

## Du bon usage de l'intelligence artificielle pour l'accessibilité et la découvrabilité des archives

Léon Robichaud  
Département d'histoire  
Université de Sherbrooke

Depuis quelques années, je suis, comme d'autres historiens dont le travail repose en grande partie sur l'informatique, confronté à l'émergence de l'intelligence artificielle (ou IA). En moins de cinq ans, l'IA s'est imposée comme le principal moteur d'une nouvelle révolution informatique qui touche tous les domaines : travail, divertissement, consommation, transport, gestion, médecine, et bien d'autres encore. Une diversité étonnante d'applications est désormais disponible<sup>1</sup>, mais le paysage est toujours dominé par ChatGPT, développée par OpenAI<sup>2</sup>. Cette application profite de la notoriété médiatique acquise lors de son lancement en novembre 2022 et est devenue synonyme d'IA pour bien des gens.

Comme dans toute bulle alimentée par la nouveauté et par les médias, les développements se poursuivent à un rythme accéléré. Selon le rapport annuel de l'Université Stanford sur l'IA, presque tous les indicateurs sont en hausse : performance, intégration dans la vie quotidienne, investissements, recherche et développement, tant aux États-Unis qu'en Chine, l'autre poids lourd dans le domaine<sup>3</sup>. Le rapport signale la persistance des inégalités géographiques et sociales dans le rapport à l'IA et souligne que les questions de développement « responsable » ne sont pas encore résolues.

Cette dichotomie de l'IA teinte les discours. D'une part, elle est présentée comme la solution à tous les problèmes : détecter les cancers, embaucher la bonne personne, planifier les vacances, se laisser conduire par une voiture sans chauffeur. D'autre part, le spectre de l'IA plane sur la société : perte

---

1 TAAFT, *There's an AI for that*, Bucarest, mise à jour continue, [<https://theresanaiforthat.com/>]. En juin 2025, on comptait plus de 36 000 applications différentes pour un peu plus de 13 500 types de tâches.

2 Kayla Zhu, « Ranked: Most Popular AI Tools by Monthly Site Visits », *Visual Capitalist*, 24 mars 2025, [<https://www.visualcapitalist.com/ranked-most-popular-ai-tools-by-monthly-site-visits/>].

3 Yolanda Gil et Raymond Perrault, *Artificial Intelligence Index Report 2025*, Stanford University Human-Centered Artificial Intelligence, 2025, [<https://hai.stanford.edu/ai-index/2025-ai-index-report>].

d'emplois, reproduction des biais systémiques, destruction de la propriété intellectuelle, plagiat, éthique boiteuse, fausses réponses, fausses nouvelles. Certains textes planchent sur les deux tableaux, tel cet article du *Devoir* qui commence par : « L'intelligence artificielle fait trembler les uns et rêver les autres<sup>4</sup>. » Toute analyse de l'IA implique de prendre en considération les risques d'un usage optimiste, voire naïf, tout en explorant son potentiel pour maximiser nos efforts tout en maintenant l'intégrité de nos fonctions.

Dans le cadre de cet article, je vais d'abord vous offrir une définition de l'IA en la comparant à l'informatique conventionnelle. Par déformation professionnelle, je ferai quelques rappels concernant l'histoire de l'IA, ce qui nous permettra d'en examiner les cycles de développement. Je vais ensuite présenter quelques dérives actuelles et potentielles avant de proposer quelques pistes pour orienter notre prise en charge de la technologie. À partir du projet Nouvelle-France numérique, j'espère décrire une approche concrète avec quelques réalisations en main et des initiatives à mettre en branle prochainement. Je mettrai l'accent sur la reconnaissance des écritures manuscrites, l'aide à l'indexation et l'aide à la recherche, trois domaines pour lesquels je suis d'avis que nous pouvons faire un bon usage de l'IA.

Après avoir rédigé mon plan, j'ai posé la question suivante à deux IA généralistes (ChatGPT et Perplexity) : « Comment parler d'intelligence artificielle à des archivistes? ». Hormis un ton de type marketing chez ChatGPT, j'ai constaté plusieurs recoupements intéressants avec mon plan et je suis certain que certains thèmes proposés vous intéresseraient, que ce soit l'analyse prédictive des tendances ou la restauration de fichiers numériques endommagés. Ils se situent toutefois hors de mes compétences et je m'en tiendrai donc à mon plan de match initial. J'ai aussi choisi d'agrémenter ce texte de références relativement pérennes (du moins, espérons-le) et que vous pourrez consulter à loisir.

---

4 Roxane Léouzon, « L'intelligence artificielle, partenaire de création scénique », *Le Devoir*, 20 mai 2025, [<https://www.ledevoir.com/culture/881536/intelligence-artificielle-partenaire-creation-scenique>].

## Définir l'IA

L'expression même « intelligence artificielle » est devenue un fourre-tout pour désigner une diversité de systèmes informatiques pouvant fonctionner sans instructions prédéfinies. De tels systèmes peuvent réaliser des tâches spécifiques (IA étroite ou narrow IA) ou s'attaquer à une grande variété de tâches intellectuelles de manière créative (IA généraliste ou generalist AI). En ce moment, le G de IAG signale plutôt une IA générative, capable de répondre à une grande variété de questions posées en langage naturel et de résoudre des problèmes en mimant l'intelligence humaine à partir des données utilisées pour son entraînement. Quant à la Super IA qui pourrait réellement remplacer l'intelligence humaine, c'est un objectif de plusieurs acteurs dans le domaine<sup>5</sup>.

Le dictionnaire *Usito* définit l'intelligence artificielle comme une : « branche de l'informatique consacrée aux théories et aux techniques visant à la création de systèmes ou de machines exploitant des fonctions simulant l'intelligence humaine, notamment grâce à des modèles et des algorithmes<sup>6</sup>. » Comment cela diffère-t-il de l'informatique traditionnelle? De ses origines en cryptographie et en calculs balistiques, jusqu'aux logiciels modernes de traitement de texte, de traitement de l'image, etc., l'informatique traditionnelle exécute des instructions prédéfinies à partir d'ensembles de données fermées et ne peut s'en écarter. Les actions imprévues et les données non conformes peuvent ainsi causer des erreurs, voire faire planter un système. Les applications en IA reposent elles aussi sur du code rédigé par des programmeurs (en général de jeunes hommes avec les biais que cela implique), mais sont conçues pour que le système définisse les règles à partir des données d'entraînement et puisse ainsi résoudre les problèmes présents dans de nouvelles données qui leur sont soumises, d'où les concepts d'apprentissage machine, d'adaptabilité et de reproduction de la pensée humaine.

---

5 Sanksshep Mahendra, « Types of AI: Narrow, General, and Super AI », *Artificial Intelligence Plus*, 18 novembre 2024, [<https://www.aiplusinfo.com/blog/types-of-ai-narrow-general-and-super-ai/>].

6 « Intelligence », *Le Dictionnaire Usito*, Sherbrooke, Université de Sherbrooke, [<https://usito.usherbrooke.ca/d%C3%A9finitions/intelligence>].

Si je prends un exemple, dans un logiciel traditionnel ayant pour objectif de reconnaître des tables de cuisine dans des illustrations, un programmeur pourrait définir une règle simple : une surface rectangulaire et plane posée sur quatre pattes<sup>7</sup>. Cette définition n'inclue pas les tables rondes ou ovales et n'exclut pas les bancs sans dossier. Des règles additionnelles seraient nécessaires pour que les tables posées sur des tréteaux, des tables placées à l'envers, etc., soient correctement identifiées. L'apprentissage machine apprend autrement. Les règles sont définies en accumulant les exemples et en appliquant ces critères aux autres cas. À partir d'une série d'images contenant des tables, l'algorithme passe et repasse sur l'échantillon jusqu'à ce qu'un très grand pourcentage des tables soient détectées. C'est cette étape d'entraînement qui nécessite énormément de ressources en puissance computationnelle et évidemment en électricité. Au lieu de prendre le temps de rédiger toutes les règles possibles qui permettent d'identifier une table à partir d'une définition formelle, le système les définit en voyant quelles règles permettent d'identifier le maximum de tables. Les règles utilisées ne sont toutefois pas explicites et peuvent parfois mettre de l'avant un élément inattendu. Ainsi, si toutes les tables de l'échantillon d'entraînement ont un pot de fleurs, toute structure sur laquelle on trouve un pot de fleurs pourrait être identifiée comme une table<sup>8</sup>.

Cette démarche n'est pas encore identique à la pensée humaine, comme les chercheurs l'ont constaté en entraînant des IA à lire l'heure sur une horloge analogique<sup>9</sup>. L'IA peut développer des règles pour calculer l'heure, mais s'adapte difficilement aux changements de style de face. Les humains sont donc

---

7 Pour un exemple du fonctionnement de l'IA, voir Ben Brubaker et Mark Belan, « How Can AI ID a Cat? An Illustrated Guide », *Quanta Magazine*, 30 avril 2025, [<https://www.quantamagazine.org/how-can-ai-id-a-cat-an-illustrated-guide-20250430/>].

8 Un algorithme entraîné à différencier les loups et les chiens avait plus de chance de classer les chiens parmi les loups lorsqu'il y avait de la neige, car les images d'entraînement présentaient les loups dans un contexte hivernal. Marco Tulio Ribiero, Sameer Singh et Carlos Guestin, « “Why Should I Trust You?”. Explaining the Predictions of Any Classifier », *KDD '16: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (2016), p. 1142-1143.

9 Martin Anderson, « AI's Struggle to Read Analogue Clocks May Have Deeper Significance », *Unite AI*, 19 mai 2025, [<https://www.unite.ai/ais-struggle-to-read-analogue-clocks-may-have-deeper-significance/>].

en mesure de définir un modèle abstrait de l'horloge et des aiguilles, ce qui nous permet de lire l'heure même lorsque la face est complexe ou déformée. Les bases sur lesquelles une IA s'appuie pour résoudre les problèmes sont simplement celles qui lui permettent d'obtenir la bonne réponse la plupart du temps.

Il est pertinent de rappeler qu'en IA, l'identification d'une table ou la lecture de l'heure n'est pas, comme ce serait la pratique dans une application traditionnelle, vrai ou fausse (0 ou 1), mais une probabilité. Le système peut donc toujours donner une réponse, même si celle-ci est fautive selon des critères objectifs. Lorsque l'ordinateur Watson d'IBM a vaincu deux concurrents humains au jeu télévisé *Jeopardy* en 2011, une réponse complètement erronée a fait les manchettes. Dans la catégorie « U.S. Cities », l'expression était « Its largest airport is named for a World War II hero; its second largest, for a World War II battle »; Watson répondit Toronto au lieu de Chicago<sup>10</sup>. Watson avait néanmoins signalé son faible niveau de confiance (14%) en ajoutant 5 points d'interrogation à sa réponse : « Toronto????? ». En dépit de cette bourde, IBM avait démontré au grand public qu'il était possible de fournir à l'ordinateur une grande quantité de faits, de les mettre en relation et de régurgiter très rapidement une réponse selon la syntaxe spécifique à *Jeopardy*. L'IA avait déjà parcouru un très long chemin depuis ses débuts 60 ans plus tôt.

## Retour historique<sup>11</sup>

Au cours de la Seconde Guerre mondiale, la science a bénéficié d'investissements massifs dans chacun des camps dans le but d'obtenir la victoire. Du côté des Alliés, ces investissements ont notamment permis

---

10 Denyse O'Leary, « Why did Watson think Toronto was in the U.S.A.? », *Mind Matters*, 3 août 2019, <https://mindmatters.ai/2019/08/why-did-watson-think-toronto-was-in-the-u-s-a/>; Associated Press, « IBM's Watson trounces humans at Jeopardy », *CBC News*, 16 février 2011, [<https://www.cbc.ca/news/entertainment/ibm-s-watson-trounces-humans-at-jeopardy-1.1028918>]. Le lieutenant-commandant Edward O'Hare est devenu le premier as de l'aviation américaine pendant la Seconde Guerre mondiale lorsqu'il a attaqué une formation de neuf bombardiers japonais et en a descendu cinq. Midway est une grande bataille navale entre les États-Unis et le Japon, un tournant dans la Guerre du Pacifique.

11 La trame chronologique s'appuie en grande partie sur Andrew Spoeth, « AI Timeline : News and Events from 1950 - 2024 that have defined Artificial Intelligence », *Time. Graphics*, [<https://time.graphics/line/809611>].

d'améliorer les trajectoires balistiques des tirs en mer et à décrypter les messages encodés par la machine Enigma développée en Allemagne nazie<sup>12</sup>. L'un des mathématiciens les plus brillants de son époque, Alan Turing, avait justement travaillé à ce projet de cryptographie. Dès 1936, il avait conçu une « Universal Machine », démonstration théorique et logique des concepts nécessaires à l'informatique moderne. Quelques mois après la fin de la guerre, il définit l'architecture d'un ordinateur, le Automatic Computing Machine, mais les Américains vont finaliser plus rapidement le premier ordinateur fonctionnel dès décembre 1945, le « Electronic Numerical Integrator And Computer » ou ENIAC. Turing ne s'arrête pas aux machines qui exécutent des programmes chargés en mémoire. En 1950, il publie un article<sup>13</sup> qui commence avec la phrase : « I propose to consider the question, 'Can machines think?'<sup>14</sup> ». C'est aussi dans cet article qu'il propose le « célèbre » test de Turing, soit des questions qui devaient permettre de déterminer si nous avons affaire à un humain ou à un ordinateur.

La carrière de Turing est toutefois interrompue par les politiques de répression de l'homosexualité. Condamné au criminel, il voit son accès aux agences de sécurité britanniques révoqué. Il met fin à ses jours en 1954 et c'est sans lui que d'autres chercheurs se rassemblent à Dartmouth en 1956. C'est dans l'invitation à cet événement, transmise en 1955, que les organisateurs choisissent de ne pas utiliser les expressions alors en vogue, telle que « thinking machine », pour privilégier « artificial intelligence »<sup>15</sup>. Après un été à discuter d'enjeux théoriques, les chercheurs tentent de les appliquer grâce aux ordinateurs disponibles à l'époque. Après deux décennies d'avancées, notamment en traitement du langage naturel avec l'agent conversationnel Eliza, le mouvement s'essouffle et vit un « premier hiver » de l'IA (1974-1980). Lors d'une brève résurgence (1980-1987), la reconnaissance du langage et les premiers véhicules

---

12 André Mondoux, *Histoire sociale des technologies numériques de 1945 à nos jours*, Québec, Éditions Nota bene, 2011.

13 Alan M. Turing, « Computing Machinery and Intelligence », *Mind, A Quarterly Review of Psychology and Philosophy*, vol. LIX, no 236 (oct. 1950), p. 433-460, [<https://archive.org/details/MIND—COMPUTING-MACHINERY-AND-INTELLIGENCE>].

14 Turing, « Computing Machinery and Intelligence », p. 433.

15 « Dartmouth Workshop », Wikipedia, [[https://en.wikipedia.org/wiki/Dartmouth\\_workshop](https://en.wikipedia.org/wiki/Dartmouth_workshop)].

autonomes font émerger des applications pratiques. Les développements sont néanmoins limités par les contraintes de la puissance de calcul disponible et par le peu de données d'entraînement à la disposition des chercheurs.

Pendant un 2<sup>e</sup> hiver (1987-1995), des chercheurs se tournent vers des applications démontrant la puissance et le potentiel de l'IA. Le jeu d'échecs, dont la maîtrise est associée à un niveau intellectuel élevé, devient la cible à atteindre pour démontrer que l'ordinateur peut être plus efficace que l'humain. En 1985, des étudiants de l'université Carnegie-Mellon se lancent dans la création d'un ordinateur pouvant battre des grands maîtres<sup>16</sup>. Intégrés par IBM, leurs travaux mènent à l'ordinateur DeepBlue et le défi est lancé au champion du monde Gary Kasparov. Le premier duel entre la machine et le grand maître est remporté par ce dernier en 1996. IBM revient à la charge avec une version améliorée, « Deeper Blue » et l'année suivante, l'ordinateur éblouit le monde en remportant la ronde de 6 matchs par 3 victoires contre 2 défaites et 1 nulle. Malgré des controverses<sup>17</sup>, la possibilité pour un ordinateur de battre un humain à son propre jeu nous lance dans une nouvelle ère. Les chercheurs demeurent toutefois confrontés à un défi de taille. Le jeu d'échecs demeure un ensemble fermé de possibilités avec des options bien connues et documentées, comment obtenir des données d'entraînement pour des applications généralistes?

La solution à ce problème provient en grande partie de l'ouverture d'Internet (jusque là réservé à une constellation d'agences gouvernementales, des universités et de grandes entreprises), aux applications commerciales et au grand public. Un nouveau protocole de publication (hypertext transfer protocol ou http) facilite aussi la diffusion de contenu sur le « world-wide web ». Contrairement aux environnements

---

16 « Matches Deep Blue contre Kasparov », *Wikipedia*,  
[[https://fr.wikipedia.org/wiki/Matches\\_Deep\\_Blue\\_contre\\_Kasparov](https://fr.wikipedia.org/wiki/Matches_Deep_Blue_contre_Kasparov)].

17 Le coup déterminant a-t-il été suggéré par un grand maître qui assistait l'ordinateur? S'agissait-il d'un bogue et donc d'un coup imprévisible? Le refus par IBM de rendre public les journaux de traitement n'a pas permis de mettre fin aux rumeurs.

fermés privilégiés par AOL, CompuServe et Prodigy<sup>18</sup>, de multiples startups misent sur l'ouverture du protocole que son créateur, Tim Berners-Lee, avait délibérément laissé dans le domaine public. « Nous » (au sens très large du terme) nous sommes alors lancés dans la création de contenu, d'abord sous forme de texte, puis d'images et enfin d'enregistrements audio et vidéo. Cette pléthore de blogues et de photos de voyage a créé une première récolte à moissonner pour entraîner des IA, mais ce n'est rien par rapport aux quantités astronomiques disponibles à partir de 2004 avec le projet Google Print, devenu Google Books, le lancement de YouTube (2005), celui de Twitter (2006), de l'ouverture au public de Facebook la même année et enfin l'apparition du iPhone (2007).

La quantité de matériel « né numérique » explose et est hébergée sur les plateformes d'entreprises qui ont pu s'en servir pour développer des IA. Nous avons participé à l'entraînement d'algorithmes chaque fois que nous avons rédigé un texte (même de 140 caractères), identifié une personne sur une photo, posé une question à Siri ou à Alexa. Nous contribuons à améliorer la reconnaissance d'images chaque fois que nous identifions les feux de circulation, les motocyclettes, les traverses de piétons, les autobus, etc. sur les applications de reCAPTCHA. Les ouvrages et les articles mis en ligne, les encyclopédies ouvertes, les sites de référence, tout est mis à profit pour enrichir et entraîner les modèles de langage naturel. Et c'est sans compter tous les ouvrages, les images, la musique et les films sous droit d'auteur qui ont été piratés par ces entreprises<sup>19</sup>. Selon la logique qui anime les principaux acteurs du domaine, les services reçus gratuitement sont une récompense suffisante pour les données que nous avons fournies et les tâches de reconnaissance que nous avons réalisées.

---

18 Alina Selyukh, « The Big Internet Brands Of The '90s — Where Are They Now? », *NPR. All Tech Considered. Tech, Culture and Connection*, [<https://www.npr.org/sections/alltechconsidered/2016/07/25/487097344/the-big-internet-brands-of-the-90s-where-are-they-now>].

19 Gil Appel, Juliana Neelbauer et David A. Schweidel, « Generative AI Has an Intellectual Property Problem », *Harvard Business Review*, 7 avril 2023, [<https://hbr.org/2023/04/generative-ai-has-an-intellectual-property-problem>].

Les obstacles qui ont ralenti le développement de l'IA sont, pour l'instant, levés. La puissance de calcul, l'efficacité des algorithmes et la masse de données d'entraînement ont permis des avancées plus rapides que prévu. La nature même de ces mégadonnées, ou *big data*, pose toutefois des défis qui ne sont pas toujours pris en compte. Celles-ci peuvent donner l'impression d'être neutres, car elles comprennent l'ensemble des connaissances humaines téléchargeables. Elles sont en fait le reflet de nos forces, de nos faiblesses, et de nos travers. L'IA conversationnelle de Microsoft lancée en 2016 a appris en moins de 24 heures d'interaction avec les usagers de Twitter à émettre des commentaires racistes, sexistes et antisémites<sup>20</sup>. Une série d'autres applications visant à prédire la criminalité, à embaucher des employés et à évaluer les problèmes de santé ont démontré qu'elles reproduisaient les biais inhérents à leurs données d'entraînement ainsi que ceux de leurs développeurs<sup>21</sup>. Certains correctifs ont été mis en place, mais la séquence de propositions ridicules, aberrantes, inquiétantes ou potentiellement fatales se poursuit.

## Les dérives

Une IA peut difficilement développer un sens de l'éthique à partir des données et les chercheurs ont dû ajouter des garde-fous pour éviter les dérives. La solution la plus simple avec la plus haute probabilité de « succès rapide » est rarement éthique. Des choix doivent être faits et ces choix vont refléter les valeurs d'une époque.

N'oublions pas que nous ne comprenons pas vraiment comment les machines s'entraînent elles-mêmes. Nous ne voyons que les résultats. Ce que les spécialistes aiment appeler des « hallucinations »

---

20 Elle Hunt, « Tay, Microsoft's AI chatbot, gets a crash course in racism from Twitter », *The Guardian*, 24 mars 2016, [<https://www.theguardian.com/technology/2016/mar/24/tay-microsofts-ai-chatbot-gets-a-crash-course-in-racism-from-twitter>].

21 Tim Lau, « Predictive Policing Explained », *Brennan Center for Justice*, 1<sup>er</sup> avril 2020, [<https://www.brennancenter.org/our-work/research-reports/predictive-policing-explained>]; Midhat Tilawat, « Rapport sur les biais de l'IA 2025 : la discrimination des LLM est pire que vous ne le pensez! », *AllAboutAI*, 1<sup>er</sup> mai 2025, [<https://www.allaboutai.com/fr-fr/ressources/statistiques-ia/biais-intelligence-artificielle/>].

risquent d'empirer, car les IA s'entraînent désormais aussi sur du contenu qu'elles génèrent<sup>22</sup>. Des mécanismes pourraient-ils être mis en place pour réduire ou jeter un peu plus de lumière sur la boîte noire? Mais si on tourne la question autrement, les usagers veulent-ils vraiment connaître les détails ayant mené à leur réponse?

À défaut de révéler le niveau de confiance d'une réponse, l'affichage des sources d'informations serait déjà une pratique recommandable. L'application Perplexity, que m'a fait découvrir Sophie St-Cyr, de l'Université de Sherbrooke, présente sa liste principale de références avec d'autres ressources « pour en savoir plus », ce qui en fait un outil à privilégier. On peut aussi demander à ChatGPT d'ajouter des références. Celles-ci comprennent des ouvrages réels, mais aussi de fausses références<sup>23</sup>. Pour les non-initiés, la référence imaginaire semble valide. Est-elle ajoutée comme indice de plagiat à l'intention des professeurs ou s'agit-il d'un amalgame d'informations éparses regroupées en un titre d'ouvrage? L'opacité règne comme d'habitude.

L'omniprésence des IA remet même en cause les fondements du web tel que nous le connaissons depuis les années 1990. Lorsque Tim Berners-Lee conçoit son protocole hypertexte, il espère que les usagers pourront naviguer de page en page et de site en site pour mettre en lien les différentes informations. Avec l'agrégation des contenus du web dans de grands modèles de langage, les moteurs de recherche répondent à nos requêtes à partir d'une IA qui a déjà fait ce travail de mise en relation des informations pour nous. Le trafic est alors limité aux résultats offerts par les agents conversationnels. Les usagers peuvent très bien se contenter de cette manière d'obtenir l'information, mais le phénomène peut mener au « Zero-click internet ». Le chef de la direction de Cloudflare a lancé l'alerte, remarquant que

---

22 Cade Metz et Karen Weise, « L'IA s'améliore, mais hallucine plus fort », *Le Devoir*, 7 mai 2025 (publié originalement dans le *New York Times*), [<https://www.ledevoir.com/societe/876925/ia-ameliore-mais-hallucine-plus-fort?>].

23 ChatGPT a généré une référence à une publication inexistante de Bertrand Desjardins sur les Filles du roi en Nouvelle-France. Le recours à l'IA dans la formation universitaire soulève de nombreux enjeux qui dépassent le cadre de cet article. Il s'agit d'un défi majeur pour l'enseignement universitaire, tant pour définir la relation entre compétences et connaissances, que pour établir la nature du raisonnement, de l'analyse et de la rédaction originale.

les clics additionnels ont été réduits de moitié sur Google. Le modèle d'affaires fondé sur le trafic et sur la publicité pourrait alors s'effondrer<sup>24</sup>. Peu de gens pleureront la diminution des pièges à clics (*clickbait*) et de la désinformation, mais la fiabilité des IA génératives n'est pas assurée. Peu après la sortie de DeepSeek, la principale IA développée en Chine, les usagers se sont précipités pour vérifier comment seraient gérées des questions relatives à la place Tiananmen, Taïwan, etc. Sans surprise, la version en ligne suivait essentiellement la ligne du Parti communiste chinois<sup>25</sup>. Récemment, à la suite d'une intervention inexplicée, l'IA contrôlée par Elon Musk (Grok), a inséré pendant plusieurs heures des références au prétendu génocide des fermiers blancs en Afrique du Sud dans des réponses sans lien avec le sujet<sup>26</sup>. De telles occurrences sont bien plus sérieuses que les atermoiements de Christian Rioux à propos des IA qui refusent de générer des blagues à propos des femmes et des minorités<sup>27</sup>.

À terme, la manipulation des résultats n'est qu'un problème parmi d'autres. Dès 2023, Yoshua Bengio et quelques autres spécialistes ont lancé l'alerte quant aux dérives possibles d'une IA qui atteindrait les capacités de l'intelligence humaine quelque part entre 2028 et 2043. Bengio a notamment insisté pour que les développements soient encadrés et que les systèmes soient analysés afin de mieux comprendre les risques pour la sécurité du public. Une telle démarche est, à son avis, nécessaire « autant pour tirer profit des avantages de l'IA que pour protéger l'humanité<sup>28</sup>. » Aux États-Unis, la

---

24 Rob Thubron, « Cloudflare CEO warns AI and zero-click internet are killing the web's business model », *Techspot*, 9 mai 2025, [<https://www.techspot.com/news/107859-cloudflare-ceo-warns-ai-zero-click-internet-killing.html>].

25 Antoine Gautherie, « Ces questions politiques sensibles auxquelles DeepSeek refuse de répondre », *Journal du Geek*, 31 janvier 2025, [<https://www.journaldugeek.com/2025/01/31/ces-questions-politiques-sensibles-auxquelles-deepseek-refuse-de-repondre/>]. En faisant quelques tests sur la version en ligne, j'ai pu constater que l'outil affiche parfois une réponse fondée sur des faits courants avant de la remplacer par un message indiquant que la question est hors de son champ d'action. Selon J.E. Ramos, l'installation de la version locale et l'ajustement des questions permettrait d'obtenir des réponses non-censurées. Voir « DeepSeek Model DOES NOT Censor Tiananmen Square », *Dev*, 29 janvier 2025, [<https://dev.to/jeramos/deepseek-model-does-not-censor-tiananmen-square-2kcb>].

26 Ali Breland et Matteo Wong, « The Day Grok Told Everyone About 'White Genocide' », *The Atlantic*, 15 mai 2025, [<https://www.theatlantic.com/technology/archive/2025/05/elon-musk-grok-white-genocide/682817/>].

27 Christian Rioux, « Bienvenue dans le merveilleux monde de l'IA », *Le Devoir*, 15 mai 2025, p. B4-B5.

28 « L'IA atteindrait le niveau de l'intelligence humaine d'ici 5 à 20 ans, croit Bengio », *Radio-Canada*, 25 juillet 2023, [<https://ici.radio-canada.ca/nouvelle/1998915/ia-conscience-evolution-bengio>].

réglementation s'est resserrée jusqu'en 2024<sup>29</sup>, mais la nouvelle administration privilégie plutôt le laissez-faire sauf pour une loi visant à limiter la prolifération des hypertrucages (*deepfakes*) à caractère pornographique<sup>30</sup>. Ces lacunes en matière de réglementation, les défaillances du point de vue de la sécurité ainsi que le rattrapage en cours par la Chine sont des facteurs qui inquiètent les auteurs du scénario AI-2027<sup>31</sup>. La programmation par les IA étant utilisée pour accélérer le développement de nouveaux modèles plus puissants et les applications ayant déjà commencé à mentir « sciemment » aux programmeurs, la perte de contrôle pourrait se produire dès l'automne 2027, d'où le titre de leur rapport. Un coup de barre sera bientôt nécessaire pour éviter les pires dérives.

Les gestes à poser pour influencer les politiques gouvernementales dépassent le cadre de cet article. Dans les domaines de l'histoire et de l'archivistique, nous pouvons faire la promotion d'un usage responsable et éthique de l'IA en poursuivant deux objectifs : tirer profit de la puissance des algorithmes pour faciliter l'accessibilité et l'analyse, d'une part, et conserver la maîtrise d'œuvre sur les processus et les données d'autre part, le tout en respectant la mission des institutions. Tel qu'annoncé ci-dessus, je me concentrerai sur trois volets intégrés au projet Nouvelle-France numérique : la reconnaissance des écritures manuscrites, l'aide à l'indexation et l'aide à la recherche.

## Du caractère analogique au caractère numérique

La transposition de caractères alignés sur une feuille de papier vers des caractères numériques que l'on peut traiter avec les ordinateurs fait déjà partie du paysage grâce à la reconnaissance optique de

---

29 Gil et Perreault, *Artificial Intelligence Index Report 2025*.

30 Andrew R. Chow, « Inside the First Major U.S. Bill Tackling AI Harms—and Deepfake Abuse », *Time*, 29 avril 2025, [<https://time.com/7277746/ai-deepfakes-take-it-down-act-2025/>]. La loi était prête à adopter sous l'administration Biden, mais des législateurs républicains ont fait obstruction pour que le président élu Trump puisse en avoir le crédit. Le projet de loi budgétaire omnibus déposé en mai 2025 comprend un article visant à empêcher les états américains d'appliquer toute réglementation relative à l'IA pendant 10 ans. Cette proposition pourrait disparaître du bill lors du processus de réconciliation entre le Sénat et la Chambre des représentants. Voir Jordain Carney, Ben Leonard, Grace Yarrow et James Bikales, « The 7 pieces of the House megabill that could succumb to Senate rules », *Politico*, 30 mai 2025, [<https://www.politico.com/news/2025/05/30/megabill-policies-senate-parliamentarian-byrd-rule-00375507>].

31 Daniel Kokotajlo, Scott Alexander, Thomas Larsen, Eli Lifland, Romeo Dean, « AI 2027 », *AI Futures Project*, 3 avril 2025, [<https://ai-2027.com/>].

caractères (ROC) ou Optical Character Recognition (OCR) de l'imprimé. Je me permets quelques rappels concernant cette technologie afin de la distinguer de la reconnaissance des écritures manuscrites ou Handwriting Text Recognition (HTR). L'océrisation est appliquée sur des pages converties en noir et blanc pour avoir le meilleur contraste possible entre les caractères et le fond. Après avoir segmenté la page en lignes, mots et caractères, la forme de chaque signe est comparée à celles du dictionnaire de caractères. Des modèles de langage sont maintenant utilisés pour améliorer les performances en intégrant les probabilités grammaticales et contextuelles, mais l'OCR repose néanmoins sur des ensembles fermés de caractères à interpréter. L'application de l'OCR à des ouvrages plus anciens donne parfois des résultats incongrus lorsque le « s long » (ſ), utilisé au XVIII<sup>e</sup> siècle, est mal interprété. Les ouvrages océrisés donnent alors donc l'impression que le mot « suck » commence par un « f », ce qui génère des fréquences de mots incorrectes et appelle à la plus grande prudence dans l'usage des ngrams de Google pour l'analyse du langage<sup>32</sup>.

Malgré ces limites, l'océrisation à grande échelle ouvre de vastes corpus de journaux et d'ouvrages imprimés à la consultation, déverrouillant ce qui était autrefois consulté péniblement sur microfilm. L'impact est majeur, mais rarement analysé, bien qu'Ian Milligan ait fait quelques mises en garde dès 2013<sup>33</sup>. La recherche est-elle orientée (*skewed*) vers les documents les plus accessibles? Comprendons-nous réellement comment les documents ont été traités et comment les moteurs de recherche sont construits? Ces précautions émises il y a une douzaine d'années nous rappellent que la puissance de la technologie ne nous libère pas des règles de base de la discipline en matière de méthodologie et de critique de sources. Elles sont d'autant plus importantes à retenir alors que nous entreprenons de grands chantiers de transcriptions de documents manuscrits.

---

32 Étienne [Ollion], « Fuck / NGram », *Data Sciences Sociales*, 26 octobre 2015, [<https://doi.org/10.58079/ngy7>].

33 Ian Milligan, « Illusionary Order: Online Databases, Optical Character Recognition, and Canadian History, 1997–2010 », *Canadian Historical Review*, vol. 94, no 4 (décembre 2013), p. 540-569, [<https://doi.org/10.3138/chr.694>].

La variabilité des écritures présente dans ces manuscrits a limité leur traitement par les approches de l'informatique conventionnelle utilisées par l'océrisation. La solution, comme dans tant d'autres domaines, est de recourir à l'apprentissage machine supervisé à l'instar des paléographes humains qui apprennent par la pratique. La reconnaissance des écritures manuscrites commence par un échantillon de pages segmentées en lignes puis transcrites par des humains avec le niveau de précision le plus parfait possible. Le nombre de pages peut varier, avec un minimum de 25 pages pour une écriture régulière et un vocabulaire récurrent, en allant jusqu'à 200 pages pour des corpus plus volumineux pour lesquels le scribe est moins régulier et le vocabulaire plus varié. Au-delà de 200 pages de données d'entraînement, le gain de précision n'est plus significatif lorsqu'on prend en compte le temps requis pour préparer l'échantillon et les ressources informatiques nécessaires à la création du modèle<sup>34</sup>.

À l'instar de l'OCR, le HTR segmente d'abord la page et distingue les lignes. Par la suite, chaque ligne est « lue » pour identifier chaque caractère. Le résultat est comparé avec la version transcrite par les humains. Quand on analyse la courbe d'apprentissage, on constate que lors des premiers passages, le taux d'erreur avoisine les 30% (parfois plus). Le système repasse sur les mêmes pages plusieurs fois (100, 200, 250 fois) pour « apprendre » à mieux détecter les caractères. La courbe d'apprentissage indique que le taux d'erreur diminue rapidement puis les gains de précision sont moins en moins grands jusqu'à ce qu'on atteigne le nombre d'« epochs » ou de cycles prédéterminé. Si le système détecte qu'il n'apprend plus (le taux d'erreur reste le même et il n'arrive pas à s'améliorer), l'entraînement peut aussi s'arrêter automatiquement pour éviter de gaspiller des ressources informatiques. Le système vérifie aussi son

---

34 Ces constats découlent des recommandations de READ-COOP, « 4. Model Setup and Training », *Transkribus*, 2025, [<https://help.transkribus.org/model-setup-and-training>] et des expérimentations réalisées par les Gardenotes dans le cadre du projet Nouvelle-France numérique. Un article de Pierre Dubois et Maxime Gohier est en préparation sur le sujet.

niveau de précision sur un autre échantillon de validation pour voir s'il peut prédire correctement les caractères.

Les résultats sont impressionnants, sans être parfaits. Le taux d'erreur, de 2 à 5%, permet de réaliser des recherches générales ou pointues permettant d'obtenir la très grande majorité des documents pertinents. Dans le cas de notaires qui font même rager les paléographes humains, le taux d'erreur augmente à 6 ou 8%, ce qui est compréhensible quand on a déjà tenté de déchiffrer leur écriture. Si on prend un nombre moyen de 6 caractères par mot et 300 mots par page, on peut s'attendre, au pire, à 90 caractères erronés sur 1 800. Nous sommes loin des 30 à 40% d'erreurs (720 caractères par page de 300 mots) sur les océrisations des années 1980 et 1990.

Signalons qu'avec un logiciel HTR tel que Transkribus ou eScriptorium, les caractères erronés sont présentement laissés en place. Un humain peut facilement repérer les erreurs et procéder à des corrections manuelles, ou, s'il est impossible de lire un caractère pour diverses raisons, signaler que le mot est illisible. ChatGPT peut aussi faire de la reconnaissance de manuscrit et le résultat, au premier coup d'oeil, semble meilleur. Si des caractères sont illisibles, le système inscrit des mots selon le contexte. Nous ne connaissons toutefois ni le taux d'erreur par caractère, ni si un mot provient de l'analyse des caractères ou d'un calcul probabiliste du discours. La transcription reposant sur des IA généralistes opaques soulève alors de grandes inquiétudes. Des corpus complets pourraient être truffés d'erreurs sans que leur existence soit clairement indiquée. La fiabilité de corpus transcrits avec de tels outils, tout comme la qualité des études qui en découleront, ne peut pas être vérifiée, nous ramenant aux mises en garde méthodologiques de Milligan.

### **Nouvelle-France numérique. Transcrire dans le respect et la transparence**

J'ai fait référence au projet Nouvelle-France numérique à quelques reprises. Il s'agit d'un projet de recherche en partenariat dirigé par Maxime Gohier et moi-même, regroupant plus de 50 chercheurs

répartis au Québec, au Canada, aux États-Unis et en France, des partenaires dans les centres d'archives et dans les communautés autochtones, ainsi que de précieux citoyens bénévoles regroupés au sein des Gardenotes. Bien que dirigé par des historiens qui ont une connaissance très fine des documents et de leur contexte de production, le projet comprend aussi des spécialistes de l'IA, de la linguistique, de la gestion des données de recherche, et j'en passe. Nous avons initié un chantier sans précédent de transcription et de balisage de documents de la Nouvelle-France visant à faciliter les recherches très pointues ainsi que les analyses transversales. Nous souhaitons en même temps valider et améliorer les protocoles de travail collaboratif (numérisation, transcription, balisage des métadonnées, interopérabilité des systèmes); identifier des modalités de structuration, de mise en relation et de diffusion des données qui facilitent leur réutilisation; adapter l'IA à l'analyse du patrimoine documentaire ancien, notamment aux collections manuscrites et iconographiques, au français moderne et aux langues autochtones; créer une communauté et élaborer un modèle efficace et pérenne de science participative; et mettre à profit l'intelligence artificielle pour amorcer un renouvellement de l'historiographie sur la Nouvelle-France. Notre démarche s'appuie présentement sur l'application Transkribus, développée par la coopérative READ, basée à l'Université d'Innsbruck<sup>35</sup>. Notre architecture nous permettrait toutefois de nous découpler de cette application pour nous tourner, au besoin, vers eScriptorium, projet de code libre et ouvert<sup>36</sup>.

La structure du projet est fédérative. Bien que dirigées à partir de Rimouski et de Sherbrooke, des antennes importantes et autonomes sont basées à l'Université de Montréal et à l'UQAM. Nouvelle-France numérique agit comme un parapluie pour coordonner et éviter les redondances dans le travail. Nos partenaires sont évidemment les centres d'archives, mais aussi des organismes spécialisés dans la

---

35 *READ Co-op*, [<https://readcoop.org/>].

36 *Escriptorium Documentation*, [<https://escriptorium.readthedocs.io/en/latest/>].

diffusion du patrimoine documentaire et culturel du Québec et du Canada (RCDR et RCIP). Les documents nous rapportent les activités des différentes populations présentes sur le territoire revendiqué par la France, incluant les peuples autochtones avec lesquels nous développons des collaborations. Afin de maximiser l'alliance entre les humains et l'IA pour nous attaquer à un corpus qui comprend plusieurs millions de pages, nous nous sommes associés à des passionnés de paléographie, les Gardenotes. Animés par leur amour de l'histoire, de la généalogie et du bien commun, l'entraînement de nos modèles repose sur la qualité indéniable de leur travail.

Nouvelle-France numérique est aussi partenaire dans une demande soumise à la Fondation canadienne pour l'innovation (FCI) pilotée par l'Université d'Ottawa afin de renouveler les infrastructures et les interfaces de [Canadiana.ca](http://Canadiana.ca) et de [Heritage.ca](http://Heritage.ca). Cette collaboration nous permettra de codévelopper des éléments de nos infrastructures respectives et de partager nos expertises. À l'automne 2025, notre plateforme de recherche et de diffusion sera accessible à nos partenaires pour réaliser les tests nécessaires à son perfectionnement avant le lancement public prévu en 2026. Mais déjà, l'équipe de l'Atelier permanent d'analyse documentaire de Dominique Deslandres et celle des Gardenotes a développé des modèles qui sont disponibles gratuitement. Nous travaillons de notre côté pour obtenir le financement nécessaire pour appliquer ces modèles à l'ensemble des greffes de notaires, des archives judiciaires, de la correspondance coloniale, des archives de congrégations religieuses, etc.

Un tel projet implique une collaboration étroite avec les centres d'archives partenaires. Certains sont numérisés par les centres d'archives alors que dans d'autres cas, nous allons sur place pour les numériser. La première option est à privilégier, étant donné que le travail est alors réalisé selon vos normes et que les métadonnées appropriées y sont associées<sup>37</sup>. Présentement, les fichiers numérisés nous sont transmis

---

<sup>37</sup> La numérisation en soi est un processus déjà coûteux auquel il faut parfois ajouter des frais supplémentaires pour traiter les documents reliés ou pour procéder à la stabilisation préalable de documents fragiles. Nous sommes à la recherche de fonds pour appuyer les centres d'archives dans ces tâches.

électroniquement selon le mode le plus simple pour le centre d'archives. Les documents sont ensuite versés sur la plateforme Transkribus. Nous cherchons à l'avenir à éviter ces transferts de fichiers en ayant plutôt recours au cadre d'application IIIF. Les images restent alors sur vos serveurs et un « manifest » nous fournit le lien vers chaque image numérisée ainsi que les métadonnées qui y sont associées. Quelle que soit la méthode utilisée, les paléographes peuvent ensuite faire les transcriptions manuelles qui servent de base à l'entraînement d'un modèle qui pourra être appliqué à l'ensemble d'une collection rédigée par la même personne<sup>38</sup>. Dans le cas des ordonnances des intendants de la Nouvelle-France, les paléographes sont par la suite repassés sur certains cahiers pour corriger les erreurs laissées par l'IA. Ces cahiers sont désormais disponibles en format PDF enrichi permettant de réaliser des recherches plein texte<sup>39</sup>.

## L'indexation et les modèles de langage

L'IA est essentielle pour traiter de grands volumes de manuscrits pour en faciliter la lecture et pour que la recherche plein texte soit possible. Ce type de recherche comprend toutefois certaines limites à cause des variantes en genre, en nombre et en conjugaison, ainsi que le recours aux synonymes et aux euphémismes. C'est pour pallier ces lacunes que les grands moteurs de recherche, Google, Bing, etc., ont intégré des modèles de langage pour analyser le sens des requêtes et améliorer les résultats, bien ceux-ci puissent être mitigés. Ces modèles de type « BERT » (Bidirectional Encoder Representations from Transformers) fonctionnent maintenant très bien pour les textes en langue anglaise. Des variantes pour d'autres langues ont été développées, dont CamemBERT et FlauBERT pour le français, mais il ne sont pas très efficaces pour le français d'Ancien Régime. Nouvelle-France numérique souhaite donc participer

---

38 Des supers modèles très performants existent pour traiter des mains d'écriture différentes en anglais et en allemand, mais aucun n'a pu offrir des résultats satisfaisants en français.

39 Archives nationales à Québec, E1,S1,D1, Fonds Intendants, Ordonnances, Cahier 1. Registre des commissions et ordonnances rendues par monsieur Raudot, 7 septembre 1705 au 18 novembre 1707 [<https://advitam.banq.qc.ca/notice/989303>].

au développement du modèle d'AleMBERT, conçu pour le français de l'époque moderne. Le modèle de langage seul ne permet pas de cibler certains types de contenus. Nous préparons aussi des listes d'autorité qui pourront faciliter l'entraînement de l'IA, et aussi offrir des filtres pour la recherche.

L'identification de ces entités nommées de base (dates, lieux, personnes et organisations) peut être automatisée grâce aux outils de reconnaissances existants. Dans le cas des dates, nous devons les compléter avec des listes d'autorités relatives au calendrier religieux et aux abréviations utilisées au XVII<sup>e</sup> siècle. Des références relatives (hier, dans huitaine, après les semences) doivent aussi être signalées et être résolues à partir d'autres éléments contenus dans le document ou, dans le cas des semences, en se basant sur le calendrier des audiences judiciaires<sup>40</sup>. Les toponymes nécessitent aussi des précautions. Jean-François Palomino a relevé ceux que l'on retrouve sur les cartes géographiques avec leurs variantes, ce qui nous offre une liste d'autorité qui n'existait pas jusqu'à maintenant. D'autres noms de lieux qui font partie du discours émergeront aussi d'un corpus de millions de mots, ce qui nous permettra d'enrichir cette banque toponymique.

Du côté des personnages, les démographes ont établi depuis un demi-siècle des règles de jumelage qui facilitent l'identification des individus dans plusieurs documents. Nos textes étant moins structurés que les recensements et les actes d'état civil (baptême, mariage, sépulture), des règles additionnelles devront être développées pour résoudre les désignations partielles que l'on retrouve sans le contexte familial des documents plus courants en démographie. Les groupes et les organisations présentent aussi des embûches. Le cas des sulpiciens met en évidence les multiples manières de référer à cette communauté de prêtres (Sulpiciens, prêtres du Séminaire, prêtres de Saint-Sulpice) ainsi que l'euphémisme « Messieurs » utilisé couramment à Montréal. En nous appuyant sur des spécialistes qui

---

40 Thomas Wien, « Les travaux pressants ». Calendrier agricole, assolement et productivité au Canada au XVIII<sup>e</sup> siècle », *Revue d'histoire de l'Amérique française*, vol. 43, n° 4 (1990), p. 535-558.

peuvent déjà donner des lignes directrices pour l'apprentissage machine, nous éviterons de gaspiller des cycles de traitement qui pourraient, ou non, arriver au même résultat.

Ces mêmes spécialistes nous aideront à analyser différents champs sémantiques sur des thèmes très variés : pouvoir, commerce, maritimité, religion, genre, altérité, sexualité, justice, occupations et bien d'autres. Nous signalerons ces éléments du langage (tant sa classe que ses attributs, incluant son identifiant dans une liste d'autorité) grâce au balisage XML. Nous avons privilégié le standard de la Text Encoding Initiative (TEI), norme pour laquelle un de nos cochercheurs, Emmanuel Château-Dutier, de l'Université de Montréal, est reconnu pour son expertise. À partir des fichiers XML, nous pourrons aussi générer des jeux de données en triplets rdf afin de participer au mouvement des données ouvertes et liées. Nous pourrons nous appuyer pour le faire sur l'expertise de nos collègues du Réseau canadien d'information sur le patrimoine, du Canadian Writing Research Collaboratory (CWRC) et du ministère de la Culture et des Communications du Québec.

La collaboration entre les chercheurs et l'IA est donc essentielle. Les humains pourront orienter les modèles de langage et ensuite valider les résultats de l'entraînement avant de les appliquer à l'ensemble du corpus. Ce travail ne se terminera pas à la fin de notre cycle de recherche, comme le serait une édition critique imprimée. Notre infrastructure est conçue dans le but de l'enrichir à long terme, aussi bien en documents qu'en thématiques.

## **Désinvisibiliser les informations et favoriser la recherche**

La description des pièces judiciaires (grâce auxquelles les archives de la Juridiction royale de Montréal ont été privilégiées dans les recherches) et des actes notariés (par Archiv-Histo pour la banque Parchemin) a grandement facilité le travail des historiens. Une partie des contenus demeure toutefois invisible et ne peut être accessible que par la recherche plein texte, voire assistée par un modèle de langage. À titre d'exemple, le concept de prix des denrées est signalé dans les instruments de recherche

lorsqu'une ordonnance a pour objectif de fixer les prix lors de pénuries, mais une mention de prix ne sera pas prise en compte de manière systématique dans le descripteur d'autres documents, incluant les baux agricoles ou les procès. La recherche plein texte aurait aussi ses limites dans un tel domaine de recherche. Nous proposons donc de commencer par l'analyse du champ sémantique des aliments et de leur commerce afin de repérer un plus grand nombre de cas et de permettre des études plus complètes des variations des prix.

L'invisibilisation de contenu n'est pas limitée aux détails des échanges commerciaux. De nombreux instruments de recherche réalisés il y a plusieurs décennies laissent à désirer, tant par les termes qui sont utilisés pour désigner certaines personnes que par leur approche qui met de l'avant un certain type d'individus (généralement des hommes blancs privilégiés) et laisse de côté les femmes, les hommes des classes populaires et les gens d'origine ethnique autre qu'européenne. L'application de nouvelles approches de classement, d'indexation, de balisage, de sémantisation, etc. s'inscrit aussi dans une démarche de décolonisation et de désinvisibilisation. Est-ce que certains y verront une nouvelle manifestation du « wokisme »? Probablement. Mais notre démarche ne va pas effacer Champlain, Frontenac, d'Iberville ou Montcalm. Ceux-ci auront toujours une place prépondérante dans la documentation sur la Nouvelle-France. Ils partageront toutefois la scène avec des hommes et des femmes -- eurodescendants, autochtones ou afrodescendants – souvent négligés. Nous connaissons déjà mieux les femmes de la Nouvelle-France grâce à Dominique Deslandres<sup>41</sup>, les couturières étudiées par Suzanne Gousse<sup>42</sup> et les artisans du métal mis en évidence par Sonia Blouin<sup>43</sup>. De nouvelles études pourront être

---

41 Dominique Deslandres, « Femmes de Nouvelle-France », *Les Cahiers des Dix*, n° 75, 2021, p. 311-344, [<https://doi.org/10.7202/1088878ar>].

42 Suzanne Gousse, « Les Couturières de Montréal au XVIII<sup>e</sup> siècle (Québec, Septentrion, 2013).

43 Sonia Blouin, « L'importance des métaux dans une ville naissante: les artisans du métal dans l'espace social et urbain de Montréal de 1642 à 1701 », mémoire de maîtrise (histoire), Université de Sherbrooke, 2021.

plus facilement réalisées afin de nourrir le renouvellement de l’histoire sociale, de l’histoire des femmes, de l’histoire culturelle et même de l’histoire politique.

L’écosystème de la recherche doit aussi s’adapter à la révolution numérique déjà en cours. De bonnes pratiques<sup>44</sup> de même que des systèmes de référence uniformes sont nécessaires pour éviter la multiplication des initiatives dispersées comme celle qu’avait engendrée l’apparition de la micro-informatique il y a 40 ans. Nous souhaitons donc valider et améliorer des protocoles de travail collaboratif (numérisation, transcription, balisage des métadonnées, interopérabilité des systèmes) qui tiennent compte de la diversité des sources de la Nouvelle-France, de la variété des problématiques de recherche et de la pluralité d’agent.e.s qui interviennent dans le traitement, l’analyse et la conservation des données. Il est aussi nécessaire d’identifier les modalités de structuration, de mise en relation et de diffusion des données qui facilitent leur réutilisation dans le respect des règles de conservation du patrimoine. Le système doit promouvoir la reconnaissance de la contribution des créateurs de données, et assurer la traçabilité de l’information. Nous pourrions ainsi adapter l’IA, grâce aux données produites, à l’analyse du patrimoine documentaire ancien, notamment les manuscrits en français moderne et en langues autochtones. L’accès à de nouvelles données et la mise en relation de données anciennes, le tout avec un niveau de confiance élevé favorisera le renouvellement de l’historiographie. Cet ensemble s’appuie sur un modèle de communauté efficace et pérenne de science participative, adapté à la recherche sur le patrimoine documentaire et qui favorise la démocratisation et la coconstruction des savoirs.

---

44 D. Dennie et A. Guindon, « Résultats d’une enquête sur les pratiques et attitudes des chercheurs de l’Université Concordia en matière de gestion des données de recherche », *Documentation et bibliothèques*, vol. 63, n° 4 (2017) 59-72; S. Higgins *et al.*, « Research Data Management Support in the Humanities: Challenges and Recommendations [Draft] », *Open Scholarship Policy Observatory, Community News*, vol. 3, (2021), [<https://ospolicyobservatory.uvic.ca/draft-report-rdm/>]; W. Kurtz, « Founders Online: Early Access: Reflections on Open Access, Crowd Sourcing, and Metadata Standards », *Digital Studies / Le champ numérique*, vol. 7, n° 1 (2017) [<https://doi.org/10.16995/dscn.277>].

## Appuyer les partenaires sans leur nuire

Même avec une démarche transparente et partenariale, nous sommes bien conscients que certains centres d'archives ne sont pas prêts à se lancer dans une telle aventure. La création de PDF enrichis avec une transcription de grande qualité apporte une plus-value aux centres d'archives. Ces documents peuvent alors être intégrés aux systèmes de recherche et de visualisation des institutions. Le recours au cadre d'application IIF facilite de plus l'intégration des documents et de leur transcription dans des visionneuses telles que Mirador, Universal ou OpenSeadragon. Cependant, si les documents sont intégrés dans un métacorpus de la Nouvelle-France, perdent-ils leur identité? La structure même de l'outil doit permettre de signaler la provenance des manuscrits. L'indication « comment citer ce document » doit être une notice très visible à l'écran et accompagnant tout téléchargement<sup>45</sup>. Nos partenaires auront aussi l'option de contrôler la découvrabilité des documents en ne permettant pas au public d'accéder au document complet en ligne. La consultation nécessiterait alors une visite *in situ*, mais les documents auraient déjà été repérés.

L'infrastructure TNF donnera évidemment accès à des statistiques relatives aux documents repérés lors des recherches, à leur consultation et au téléchargement. Les listes d'autorité relatives aux entités nommées seront aussi utiles pour enrichir les inventaires de vos collections. Enfin, la mise en relation de documents traitant de mêmes thèmes ou mentionnant les mêmes individus dans différentes collections, voire dans différents centres d'archives participera au décloisonnement de la recherche.

Nous espérons donc que les gains offerts par la participation à Nouvelle-France numérique vont permettre de surmonter les inquiétudes. Notre approche d'enrichissement continu permettra d'ailleurs à de nouveaux partenaires de se joindre à nous une fois que nous aurons fait nos preuves.

---

45 Bien que ses technologies puissent sembler datées après seulement deux décennies d'existence, le portail des Archives de la Nouvelle-France avait déjà donné le ton en indiquant clairement la provenance des documents, [<https://nouvelle-france.org/fra/item>].

## Conclusion

La révolution associée à l'IA dépasse de loin le domaine de l'informatique. Les ramifications sociales, économiques et politiques sont même inquiétantes en ce moment. Les applications commerciales sont caractérisées par leur opacité, leur éthique minimaliste, et leur stratégie visant à concentrer les bénéfices, pour ne pas dire le pouvoir. La facilité offerte par ces outils masque des coûts non monétaires dont l'ampleur est difficile à mesurer.

En rassemblant chercheuses et chercheurs, centres d'archives, coopérative, bénévoles et réseaux de spécialistes, nous pouvons contribuer au bien commun et partager le savoir. À partir des pratiques du partenariat Nouvelle-France numérique, nous avons mis de l'avant quelques principes du bon usage de l'IA :

- la transparence en matière de méthodologie et de gouvernance
- le respect des sources, des partenaires, des bénévoles, des cochercheurs
- la complémentarité de l'humain et de la machine pour tirer profit des meilleures qualités de chacun.

En misant sur ces principes, nous espérons tirer profit de la puissance de la technologie tout en protégeant la mission des centres d'archives et en évitant d'enrichir des entreprises dont les objectifs s'arriment rarement avec ceux des sciences humaines et du milieu de la culture.