on peut fournir le prompt aux élèves, ou les laisser l'écrire (et prendre en compte dans la notation).

On peut aussi organiser des débats ou des panels de discussions en demandant à ChatGPT de jouer plusieurs rôles. Et l'élève peut travailler tout seul chez lui/ comme s'il était dans une classe complète !

La deuxième ressource que je vous propose est payante (mais encore une fois j'ai demandé à la bibliothèque de l'acheter). elle donne 80 façons d'utiliser ChatGPT dans la classe, mais cette fois-ci de manière un peu plus concrète, avec des exemples de prompts.

Enfin, Karine RICHOU avait trouvé des documents/outils dans la pédagogie de l'IMT, ça s'appelait Ayden et Ariane, sur ces sujets. Malheureusement, il semblerait qu'ils ne sont plus accessibles aujourd'hui.

[SLIDE 26/36]

En tant qu'enseignants, il faut évidemment qu'on aborde la question de la détection, pour éviter la triche.

Il y a des outils qui existent, vous avez sûrement entendu parler de GPTero (un des tout premiers systèmes de détection). L'outil qu'on utilise pour la détection du plagiat, Compilatio, intègre aussi maintenant des outils qui vous donnent un pourcentage de chances pour qu'une IA générative soit à l'origine du texte que l'élève a rendu.

Malheureusement, GPT est un outil qui vient du Machine *Learning*, donc il est entraîné pour minimiser une métrique (qui ici combine différentes mesures de la "qualité" du texte généré). Si on a un détecteur qui répond oui/non (le texte est généré artificiellement) ou donne un pourcentage, on peut intégrer cette valeur dans la métrique à minimiser et la nouvelle version de l'outil saura tromper parfaitement l'outil de détection.

En ce sens, je pense que c'est une course qui est perdue d'avance. Il vaut mieux considérer qu'on ne pourra jamais détecter ces outils et essayer de faire avec.

[SLIDE 27/36]

On a parlé enseignement, maintenant parlons un petit peu recherche, puisque beaucoup d'entre vous sont aussi chercheurs.

Voici tout les outils que j'ai trouvé qui concerne la recherche.

Le premier s'appelle "consensus", il y a assez rigolo. Il n'est pas forcément pour la recherche pure, mais ça peut vous intéresser. Soit pour des domaines que vous ne maîtrisez pas forcément. Soit juste à titre personnel, par curiosité scientifique. L'idée c'est, on lui pose une question, typiquement celle que je disais tout à l'heure ("qu'est-ce qui est mieux pour la planète, est-ce que c'est manger moins de gens ou manger local"), et l'outil va chercher dans des

articles scientifiques ce qu'il trouve sur le sujet et le synthétise sous la forme d'une barre de pourcentage pour montrer si le consensus scientifique penche plutôt pour l'une ou l'autre réponse.

Ensuite il y a le site qui s'appelle "elicit". On peut mettre une question scientifique, et comme un moteur de recherches classique, il va chercher des articles scientifiques qui matchent. Mais ici il ne le fait pas juste avec des mots clés, il se base vraiment sur le contenu des articles et le sens/sémantique des mots (ça marche par exemple même si on utilise des mots différents mais du même champ lexical : véhicule au lieu de voiture, etc.).

Il liste les résultats dans un tableau, comme vous pouvez voir ici, dans la petite image au milieu, et l'idée c'est qu'il vous résume l'article en une phrase, donc il n'y a même pas à lire les abstracts. C'est très intéressant pour scanner rapidement un nouveau domaine qu'on ne maîtrise pas encore trop.

Un petit bémol : durant mes essais, j'ai eu l'impression qu'il était basé principalement sur des articles scientifiques de IEEE, donc y a pas forcément tous les articles/disciplines... Les résultats sont un peu biaisés mais j'ai trouvé le concept intéressant.

Un autre système lui aussi pour trouver des articles intéressants, fonctionne avec une autre approche, avec un système de recommandation. Il s'appelle "ResearchRabbit". Là l'idée c'est qu'on lui donne une liste d'articles qu'on aime et il cherche les articles similaires. Mais encore une fois, pas juste suivant des mots clefs ou des listes d'auteurs, vraiment par rapport au contenu des papiers. En plus, on peut créer diverses collections d'articles et connecter avec Zotero.

Comme il y a pas mal de doctorants ici à l'École des Mines, j'ai mis un lien vers un article, mais il y en a des tonnes qui expliquent aux doctorants comment utiliser le ChatGPT pour mieux écrire sa thèse.

Il y a aussi une astuce que j'ai découverte récemment : si vous trouvez un article sur ArXiv, vous remplacez juste "arxiv" dans l'URL par "talk2arxiv" et

vous pouvez discuter avec l'article (vous pouvez demander au système, "tiens de quoi parle cet article", "quelle est la méthode utilisée dans cet article", "à quel conclusion arrive cet article", etc.) © C'est quelque chose qui est déjà possible dans ChatGPT d'OpenAI, mais là c'est automatique (je crois qu'il faut un compte au bout d'un moment). Attention quand même ArXiv c'est pas toujours peer-reviewed, donc il n'y a pas que des bons articles !

Enfin, une petite note pour les chercheurs : ATTENTION, certains publishers interdisent l'utilisation des IA Gen. Je vous ai mis quelques liens sur les politiques de certains publishers de mon domaine (informatique/maths) pour l'exemple. Eux l'autorise, mais donnent des instruction sur comment signaler l'usage de ces outils. Pendant un moment s'était posée la questions de savoir s'il fallait citer ChatGPT parmi les auteurs, ici tous répondent que non !

Si vous n'avez pas le droit de les utiliser, attention.

[SLIDE 29/36]

Voilà, j'en ai fini avec les utilisations, notamment pédagogiques et de recherche de ChatGPT, donc avec les aspects "positifs" de ChatGPT. Maintenant, je vais aborder plutôt le(s) revers de la médaille.

J'ai essayé de les synthétiser notamment, en m'appuyant sur les trois piliers du développement durable : l'environnement, la société et l'économie.

- Sur la partie sociétale

+ Premier point c'est la question de la *souveraineté*. Parce que tous ses

outils viennent des GAFAM (Google, Amazon, Facebook, Apple, et Microsoft. Vous avez dû entendre dans les news récemment, qu'il y a une

entreprise française qui s'appelle Mistral, qui a levé des énormes fonds,

parce qu'ils ont un outil qui marche relativement bien. Le problème, c'est

que cet outil n'est pas aussi facile à utiliser que ses concurrents américains. Il n'est pas n'est pas accessible à tout un chacun sur un siteweb évident et gratuitement comme BingChat/ChatGPT. Selon moi, il y a

très fortes chances que, comme pour le Cloud Computing, on aie très rapidement, les GAFAM qui deviennent bien meilleurs que n'importe quelle

entreprise, et que ce soit irrattrapable. En effet, en mettant dès le départ

ces outils, même imparfaits, dans les mains de tout un chacun, les gens les utilisent plus, donc les GAFAM collectent plein d'infos pour s'améliorer (on est fassé à des outils qui apprennent !) et très rapidement n'importe quelle autre entreprise qui arrive sur le marché sera complètement écrasée et dépassée par la puissance ces outils et ne pourra pas monter un modèle économique viable fassé à ces géants.

+ Ensuite, il y a la question de la vie privée, je pense que vous avez entendu beaucoup parler, là aussi. Comme je le disait, il y a deux phases à ce type de modèles : une phase où on l'entraîne, on lui donne à manger beaucoup de textes, et il apprend à un modèle de langage, et puis la phase où on l'utilise, par exemple, dans le modèle ChatGPT, c'est le moment on discute avec lui, et répond à nos questions.

Et bien, il y a des problèmes de privacy aux 2 niveaux.

Au premier, c'est que les société qui ont créé ces outils, elles leur ont donné à mangé des tas et des tas et des tas de données, à peu près tout ce qu'elles ont pu trouver : sites internet, livres, codes sur GitHub, réseaux sociaux, etc. sans regarder si elles avaient le droit d'utiliser ces données. Or vous devez savoir que ce n'est pas parce qu'un document est accessible sur Internet, qu'on a le droit d'utiliser et le réutiliser (propriété intellectuelle, licences). Il y a des gros problèmes là-dessus, d'autant que les entreprises refusent de diffuser la liste des documents qu'ils ont utilisé pour entraîner leur modèle.

Ça va même plus loin que ça : lors de la présentation de dernière version de ChatGPT (4 Turbo), OpenAI a annoncé le "privacy shield" : l'idée est que les modèles d'IA Gen étant probabilistes, il y a très peu de chances qu'ils recrachent /exactement/mot-pour-mot/ un document qu'ils ont vu passer durant leur entraînement. Or, si on recrache un document suffisamment différent, on peut (légalement) argumenter qu'il y a une nouvelle valeur, qu'il s'agit d'une œuvre nouvelle et qu'on rentre dans le domaine du FairUse. L'entreprise OpenAI est tellement sûre d'elle qu'elle propose de vous payer vos frais d'avocat si, en utilisant un document généré par leurs

systèmes, vous vous faites attaquer par un artiste, parce que l'image ressemble trop à une de ses œuvres

Concernant la phase d'utilisation, ce qu'il faut savoir, c'est que, normalement, quand vous créez un compte sur (la version gratuite de) leur

outils, vous acceptez en général les conditions d'utilisation ces outils, qui stipulent que tout ce que vous y entrez (questions, documents, etc.)

pourra être utilisé par l'entreprise pour "améliorer" l'outil (i.e. entraîner la prochaine version de l'outil). Or, imaginons que vous

soyez l'administrateur financier de MinesSE et que vous utilisiez (à votre

insu) le Copilot inclut dans la suite Office365/Excel, vous risquez de

laisser fuiter tous les comptes de l'École ! Il en va de même avec le code

source, dans lequel il y a des fois des mot de passes/clef d'accès à certains services.

+ Ensuite, il y a l'épineuse question de l'*emploi*. Je vous ai mis toute une

série d'articles que j'ai trouvé, qui vont depuis des enquêtes auprès des

travailleurs qui ont peur de perdre leur emploi, jusqu'à l'ONU, qui dit que

ce n'est pas grave car au final plus d'emplois seront créés que détruits.

Dans l'article sur les travailleurs, ils font état de très fortes craintes

du côté des employés, qui ont peur soit de se faire remplacer, soit que

l'emploi change plus vite qu'il ne sont capable de s'adapter, ou encore que

des petits jeunes qui arrivent, et qui savent utiliser ces outils soient

plus efficaces qu'eux, et donc que leurs employeurs leur demande de travailler plus pour compenser, ou d'être payé moins. J'ai entendu ce

genre de discours ici à MinesSE lors des présentations (similaires à celle-ci)

qu'on fait avec Karine.

Dans l'article sur l'Angleterre, il est dit que les IA Gen ne changeront pas

grand chose, il y aura autant de perte d'emploi que de gain et que les

emplois suffisamment peu pour que les employés s'adaptent.

Et puis, il y a des sites de développeurs informatiques et celui de l'ONU,

qui disent que de toute façons, ce n'est pas grave, parce que ça va créer

plus d'emplois que ça ne va en détruire. À mon sens ils oublient un peu au

passage, que tous ces emplois que ça va créer demanderont de nouvelles compétences que ceux qui ont perdu leur emploi ne pourront pas acquérir si vite que ça et se feront donc "piquer" ces emplois par les petits jeunes ou devront faire de longues période de réhabilitation payées (ou pas) par le chômage. Une personne qui était par exemple, je ne sais pas, moi, mécanicienne dans un garage, si on lui enlève son emploi, elle ne peut pas du jour au lendemain devenir ingénieur en informatique.

Et puis, comme vous l'avez sûrement aussi vu dans les news, on parle de création de jobs, mais pour le moment ça a surtout consisté en la créations d'entreprises au Kenya, où des gens étaient payés des misères pour regarder les contenus les plus horrible d'internet afin d'apprendre au système à les filtrer. Je ne suis personnellement pas super convaincu que si ce genre d'emplois venait à ce démultiplier grandement, ce soit un grand progrès pour la société...

Une des grandes questions qui se pose, c'est est ce qu'il y a vraiment des emplois qui vont disparaître ou est ce qu'il est "juste" que nos emplois vont se transformer et qu'on sera assistés par des IA Gen dans nos tâches de tous les jours ?

Ma réflexion personnelle sur le sujet, c'est que pour sûr, nos travaux vont changer. Il suffit de voir aujourd'hui pour nous, enseignants. Nous savons que les élèves utilisent ChatGPT, donc nous sommes obligés d'adapter nos cours ou au minimum notre façon d'évaluer pour aller rechercher de nouveau types de triche.

J'ai eu une lève de troisième, lors de rencontres pour "Réussir Aujourd'hui", qui m'a demandé si le métier de radiologue allait disparaître.

J'ai compris entre les lignes que son père ou sa mère devait être radiologue et avait très peur d'être remplacé·e, puisqu'on entend régulièrement parler de systèmes d'IA qui sont capables de détecter des les maladies mieux que les "meilleurs" médecins.

Ma réponse un peu naïve a été de lui dire que pour l'instant, on a toujours

besoin d'un être humain dans la médecine, que ce soit pour annoncer avec les bons mots à une personne ça maladie, ou pour vérifier qu'il n'y a pas eu d'erreur dans le système, ou même, juste d'un point de vue légal, pour prendre la responsabilité en cas de mauvaise décision. Et que le métier de radiologue n'allait probablement pas disparaître, mais qu'en effet, il allait serait certainement beaucoup changer, puisque les radiologues n'allaient probablement presque plus faire le travail d'analyser les radios.

D'un côté, on pourrait trouver ça intéressant puisque les IA font ce travail très rapidement, ça pourrait permettre à un même radiologue de traiter dix/vingt fois plus de patients, ce qui pourrait être intéressant dans notre monde actuel où l'hôpital est complètement surchargé et qu'il faut souvent plusieurs heures pour gérer un patients qui arrive aux urgences.

D'un autre côté, quand je dis que ma réponse a été naïve, c'est parce que, comme je le disais, tout cela tient pratiquement à un fil humain (annoncer la maladie), ou législatif (prendre la responsabilité), et que tout cela, pour des raisons bassement budgétaires, passe très souvent complètement à la trappe.

On voit bien, par exemple, dans l'aviation (où j'ai travaillé) ou des véhicules autonomes, où il y a deux plus en plus de systèmes informatisés pour remplacer les pilotes ou les contrôleurs, que finalement, on en vient à changer les lois pour que les systèmes informatiques fassent l'intégralité du travail, qui avant était fait par les humains, et que cette fameuse "barrière" de responsabilité ou de l'humanité commence à s'effacer de la loi. Donc il est tout à fait possible, que, dans quelques années, le métier de radiologue disparaisse aussi.

Il y a aussi toute la question des jobs "créatifs". Vous avez sûrement entendu parler de la fronde des acteurs d'Hollywood qui ont peur de se faire piquer leur voix ou leur image. Plus proche de nous, je vous ai cité l'exemple de générer une image de chaise en forme d'avocat ou celui du site pour créer ses propres habits. Pour sûr ce qui est généré n'est (pour

l'instant) pas du même niveau de qualité que ce que pourrait faire un vrai artiste humain. Mais pour les personnes à petit budget, c'est largement suffisant.

Je vois actuellement fleurir des tas de sites comme le suivant : un site qui a entraîné une IA Gen en associant le texte décrivant une entreprise (par exemple pris sur la page about de site web) avec le logo de l'entreprise.

Leur système est maintenant capable de générer des logos très satisfaisants à partir d'une simple description de quelques lignes de ce que fait l'entreprise. Ayant pas mal d'amis informaticiens qui sont auto-entrepreneurs je suis sûr que ce sera beaucoup plus facile et moins cher pour eux d'utiliser un site de ce genre plutôt que de faire appel à un vrai artiste. Il est donc à craindre que, dans les années à venir, toute la frange des artistes qui vivaient de petits boulots aura de plus en plus de mal à trouver des clients.

En conclusion, le scénario le plus probable soit que, dans un premier temps, nos boulots vont changer en étant de plus en plus assistés par IA Gen, mais il est fort probable que, à plus long terme, de nombreux emplois soient entièrement remplacés par des IA Gen, en commençant comme d'habitude par grignoter les emplois des plus précaires.

En gros, la grosse question qui reste, c'est l'échelle de temps de ses transformations.

[SLIDE 30/36]

+ Toujours d'un point de vue social, un autre point de vigilance c'est la cyber-sécurité.

- J'imagine que vous avez tous reçu ces emails qui vous proposent d'agrandir votre pénis ou blablabla. Jusqu'à maintenant, vous saviez que la façon la plus directe pour les détecter, c'était les fautes d'orthographe, notamment en français. Malheureusement, avec l'arrivée des outils d'IA Gen qui sont capables de générer du texte magnifique, et même carrément une pages web avec le "style de" n'importe quelle entreprise, il va être de

plus en plus difficile de détecter les attaque de phishing (vous amener à cliquer sur un lien vérolé) ou de fake news.

- Je ne sais pas si vous avez entendu parler de l'"arnaque au président" ?

L'idée c'est que le hacker se fait passer pour le président du d'entreprise. Il appelle la secrétaire de direction, qui a le pouvoir de

signature au nom du président, et lui dit : "en fait, Germaine, je suis au

Congo actuellement, tu te rappelles, il y avait le contrat avec la société

Tartampion, il fallait le signer absolument avant ce soir 19h, donc je

vais t'envoyer le PDF, il faut que tu le signes, et que tu l'envoie aux

services des financiers. Il faut savoir qu'en France, et même partout dans

le monde, d'ailleurs, un contrat qui est signé a valeur juridique, peu

importe si la personne qui l'a signé a été manipulée pour le signer. Le

contrat lui-même est valable.

Du coup, dans ce genre de cas, si l'arnaqueur arrive à faire que la secrétaire signe un document, ça peut être très grave, donc il y a eu des

très très gros détournements de fonds aux USA il y a quelques années, qui

ont été évités de justesse (arreés par le systèmes SWIFT qui trouvait

bizarre des énormes virements vers les Philippines).

Pour l'instant, cette arnaque fonctionnait via un simple appel téléphonique. Aujourd'hui avec les outils dont je vous ai parlé, où

on peut carrément générer un avatar à partir de la photo du président de

l'entreprise, et même synthétiser une fausse voix du président. On peut

donc imaginer, organiser une web conférence avec la secrétaire et un faux

président.

- Comme je vous ai dit aussi, ChatGPT est capable de générer du code. Des

gens qui se sont amusés à lui faire générer des virus, ou des malwares, à

partir de rien. Il se trouve que comme ces virus malwares sont nouveaux,

la plupart des antivirus ne sont pas capables de les détecter.

- Comme les IA sont entraînées à vous fournir l'information qui vous intéresse le plus, il y a aussi le problème d'enfermer les gens dans des

bulles informationnelles.

- Il est aussi possible (cf. Cambridge Analytica) générer du texte/vidéos très crédibles, pour faire de la désinformation ou de la manipulation ciblée. Je ne sais pas si vous aviez entendu, par exemple que, pour le cas du Brexit, en Angleterre ou l'élection de Trump aux USA, il est avéré que le nombre de personnes qui ont été manipulées sur Facebook à l'époque est plus important que la différence, entre les 2 candidats/réponses. Donc le nombre de personnes manipulées est plus grand que le nombre de personnes qui séparent le nombre de votes qui séparent les deux candidats. Ça ne veut pas dire que Trump n'aurait pas dû être élu, mais ça veut dire que la marge d'erreur était telle que, en fait, que mathématiquement, le système de vote aurait dû se déclarer incapable de définir la vainqueur.

- Encore pire que ça, j'ai vu tout un reportage récemment sur TikTok où les journalistes expliquent comment début TikTok, qui était juste une entreprise chinoise basique a été petit à petit phagocytée par le gouvernement chinois a commencé (investissement, puis a placé des gens du gouvernement à l'intérieur, et maintenant, l'entreprise est entièrement dirigée par le Parti communiste chinois). L'explication qu'ils donnent, c'est qu'en fait, ils se sont rendus compte que le principe de fonctionnement même de TikTok, c'est de connaître vos goûts pour vous proposer que des vidéos qui vous intéressent. Donc ils sont rendus compte, que quand les gens commencent à regarder un peu trop des vidéos sur les ouïgurs ou des choses un peu gênantes pour le gouvernement, il suffit d'inonder le système avec des jolies vidéos, par exemple des vidéos de chats, et que les gens vont détourner leur attention des sujets qui fâchent ; il n'y a même plus besoin de réfléchir à des campagnes complexes de manipulation les gens en générant des fausses information cohérentes/plausibles ; il suffit juste de les faire regarder autre chose. Vu la surcharge informationnelle qu'on subit actuellement, notre cerveau sature un peu et est trop content d'aller voir des vidéos qui lui font produire de la dopamine. Et comme les outils d'IA savent parfaitement (et même mieux que vous) vos

goûts il suffit d'un clic sur un bouton au "directeur de TikTok"
pour que
l'IA envoie les-videos-qui-vont-bien à chacun d'entre vous pour que
vous
arrêtiez de regarder les sujets qui fâchent.

[SLIDE 31/36]

+ Concernant *l'éducation*, pas besoin de vous faire un dessin. Vous
avez
certainement déjà eu ou entendu parler de cas ici même, à l'école des
mines,
d'élèves qui utilisent de plus en plus ces outils et à des fins de
travailler le moins possible, et pas à des fins d'améliorer leur
travail.

- Si Vous vous êtes intéressés au sujet, vous savez aussi que
certaines des
formes d'évaluation que vous utilisez jusqu'à maintenant ne
fonctionneront
plus avec ces outils indétectables, et qu'il va falloir reréfléchir à
vos
cours.

- Il y a aussi la question de *l'équité* faite à ces outils. Entre
les
élèves qui vont pouvoir se payer la dernière version de ChatGPT
ultra
performante, à 30 euros par mois, et ceux qui ne peuvent pas se le
payer,
qui n'utiliseront que la version gratuite, plus limitée.

+ En termes de recherches, on a déjà entendu parler, d'articles générés
automatiquement, qui étaient validés dans des conférences. Souvent
c'était
parce que les conférences ne mettaient pas en place de reviewing
correct. Mais désormais on peut faire écrire un article complet à
ChatGPT,
mais aussi, plus sournoisement, faire générer les données qui vont
bien pour
aller dans le sens de l'hypothèse qu'on veut illustrer, ce qui est
encore
plus dangereux.

- S'est posé aussi la question de savoir si ChatGPT était un auteur
ou
pas. J'en ai déjà parlé. En gros le consensus est que non, mais
qu'il faut
quand même par transparence signaler les usages des IA Gen.

+ Une autre grosse question est celle *environnementale*. Encore une
fois, il
faut considérer les deux phases des IA Gen : entraînement et
utilisation.

- Pour l'entraînement il faut faire tourner le système sur de gros
ordinateurs. Cela consomme une quantité incroyable d'énergie, mais
aussi

de tout un tas d'autres ressources, comme de l'eau, les terres rares, etc.

En termes d'énergie on estime que la version 3.5 de GPT a nécessité l'équivalent de 205 A/R Paris/NYork en avion.

- Pour ce qui est de l'utilisation, c'est plus compliqué à estimer, car un minimum de serveur tournent quoi qu'il arrive pour offrir le service. Le coût par requête/utilisateur dépend donc du nombre de personnes qui utilisent le service. Par ailleurs, le nombre de serveurs s'adapte avec le nombre d'utilisateurs.

Laurent Avaro, de l'école, du CIS, me disait qu'il a trouvé des chiffres récemment qui comparaient les IA Gen et le cerveau humain. En gros pour une même requête, ChatGPT consommerait l'équivalent de 1 à 10 Wh, alors que le cerveau humain, pour n'aurait besoin que 10 puissance moins 4 Wh. Il y a un donc facteur quasiment 10 000, entre la consommation énergétique qu'on peut répondre à une question par ChatGPT ou par un cerveau humain. Ça fait réfléchir quand on pense que beaucoup d'entreprises rêvent de remplacer une grande partie de leurs employés par des IA Gen pour économiser côté financier.

[SLIDE 32/36]

Il y a aussi des problèmes intrinsèques au fonctionnement même de ces modèles.

Je vous en ai illustré deux ici.

Le premier, c'est ce qu'on appelle le "model decay", qui est lié aux changements de distribution probabilités.

Je vais vous donner un exemple, juste pour illustrer, admettons que la distribution qui est en bleu ici soit, par exemple, la question que je pose à ChatGPT, c'est, est-ce que demain, je dois mettre ma robe verte, ou mon pantalon noir ? Et bien, la distribution de probabilités bleue représente les différentes réponses possibles, là où il y a le 4, ça serait par exemple que je dois mettre ma robe verte, qui serait la réponse la plus probable aujourd'hui. Mais on sait que, en fait, tout ce qui concerne la mode, ça évolue dans le temps, donc en fait, la distribution de probabilité en bleu c'est celle qui existait dans les

données collectées quand ChatGPT a été entraîné. Mais il est fort probable que dans 10 ans, ou même dans 6 mois, en fait, la mode a changé, et que la nouvelle distribution, en fait, ça serait plutôt la distribution jaune qui est ici, qui est décalée vers la droite et très resserrée. Et bien, du coup, quand je vais poser la question à ChatGPT dans 6 mois, il devrait répondre la réponse la plus probable, là où il y a le 7 sur la distribution jaune, ça serait par exemple, justement, mettre le pantalon noir. Mais en fait comme lui il est resté sur la distribution bleue, il continuera de répondre qu'il faut mettre la robe verte.

Là, je vous ai trouvé un exemple un peu bête sur la mode, parce que c'est quelque chose, qui change vite de façon évidente, mais en fait il y a beaucoup de choses dans la vie qui changent progressivement, sans vraiment qu'on s'en rende compte. Par exemple, la réponse la plus probable à un problème scientifique peut évoluer avec le temps et l'acquisition de connaissances ou la proportion de textes racistes sur internet.

Pour corriger le problème il faudrait ré-entraîner le modèle régulièrement et on a vu que les coûts économique et écologique seraient drastiques.

La deuxième chose, c'est ce qu'on a appelé avec Karine, la "consentignité" (enfin c'est surtout Karine, moi, j'appelais ça plutôt l'"entraînement en boucle fermé", mais c'est plus technique, moins parlant). L'idée est illustrée ici. Je ne sais pas si vous avez déjà vu ces images ? L'idée c'est que, il y a des gens (des chercheurs de chez NVIDIA la société qui fait les cartes graphiques très utilisées pour accélérer les calculs d'entraînement des IA) qui s'étaient amusés à donner plein de photos de stars américaines à un outil (pas ChatGPT, plutôt de style Dall-E), et après, il lui disaient : "maintenant, vas-y, génère une nouvelle image", et le système générait des visages, comme ceux qu'on voit ici, à gauche, qui sont des visages de personnes qui n'existent pas, mais qui sont très, très réalistes, parce que le système a compris comment est composé un visage : deux yeux, un nez, une bouche, etc.

Dans le présent travail, d'autres chercheurs, ce qu'ils ont essayé de faire, c'est OK, on a notre 1ère génération d'images ici à gauche, qui sont

générées. Admettons qu'on oublie toutes les images originales (les vrais photos de vrais stars). Qu'est-ce qui se passe si on utilise que ces photos générées à l'étape 1, pour entraîner un modèle sur ces images. Puis, on recommence sur les images générées à l'étape 2, on génère des images à l'étape 3, etc.

On voit qu'à l'étape 9, il y a des artefacts absolument ignobles qui apparaissent : les cheveux de la demoiselle en haut, ou des traits dans la barbe du monsieur en bas.

Si on regarde à l'étape 7 ou même l'étape 3 on voit en fait que ces artefacts existaient déjà de manière embryonnaire. La femme, dans l'image en haut, on voit qu'elle est quand même relativement ridée alors qu'elle paraît jeune.

Tout ça, en fait, c'est lié au fait que les modèles ils généralisent le concept qu'ils apprennent. Ils comprennent bien qu'un visage a 2 yeux, etc. mais aussi que la peau c'est un pattern assez compliqué avec du grain. Et ils le simplifie en y associant un certain bruit paramétrique. Et c'est ce bruit qui, de génération en génération s'amplifie (typique accumulation d'erreurs d'un modèle entraîné sur un modèle, etc.).

Dans l'exemple que je donne ici ce n'est pas très grave. C'est plutôt rigolo à voir. Mais dans la vraie vie, c'est un gros gros problème. Par exemple, il est dit qu'aujourd'hui environ 30% du contenu sur Internet a été plus ou moins généré par des outils IA Gen. Donc ça veut dire que les systèmes d'IA Gen comme ChatGPT qui seront entraînés dans 1 ou 2 mois / dans 2 ou 2 générations vont être entraînés sur une grosse partie de contenu qui aura déjà été générée par les systèmes de la génération d'avant. Et cela va évidemment renforcer les biais qui existaient déjà dans les modèles d'avant. Or, on sait que sur Internet (notamment, sur Twitter) , une grosse proportion des textes est fortement raciste. Donc ce genre de biais ne va faire qu'empirer au fur et à mesure des générations des outils d'IA Gen.

[SLIDE 34/36]

J'en ai fini avec le cœur de la présentation.

Pour ceux que ce sujet intéresse et qui voudraient poursuivre les discussion, il

faut savoir que j'ai créé une mailing list (elle n'est pas limitée à Mines SE, elle comprend aussi des gens de l'IMT que j'ai rencontré à MÉDIANE en 2023.

Pour s'inscrire, il suffit d'envoyer un mail à l'adresse indiquée avec pour sujet "Subscribe".

[SLIDE 35/36]

Enfin sur ce dernier transparent, je vous ai mis quelques liens bibliographiques:

- GPT France pour des news en français sur les IA Gen.
- Un article qui discute la notion de plagiat face à ces nouveaux outils d'IA Gen. Pour info, l'IMT-BS a créé une "charte éthique" spéciale IA Gen. Elle a circulé un peu à Mines SE et on s'est dit qu'elle était trop vague et pas assez contraignante pour qu'on en fasse une pareil.

La DF s'est surtout concentrée à ré-écrire les articles du "code de l'enseignement" pour qu'il couvre les IA Gen.

- Enfin, il y a deux sites qu'on vous a mis ici, sur ChatGPT et les enseignants. Il y a des versions françaises et anglaises pour les 2 documents.