

La Route Moins Empruntee

Repenser les couts, la securite et la confiance lies a l'IA generative dans les institutions reglementees

POURQUOI CETTE CONVERSATION EST IMPORTANTE

Les institutions se font dire que l'intelligence artificielle est inevitable. Pas seulement utile - inevitable. Le message est constant : l'IA va transformer le fonctionnement des organisations, mais il vous faudra les bonnes protections.

Ce qu'on remet rarement en question, c'est de savoir si les institutions sont invitees a s'adapter a un modele industriel qui n'a jamais ete concu pour elles. L'objectif ici n'est pas de critiquer les outils et produits d'IA standard, ni de soutenir qu'une approche est universellement superieure. L'objectif est d'expliquer pourquoi les organisations reglementees pourraient avoir besoin d'un type fondamentalement different de systeme d'IA - et pourquoi cette difference compte davantage que les fonctionnalites, les modeles, les arguments commerciaux ou le prestige des fournisseurs.

Pour mieux illustrer le probleme, j'utiliserai deux metaphores simples que j'ai explorees tout au long du developpement de COMPAiSS, un assistant IA concu pour les environnements institutionnels :

La Matrice IA

l'ensemble des hypotheses cachees qui faconnent la conception de la plupart des systemes d'IA.

La Route (et la Route Moins Empruntee)

la maniere dont ces hypotheses se traduisent en couts, en risques, en correctifs recommandes, en arguments de vente et en adequation institutionnelle.

LA MATRICE IA : DES HYPOTHESES QUE TOUT LE MONDE ACCEPTE (SOUVENT SANS LE REMARQUER)

La plupart des systemes d'IA modernes sont construits a l'interieur de ce qu'on peut raisonnablement appeler un seul cadre de pensee dominant. Dans ce cadre, on suppose que l'IA reflechit en permanence et est toujours 'en marche'. Elle peut raisonner sur le monde entier - l'internet, les connaissances generales, les faits implicites, les analogies et les suppositions. La securite, c'est quelque chose que l'on gere apres que la reflexion a commence, pas avant, et tous les risques correspondants sont acceptables du moment qu'ils peuvent etre filtres, moderes ou refuses.

Ce cadre de pensee n'est pas malveillant - c'est le principe fondateur accepte de l'IA grand public, de l'IA creative et de l'IA generaliste. Et cela est logique si votre objectif est de vendre, de fideliser et de renforcer une intelligence ouverte. Une fois entierement immerge dans cette matrice, certaines choses deviennent incontestables :

bien sur que l'IA doit reflechir d'abord.

bien sur que des erreurs se produiront.

bien sur que les hallucinations sont inevitables.

bien sur que vous avez besoin de couches de controles pour gerer les risques.

bien sur que les couts augmentent avec la complexite.

LA QUESTION QUI N'EST PAS POSEE

Et si l'ensemble de ce cadrage etait inapproprié pour les organisations réglementées ayant l'obligation de fournir des informations précises et faisant autorité aux personnes qu'elles servent ?

LA ROUTE PRINCIPALE : COMMENT L'IA D'ENTREPRISE EST CONSTRuite (ET POURQUOI)

La plupart des produits d'IA d'entreprise - y compris ceux commercialisés spécifiquement pour les environnements réglementés - empruntent la route principale, qui ressemble à ceci :

Partir d'une IA généraliste pouvant répondre à presque tout ; la connecter aux documents internes grâce à des techniques comme la génération augmentée par récupération ('RAG') ; ajouter des couches de modération, des filtres de politique, des invites de refus, des tableaux de bord de conformité, etc. ; puis espérer que la plupart du temps, le système se comportera correctement.

Pour filer la métaphore de la route, c'est comme donner à tout le monde accès à un réseau autoroutier mondial, puis investir des ressources énormes dans la régulation de la circulation, les panneaux d'arrêt, les limitations de vitesse, les caméras, la police, les rapports d'accidents, les assurances et l'application continue des règles. La route est vaste, les destinations sont imprévisibles, et l'application des règles ne s'arrête donc jamais - c'est pourquoi ces systèmes sont souvent coûteux à déployer, coûteux à maintenir, et peuvent être difficiles à faire pleinement confiance même avec une exploitation soignée.

POURQUOI LE RAG SEMBLE JUSTE DANS LA MATRICE - ET POURQUOI IL LAISSE ENCORE PASSER DES RISQUES

La génération augmentée par récupération (RAG) semble rassurante. L'idée est simple : si l'IA ne consulte que des documents approuvés, les hallucinations peuvent être maîtrisées. Mais voici ce qu'on laisse facilement passer quand on est entièrement immergé dans le paradigme accepté : le RAG contrôle ce que l'IA consulte - pas si, ni comment, l'IA est autorisée à réfléchir en premier lieu.

Même avec des documents parfaits, l'IA raisonne quand même, suppose les pièces manquantes, tire des conclusions qui n'ont jamais été explicitement énoncées, devine la façon dont les choses s'articulent et relie des points qui pourraient ne pas exister.

Le RAG réduit les entrées, mais il ne modifie pas la nature du raisonnement. C'est pourquoi les hallucinations persistent, sous des formes plus subtiles et plus dangereuses - sonnantes comme faisant autorité, mélangeant politique et inférence, remplissant avec assurance ce qui n'a jamais été écrit. Pour les organisations réglementées, ce sont les pires types d'erreurs : pas des inepties flagrantes, mais des réponses presque correctes qui semblent officielles.

Des recherches évaluées par les pairs confirment ce schéma. Des études documentent des taux d'hallucination résiduels d'environ 6 % même dans les conditions RAG optimales (Nishisako, Higashi & Wakao, 2025), et des tests empiriques indépendants de la plateforme juridique phare de Thomson Reuters ont révélé des taux d'hallucination de 17 à 33 % en conditions réelles d'utilisation (Stanford / Journal of Empirical Legal Studies, 2025). Thomson Reuters lui-même inclut une mise en garde écrite dans son produit conseillant aux utilisateurs de toujours lire les sources et de vérifier les résultats. La raison sous-jacente est architecturale : la formation d'un modèle peut surpasser le contenu récupéré au niveau de l'inférence, ce qui signifie qu'une meilleure récupération améliore ce que l'IA lit, mais ne contrôle pas ce que l'IA dit (Sun et al., ReDeEP, ICLR 2025).

LE PROBLEME DES COUTS CACHES : POURQUOI GERER LES RISQUES EST SI COUTEUX

Cela explique pourquoi les solutions d'IA d'entreprise deviennent si coûteuses avec le temps. Elles nécessitent davantage de surveillance, davantage d'outils de conformité, davantage de révision humaine, davantage de mises à jour de politiques, davantage de gestion des exceptions.

Ces coûts sont des conséquences structurelles de l'architecture d'inférence en premier lieu, plutôt que des choix d'implémentation fortuits. Parce que l'inférence est toujours autorisée à s'exécuter, le système doit gérer en permanence ce qu'il produit - par le biais de couches de modération, d'outils de validation, de flux de travail de révision humaine et de

complements de gouvernance. Ce sont de véritables exigences opérationnelles, et non des complexités artificielles. L'économie de la route principale reflète le coût réel de la gestion des risques dans un système où la réflexion se produit toujours en premier.

De l'intérieur du paradigme généralement accepté de l'IA, cette architecture semble naturelle et bien étayée. La question qui mérite d'être posée est de savoir si elle convient aux organisations réglementées dont les obligations nécessitent un point de départ entièrement différent.

LA ROUTE MOINS EMPRUNTEE : UN MONDE DE VERITE PLUS PETIT ET PLUS SUR

En sortant de cette matrice, le chemin alternatif commence par une question très différente : et si la réflexion elle-même était conditionnelle ? Non pas : 'Comment nettoyer les sorties ?' ou 'Comment détecter les hallucinations ?' ou 'Comment localiser et refuser les mauvaises réponses ?', mais plutôt : le système devrait-il être autorisé à réfléchir du tout - à moins qu'il ne se trouve déjà à l'intérieur d'un monde de confiance ?

COMPAiSS est construit sur une idée simple mais radicale : établir le monde d'abord, puis libérer l'IA pour qu'elle fonctionne uniquement à l'intérieur de ce monde. Ce cadrage a émergé de la confrontation répétée aux mêmes contraintes institutionnelles - la confiance, la précision et la responsabilité.

Plutôt que le monde de tout, le système opère à l'intérieur d'un monde de vérité pré-autorisé, défini par l'institution elle-même. Dans ce monde, seules les sources officielles existent, seuls les domaines approuvés sont accessibles et seules les connaissances institutionnelles sont présentes. L'IA n'a pas besoin d'être empêchée d'aller ailleurs - ailleurs n'existe pas.

Pour filer la métaphore de la route, COMPAiSS est conçu comme un système de transport dédié où tout le monde conduit le même type de véhicule sur la même route. Imaginez maintenant une route unique, bien conçue, où chaque véhicule autopilote roule exactement à la limite de vitesse, les voies sont uniformes, les intersections sont rares et tout le monde se dirige vers la même ville - même s'ils font des achats différents une fois arrivés. Dans cet environnement, on n'a pas besoin d'agents de la circulation à chaque coin de rue, de systèmes d'application élaborés ou d'une surveillance constante des comportements imprudents. La structure de la route elle-même fait l'essentiel du travail. C'est ce que COMPAiSS fait pour le raisonnement de l'IA.

En contraignant le monde dans lequel l'IA opère, on supprime le besoin de nombreux contrôles très coûteux sur lesquels s'appuient les systèmes d'IA d'entreprise. Il n'est pas nécessaire de 'tirer continuellement sur le volant' pour éviter des directions dangereuses, car ces directions n'ont jamais été asphaltées.

Tout le monde pose encore des questions différentes, et tout le monde cherche encore des réponses détaillées, précises et de haute qualité - mais ils le font dans le même environnement partagé et de confiance, régi par la propre compréhension de l'institution de ce qui est vrai et pertinent. C'est pourquoi COMPAiSS semble plus simple, et pourquoi il est beaucoup moins coûteux à exploiter.

POURQUOI CETTE ARCHITECTURE GERE LES HALLUCINATIONS A UN NIVEAU STRUCTUREL

Les hallucinations ne viennent pas de la malveillance ; elles émergent naturellement lorsqu'un système a une grande liberté de raisonnement. Le modèle standard produit donc deux types d'erreurs : les faux positifs, où quelque chose d'incorrect ou d'inapproprié passe quand même, et les faux négatifs, où une réponse légitime est bloquée par excès de prudence. Les deux sont inévitables lorsque le système réfléchit toujours d'abord et est corrigé ensuite.

Avec COMPAiSS, des categories entieres d'hallucinations disparaissent - non pas parce que l'IA se comporte mieux, mais parce qu'il n'y a rien sur quoi halluciner. Le raisonnement n'est autorise qu'a l'interieur d'un monde de verite preautorise, defini par l'institution. Si une reponse ne peut pas y etre fondee, le systeme n'entre jamais dans un etat ou les suppositions sont autorisees, eliminant ainsi la classe la plus dangereuse de reponses confiantes mais presque correctes.

Bien sur, COMPAiSS peut encore faire des erreurs en raison d'une ambiguite dans les documents sources, de politiques peu claires ou de liens manquants. Mais ce sont des erreurs plus sures. Elles sont bornees, explicables et ancrees institutionnellement. En controlant ou la reflexion est autorisee a exister - plutot que d'essayer de surveiller chaque resultat apres coup - COMPAiSS reduit a la fois la frequence et la gravite des erreurs.

POURQUOI CETTE ARCHITECTURE EST IMPORTANTE POUR LES INDUSTRIES REGLEMENTEES

L'argument architectural decrit ici s'etend naturellement au-dela d'une seule organisation. Les universites, les hopitaux et les systemes de sante, ainsi que les unites gouvernementales en contact avec le public, fonctionnent toutes sous des pressions similaires : elles sont responsables de la precision des informations qu'elles fournissent, elles doivent preserver la confiance institutionnelle et sont de plus en plus exposees aux risques de traduction et d'interpretation a mesure que les services s'elargissent a differentes langues et plateformes.

Dans ces environnements, les reponses 'presque correctes' sont souvent plus dangereuses que les erreurs evidentes. Un leger decalage d'interpretation - introduit lors de la recuperation, de la traduction ou de la generation - peut entrainer des consequences juridiques, cliniques ou politiques qui ne peuvent pas facilement etre annulees apres coup.

Par ailleurs, ces institutions font face a une pression croissante en matiere de couts et de gouvernance. Les systemes necessitant une surveillance, une correction et une supervision continues pour rester surs introduisent des charges operationnelles a long terme qui s'accumulent dans le temps. Une architecture qui reduit les risques de facon structurelle, plutot que de les gerer indefiniment, s'aligne mieux sur la maniere dont les institutions reglementees sont censees operer et etre tenues responsables.

POURQUOI CETTE ROUTE N'EST PAS POUR TOUT LE MONDE - ET C'EST BIEN AINSI

Le modele COMPAiSS n'est pas concu pour tout repondre. Il est concu pour repondre aux bonnes choses. Les organisations reglementees n'ont pas besoin d'une IA pouvant expliquer l'univers ; elles ont besoin d'une IA pouvant expliquer leurs regles, leurs politiques, leurs services, leurs obligations et leurs decisions. COMPAiSS decide ce que la reflexion est autorisee a etre avant qu'elle commence. Cette difference change les couts, la securite, la confiance et l'adequation institutionnelle. La route moins empruntee n'est pas clinquante, mais elle est plus silencieuse, moins couteuse, plus sure et plus honnete. Et pour les institutions fondees sur la confiance, c'est precisement ce qui compte.

L'ARGUMENT FINAL

Les systemes a generation prioritaire demandent aux institutions d'accepter un ecart structurel permanent entre ce que l'IA peut produire et ce que l'institution peut cautionner, puis de gerer cet ecart indefiniment par la surveillance, la correction et la supervision. L'inference a execution conditionnelle ferme cet ecart par conception. L'institution definit le monde. L'IA opere a l'interieur. Tout ce que le systeme dit est tracable, borne et autorise par l'institution.

La route moins empruntee demande aux institutions de partir de leurs obligations plutot que de s'adapter a un modele industriel concu a d'autres fins. C'est une conversation plus difficile a avoir avec un fournisseur. C'est une conversation beaucoup plus facile a avoir avec les personnes que l'institution sert.

