

Confrontation de la grille d'analyse

Critères ALTAIL vs Critère cadre théorique

Collecte de données ciblées

Critères ALTAIL	Critères		
	Approche réglementaire	Approche éthique by design	Approche interactionnelle
Analyse d'impact sur les droits fondamentaux			
Documentation des arbitrages éthiques			
Interaction avec la prise de décision humaine			
Préservation de l'autonomie décisionnelle			

Signalement des décisions algorithmiques			<input checked="" type="checkbox"/> Communication sur les biais
Identification des agents virtuels			
Répartition des tâches entre IA et humains			
Renforcement des capacités humaines			
Prévention de la dépendance excessive			
Définition du niveau de contrôle humain	<input checked="" type="checkbox"/> Contrôle humain : Interface permettant une supervision humaine continue et une intervention manuelle en cas de besoin.		
Identification de l'acteur humain responsable			
Mécanismes de supervision humaine	<input checked="" type="checkbox"/> Contrôle humain : Interface permettant une supervision humaine continue et une intervention manuelle en cas de besoin.		
Auditabilité et gouvernance de l'autonomie			
Contrôle renforcé pour les systèmes autonomes ou auto-apprenants			

Détection des anomalies et réponse aux risques			
Présence d'un mécanisme d'arrêt d'urgence			
Évaluation des vulnérabilités potentielles			
Identification des types de vulnérabilités			
Mesures de résilience et d'intégrité			
Comportement en environnement imprévu			
Gestion du double usage			
Existence de solutions de secours			
Évaluation du niveau de risque contextuel			
Processus d'évaluation des risques et de la sécurité			
Communication des risques pour l'intégrité physique			
Politique d'assurance			

Plan de gestion des usages abusifs ou détournés			
Évaluation des dommages potentiels			
Règles de responsabilité et protection des consommateurs			
Évaluation des impacts environnementaux et sur les animaux			
Analyse des menaces de cybersécurité			
Évaluation des défaillances critiques			
Définition de seuils de déclenchement			
Test des solutions de secours			
Définition du niveau de précision attendu			
Méthodes de mesure de la précision			
Qualité et actualisation des données			
Évaluation des besoins en données supplémentaires			
Évaluation des préjudices liés aux erreurs			

Suivi du taux d'erreurs			
Plan de correction des erreurs			
Stratégie de contrôle du système			✓ Contrôle utilisateur
Prise en compte des contextes spécifiques			
Méthodes de vérification de la fiabilité et de la reproductibilité			
Documentation des réglages critiques			
Application des processus de test			
Communication de la fiabilité aux utilisateurs			
Mécanisme de signalement des atteintes à la vie privée			
Évaluation du type et de la portée des données			
Minimisation de l'usage de données sensibles			
Mécanismes de consentement et de contrôle			
Mesures de protection des données			
Mobilisation du DPO (Data Protection Officer)			

Alignement sur les normes reconnues			✓ Documentation technique intégrée
Mécanismes de contrôle des données			
Évaluation de la qualité des sources externes			
Garantie de la qualité et de l'intégrité des données			
Vérification de la sécurité des ensembles de données			
Documentation des protocoles de gouvernance des données			
Analyse des accès aux données			
Qualification et légitimité des personnes ayant accès aux données			
Traçabilité des accès aux données	✓ Journalisation des événements : Enregistrement systématique des événements critiques pour assurer la traçabilité.		
Mise en place de mesures de traçabilité			

Documentation des méthodes de conception	<input checked="" type="checkbox"/> Documentation technique : 1 : Description des algorithmes, sources de données, méthodes de traitement (ex. : anonymisation, réduction des biais). 2 : Preuves des conformité aux normes de sécurité et d'éthique. 3 : Plans d'atténuation des risques et d'intervention en cas de défaillance. 4 : Mécanismes d'audit et d'explicabilité.		
Documentation de l'entraînement des modèles (IA fondée sur l'apprentissage)			
Documentation des méthodes de test et de validation			<input checked="" type="checkbox"/> Documentation technique intégrée
Documentation des résultats et des décisions	<input checked="" type="checkbox"/> Transparence envers les déployeurs : Fourniture d'une notice d'information claire sur le fonctionnement, les risques, les limites et les conditions d'utilisation du système. <input checked="" type="checkbox"/> Documentation technique : 1 : Description des algorithmes, sources de données, méthodes de traitement (ex. : anonymisation, réduction des biais). 2 : Preuves des conformité aux normes de sécurité et d'éthique. 3 : Plans d'atténuation des risques et d'intervention en cas de défaillance. 4 : Mécanismes d'audit et d'explicabilité.		
Compréhensibilité des décisions du système			<input checked="" type="checkbox"/> Compréhension utilisateur
Explication des choix algorithmiques			<input checked="" type="checkbox"/> Explication contextualisée
Évaluation de l'influence sur les processus décisionnels			

Justification du déploiement dans un domaine spécifique			
Évaluation du modèle économique			
Conception orientée vers l'interprétabilité		<input checked="" type="checkbox"/> Design centré utilisateur (VSD) : 1 : Intégration d'experts en valeurs dans les équipes de conception. 2 : Prise en compte du contexte cognitif et social d'usage. 3 : Conception d'interfaces favorisant la compréhension et l'appropriation.	<input checked="" type="checkbox"/> Explicabilité graduée
Choix de modèles interprétables			<input checked="" type="checkbox"/> Transparence située
Analyse des données d'entraînement et de test			
Accès à la logique interne du modèle			
Identification claire du système comme étant une IA			
Explication des critères de décision			
Accessibilité de l'information pour tous les profils d'utilisateurs		<input checked="" type="checkbox"/> Design centré utilisateur (VSD) : 1 : Intégration d'experts en valeurs dans les équipes de conception. 2 : Prise en compte du contexte cognitif et social d'usage. 3 : Conception d'interfaces favorisant la compréhension et l'appropriation.	<input checked="" type="checkbox"/> Accessibilité cognitive
Prise en compte des retours utilisateurs			<input checked="" type="checkbox"/> Recours utilisateur
Communication des risques et biais			<input checked="" type="checkbox"/> Visualisation pédagogique
Communication élargie selon les publics			
Indication de la finalité du système			

Définition et explication des scénarios d'usage			
Prise en compte des limites cognitives humaines			
Communication différenciée selon les rôles	<input checked="" type="checkbox"/> Transparence envers les déployeurs : Fourniture d'une notice d'information claire sur le fonctionnement, les risques, les limites et les conditions d'utilisation du système.		<input checked="" type="checkbox"/> Feedback formateur
Stratégie de prévention des biais			
Évaluation des limites des ensembles de données			
Prise en compte de la diversité des utilisateurs			
Utilisation d'outils d'analyse des biais			
Processus de détection et de contrôle des biais			
Mécanisme de signalement des biais			
Communication sur les biais et les recours			
Prise en compte des effets indirects		<input checked="" type="checkbox"/> Responsabilité sociale : 1 : Anticipation des médiations (morales, comportementales, cognitives).	

		2 : Analyse des effets négatifs (biais, perte d'autonomie, surveillance). 3 : intégration de médiations souhaitables (durabilité, surveillance). 4 : intégration de médiations souhaitables (durabilité, transparence). 5 : Partage de la responsabilité entre concepteurs, utilisateurs et systèmes.	
Évaluation de la variabilité des décisions			
Analyse des causes de variabilité			
Évaluation de l'impact de la variabilité sur les droits fondamentaux			
Définition de l'équité appliquée			
Mesure de l'équité			
Mécanismes de garantie de l'équité			
Adaptation aux préférences et capacités individuelles			
Prise en compte des besoins spécifiques et du handicap			
Compatibilité avec les technologies d'assistance			
Consultation des communautés concernées			

Évaluation de l'impact sur les groupes d'utilisateurs			
Diversité de l'équipe de conception			
Évaluation des effets disproportionnés		<p>✅ Responsabilité sociale :</p> <p>1 : Anticipation des médiations (morales, comportementales, cognitives).</p> <p>2 : Analyse des effets négatifs (biais, perte d'autonomie, surveillance).</p> <p>3 : intégration de médiations souhaitables (durabilité, surveillance).</p> <p>4 : intégration de médiations souhaitables (durabilité, transparence).</p> <p>5 : Partage de la responsabilité entre concepteurs, utilisateurs et systèmes.</p>	
Prise en compte des retours de groupes externes			
Participation des parties prenantes		<p>✅ Gouvernance éthique :</p> <p>1 : Cartographie des parties prenantes.</p> <p>2 : Elicitation des valeurs et traduction en exigences techniques.</p> <p>3 : Evaluation continue via audits et retours d'expérience</p>	

Mobilisation des acteurs internes			
Mesure de l'impact environnemental			
Réduction de l'empreinte écologique			
Évaluation de l'attachement émotionnel		<p>✅ Responsabilité sociale :</p> <p>1 : Anticipation des médiations (morales, comportementales, cognitives).</p> <p>2 : Analyse des effets négatifs (biais, perte d'autonomie, surveillance).</p> <p>3 : intégration de médiations souhaitables (durabilité, surveillance).</p> <p>4 : intégration de médiations souhaitables (durabilité, transparence).</p> <p>5 : Partage de la responsabilité entre concepteurs, utilisateurs et systèmes.</p>	
Clarté sur la nature simulée de l'interaction			
Analyse des incidences sociales		<p>✅ Responsabilité sociale :</p> <p>1 : Anticipation des médiations (morales, comportementales, cognitives).</p> <p>2 : Analyse des effets négatifs (biais, perte d'autonomie, surveillance).</p> <p>3 : intégration de médiations souhaitables (durabilité, surveillance).</p> <p>4 : intégration de médiations souhaitables (durabilité, transparence).</p>	

		<input checked="" type="checkbox"/> Partage de la responsabilité entre concepteurs, utilisateurs et systèmes.	
Évaluation de l'impact sociétal élargi			
Auditabilité du système	<input checked="" type="checkbox"/> journalisation des événements : Enregistrement systématique des événements critiques pour assurer la traçabilité.		
Analyse des risques multi-acteurs			
Mise en place de formations à la responsabilité			
Identification des publics concernés par la formation			
Intégration du cadre juridique dans les formations			
Mise en place d'un comité d'éthique ou d'un mécanisme équivalent		<input checked="" type="checkbox"/> Gouvernance éthique : 1 : Cartographie des parties prenantes. 2 : Elicitation des valeurs et traduction en exigences techniques. 3 : Evaluation continue via audits et retours d'expérience	
Recours à des audits ou orientations externes			
Mécanisme de signalement des vulnérabilités			
Mécanisme de recensement des valeurs et intérêts		<input checked="" type="checkbox"/> Méthodologie éthique intégrée : 1 : Identification explicite des valeurs humaines fondamentales. 2 : Co-construction de solutions techniques respectueuses de ces valeurs. 3 : Acceptation de compromis fonctionnels pour respecter l'éthique.	
Processus de décision pour les arbitrages		<input checked="" type="checkbox"/> Méthodologie éthique intégrée : 1 : Identification explicite des valeurs humaines fondamentales. 2 : Co-construction de solutions techniques respectueuses de ces valeurs. 3 : Acceptation de compromis fonctionnels pour respecter l'éthique.	

Documentation des arbitrages	<input checked="" type="checkbox"/> Documentation technique : 1 : Description des algorithmes, sources de données, méthodes de traitement (ex. : anonymisation, réduction des biais). 2 : Preuves des conformité aux normes de sécurité et d'éthique. 3 : Plans d'atténuation des risques et d'intervention en cas de défaillance. 4 : Mécanismes d'audit et d'explicabilité.		<input checked="" type="checkbox"/> Ethique by design
Mécanismes de recours en cas de préjudice			
Information sur les voies de recours			

Nous n'avons rien trouver concernant l'inscription au registre public européen, donc nous n'avons pas pu l'intégrer dans la grille.

Résultats

Notre grille d'analyse refléter 24 critères de transparence sur 129 critères aux total. Après confrontation avec notre grille d'analyse, nous constatons que les critères de transparence serait de 27 mais qu'elle serait diviser dans toutes les dimensions proposé par l'outil ALTAI et pas seulement celle définit dans la dimension de transparence (traçabilité, explicabilité, communication). Dans le tableau suivant, ont peut voir les critères concordants avec notre grille d'analyse, ceux étant non concordant ainsi que ceux relative à la dimension de transparence de l'outil ALTAI mais que nous n'avons pas retrouver.




Critères concordants (27 critères)	Critères non concordants (hors transparence)	Critères non concordants (transparence)
Signalement des décisions algorithmiques Définition du niveau de contrôle humain Mécanismes de supervision humaine Stratégie de contrôle du système Alignement sur les normes reconnues Traçabilité des accès aux données Documentation des méthodes de conception Documentation des méthodes de test et de validation Documentation des résultats et des décisions Compréhensibilité des décisions du système Explication des choix algorithmiques Conception orientée vers l'interprétabilité Choix de modèles interprétables Accessibilité de l'information pour tous les profils d'utilisateurs Prise en compte des retours utilisateurs Communication des risques et biais Communication différenciée selon les rôles Prise en compte des effets indirects Évaluation des effets disproportionnés Participation des parties prenantes Évaluation de l'attachement émotionnel Analyse des incidences sociales Auditabilité du système Mise en place d'un comité d'éthique ou d'un mécanisme équivalent Mécanisme de recensement des valeurs et intérêts Processus de décision pour les arbitrages Documentation des arbitrages	Analyse d'impact sur les droits fondamentaux Documentation des arbitrages éthiques Interaction avec la prise de décision humaine Préservation de l'autonomie décisionnelle Identification des agents virtuels Répartition des tâches entre IA et humains Renforcement des capacités humaines Prévention de la dépendance excessive Identification de l'acteur humain responsable Auditabilité et gouvernance de l'autonomie Contrôle renforcé pour les systèmes autonomes ou auto-apprenants Détection des anomalies et réponse aux risques Présence d'un mécanisme d'arrêt d'urgence Évaluation des vulnérabilités potentielles Identification des types de vulnérabilités Mesures de résilience et d'intégrité Comportement en environnement imprévu Gestion du double usage Existence de solutions de secours Évaluation du niveau de risque contextuel Processus d'évaluation des risques et de la sécurité Communication des risques pour l'intégrité physique Politique d'assurance Plan de gestion des usages abusifs ou détournés Évaluation des dommages potentiels Règles de responsabilité et protection des consommateurs Évaluation des impacts environnementaux et sur les animaux Analyse des menaces de cybersécurité Évaluation des défaillances critiques Définition de seuils de déclenchement Test des solutions de secours Définition du niveau de précision attendu Méthodes de mesure de la précision Qualité et actualisation des données Évaluation des besoins en données supplémentaires Évaluation des préjudices liés aux erreurs Suivi du taux d'erreurs Plan de correction des erreurs	Mise en place de mesures de traçabilité Documentation de l'entraînement des modèles (IA fondée sur l'apprentissage) Évaluation de l'influence sur les processus décisionnels Justification du déploiement dans un domaine spécifique Évaluation du modèle économique Analyse des données d'entraînement et de test Accès à la logique interne du modèle Identification claire du système comme étant une IA Explication des critères de décision Communication élargie selon les publics Indication de la finalité du système Définition et explication des scénarios d'usage Prise en compte des limites cognitives humaines Évaluation de l'impact sociétal élargi



	Prise en compte des contextes spécifiques Méthodes de vérification de la fiabilité et de la reproductibilité Documentation des réglages critiques Application des processus de test Communication de la fiabilité aux utilisateurs Mécanisme de signalement des atteintes à la vie privée Évaluation du type et de la portée des données Minimisation de l'usage de données sensibles Mécanismes de consentement et de contrôle Mesures de protection des données Mobilisation du DPO (Data Protection Officer) Mécanismes de contrôle des données Évaluation de la qualité des sources externes Garantie de la qualité et de l'intégrité des données Vérification de la sécurité des ensembles de données Documentation des protocoles de gouvernance des données Analyse des accès aux données Qualification et légitimité des personnes ayant accès aux données Stratégie de prévention des biais Évaluation des limites des ensembles de données Prise en compte de la diversité des utilisateurs Utilisation d'outils d'analyse des biais Processus de détection et de contrôle des biais Mécanisme de signalement des biais Communication sur les biais et les recours Évaluation de la variabilité des décisions Analyse des causes de variabilité Évaluation de l'impact de la variabilité sur les droits fondamentaux Définition de l'équité appliquée Mesure de l'équité Mécanismes de garantie de l'équité Adaptation aux préférences et capacités individuelles Prise en compte des besoins spécifiques et du handicap Compatibilité avec les technologies d'assistance Consultation des communautés concernées Évaluation de l'impact sur les groupes d'utilisateurs Diversité de l'équipe de conception Prise en compte des retours de groupes externes Mobilisation des acteurs internes Mesure de l'impact environnemental Réduction de l'empreinte écologique Clarté sur la nature simulée de l'interaction Analyse des risques multi-acteurs Mise en place de formations à la responsabilité Identification des publics concernés par la formation Intégration du cadre juridique dans les formations	
--	--	--

	Recours à des audits ou orientations externes Mécanisme de signalement des vulnérabilités Mécanismes de recours en cas de préjudice Information sur les voies de recours	
--	---	--



Critères ALTAIL vs critères de transparence des entreprises

Collecte de données ciblées

128 critères	Entreprises Salesforce	Entreprise Grammarly	Entreprise Turnitin	Copilot (lien)	Gemini (lien 1) (lien 2)	Claude AI (Anthropic) (lien)
Analyse d'impact sur les droits fondamentaux						
Documentation des arbitrages éthiques						
Interaction avec la prise de décision humaine		 Contrôle utilisateur sur les suggestions : « Les utilisateurs peuvent accepter ou refuser les suggestions, ajuster les types de suggestions, et reçoivent des explications pour prendre des décisions éclairées ».	 L'IA doit soutenir les objectifs pédagogiques et protéger l'intégrité académique : « L'IA est conçue pour renforcer l'apprentissage, non pour le remplacer ou le biaiser ».			
Préservation de l'autonomie décisionnelle		 Contrôle utilisateur sur les suggestions : « Les utilisateurs peuvent accepter ou refuser les suggestions, ajuster les types de suggestions, et reçoivent des explications pour prendre des décisions éclairées ».				

Signalement des décisions algorithmiques		 Sensibilisation des utilisateurs : « Les utilisateurs doivent savoir clairement quand ils interagissent avec l'IA, identifier le contenu généré par l'IA, et comprendre comment les décisions sont prises ».				
Identification des agents virtuels						
Répartition des tâches entre IA et humains						
Renforcement des capacités humaines			 L'IA doit soutenir les objectifs pédagogiques et protéger l'intégrité académique : « L'IA est conçue pour renforcer l'apprentissage, non pour le remplacer ou le biaiser ».			
Prévention de la dépendance excessive						
Définition du niveau de contrôle humain						
Identification de l'acteur humain responsable						
Mécanismes de supervision humaine						
Auditabilité et gouvernance de l'autonomie						









Contrôle renforcé pour les systèmes autonomes ou auto-apprenants						
Détection des anomalies et réponse aux risques						
Présence d'un mécanisme d'arrêt d'urgence					<input checked="" type="checkbox"/> Sécurité des données et stockage à froid : « Les clés privées sont stockées hors ligne dans des modules de sécurité matériels, jamais connectés à Internet, avec des contrôles d'accès biométriques »	
Évaluation des vulnérabilités potentielles						<input checked="" type="checkbox"/> Sécurité, auditabilité et conformité : « Anthropic décrit les méthodes de test, les résultats de sécurité, et les engagements de conformité (ex. : ASL-2)
Identification des types de vulnérabilités						
Mesures de résilience et d'intégrité						
Comportement en environnement imprévu						
Gestion du double usage						

Existence de solutions de secours						
Évaluation du niveau de risque contextuel						
Processus d'évaluation des risques et de la sécurité						
Communication des risques pour l'intégrité physique						
Politique d'assurance					 Règlementation et responsabilité : « Gemini est un dépositaire fiduciaire qualifié, réglementé par l'Etat de new York, avec une couverture d'assurance pour les pertes liées aux cryptomonnaies »	
Plan de gestion des usages abusifs ou détournés						
Évaluation des dommages potentiels						
Règles de responsabilité et protection des consommateurs					 Règlementation et responsabilité : « Gemini est un dépositaire fiduciaire qualifié, réglementé par l'Etat de new York, avec une couverture d'assurance pour les pertes liées aux cryptomonnaies »	
Évaluation des impacts environnementaux et sur les animaux						

Analyse des menaces de cybersécurité						
Évaluation des défaillances critiques						
Définition de seuils de déclenchement						
Test des solutions de secours						
Définition du niveau de précision attendu						
Méthodes de mesure de la précision						
Qualité et actualisation des données						
Évaluation des besoins en données supplémentaires						
Évaluation des préjudices liés aux erreurs						
Suivi du taux d'erreurs						
Plan de correction des erreurs						
Stratégie de contrôle du système						
Prise en compte des contextes spécifiques				<input checked="" type="checkbox"/> Importance du contexte d'utilisation (technologie, personnes, environnement) : « Microsoft souligne que l'IA ne se		





				limite pas à la technologie, mais inclut les utilisateurs et leur environnement ».		
Méthodes de vérification de la fiabilité et de la reproductibilité			✓ Respect de critères rigoureux d'analyse statistique et d'apprentissage automatique : « L'IA doit être fondée sur des méthodes robustes et vérifiables »			
Documentation des réglages critiques						
Application des processus de test			✓ Respect de critères rigoureux d'analyse statistique et d'apprentissage automatique : « L'IA doit être fondée sur des méthodes robustes et vérifiables »			
Communication de la fiabilité aux utilisateurs						
Mécanisme de signalement des atteintes à la vie privée						
Évaluation du type et de la portée des données		✓ Documentation détaillée : « Les capacités, limites, données utilisées, méthodologies, risques et directives doivent être documentées. Les utilisateurs doivent savoir comment les données sont utilisées ».				✓ Transparence sur les données d'entraînement : « Claude n'est pas entraîné sur les données des utilisateurs. Les données proviennent de sources publiques, de partenaires tiers, ou sont générées en internes »
Minimisation de l'usage de données sensibles						✓ Transparence sur les données d'entraînement : « Claude n'est pas entraîné sur les données des utilisateurs. Les données proviennent de sources publiques, de partenaires

						tiers, ou sont générées en internes »
Mécanismes de consentement et de contrôle	<input checked="" type="checkbox"/> « Les données et les modèles de nos clients leur appartiennent, c'est pourquoi nous nous assurons qu'ils en gardent toujours le contrôle »	<input checked="" type="checkbox"/> Documentation détaillée : « Les capacités, limites, données utilisées, méthodologies, risques et directives doivent être documentées. Les utilisateurs doivent savoir comment les données sont utilisées ».	<input checked="" type="checkbox"/> Conception selon des critères exigeants de confidentialité et de propriété des données : « Les données des utilisateurs doivent être protégées et rester sous leur contrôle »	<input checked="" type="checkbox"/> Traitement responsable des données : « Microsoft décrit les étapes de traitement des requêtes, incluant des vérifications de conformité à l'IA responsable »		
Mesures de protection des données		<input checked="" type="checkbox"/> Documentation détaillée : « Les capacités, limites, données utilisées, méthodologies, risques et directives doivent être documentées. Les utilisateurs doivent savoir comment les données sont utilisées ».	<input checked="" type="checkbox"/> Conception selon des critères exigeants de confidentialité et de propriété des données : « Les données des utilisateurs doivent être protégées et rester sous leur contrôle »	<input checked="" type="checkbox"/> Traitement responsable des données : « Microsoft décrit les étapes de traitement des requêtes, incluant des vérifications de conformité à l'IA responsable »	<input checked="" type="checkbox"/> Sécurité des données et stockage à froid : « Les clés privées sont stockées hors ligne dans des modules de sécurité matériels, jamais connectés à Internet, avec des contrôles d'accès biométriques »	<input checked="" type="checkbox"/> Transparence sur les données d'entraînement : « Claude n'est pas entraîné sur les données des utilisateurs. Les données proviennent de sources publiques, de partenaires tiers, ou sont générées en internes »
Mobilisation du DPO (Data Protection Officer)						
Alignement sur les normes reconnues						
Mécanismes de contrôle des données						
Évaluation de la qualité des sources externes						
Garantie de la qualité et de l'intégrité des données						

Vérification de la sécurité des ensembles de données					 Sécurité des données et stockage à froid : « Les clés privées sont stockées hors ligne dans des modules de sécurité matériels, jamais connectés à Internet, avec des contrôles d'accès biométriques »	
Documentation des protocoles de gouvernance des données					 Gouvernance des accès et séparation des fonds : « Les fonds des clients sont détenus à un ratio de 1:1 séparés des fonds de l'entreprise, et peuvent être retirés à tout moment »	
Analyse des accès aux données					 Gouvernance des accès et séparation des fonds : « Les fonds des clients sont détenus à un ratio de 1:1 séparés des fonds de l'entreprise, et peuvent être retirés à tout moment »	
Qualification et légitimité des personnes ayant accès aux données	 « Les données et les modèles de nos clients leur appartiennent, c'est pourquoi nous nous assurons qu'ils en gardent toujours le contrôle »		 Conception selon des critères exigeants de confidentialité et de propriété des données : « Les données des utilisateurs doivent être protégées et rester sous leur contrôle »		 Gouvernance des accès et séparation des fonds : « Les fonds des clients sont détenus à un ratio de 1:1 séparés des fonds de l'entreprise, et peuvent être retirés à tout moment »	
Traçabilité des accès aux données	 « Les données et les modèles de nos clients leur appartiennent, c'est pourquoi nous nous assurons qu'ils en gardent toujours le contrôle »				 Conformité et transparence des actifs : « Tous les actifs des clients sont séparés à l'aide d'adresses numériques uniques, vérifiables sur la blockchain. Les auditeurs ont un accès en lecture seule aux soldes, transactions et activités »	

Mise en place de mesures de traçabilité						
Documentation des méthodes de conception	<p>✓ « Nous publions des cartes de modèles qui décrivent la manière dont ils ont été créés, les cas d'usage potentiels, les implications éthiques ou sociétales connues ainsi que les indices de performance »</p>	<p>✓ Documentation détaillée : « Les capacités, limites, données utilisées, méthodologies, risques et directives doivent être documentées. Les utilisateurs doivent savoir comment les données sont utilisées ».</p>	<p>✓ Respect de critères rigoureux d'analyse statistique et d'apprentissage automatique : « L'IA doit être fondée sur des méthodes robustes et vérifiables »</p>	<p>✓ Comprendre le fonctionnement de l'IA, ses capacités et ses limites : « Microsoft insiste sur la nécessité de comprendre comment fonctionne Copilot, ce qu'il peut faire, ce qu'il ne peut pas faire, et dans quel contexte il est utilisé »</p>		<p>✓ Compréhension du fonctionnement du modèle (Claude 3.7 Sonnet) : « Anthropic fournit un résumé clair des capacités, des limites, des méthodes d'entraînement, et des techniques de sécurité de Claude ».</p>
Documentation de l'entraînement des modèles (IA fondée sur l'apprentissage)						
Documentation des méthodes de test et de validation						<p>✓ Sécurité, auditabilité et conformité : « Anthropic décrit les méthodes de test, les résultats de sécurité, et les engagements de conformité (ex. : ASL-2)</p>
Documentation des résultats et des décisions	<p>✓ « Nous publions des cartes de modèles qui décrivent la manière dont ils ont été créés, les cas d'usage potentiels, les implications éthiques ou sociétales connues ainsi que les indices de performance »</p>	<p>✓ Documentation détaillée : « Les capacités, limites, données utilisées, méthodologies, risques et directives doivent être documentées. Les utilisateurs doivent savoir comment les données sont utilisées ».</p>		<p>✓ Comprendre le fonctionnement de l'IA, ses capacités et ses limites : « Microsoft insiste sur la nécessité de comprendre comment fonctionne Copilot, ce qu'il peut faire, ce qu'il ne peut pas faire, et dans quel contexte il est utilisé »</p> <p>✓ Grounding (ancrage contextuel des réponses) : « Copilot utilise des sources</p>	<p>✓ Conformité et transparence des actifs : « Tous les actifs des clients sont séparés à l'aide d'adresses numériques uniques, vérifiables sur la blockchain. Les auditeurs ont un accès en lecture seule aux soldes, transactions et activités »</p>	<p>✓ Compréhension du fonctionnement du modèle (Claude 3.7 Sonnet) : « Anthropic fournit un résumé clair des capacités, des limites, des méthodes d'entraînement, et des techniques de sécurité de Claude ».</p>

				comme Microsoft Graph pour contextualiser les réponses et améliorer leur pertinence »		
Compréhensibilité des décisions du système	<input checked="" type="checkbox"/> « Nous souhaitons que nos clients comprennent les tenants et aboutissants de chaque suggestion générées par l'IA »	<input checked="" type="checkbox"/> Sensibilisation des utilisateurs « Les utilisateurs doivent savoir clairement quand ils interagissent avec l'IA, identifier le contenu généré par l'IA, et comprendre comment les décisions sont prises ».		<input checked="" type="checkbox"/> Grounding (ancrage contextuel des réponses) : « Copilot utilise des sources comme Microsoft Graph pour contextualiser les réponses et améliorer leur pertinence »		
Explication des choix algorithmiques	<input checked="" type="checkbox"/> « Nous souhaitons que nos clients comprennent les tenants et aboutissants de chaque suggestion générées par l'IA »	<input checked="" type="checkbox"/> Sensibilisation des utilisateurs « Les utilisateurs doivent savoir clairement quand ils interagissent avec l'IA, identifier le contenu généré par l'IA, et comprendre comment les décisions sont prises ».				
Évaluation de l'influence sur les processus décisionnels						
Justification du déploiement dans un domaine spécifique						
Évaluation du modèle économique	<input checked="" type="checkbox"/> « Nous publions des cartes de modèles qui décrivent la manière dont ils ont été créés, les cas d'usage potentiels, les implications éthiques ou sociétales connues	<input checked="" type="checkbox"/> Développement et limites du système : « Les utilisateurs doivent comprendre les risques, les conflits d'intérêts, et les motivations commerciales. Cela renforce la confiance dans l'IA ».				

	ainsi que les indices de performance »					
Conception orientée vers l'interprétabilité						 Compréhension du fonctionnement du modèle (Claude 3.7 Sonnet) : « Anthropic fournit un résumé clair des capacités, des limites, des méthodes d'entraînement, et des techniques de sécurité de Claude ».
Choix de modèles interprétables				 Transparence sur les choix de conception et d'implémentation : « Microsoft décrit les choix faits dans la conception du système, notamment le recours à différents modèles selon les besoins (vitesse, créativité, etc.) »		
Analyse des données d'entraînement et de test						
Accès à la logique interne du modèle						
Identification claire du système comme étant une IA		 Sensibilisation des utilisateurs « Les utilisateurs doivent savoir clairement quand ils interagissent avec l'IA, identifier le contenu généré par l'IA, et comprendre comment les décisions sont prises ».				
Explication des critères de décision	 « Nous souhaitons que nos clients comprennent les tenants et aboutissants de chaque suggestion générées par l'IA »			 Grounding (ancrage contextuel des réponses) : « Copilot utilise des sources comme Microsoft Graph pour contextualiser les réponses et améliorer leur pertinence »		

Accessibilité de l'information pour tous les profils d'utilisateurs		<input checked="" type="checkbox"/> Contrôle utilisateur sur les suggestions : « Les utilisateurs peuvent accepter ou refuser les suggestions, ajuster les types de suggestions, et reçoivent des explications pour prendre des décisions éclairées ».	<input checked="" type="checkbox"/> Amélioration continue pour plus d'accessibilité, d'équité et de bénéfices : « L'IA doit évoluer pour mieux servir tous les profils d'utilisateurs »			
Prise en compte des retours utilisateurs						
Communication des risques et biais		<input checked="" type="checkbox"/> Développement et limites du système : « Les utilisateurs doivent comprendre les risques, les conflits d'intérêts, et les motivations commerciales. Cela renforce la confiance dans l'IA ».				
Communication élargie selon les publics						
Indication de la finalité du système	<input checked="" type="checkbox"/> Nous leur fournissons également une charte clarifiant les conditions d'utilisation et les applications des solutions d'IA »					
Définition et explication des scénarios d'usage	<input checked="" type="checkbox"/> Nous leur fournissons également une charte clarifiant les conditions d'utilisation et les applications des solutions d'IA »					
Prise en compte des limites cognitives humaines						







Communication différenciée selon les rôles	✓ Nous leurs fournissons également une charte clarifiant les conditions d'utilisation et les applications des solutions d'IA »			✓ Comprendre le fonctionnement de l'IA, ses capacité et ses limites : « Microsoft insiste sur la nécessité de comprendre comment fonctionne Copilot, ce qu'il peut faire, ce qu'il ne peut pas faire, et dans quel contexte il est utilisé »		✓ Compréhension du fonctionnement du modèle (Claude 3.7 Sonnet) : « Anthropic fournit un résumé clair des capacités, des limites, des méthodes d'entraînement, et des techniques de sécurité de Claude ».
Stratégie de prévention des biais			✓ L'IA doit atténuer la subjectivité : « Elle doit être conçue pour réduire les biais et favoriser l'équité »			✓ Constitutionnel AI et alignement sur des principes éthiques : « Claude est entraîné avec une méthode appelée « Conditionnel AI » pour garantir des réponses utiles, inoffensives et honnêtes ».
Évaluation des limites des ensembles de données				✓ Comprendre le fonctionnement de l'IA, ses capacité et ses limites : « Microsoft insiste sur la nécessité de comprendre comment fonctionne Copilot, ce qu'il peut faire, ce qu'il ne peut pas faire, et dans quel contexte il est utilisé »		
Prise en compte de la diversité des utilisateurs			✓ Conception avec un groupe diversifié d'étudiants et d'enseignants : « L'IA doit être codéveloppée avec les utilisateurs finaux pour refléter leurs besoins »			
Utilisation d'outils d'analyse des biais						
Processus de détection et de contrôle des biais			✓ L'IA doit atténuer la subjectivité : « Elle doit être conçue pour réduire les biais et favoriser l'équité »			

Mécanisme de signalement des biais						
Communication sur les biais et les recours						
Prise en compte des effets indirects						
Évaluation de la variabilité des décisions						
Analyse des causes de variabilité						
Évaluation de l'impact de la variabilité sur les droits fondamentaux						
Définition de l'équité appliquée						
Mesure de l'équité			<input checked="" type="checkbox"/> L'IA doit atténuer la subjectivité : « Elle doit être conçue pour réduire les biais et favoriser l'équité »			
Mécanismes de garantie de l'équité						
Adaptation aux préférences et capacités individuelles			<input checked="" type="checkbox"/> Amélioration continue pour plus d'accessibilité, d'équité et de bénéfices : « L'IA doit évoluer pour mieux servir tous les profils d'utilisateurs »			

Prise en compte des besoins spécifiques et du handicap						
Compatibilité avec les technologies d'assistance						
Consultation des communautés concernées			<input checked="" type="checkbox"/> Conception avec un groupe diversifié d'étudiants et d'enseignants : « L'IA doit être codéveloppée avec les utilisateurs finaux pour refléter leurs besoins »			
Évaluation de l'impact sur les groupes d'utilisateurs						
Diversité de l'équipe de conception						
Évaluation des effets disproportionnés			<input checked="" type="checkbox"/> Amélioration continue pour plus d'accessibilité, d'équité et de bénéfices : « L'IA doit évoluer pour mieux servir tous les profils d'utilisateurs »			
Prise en compte des retours de groupes externes						

Participation des parties prenantes			<input checked="" type="checkbox"/> Conception avec un groupe diversifié d'étudiants et d'enseignants : « L'IA doit être codéveloppée avec les utilisateurs finaux pour refléter leurs besoins »	<input checked="" type="checkbox"/> Importance du contexte d'utilisation (technologie, personnes, environnement) : « Microsoft souligne que l'IA ne se limite pas à la technologie, mais inclut les utilisateurs et leur environnement ».		
Mobilisation des acteurs internes						
Mesure de l'impact environnemental						
Réduction de l'empreinte écologique						
Évaluation de l'attachement émotionnel						
Clarté sur la nature simulée de l'interaction						

Analyse des incidences sociales	<p>✔ « Nous publions des cartes de modèles qui décrivent la manière dont ils ont été créés, les cas d'usage potentiels, les implications éthiques ou sociétales connues ainsi que les indices de performance »</p>	<p>✔ Développement et limites du système : « Les utilisateurs doivent comprendre les risques, les conflits d'intérêts, et les motivations commerciales. Cela renforce la confiance dans l'IA ».</p>	<p>✔ L'IA doit soutenir les objectifs pédagogiques et protéger l'intégrité académique : « L'IA est conçue pour renforcer l'apprentissage, non pour le remplacer ou le biaiser ».</p>	<p>✔ Importance du contexte d'utilisation (technologie, personnes, environnement) : « Microsoft souligne que l'IA ne se limite pas à la technologie, mais inclut les utilisateurs et leur environnement ».</p>		
Évaluation de l'impact sociétal élargi						
Auditabilité du système				<p>✔ Traitement responsable des données : « Microsoft décrit les étapes de traitement des requêtes, incluant des vérifications de conformité à l'IA responsable »</p>	<p>✔ Conformité et transparence des actifs : « Tous les actifs des clients sont séparés à l'aide d'adresses numériques uniques, vérifiables sur la blockchain. Les auditeurs ont un accès en lecture seule aux soldes, transactions et activités »</p>	<p>✔ Sécurité, auditabilité et conformité : « Anthropic décrit les méthodes de test, les résultats de sécurité, et les engagements de conformité (ex. : ASL-2)</p>
Analyse des risques multi-acteurs						
Mise en place de formations à la responsabilité						
Identification des publics concernés par la formation						
Intégration du cadre juridique dans les formations						
Mise en place d'un comité d'éthique ou d'un mécanisme équivalent						

Recours à des audits ou orientations externes					 Réglementation et responsabilité : « Gemini est un dépositaire fiduciaire qualifié, réglementé par l'Etat de New York, avec une couverture d'assurance pour les pertes liées aux cryptomonnaies »	
Mécanisme de signalement des vulnérabilités						
Mécanisme de recensement des valeurs et intérêts						 Constitutionnal AI et alignement sur des principes éthiques : « Claude est entraîné avec une méthode appelée « Conditionnal AI » pour garantir des réponses utiles, inoffensives et honnêtes ».
Processus de décision pour les arbitrages				 Transparence sur les choix de conception et d'implémentation : « Microsoft décrit les choix faits dans la conception du système, notamment le recours à différents modèles selon les besoins (vitesse, créativité, etc.)		
Documentation des arbitrages		 Développement et limites du système : « Les utilisateurs doivent comprendre les risques, les conflits d'intérêts, et les motivations commerciales. Cela renforce la confiance dans l'IA ».		 Transparence sur les choix de conception et d'implémentation : « Microsoft décrit les choix faits dans la conception du système, notamment le recours à différents modèles selon les besoins (vitesse, créativité, etc.)		 Constitutionnal AI et alignement sur des principes éthiques : « Claude est entraîné avec une méthode appelée « Conditionnal AI » pour garantir des réponses utiles, inoffensives et honnêtes ».
Mécanismes de recours en cas de préjudice						
Information sur les voies de recours						

Nous n'avons pas pu analyser l'application ChatGPT car aucune page dédiée à la transparence n'était disponible sur leur site. En faisant des recherches sur internet nous avons pu constater de fort débat concernant la transparence de cette entreprise.

Résultats

Notre grille d'analyse refléter 24 critères de transparence sur 129 critères aux total. Après confrontation avec notre grille d'analyse, nous constatons que les critères de transparence serait de 50 mais qu'elle serait diviser dans toutes les dimensions proposé par l'outil ALTAI et pas seulement celle définit dans la dimension de transparence (traçabilité, explicabilité, communication). Dans le tableau suivant, on peut voir les critères concordants avec notre grille d'analyse, ceux étant non concordant ainsi que ceux relative à la dimension de transparence de l'outil ALTAI mais que nous n'avons pas retrouver.

Critères concordants (50 critères)	Critères non concordants (hors transparence) (70 critères)	Critères non concordants (transparence) (9 critères)
Interaction avec la prise de décision humaine Préservation de l'autonomie décisionnelle Signalement des décisions algorithmiques Renforcement des capacités humaines Présence d'un mécanisme d'arrêt d'urgence Évaluation des vulnérabilités potentielles Politique d'assurance Règles de responsabilité et protection des consommateurs Prise en compte des contextes spécifiques Méthodes de vérification de la fiabilité et de la reproductibilité Application des processus de test Évaluation du type et de la portée des données Minimisation de l'usage de données sensibles Mécanismes de consentement et de contrôle Mesures de protection des données Vérification de la sécurité des ensembles de données Documentation des protocoles de gouvernance des données Analyse des accès aux données Qualification et légitimité des personnes ayant accès aux données Traçabilité des accès aux données Documentation des méthodes de conception Documentation des méthodes de test et de validation Documentation des résultats et des décisions Compréhensibilité des décisions du système Explication des choix algorithmiques Évaluation du modèle économique Conception orientée vers l'interprétabilité Choix de modèles interprétables Identification claire du système comme étant une IA Explication des critères de décision Accessibilité de l'information pour tous les profils d'utilisateurs Communication des risques et biais Indication de la finalité du système Définition et explication des scénarios d'usage Communication différenciée selon les rôles Stratégie de prévention des biais Évaluation des limites des ensembles de données Prise en compte de la diversité des utilisateurs	Analyse d'impact sur les droits fondamentaux Documentation des arbitrages éthiques Identification des agents virtuels Répartition des tâches entre IA et humains Prévention de la dépendance excessive Définition du niveau de contrôle humain Identification de l'acteur humain responsable Mécanismes de supervision humaine Auditabilité et gouvernance de l'autonomie Contrôle renforcé pour les systèmes autonomes ou auto-apprenants Détection des anomalies et réponse aux risques Identification des types de vulnérabilités Mesures de résilience et d'intégrité Comportement en environnement imprévu Gestion du double usage Existence de solutions de secours Évaluation du niveau de risque contextuel Processus d'évaluation des risques et de la sécurité Communication des risques pour l'intégrité physique Plan de gestion des usages abusifs ou détournés Évaluation des dommages potentiels Évaluation des impacts environnementaux et sur les animaux Analyse des menaces de cybersécurité Évaluation des défaillances critiques Définition de seuils de déclenchement Test des solutions de secours Définition du niveau de précision attendu Méthodes de mesure de la précision Qualité et actualisation des données Évaluation des besoins en données supplémentaires Évaluation des préjudices liés aux erreurs Suivi du taux d'erreurs Plan de correction des erreurs Stratégie de contrôle du système Documentation des réglages critiques Communication de la fiabilité aux utilisateurs Mécanisme de signalement des atteintes à la vie privée Mobilisation du DPO (Data Protection Officer)	Mise en place de mesures de traçabilité Documentation de l'entraînement des modèles (IA fondée sur l'apprentissage) Évaluation de l'influence sur les processus décisionnels Justification du déploiement dans un domaine spécifique Analyse des données d'entraînement et de test Accès à la logique interne du modèle Prise en compte des retours utilisateurs Communication élargie selon les publics Prise en compte des limites cognitives humaines

Processus de détection et de contrôle des biais Mesure de l'équité Adaptation aux préférences et capacités individuelles Consultation des communautés concernées Évaluation des effets disproportionnés Participation des parties prenantes Analyse des incidences sociales Auditabilité du système Recours à des audits ou orientations externes Mécanisme de recensement des valeurs et intérêts Processus de décision pour les arbitrages Documentation des arbitrages	Alignement sur les normes reconnues Mécanismes de contrôle des données Évaluation de la qualité des sources externes Garantie de la qualité et de l'intégrité des données Utilisation d'outils d'analyse des biais Mécanisme de signalement des biais Communication sur les biais et les recours Prise en compte des effets indirects Évaluation de la variabilité des décisions Analyse des causes de variabilité Évaluation de l'impact de la variabilité sur les droits fondamentaux Définition de l'équité appliquée Mécanismes de garantie de l'équité Prise en compte des besoins spécifiques et du handicap Compatibilité avec les technologies d'assistance Évaluation de l'impact sur les groupes d'utilisateurs Diversité de l'équipe de conception Prise en compte des retours de groupes externes Mobilisation des acteurs internes Mesure de l'impact environnemental Réduction de l'empreinte écologique Évaluation de l'attachement émotionnel Clarté sur la nature simulée de l'interaction Évaluation de l'impact sociétal élargi Analyse des risques multi-acteurs Mise en place de formations à la responsabilité Identification des publics concernés par la formation Intégration du cadre juridique dans les formations Mise en place d'un comité d'éthique ou d'un mécanisme équivalent Mécanisme de signalement des vulnérabilités Mécanismes de recours en cas de préjudice Information sur les voies de recours	
--	--	--

Comparaisons critères de transparence entre : outil ALTAIL, cadre théorique, entreprises

Critères de transparence outil ALTAIL (24 critères)	Critères de transparence cadre théorique (27 critères)	Critères de transparence des entreprises (50 critères)
---	--	--

<p>Mise en place de mesures de traçabilité</p> <p>Documentation des méthodes de conception</p> <p>Documentation de l'entraînement des modèles (IA fondée sur l'apprentissage)</p> <p>Documentation des méthodes de test et de validation</p> <p>Documentation des résultats et des décisions</p> <p>Compréhensibilité des décisions du système</p> <p>Explication des choix algorithmiques</p> <p>Évaluation de l'influence sur les processus décisionnels</p> <p>Justification du déploiement dans un domaine spécifique</p> <p>Évaluation du modèle économique</p> <p>Conception orientée vers l'interprétabilité</p> <p>Choix de modèles interprétables</p> <p>Analyse des données d'entraînement et de test</p> <p>Accès à la logique interne du modèle</p> <p>Identification claire du système comme étant une IA</p> <p>Explication des critères de décision</p> <p>Accessibilité de l'information pour tous les profils d'utilisateurs</p> <p>Prise en compte des retours utilisateurs</p> <p>Communication des risques et biais</p> <p>Communication élargie selon les publics</p> <p>Indication de la finalité du système</p> <p>Définition et explication des scénarios d'usage</p> <p>Prise en compte des limites cognitives humaines</p> <p>Communication différenciée selon les rôles</p>	<p>Signalement des décisions algorithmiques</p> <p>Définition du niveau de contrôle humain</p> <p>Mécanismes de supervision humaine</p> <p>Stratégie de contrôle du système</p> <p>Alignement sur les normes reconnues</p> <p>Traçabilité des accès aux données</p> <p>Documentation des méthodes de conception</p> <p>Documentation des méthodes de test et de validation</p> <p>Documentation des résultats et des décisions</p> <p>Compréhensibilité des décisions du système</p> <p>Explication des choix algorithmiques</p> <p>Conception orientée vers l'interprétabilité</p> <p>Choix de modèles interprétables</p> <p>Accessibilité de l'information pour tous les profils d'utilisateurs</p> <p>Prise en compte des retours utilisateurs</p> <p>Communication des risques et biais</p> <p>Communication différenciée selon les rôles</p> <p>Prise en compte des effets indirects</p> <p>Évaluation des effets disproportionnés</p> <p>Participation des parties prenantes</p> <p>Évaluation de l'attachement émotionnel</p> <p>Analyse des incidences sociales</p> <p>Auditabilité du système</p> <p>Mise en place d'un comité d'éthique ou d'un mécanisme équivalent</p> <p>Mécanisme de recensement des valeurs et intérêts</p> <p>Processus de décision pour les arbitrages</p> <p>Documentation des arbitrages</p>	<p>Interaction avec la prise de décision humaine</p> <p>Préservation de l'autonomie décisionnelle</p> <p>Signalement des décisions algorithmiques</p> <p>Renforcement des capacités humaines</p> <p>Présence d'un mécanisme d'arrêt d'urgence</p> <p>Évaluation des vulnérabilités potentielles</p> <p>Politique d'assurance</p> <p>Règles de responsabilité et protection des consommateurs</p> <p>Prise en compte des contextes spécifiques</p> <p>Méthodes de vérification de la fiabilité et de la reproductibilité</p> <p>Application des processus de test</p> <p>Évaluation du type et de la portée des données</p> <p>Minimisation de l'usage de données sensibles</p> <p>Mécanismes de consentement et de contrôle</p> <p>Mesures de protection des données</p> <p>Vérification de la sécurité des ensembles de données</p> <p>Documentation des protocoles de gouvernance des données</p> <p>Analyse des accès aux données</p> <p>Qualification et légitimité des personnes ayant accès aux données</p> <p>Traçabilité des accès aux données</p> <p>Documentation des méthodes de conception</p> <p>Documentation des méthodes de test et de validation</p> <p>Documentation des résultats et des décisions</p> <p>Compréhensibilité des décisions du système</p> <p>Explication des choix algorithmiques</p> <p>Évaluation du modèle économique</p> <p>Conception orientée vers l'interprétabilité</p> <p>Choix de modèles interprétables</p> <p>Identification claire du système comme étant une IA</p> <p>Explication des critères de décision</p> <p>Accessibilité de l'information pour tous les profils d'utilisateurs</p> <p>Communication des risques et biais</p> <p>Indication de la finalité du système</p> <p>Définition et explication des scénarios d'usage</p> <p>Communication différenciée selon les rôles</p> <p>Stratégie de prévention des biais</p> <p>Évaluation des limites des ensembles de données</p> <p>Prise en compte de la diversité des utilisateurs</p> <p>Processus de détection et de contrôle des biais</p> <p>Mesure de l'équité</p> <p>Adaptation aux préférences et capacités individuelles</p> <p>Consultation des communautés concernées</p> <p>Évaluation des effets disproportionnés</p> <p>Participation des parties prenantes</p> <p>Analyse des incidences sociales</p> <p>Auditabilité du système</p>
--	---	---

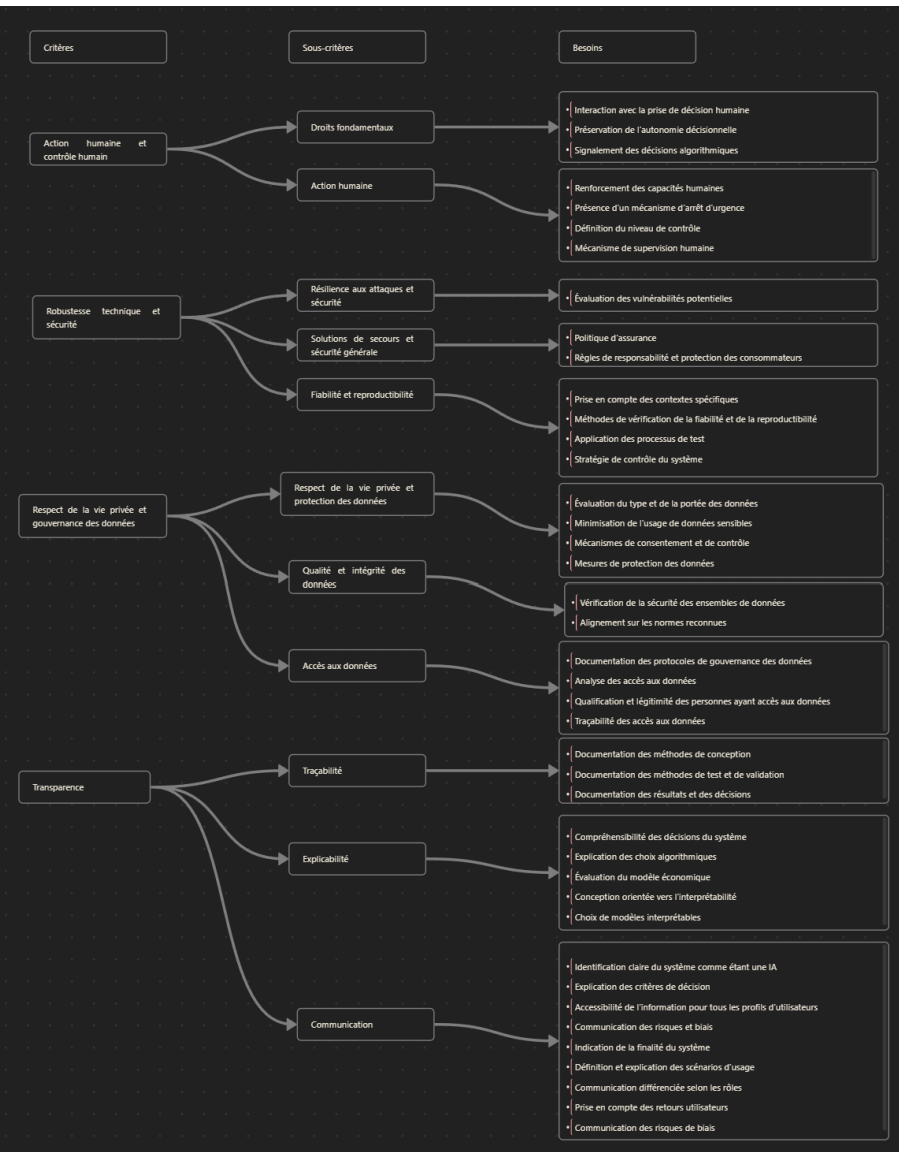
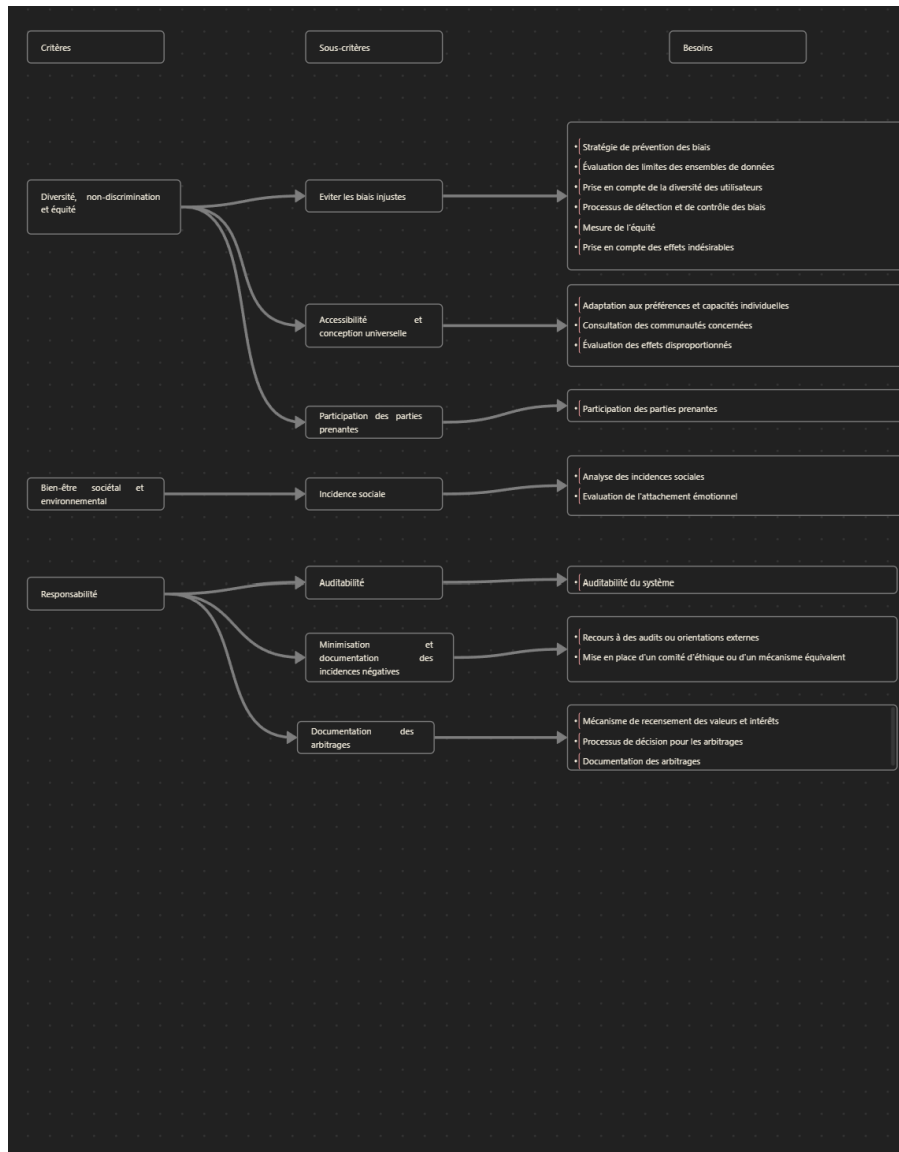
		Recours à des audits ou orientations externes Mécanisme de recensement des valeurs et intérêts Processus de décision pour les arbitrages Documentation des arbitrages
--	--	--

Nous pouvons analyser que dans ce tableau récapitulatif, une quantité et une richesse de critères sont présent :

- Le référentiel de l'outil ALTAI comprend 24 critères qui reflète du cadre normatif européen, structuré autour de la traçabilité, l'explicabilité et la communication.
- Le référentiel de notre cadre théorique comprend 27 critères qui enrichissent le référentiel ALTAI avec des dimensions comme l'équité, les effets indirects ou la participation.
- Le référentiel des entreprises comprend 50 critères d'une très grande granularité, mais hétérogène avec des critères très détaillées et d'autres absent ou flous.

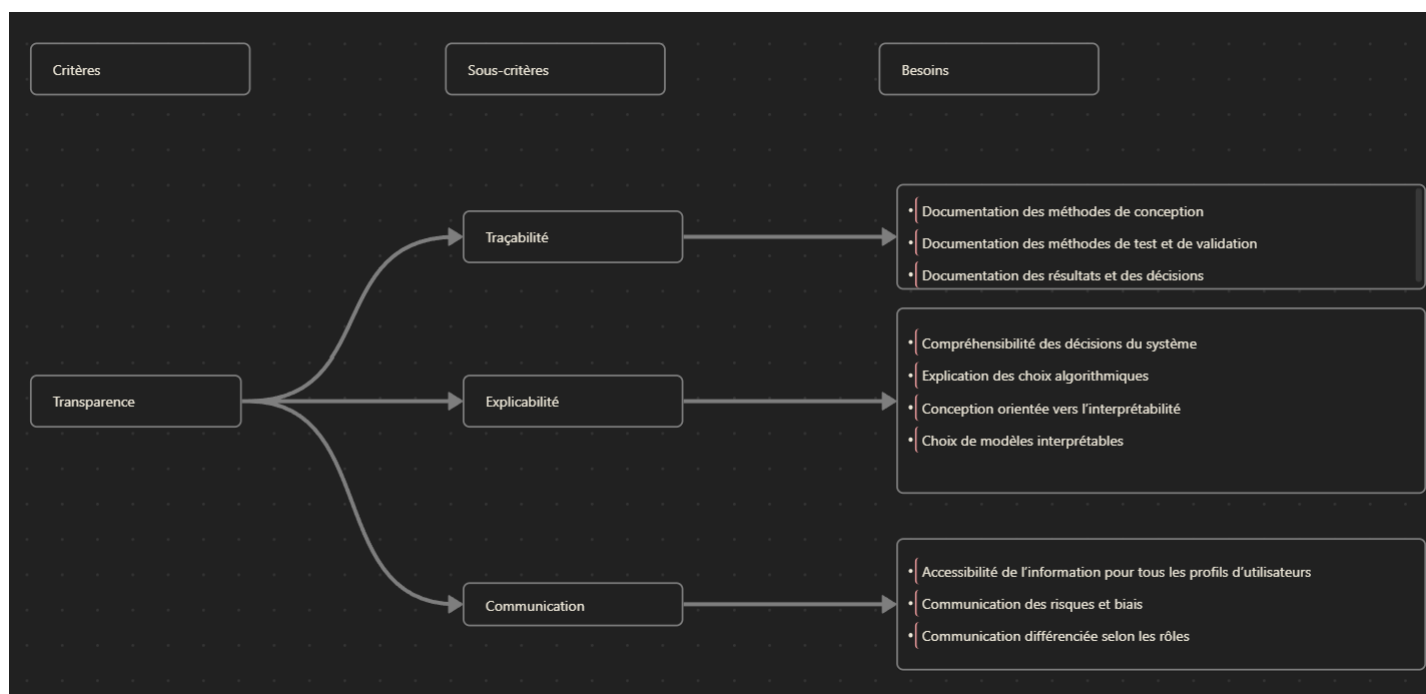
Les entreprises mobilisent un plus grand nombre de critères, mais cela ne signifie pas nécessairement une meilleure couverture. Il s'agit souvent de redondances, de reformulations ou de critères spécifiques à leur contexte commercial.

Les critères sont présents à la fois dans l'outil ALTAI, la grille d'analyse et les pratiques des entreprises sont présenté ci-dessous :



Ces critères représentent le socle commun de la gouvernance éthique de l'IA. Nous pouvons voir qu'il y a un équilibre entre droits fondamentaux, action humaine, résilience aux attaques et sécurité, solutions de secours et sécurité générale, fiabilité et reproductibilité, responsabilité de la vie privée et protection des données, qualité et intégrité des données, accès aux données, traçabilité, explicabilité, communication, éviter les biais injustes, accessibilité et conception universelle, participation des parties prenantes, incidences sociales, auditabilité, minimisation et documentation des incidences négatives, documentation des arbitrages.

Les critères suivants sont présents dans les trois référentiels, ce qui montre un noyau dur de la transparence :



Ces critères constituent les fondamentaux de la transparence algorithmique. Leur présence dans les trois référentiels montre leur légitimité et leur applicabilité.

Certains critères sont recommandés par l'ALTAI et intégrés dans le cadre théorique, mais peu ou pas repris par les entreprises, tel que :

- Définition du niveau de contrôle humain
- Mécanisme de supervision humaine
- Signalement des décisions algorithmiques
- Stratégie de contrôle du système
- Alignement sur les normes reconnues.

Ces critères relèvent d'une transparence structurelle ou organisationnelle, souvent négligée par les entreprises au profit d'une transparence plus marketing ou centrée sur l'interface utilisateur.

Le cadre théorique introduit des critères absents d'ALTAI, mais essentiels pour une approche éthique approfondie :

- Prise en compte des effets indirects
- Evaluation des effets disproportionnés
- Participation des parties prenantes
- Evaluation de l'attachement émotionnel
- Analyse des incidences sociales
- Mise en place d'un comité d'éthique
- Documentation des arbitrages

Ces critères traduisent une volonté de gouvernance éthique et de responsabilité sociale, encore peu intégrée dans les pratiques industrielles.

Les entreprises introduisent des critères non présents dans ALTAI ni dans notre cadre théorique, souvent liées à la sécurité des données (ex. : stockage à froid, séparation des fonds), la conformité réglementaire (ex. : ASL-2, SOC2), la contextualisation des réponses (ex. : grounding via Microsoft Graph), la gestion des accès (ex. : journalisation blockchain). Ces critères montent une adaptation aux contraintes opérationnelles, mais peuvent aussi refléter une stratégie de communication plus qu'un engagement éthique profond.

Confrontation réglementation IA Act avec les critères de transparence des entreprises

Critères des entreprises (Approche réglementaire en gras) (5/50 critères)	Critères de notre grille d'analyse (Approche réglementaire en gras) (7/27 critères)
<p>Interaction avec la prise de décision humaine</p> <p>Préservation de l'autonomie décisionnelle</p> <p>Signalement des décisions algorithmiques</p> <p>Renforcement des capacités humaines</p> <p>Présence d'un mécanisme d'arrêt d'urgence</p> <p>Évaluation des vulnérabilités potentielles</p> <p>Politique d'assurance</p> <p>Règles de responsabilité et protection des consommateurs</p> <p>Prise en compte des contextes spécifiques</p> <p>Méthodes de vérification de la fiabilité et de la reproductibilité</p> <p>Application des processus de test</p> <p>Évaluation du type et de la portée des données</p> <p>Minimisation de l'usage de données sensibles</p> <p>Mécanismes de consentement et de contrôle</p> <p>Mesures de protection des données</p> <p>Vérification de la sécurité des ensembles de données</p> <p>Documentation des protocoles de gouvernance des données</p> <p>Analyse des accès aux données</p> <p>Qualification et légitimité des personnes ayant accès aux données</p> <p>Traçabilité des accès aux données</p> <p>Documentation des méthodes de conception</p> <p>Documentation des méthodes de test et de validation</p> <p>Documentation des résultats et des décisions</p> <p>Compréhensibilité des décisions du système</p> <p>Explication des choix algorithmiques</p> <p>Évaluation du modèle économique</p> <p>Conception orientée vers l'interprétabilité</p> <p>Choix de modèles interprétables</p> <p>Identification claire du système comme étant une IA</p> <p>Explication des critères de décision</p> <p>Accessibilité de l'information pour tous les profils d'utilisateurs</p> <p>Communication des risques et biais</p> <p>Indication de la finalité du système</p> <p>Définition et explication des scénarios d'usage</p> <p>Communication différenciée selon les rôles</p> <p>Stratégie de prévention des biais</p> <p>Évaluation des limites des ensembles de données</p> <p>Prise en compte de la diversité des utilisateurs</p> <p>Processus de détection et de contrôle des biais</p> <p>Mesure de l'équité</p> <p>Adaptation aux préférences et capacités individuelles</p>	<p>Signalement des décisions algorithmiques</p> <p>Définition du niveau de contrôle humain</p> <p>Mécanismes de supervision humaine</p> <p>Stratégie de contrôle du système</p> <p>Alignement sur les normes reconnues</p> <p>Traçabilité des accès aux données</p> <p>Documentation des méthodes de conception</p> <p>Documentation des méthodes de test et de validation</p> <p>Documentation des résultats et des décisions</p> <p>Compréhensibilité des décisions du système</p> <p>Explication des choix algorithmiques</p> <p>Conception orientée vers l'interprétabilité</p> <p>Choix de modèles interprétables</p> <p>Accessibilité de l'information pour tous les profils d'utilisateurs</p> <p>Prise en compte des retours utilisateurs</p> <p>Communication des risques et biais</p> <p>Communication différenciée selon les rôles</p> <p>Prise en compte des effets indirects</p> <p>Évaluation des effets disproportionnés</p> <p>Participation des parties prenantes</p> <p>Évaluation de l'attachement émotionnel</p> <p>Analyse des incidences sociales</p> <p>Auditabilité du système</p> <p>Mise en place d'un comité d'éthique ou d'un mécanisme équivalent</p> <p>Mécanisme de recensement des valeurs et intérêts</p> <p>Processus de décision pour les arbitrages</p> <p>Documentation des arbitrages</p>

Consultation des communautés concernées Évaluation des effets disproportionnés Participation des parties prenantes Analyse des incidences sociales Auditabilité du système Recours à des audits ou orientations externes Mécanisme de recensement des valeurs et intérêts Processus de décision pour les arbitrages Documentation des arbitrages	
---	--

L'analyse des pratiques des entreprises à travers la grille révèle des écarts significatifs entre les critères théoriques issus de l'outil ALTAI et leur mise en œuvre concrète. Plusieurs critères clés ont été examinés, notamment en matières de contrôle humain, de transparence, de documentation et d'audibilité.

Contrôle humain

Deux critères fondamentaux liés au contrôle humains sont absents chez l'ensemble des entreprises analysées :

- Définition du niveau de contrôle humain : aucune entreprise ne précise clairement le niveau de contrôle attendu de la part des utilisateurs humains dans les processus décisionnels.
- Mécanismes de supervision humaine : aucun dispositif explicite de supervision ou de validation humaine n'est documenté, ce qui soulève des questions sur la capacité des utilisateurs à reprendre la main en cas de dérive du système.

Traçabilité des accès aux données

Seules deux entreprises mettent en œuvre des mécanismes de journalisation des accès aux données :

- L'une affirme que « les données et les modèles de nos clients leur appartiennent », garantissant ainsi un contrôle utilisateur.
- L'autre met en place une traçabilité blockchain avec accès en lecture seule pour les auditeurs, ce qui constitue une bonne pratique en matière de transparence et de sécurité.

Documentation des méthodes de conception

Cinq entreprises se distinguent par une documentation technique approfondie :

- Elles publient des cartes de modèles, décrivent les méthodologies d'entraînement, les limites, les risques et les cas d'usage.
- Certaines insistent sur la compréhension du fonctionnement de l'IA par les utilisateurs, ce qui renforce la transparence et l'appropriation.

Documentation des résultats et décisions

Là encore, cinq entreprises fournissent une documentation claire sur les résultats produits par l'IA :

- Elles expliquent les choix algorithmiques, les sources de données, et les mécanismes de contextualisation (ex. : grounding via Microsoft Graph).
- Certaines utilisent des adresses numériques vérifiables pour assurer la transparence des transactions et décisions.

Communication différenciée selon les rôles

Trois entreprises adoptent une communication différenciée selon les profils d'utilisateurs :

- Elles fournissent des chartes d'utilisation, des notices explicatives et des résumés techniques adaptés aux administrateurs ou aux utilisateurs finaux.
- Cette approche est cohérente avec l'idée de transparence située, qui adapte l'information au contexte et au rôle de l'utilisateur.

Auditabilité du système

Trois entreprises mettent en œuvre des mécanismes d'auditabilité :

- Cela inclut la journalisation des événements, des engagements de conformité (ex. : ASL-2), et des accès en lecture seule pour les auditeurs.
- Ces pratiques permettent une vérification externe du fonctionnement du système, mais restent encore limitées à certaines entreprises.

Documentation des arbitrages

Enfin, trois entreprises documentent les arbitrages éthiques réalisés lors de la conception :

- Elles évoquent les conflits d'intérêts, les motivations commerciales, ou encore les choix de modèles selon les besoins.
- L'approche « Constitutionnal AI » d'Anthropic est un exemple de formalisation des valeurs dans la conception algorithmique.