



**MINISTÈRE
DE LA CULTURE**

*Liberté
Égalité
Fraternité*

Rapport de mission relative
à la mise en œuvre du règlement
européen établissant des règles
harmonisées sur l'intelligence
artificielle (« template »)



**PRÉSENTÉ AU CONSEIL SUPÉRIEUR
DE LA PROPRIÉTÉ LITTÉRAIRE ET ARTISTIQUE**

Présidente de la mission : Alexandra Bensamoun

Rapporteur : Lionel Ferreira

Avec le soutien de Frédéric Pascal

Présidente de la mission

Alexandra Bensamoun

Professeur des universités (droit)

Personnalité qualifiée du CSPLA

Rapporteur

Lionel Ferreira

Maître des requêtes au Conseil d'Etat

Avec le soutien de

Frédéric Pascal

Professeur des universités (mathématiques appliquées)

Personnalité qualifiée du CSPLA

Rapport présenté à la réunion plénière du CSPLA du 9 décembre 2024

Son contenu n'engage que ses auteurs

Image de couverture générée par IA avec MidJourney

Prompt : « *Abstract European flag made of glowing digital data and network connections, on a blue background with yellow stars, in a wide banner design. --ar 121:62 --v 6.1* »

Résumé décideurs	5
RAPPORT	8
I. Etat des lieux	9
1. La collecte et l'exploitation des données revêtent une importance majeure, mais s'effectuent dans des conditions qui ne garantissent pas le respect des valeurs et du droit de l'UE.	9
a. La collecte et l'exploitation de données générées par des humains sont devenues un enjeu stratégique pour les fournisseurs de modèles d'IA.....	9
b. L'encadrement des conditions de récupération (par moissonnage ou autrement) et d'exploitation de ces données est insatisfaisant.	11
c. Adoptés alors que l'utilisation massive des contenus protégés par les modèles d'IA générative n'avait pas été anticipée, le règlement sur la protection des données à caractère personnel et la directive relative au droit d'auteur dans le marché unique numérique n'étaient plus à même de garantir de façon satisfaisante le respect des droits des citoyens européens.	13
2. Pour mettre fin à une situation préjudiciable à l'innovation et aux citoyens, l'Union européenne a adopté un règlement relatif à l'intelligence artificielle qui prévoit notamment une obligation de transparence dont la mission doit préciser la portée dans la perspective des négociations entre Etats membres.....	15
a. Cette situation préjudicie à l'innovation et aux citoyens.	15
b. Le règlement sur l'intelligence artificielle (RIA) a pour ambition de constituer un cadre à la fois favorable à l'innovation et respectueux des valeurs de l'UE.....	16
c. La mission créée par la ministre de la Culture a pour objet de préciser la portée des dispositions de l'article 53, 1, d, et de proposer un modèle de résumé qui puisse être porté au nom de la France au niveau européen.	18
II. Analyse.....	20
1. L'obligation de mettre en place une politique de conformité et celle de mettre à disposition du public un résumé suffisamment détaillé participent d'un même objectif : améliorer la transparence.	20
a. Le RIA semble considérer qu'il s'agit de deux obligations à traiter en silo.	20
b. Pour la mission, les deux obligations sont indissociables.....	20
c. Le modèle de résumé doit intégrer des éléments relatifs à la conformité, et notamment au respect de la réserve de droits.....	21
2. La transparence ne consiste pas à laisser les acteurs s'autoréguler et peut aller jusqu'à exiger la production d'une liste des contenus utilisés.	22
a. Le résumé ne saurait se borner à lister les principales sources de données dans l'attente de la création d'un marché de la donnée.....	22
b. Le texte n'exclut pas de lister les contenus protégés utilisés pour l'entraînement des modèles.	24
c. La portée normative du résumé doit être proportionnée à l'objectif poursuivi : aider les intéressés à faire valoir leurs droits.....	24

d.	L’effort doit être poursuivi pour donner à la transparence les conséquences attendues, à savoir créer un marché et permettre la rémunération des contenus.	25
III.	Lignes directrices pour le modèle de résumé.	28
1.	Le modèle doit être « simple et utile » pour permettre au fournisseur d’IA d’élaborer son résumé.	28
2.	Les principaux éléments de la politique de conformité doivent figurer en amont, puisqu’ils justifient, en aval, la présence ou l’absence de certains éléments.	28
3.	S’agissant, ensuite, des informations relatives aux contenus, le degré de détail attendu est fonction du degré de fiabilité des sources.	28
4.	Le modèle de résumé doit enfin requérir en amont des informations contextuelles importantes.	29
IV.	Modèle de résumé	30
	Annexes	33
	Liste des contributeurs et des personnes auditionnées	34
	Lettre de mission	37
	Règlement (UE) 2024/1689 du Parlement européen et du Conseil du 13 juin 2024 établissant des règles harmonisées concernant l’intelligence artificielle (extraits)	39
	Loi AB 2013 (Californie): Generative artificial intelligence: training data transparency	46

Résumé décideurs

La collecte et l'exploitation des données de qualité, en particulier de données culturelles, revêtent une importance stratégique pour les fournisseurs de modèles d'intelligence artificielle (IA). Pourtant, **les données humaines s'appauvrissent** sur le web et l'entraînement d'un modèle d'IA sur des **données synthétiques** conduit à la **dégénérescence** de celui-ci. En dépit de ces constats, les données sont, paradoxalement, les seuls « intrants » de la chaîne dont la **valeur commerciale** est remise en cause.

Destiné à créer un cadre favorable à l'innovation et protecteur des droits et valeurs de l'Union européenne, le règlement relatif à l'intelligence artificielle (RIA) du 13 juin 2024 complète le paysage normatif, notamment le règlement général sur la protection des données (RGPD) du 27 avril 2016 et la directive relative au droit d'auteur et aux droits voisins dans le marché unique numérique du 17 avril 2019, ces deux derniers textes ayant été adoptés avant l'émergence de l'IA générative « grand public ».

En particulier, l'article 53 du RIA crée une **obligation de transparence** qui impose aux fournisseurs d'IA à usage général, y compris lorsque les modèles sont publiés dans le cadre d'une licence libre et ouverte, de mettre en place une **politique visant à se conformer à la législation de l'Union** en matière de droit d'auteur et de droits voisins (art 53, 1, c) et de mettre à la disposition du public « **un résumé suffisamment détaillé du contenu utilisé pour entraîner le modèle d'IA à usage général** » (art 53, 1, d). Ce résumé doit être conforme à un modèle (« *template* ») fourni par le Bureau de l'intelligence artificielle, service de la Commission européenne créé par le RIA. La première version est attendue pour janvier 2025.

La présente mission a pour objet de **préciser la portée des dispositions** de l'article 53, 1, d, et de **proposer un modèle de résumé** qui nourrira les positions de la France au niveau européen.

Le champ de cette mission flash étant **le droit d'auteur et les droits voisins**, ne sont pas expertisés la question de l'entraînement à partir de données à caractère personnel, l'articulation avec d'autres champs du droit, en particulier le droit de la concurrence, ainsi que l'enjeu de la diversité des données nécessaire pour éviter des biais et assurer le rayonnement de la culture française. Ces sujets mériteraient toutefois d'être précisément expertisés.

Selon la mission, la politique de conformité exigée à l'article 53, 1, c du RIA et la mise à la disposition du public d'un résumé suffisamment détaillé imposée à l'article 53, 1, d sont **indissociables**. La politique de conformité est le négatif du résumé détaillé : ce que le second dit en plein, la première le dit en creux. Ils composent ainsi les deux faces d'une même obligation : l'**obligation de transparence**. Le modèle de résumé doit en conséquence **intégrer les éléments pertinents de la politique de conformité**, et notamment ceux relatifs à la clause de réserve de droits (« *opt out* ») prévue par l'article 4 de la directive de 2019 relative au droit d'auteur et aux droits voisins.

La finalité du résumé est, ainsi que le précisent les considérants du RIA, d'« *aider les parties ayant des intérêts légitimes, y compris les titulaires de droits d'auteur, à exercer et à faire respecter les droits que leur confère la législation de l'Union* ». Le contenu du résumé ne doit cependant pas porter atteinte au **secret des affaires**. Le degré de détail du résumé doit donc s'apprécier au regard de cet objectif et en tenant compte de cette limite. Dans ce cadre, les exigences doivent être appréciées en contemplation les unes des autres, impliquant une **lecture finaliste et aussi globale** de l'obligation. Les dispositions européennes doivent en effet avoir un « **effet utile** », comme le rappelle régulièrement la CJUE.

Cette orientation permet d'être **exigeant sur l'identification des contenus utilisés**. Contrairement à ce que soutiennent certains fournisseurs de modèles d'IA, le résumé n'a pas à se contenter de lister les « principales » sources de données. Il est notamment indispensable de requérir une liste de noms de domaine, et même d'URLs datés. Comme l'indique le considérant 107 du RIA, le résumé doit être « **complet en termes de contenu** ».

Toutefois, les informations techniques, par nature susceptibles de porter atteinte au secret des affaires, doivent être limitées (même considérant). Il s'en suit que ce résumé public doit permettre d'**identifier l'utilisation potentielle** d'une œuvre ou d'un contenu protégé, mais **pas de détailler comment ce contenu a été utilisé**. Les informations techniques relatives à la tokenisation, au processus de filtrage, n'ont pas à figurer dans ce résumé.

Autrement dit, la liste précise d'ingrédients peut être rendue publique, mais pas la recette. S'arrêter sur le seul terme « résumé » pour réduire au minimum l'information relative aux ingrédients conduirait à ignorer l'ordre législatif. L'absence de complétude (qui justifie l'usage du terme « résumé ») vise donc la recette – les techniques –, non les ingrédients – le contenu.

Le résumé constitue alors un premier pas sur le chemin qui permettra à l'intéressé de faire respecter ses droits, mais pas le dernier. L'obligation de transparence crée un pont vers le respect des droits et la création d'un marché respectueux de la chaîne de valeur. Mais il reste des pas à franchir pour construire un écosystème éthique. Comment, concrètement, exercer et faire respecter ses droits, ce qui suppose aussi l'accès à des informations techniques protégées par le secret des affaires ? Le tout **à droit constant**, puisque la mission n'est pas prospective.

Deux voies s'ouvrent à ce stade.

La première implique un **dialogue direct** entre les titulaires de droits (ou leurs représentants) et les fournisseurs d'IA, dans lequel pourront le cas échéant être échangées des informations permettant la négociation. Il suffit pour cela d'intégrer au résumé un point de contact pour faciliter la mise en relation.

La seconde, envisagée par le RIA, passe par l'intermédiaire d'une autorité administrative (Bureau de l'IA en première intention), saisie d'une **réclamation** (sans préjudice d'une éventuelle action contentieuse), avec l'objectif d'éviter la judiciarisation des affaires (rappelons que plus de 30 actions judiciaires sont engagées aux Etats-Unis et que la GEMA en Allemagne a intenté une action en contrefaçon contre OpenAI). Elle ouvrirait aussi un **espace d'échange encadré** et permettrait peut-être une **médiation en facilitant le dialogue sur la preuve**.

En l'état actuel, il est presque impossible pour un titulaire de droits de rapporter la preuve de l'utilisation de son contenu. La complétude exigée du résumé pourrait l'y aider, mais cela dépendra beaucoup du modèle qui sera retenu par le Bureau de l'IA. En outre, cette exigence, couplée à l'impossibilité d'intégrer dans le résumé des informations relevant du secret des affaires (méthodes de filtrage...), fait que certains contenus pourtant mentionnés dans le résumé pourraient finalement ne pas avoir été utilisés pour l'entraînement du modèle. Le fournisseur aurait dans ce cadre la possibilité de prouver qu'il respecte le droit.

Fort de ces précisions, la mission propose de retenir pour le modèle de résumé une **approche par type de contenus, avec un degré de détail croissant en fonction de leur sensibilité au droit** :

- Pour les contenus libres de droits (domaine public ou utilisation expressément autorisée par le titulaire en « licence libre »), il est possible de se contenter d'informations générales. En revanche, si des identifiants sont disponibles, il est important de les mentionner ;
- Pour les autres données, il est essentiel d'exiger des précisions sur les méthodes de recueil utilisées pour s'assurer que les données ont été collectées en conformité avec le droit de

l'Union, en précisant notamment la base légale de la collecte. Pour les données moissonnées sur internet, les URLs et la date de moissonnage doivent être révélées. Les bases d'entraînement utilisées doivent être documentées. Notamment, les identifiants uniques doivent être mentionnés lorsqu'ils sont disponibles.

Le résumé doit également contenir certaines informations essentielles, notamment le point de contact ou l'existence d'accords commerciaux ou partenariaux le cas échéant.

La mission rappelle que, de manière générale, la transparence est une condition de l'effectivité des droits et que l'opacité emporte nécessairement des conséquences dysfonctionnelles sur le marché. Ici, la transparence est le préalable à l'émergence d'un marché éthique et compétitif, respectueux de la chaîne de valeur et rémunérant à ce titre les contenus sous droits.

RAPPORT

Prolégomènes

Cette « mission flash » a été réalisée en six mois, dans un calendrier contraint, notamment par celui des négociations conduites au niveau européen pour l'élaboration d'un « *template* » par le Bureau de l'intelligence artificielle (IA). Les auditions, lorsque c'était pertinent pour une meilleure compréhension des positions, ont été conduites sur la base des contributions produites par les différentes parties prenantes.

Les acteurs mis à contribution sont tant **nationaux**, qu'**européens**, voire **internationaux**, représentant des **intérêts variés** (titulaires de droits – tous secteurs –, fournisseurs d'IA, acteurs institutionnels).

Les positions, qui n'engagent que la mission, s'inscrivent dans le prolongement des réflexions, notamment sur la transparence, menées au sein de la **Commission interministérielle de l'IA**, qui a remis son rapport, *IA : notre ambition pour la France*, au Président de la République en mars 2024¹. Conformément au droit applicable, le rapport soutient aussi que l'utilisation de contenus protégés pour l'entraînement des modèles d'IA ne peut se faire que « *dans le respect des droits de propriété intellectuelle* » (v. « Recommandations clés »).

A titre liminaire, il convient en effet de rappeler que le droit de propriété intellectuelle est un **droit fondamental**, spécifiquement mentionné à l'article 17 de la Charte des droits fondamentaux de l'Union européenne. En droit interne, la propriété intellectuelle est rattachée au droit de propriété, qui fait l'objet d'une protection constitutionnelle.

Aussi, ce droit implique un **monopole**, qui se traduit par des droits exclusifs sur un objet, accordant au titulaire le **pouvoir de dire oui ou non**, c'est-à-dire d'accepter ou de refuser l'utilisation de son objet par un tiers, et le cas échéant de demander une rémunération. Lorsqu'une exception/limitation confisque cette exclusivité, elle ne peut le faire qu'avec **mesure**, au risque que le principe devienne l'exception et que l'atteinte au droit fondamental soit en conséquence **disproportionnée**.

Dans ce cadre, la construction d'un marché est essentielle. L'obligation de transparence imposée par le législateur européen a pour ambition de servir de **levier** à l'établissement de ce marché, en permettant la négociation.

¹ <https://www.elysee.fr/emmanuel-macron/2024/03/13/25-recommandations-pour-lia-en-france>

L'émergence de l'IA générative est porteuse d'innovation, sans aucun doute de progrès, mais également de risques, notamment économiques et culturels. Une étude du MIT publiée en août 2024 relève ainsi que des données protégées par le droit d'auteur obtenues sans autorisation sont fréquemment utilisées pour l'entraînement de modèles d'IA². Le sondage réalisé par le Comité consultatif sur l'IA des Nations Unies auprès de 348 experts dans 68 pays révèle encore que la violation de la propriété intellectuelle figure en bonne place des risques que fait courir l'IA et constitue une préoccupation pour plus de la moitié d'entre eux³.

Alors que les solutions technologiques et le cadre juridique s'avèrent inadaptés aux défis que pose l'IA, et notamment l'utilisation sans consentement des œuvres pour entraîner les modèles d'IA générative, l'Union européenne a adopté un règlement sur l'IA, imposant une logique de transparence sur les données d'entraînement, consistant dans la mise en place d'une politique de conformité et la mise à disposition du public d'un résumé suffisamment détaillé relatif à ces données (I).

Pour que les dispositions de ce règlement aient une portée réelle (la CJUE dirait un « effet utile »), le contenu du résumé doit lui-même prendre en compte la politique de conformité et contenir les informations nécessaires pour aider les ayants droit à exercer et faire respecter leur droits (II).

Il s'en déduit des lignes directrices (III) pour l'élaboration du modèle de résumé (IV).

I. Etat des lieux

1. La collecte et l'exploitation des données revêtent une importance majeure, mais s'effectuent dans des conditions qui ne garantissent pas le respect des valeurs et du droit de l'UE.

a. La collecte et l'exploitation de données générées par des humains sont devenues un enjeu stratégique pour les fournisseurs de modèles d'IA.

- *Pour s'entraîner, les grands modèles, par exemple linguistiques (Large Language Models ou LLM), ont besoin de grandes quantités de données.*

Les fournisseurs de modèles d'IA se situent au milieu d'une chaîne de valeur qui comprend, en amont, des équipements (processeurs graphiques – GPU – pour la puissance de calcul, serveurs pour stocker des données), des plateformes de données et de puissance de calcul type Amazon Web Services, de l'électricité, des talents humains et des données⁴, et, en aval, des entreprises,

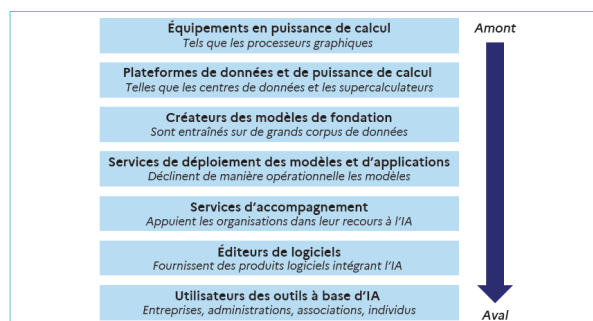
² Peter Slattery, Alexander K. Saeri, Emily A. C. Grundy, Jess Graham, Michael Noetel, Risto Uuk, James Dao, Soroush Pour, Stephen Casper et Neil Thompson, [The AI Risk Repository: A Comprehensive Meta-Review, Database, and Taxonomy of Risks From Artificial Intelligence](#), MIT, 13 août 2024, p. 40.

³ United Nation, AI advisory board, [Governing AI for Humanity](#), septembre 2024, p. 29.

V. également : Hagendorff, *Mapping the Ethics of Generative AI: A comprehensive scoping Review*, université de Stuttgart, 2024 ; Gabriel, I., Manzini, A., Keeling, G., Hendricks, L. A., Rieser, V., Iqbal, H., Tomašev, N., Ktena, I., Kenton, Z., Rodriguez, M., El-Sayed, S., Brown, S., Akbulut, C., Trask, A., Hughes, E., Stevie Bergman, A., Shelby, R., Marchal, N., Griffin, C., Mateos-García, J., Weidinger, L., Street, W., Lange, B., Ingerman, A., Lentz, A., Enger, R., Barakat, A., Krakovna, V., Siy, J. O., Kurth-Nelson, Z., McCroskery, A., Bolina, V., Law, H., Shanahan, M., Alberts, L., Balle, B., de Haas, S., Ibitoye, Y., Dafoe, A., Goldberg, B., Krier, S., Reese, A., Witherspoon, S., Hawkins, W., Rauh, M., Wallace, D., Franklin, M., Goldstein, J. A., Lehman, J., Klenk, M., Vallor, S., Biles, C., Ringel Morris, M., King, H., Agüera y Arcas, B., Isaac, W., Manyika, J., [The Ethics of Advanced AI Assistants](#), Google DeepMind, 2024.

⁴ V. Commission de l'intelligence artificielle, [IA : notre ambition pour la France](#), mars 2024, p. 21 et suivantes.

administrations, associations et individus qui retirent un avantage économique (gain de productivité, qualité...) des éléments produits par les modèles d'IA à partir des éléments d'amont.



Rapport IA : notre ambition pour la France, mars 2024

Dans cette chaîne, les données entrantes (*inputs*) constituent les ingrédients sans lesquels rien n'est possible. Ces ingrédients sont acquis (récupération de bases de données, moissonnage de données), nettoyés (filtrage et structuration) puis préparés (tokenisation et vectorisation) pour obtenir en « sortant » (*output*) un met (texte, image, musique, etc.).

Ces données, et notamment les données culturelles, constituent les seuls intrants de la chaîne dont la valeur commerciale, bien qu'évidente, est remise en cause⁵.

- *Alors que le besoin de données de qualité croît, les gisements exploitables se tarissent.*

Il est aujourd'hui scientifiquement démontré que l'entraînement à partir de **données synthétiques**, générées par des modèles d'IA, dégrade les performances du modèle et conduit, à terme, à sa **dégénérescence**⁶. Or une part croissante des données disponibles sur internet est constituée de données synthétiques, ce qui détériore la qualité des données qui y sont collectées et donc des modèles. Au-delà de la qualité même du modèle, la qualité et la spécialisation des données d'entraînement limitent également le risque d'hallucinations.

Le besoin de données n'a pourtant pas diminué : à mesure que les modèles se perfectionnent et se multiplient, il augmente. Certains spécialistes de l'IA pointent même la **saturation des performances** des grands modèles de langage alors même qu'une grande partie des données « disponibles » sur internet ont déjà été utilisées.

Alors que la proportion de données de qualité, issues d'interactions humaines, diminue, le **besoin de données** croît.

- *La constitution de jeux de données de qualité, notamment de données générées par des interactions humaines, devient cruciale.*

Les données générées par des interactions humaines constituent une ressource rare et précieuse. C'est le cas, en particulier, parmi elles, des créations humaines contemporaines, indispensables pour entraîner des modèles en phase avec leur époque. Dans sa contribution au Comité spécial des communications et du numérique de la Chambre des Lords du Parlement britannique, OpenAI relevait en ce sens que limiter l'entraînement des données aux livres et tableaux du

⁵ Adi Robertson, [Mark Zuckerberg: creators and publishers 'overestimate the value' of their work for training AI](#), revue TheVerge, 25 septembre 2024.

⁶ Iliia Shumailov, Zakhar Shmaylov, Yiren Zhao, Nicolas Papernot, Ross Anderson et Yarin Gal, [AI models collapse when trained on recursively generated data](#), revue Nature, 24 juillet 2024.

domaine public créés il y a plus d'un siècle constituerait une expérience intéressante, mais ne permettant pas de fournir des systèmes d'IA adaptés aux **besoins des citoyens d'aujourd'hui**⁷. Or ces créations contemporaines sont, par construction, les plus susceptibles d'être protégées par le droit d'auteur et/ou les droits voisins.

Ces données peuvent être collectées à partir d'internet. Il peut s'agir de jeux de données non structurées (URLs ou contenus), comme celles fournies par l'association Common Crawl, qui aspire régulièrement, depuis plusieurs années, l'ensemble des données disponibles sur internet à l'aide de robots appelés « web crawlers » ou moissonneurs du web. Ces bases de données peuvent ensuite être retravaillées pour créer des jeux de données (datasets) structurés. Par exemple, les bases d'entraînement, comme LAION ou BOOKS3, constituent une autre source d'approvisionnement pour l'entraînement des modèles d'IA. Il a cependant été dénoncé que des contenus illicites y étaient intégrés.

S'agissant enfin de la nature des données, plusieurs types sont exploités. En plus des données protégées par le droit d'auteur, on peut mentionner les données à caractère personnel des utilisateurs de réseaux sociaux, type X ou Meta (vidéos, sons, textes publiés par les utilisateurs), ou encore les instructions données à une IA générative par un utilisateur (« *prompts* ») qui peuvent elles-mêmes contenir des données à caractère personnel ou des contenus protégés par le droit d'auteur.

b. L'encadrement des conditions de récupération (par moissonnage ou autrement) et d'exploitation de ces données est insatisfaisant.

- *En amont, il n'existe pas, en l'état de l'art, de standardisation.*

De nombreuses technologies existent pour permettre aux titulaires de droits d'indiquer aux moissonneurs de données si et comment des contenus peuvent être utilisés⁸. Elles se répartissent en deux grands ensembles selon l'approche retenue pour identifier ces contenus. La première approche consiste à utiliser des identificateurs fondés sur la localisation du contenu (site web ou domaine : *location-based identifiers*). L'ensemble des contenus présents sur ce lieu virtuel sont concernés. Dans ce premier bloc figurent notamment les protocoles robots.txt, ai.txt, DeviantArt's noai meta-tags, l'utilisation d'en-têtes http, ou encore l'enregistrement d'un domaine dans un registre type « do-not-train ». La seconde approche consiste à permettre aux titulaires de droits d'indiquer comment un contenu protégé peut ou non être utilisé (*unit-based identifiers*). Cela peut passer par la création d'un identificateur permettant de lier les métadonnées à l'œuvre (voir, notamment, l'*International Standard Content Code*). Cela peut aussi consister à établir des normes de référence d'intégration des métadonnées dans les contenus numériques afin d'en retracer la provenance : c'est l'approche de la coalition pour la provenance des contenus et l'authenticité (C2PA : *Coalition for Content Provenance and Authenticity*), fondée par des entreprises comme Microsoft et Adobe. Il peut enfin s'agir d'inscrire une œuvre dans un registre (par exemple haveibeentrained.com).

On peut encore mentionner la technique d'*opt out* TDMRep, conçue par les ayants droit, à la fois *location* et *unit-based*.

Toutefois, ces technologies sont plus ou moins efficaces selon le type de contenu à identifier. L'approche par localisation de contenus est adaptée aux contenus textuels. L'approche par identification des œuvres et autres objets protégés est plus pertinente pour les contenus qui se présentent sous d'autres formes de fichiers. Il est donc difficile d'imaginer un modèle de vérification qui serait approprié pour tous les objets. Et à supposer même que l'on puisse créer

⁷ [OpenAI written evidence \(LLM0113\)](#), House of Lords communications and digital select committee inquiry: Large language models (p. 4).

⁸ Paul Keller, *Considerations for opt-out compliance policies by AI model developers*, Open_Future, 16 mai 2024.

un tel système, son efficacité réelle demeurerait douteuse, puisque les titulaires de droits ne sont pas systématiquement à l'origine de toutes les publications et mises en ligne de leurs contenus.

En outre, ces technologies sont parfois volontairement ignorées, même lorsque, à l'instar du protocole « robots.txt », elles sont largement diffusées⁹.

Dans ce contexte, et alors qu'il n'apparaît ni opérationnel au regard des évolutions rapides du secteur, ni même souhaitable d'imposer une solution unique, la Commission européenne, actant qu'il est difficile d'imaginer un système qui fonctionnerait pour tout type de contenus, explore la piste de la création d'un registre centralisé de réserve de droits, solution de type « *unit-based identifiers* », qui viendrait **en complément** des autres outils techniques, ainsi que l'a précisé Renate Nikolay, directrice générale de la DG Connect, le 9 septembre 2024¹⁰.

Selon la Commission, ce registre pourrait servir de base au futur **marché des licences** pour l'entraînement des modèles d'IA. Certains ayants droit font cependant valoir, non sans force, que l'absence de référencement dans un tel registre ne saurait valoir présomption d'utilisation libre¹¹ et, qu'en tout état de cause, un tel registre nécessiterait des moyens conséquents pour être mis en œuvre et mis à jour, tout en posant des questions redoutables en termes de responsabilité en cas d'oubli et / ou d'erreur. Des éléments complémentaires de la Commission permettront sans doute de répondre aux interrogations légitimes des ayants droit.

- *En aval, les méthodes de désapprentissage ne sont pas opérationnelles.*

Le désapprentissage machine a pour objet d'effacer une information des connaissances apprises par un modèle d'IA.

Le désapprentissage exact consiste à entraîner à nouveau un modèle à partir du même ensemble de données duquel sont retirées les données litigieuses. Le coût d'un tel réentraînement en fait une solution peu réaliste.

Le désapprentissage approximatif peut être mis en œuvre par une multitude de techniques réparties en trois groupes, selon que le désapprentissage s'effectue par une modification des données¹², du protocole d'apprentissage¹³ ou du modèle entraîné¹⁴.

Cependant, en raison de la nature probabiliste et opaque du procédé d'apprentissage, il n'existe aujourd'hui aucun indicateur fiable pour mesurer l'efficacité du désapprentissage approximatif qui, en théorie, ne doit pas dégrader la performance du modèle. En outre, ces méthodes n'apportent pas systématiquement de garanties vérifiables¹⁵.

⁹ Katie Paul, [Multiple AI companies bypassing web standard to scrape publisher sites, licensing firm says](#), Reuters, 21 juin 2024.

¹⁰ [UE La DG Connect tient à son registre de l'« opt-out » de l'IA](#), Briefing Médias, revue Contexte, 12 septembre 2024.

¹¹ Ce qui reviendrait à imposer des formalités préalables à la protection – modèle prohibé par la Convention de Berne.

¹² Hanxun Huang, Xingjun Ma, Sarah Monazam Erfani, James Bailey, Yisen Wang, [Unlearnable examples: Making personal data unexploitable](#), conference paper ICLR 2021, 13 janvier 2021 ; Ayush K Tarun, Vikram S Chundawat, Murari Mandal, Mohan Kankanhalli, [Fast yet effective machine unlearning](#), 31 mai 2023.

¹³ Lucas Bourtole, Varun Chandrasekaran, Christopher A. Choquette-Choo, Hengrui Jia, Adelin Travers, Baiwu Zhang, David Lie, Nicolas Papernot, [Machine Unlearning](#), 42^e IEEE symposium of security and privacy, 15 décembre 2020 ; Yinzhi Cao et Junfeng Yang, [Towards making systems forget with machine unlearning](#), IEEE symposium of security and privacy, 2015.

¹⁴ Vikram S Chundawat, Ayush K Tarun, Murari Mandal, Mohan Kankanhalli, [Can bad teaching induce forgetting? Unlearning in deep networks using an incompetent teacher](#), 31 mai 2023.

Aditya Golatkar, Alessandro Achille, Stefano Soatto, [External sunshine of the spotless net: selective forgetting in deep networks](#), 31 mars 2020.

¹⁵ Alexis Léautier, [Comprendre le désapprentissage machine : anatomie du poisson rouge](#), CNIL, 26 mai 2023.

c. Adoptés alors que l'utilisation massive des contenus protégés par les modèles d'IA générative n'avait pas été anticipée, le règlement sur la protection des données à caractère personnel et la directive relative au droit d'auteur dans le marché unique numérique n'étaient plus à même de garantir de façon satisfaisante le respect des droits des citoyens européens.

- *Pour les données personnelles.*

Le Règlement (UE) 2016/679 du Parlement européen et du Conseil du 27 avril 2016, relatif à la protection des personnes physiques à l'égard du traitement des données à caractère personnel et à la libre circulation de ces données, dit « règlement général sur la protection des données » (RGPD), a un champ d'application particulièrement étendu. Il encadre le traitement des données à caractère personnel, entendu en un sens large, effectué par des organisations établies sur le territoire de l'Union européenne ou par des organisations qui, quel que soit leur lieu d'établissement, exercent une activité visant à cibler ou fournir des biens et services à des résidents européens.

Chaque Etat membre dispose d'une autorité de protection des données – en France la CNIL –, qui accompagne et contrôle les acteurs. Si le règlement repose sur une **logique de conformité** (« *compliance* ») et encourage la protection des données dès la conception et par défaut (« *privacy by design* » et « *privacy by default* »), dans une démarche de **responsabilisation** des acteurs (principe de responsabilité ou « *accountability* »), des sanctions graduées peuvent être infligées en cas de non-respect des obligations qu'il prévoit.

Des litiges récents illustrent les interrogations que suscite l'utilisation des données personnelles par des fournisseurs de modèles d'IA. Après que l'autorité irlandaise pour la protection des données a introduit une procédure en justice, la plateforme X s'est engagée en septembre 2024 à ne plus utiliser les données personnelles de ses utilisateurs européens pour entraîner son programme d'intelligence artificielle. Le groupe Meta, visé par des plaintes déposées dans onze pays européens par l'association *None of Your Business* (NOYB), a annoncé en juin 2024 qu'il renonçait, pour le moment, à utiliser les *posts* Facebook et Instagram de ses utilisateurs en Europe pour entraîner ses modèles d'IA. En avril dernier, l'association NOYB a déposé une plainte contre OpenAI qui, selon elle, ne permet pas de rectifier les données personnelles inexacts que produit son service ChaptGPT.

Beaucoup s'interrogent en particulier sur la pertinence de l'**intérêt légitime** comme base légale du traitement¹⁶. Des conditions sont requises pour pouvoir bénéficier de ce cas d'ouverture à traitement, ce que rappelle d'ailleurs le Conseil européen de la protection des données dans ses lignes directrices d'octobre 2024¹⁷, dans la lignée de la jurisprudence de la CJUE et des avis du G29.

- *Pour le droit d'auteur et les droits voisins.*

La directive (UE) 2019/790 du Parlement européen et du Conseil du 17 avril 2019 sur le droit d'auteur et les droits voisins dans le marché unique numérique définit notamment les règles applicables à la fouille de textes et de données, c'est-à-dire à « *toute technique d'analyse automatisée visant à analyser des textes et des données sous une forme numérique afin d'en dégager des informations, ce qui comprend, à titre non exhaustif, des constantes, des tendances et des corrélations* » (article 2, 2).

L'article 3 prévoit une exception au monopole sur les contenus protégés, en autorisant les organismes de recherche et institutions du patrimoine à procéder à une fouille de textes et de

¹⁶ V. not. la 2^e série de fiches CNIL sur l'IA, soumises à consultation jusqu'en oct. 2024.

¹⁷ EDPB (CEPD), [Guidelines 1/2024 on processing of personal data based on Article 6\(1\) \(f\) GDPR](#), 8 oct. 2024.

données sur des œuvres ou autres objets protégés, sous deux conditions. La première, objective, tient à l'accès licite à ces objets protégés. La seconde, subjective, tient à la finalité poursuivie : la fouille doit être effectuée à des fins de recherche scientifique. L'exception est encadrée et finalisée. Elle ne saurait bénéficier à des acteurs économiques.

L'article 4 prévoit une seconde exception permettant l'utilisation de contenus protégés sans finalité assignée. Cette dérogation au monopole a un champ plus large, puisqu'elle autorise tous les acteurs à procéder à la fouille de textes et données, pour **toute finalité, y compris commerciale** donc. L'article 4 maintient la condition objective de l'**accès licite** au contenu et ajoute une clause à la main des titulaires de droits tenant à ce que l'utilisation des objets protégés « *n'ait pas été expressément réservée par [eux] de manière appropriée, notamment par des procédés lisibles par machine pour les contenus mis à la disposition du public en ligne* ». Il s'agit de la clause de réserve de droits, dite clause d'**opt out**, qui implique alors un retour au monopole (principe d'autorisation et, le cas échéant, rémunération).

Cependant, la condition d'accès licite ne paraît pas systématiquement respectée, des éléments dans la presse donnant à penser que, par exemple, des romans protégés par le droit d'auteur, ou encore des contenus de presse protégés par les droits voisins, figurent dans des bases de données ayant servi à l'entraînement de grands modèles¹⁸.

En outre, comme on l'a vu, l'effectivité et l'efficacité de l'**opt out** sont pour l'heure douteuses. Enfin, on peut s'interroger sur ce que recouvre l'expression « *lisibles par machine* »¹⁹. Les interprétations peuvent diverger, comme l'illustre l'affaire qui oppose l'association LAION e.v à l'artiste Robert Kneschke, dont a été saisie la Cour de district de Hambourg. Sous quelle forme la clause de réserve de droits doit-elle être formulée pour être considérée comme « lisible » par les robots moissonneurs de données sur internet ?²⁰ Est-ce à l'utilisateur d'utiliser un protocole répandu type robots.txt, ou bien au moissonneur de se mettre en capacité de lire tout type d'instructions, y compris celles écrites « en clair » dans le code html ? En l'absence de standardisation, et à défaut d'un vocabulaire commun permettant de définir de façon univoque les usages autorisés, on pourrait être tenté de considérer qu'imposer au fournisseur d'IA de prendre en compte la multitude de solutions existantes revient à faire peser sur lui une charge déraisonnable. Toutefois, le considérant 18 de la directive de 2019 précise que sont au nombre des procédés lisibles par machine « *les conditions générales d'utilisation d'un site internet ou d'un service* ». Cette approche ouverte donne à penser que c'est au moissonneur de se mettre en capacité de lire une information disponible, et non à l'ayant droit d'utiliser une technologie particulière qui serait imposée. Sans compter que la charge serait alors insurmontable pour les titulaires, qui ne sont pas à l'origine de toutes les mises en ligne.

On relèvera donc avec intérêt que dans l'affaire précitée, la Cour de district de Hambourg, dont la position devra certes être confirmée en appel, a également retenu qu'une réserve d'utilisation rédigée en langage naturel était « *lisible par une machine* » au sens de la directive (UE) 2019/790²¹.

¹⁸ Alex Reisner, [Revealed: The Authors Whose Pirated Books Are Powering Generative AI](#), TheAtlantic, 19 août 2023.

¹⁹ V. not. CSPLA, [Rapport Transposition des exceptions de fouille de textes et de données](#), décembre 2020.

²⁰ V. Paul Keller, *Machine readable or not? – notes on the hearing in LAION e.v. vs Kneschke*, 22 juillet 2024, Institute for Information Law.

²¹ [Cour de district de Hambourg, 27 septembre 2024](#), p 15 : « *Die Kammer neigt allerdings dazu, als „maschinenverständlich“ auch einen allein in „natürlicher Sprache“ verfassten Nutzungsvorbehalt anzusehen* ». La demande de M. Kneschke a cependant été rejetée au fond car la cour a considéré que l'association bénéficiait d'une disposition spécifique du droit allemand (transposition nationale de l'art. 3 de la dir. 2021/790) qui autorise les reproductions de texte et l'exploration de données à des fins de recherche scientifique et que M. Kneschke

Ce sujet suscite l'intérêt de la Commission européenne qui a organisé des réunions sur la clause de réserve de droits. Il figure également dans le questionnaire de la présidence hongroise du Conseil adressé aux Etats membres²².

2. Pour mettre fin à une situation préjudiciable à l'innovation et aux citoyens, l'Union européenne a adopté un règlement relatif à l'intelligence artificielle qui prévoit notamment une obligation de transparence dont la mission doit préciser la portée dans la perspective des négociations entre Etats membres.

a. Cette situation préjudicie à l'innovation et aux citoyens.

- *Elle est source d'incertitude pour les entreprises et les citoyens.*

Cette situation d'incertitude sur l'application de l'exception de fouille de textes et de données, précisément sur l'effectivité de ses conditions, est catastrophique pour les petites entreprises, mais également pour les plus grandes qui soulignent le besoin de sécurité juridique pour permettre l'éclosion du marché.

L'incertitude juridique ne peut que conduire, comme c'est déjà le cas aux Etats-Unis²³, à la multiplication des contentieux ou à des transactions défavorables au développement du marché.

Seule la fixation de conditions assurées permettra un développement serein et pérenne du marché.

- *Elle préjudicie aux citoyens, car elle favorise la course au moins-disant.*

Au cours des auditions, la mission a entendu que refuser d'exploiter des données d'origine douteuse, c'était prendre le risque de se faire distancer dans la course à l'innovation. Mais l'on ne saurait sacrifier les droits sur l'autel de l'innovation.

Attendre ne conduit qu'à dégrader un peu plus la situation car il n'existe pas, du moins en ce domaine, de main invisible conduisant à un marché autorégulé.

Enfin, ce statu quo ne durera qu'un temps et il faudra bien, à un moment, revenir à la raison. Le temps qui passe ne fera que renforcer la rupture d'égalité entre les acteurs, avec les risques concurrentiels liés.

- *Elle entrave l'innovation.*

Cette situation renforce les sociétés dominantes, en capacité d'assumer des contentieux longs et coûteux, ou de signer des accords avantageux avec les ayants droit qui pourraient même se voir contraints de négocier en exclusivité, alors qu'au contraire leur modèle économique est bien de céder leurs droits à une multitude d'acteurs.

n'établissait pas que l'association poursuivait, en outre, des fins commerciales. Sur ce point, la conformité de la décision avec le droit européen est douteuse.

²² [Hungarian Presidency policy questionnaire on the relationship between generative Artificial Intelligence and copyright and related rights](#), 27 juin 2024.

²³ Pour le droit d'auteur aux Etats-Unis, v. par exemple 9 affaires récentes recensées par Luiza Jarovsky (LinkedIn) : UMG Recordings, Capitol Records, Sony Music Entertainment, Atlantic Recording Corporation, Atlantic Records, Rhino Entertainment, The All Blacks, Warner Music International & Warner Records vs. Suno (24/06/2024) ; Andre Dubus III & Susan Orlean vs. NVIDIA (02/05/2024) ; The Intercept Media vs. OpenAI & Microsoft (28/02/2024) ; The NY Times vs. Microsoft & OpenAI (27/12/2023) ; Mike Huckabee, Relevate Group, David Kinnaman, Tsh Oxenreider, Lysa TerKeurst & John Blase vs. Meta Platforms, Microsoft, Bloomberg, EleutherAI (17/10/2023) ; Author's Guild and others vs. Open AI (19/09/2023) ; J.L., C.B., K.S., P.M., N.G., R.F., J.D. & G.R vs. Google (11/07/2023) ; Kadrey, Silverman & Golden vs. Meta (07/07/2023) ; Paul Tremblay & Mona Awad vs. Open AI (28/06/2023). – Et le recensement fait sur ce site : <https://chatgptiseatingtheworld.com/2024/08/27/master-list-of-lawsuits-v-ai-chatgpt-openai-microsoft-meta-midjourney-other-ai-cos/>

L'exploitation sans encadrement de l'ensemble des données personnelles et culturelles ne pouvant, à l'évidence, se poursuivre indéfiniment, la mise en place d'une régulation est inévitable. Retarder son émergence ne fait que hisser un peu plus haut la barrière à l'entrée pour les futurs acteurs de ce secteur, puisque dans l'intervalle, les acteurs dominants peuvent tirer profit des failles existantes pour accroître leur avance.

Cette situation illustre bien l'**opposition stérile entre innovation et régulation**. Une régulation adaptée est nécessaire pour permettre l'innovation²⁴, pour que l'innovation soit synonyme de progrès.

L'UE, qui accuse un retard dans le secteur, pourrait bien se distinguer en développant un modèle d'**IA de confiance**, respectant des critères éthiques. D'où l'importance de modifier le cadre, ce qu'a souhaité le RIA.

b. Le règlement sur l'intelligence artificielle (RIA) a pour ambition de constituer un cadre à la fois favorable à l'innovation et respectueux des valeurs de l'UE.

- *Le RIA a pour objet d'améliorer le fonctionnement du marché intérieur.*

Ce règlement²⁵ crée un cadre garantissant la libre circulation des biens et services fondés sur l'IA, favorable à l'innovation, dans le respect des valeurs de l'Union européenne²⁶.

Dans une lettre ouverte parue le 19 septembre 2024, une trentaine d'entreprises, dont Meta et Spotify, ont mis en garde contre une fragmentation et des inconsistances de la régulation de l'Union qui affecteraient la capacité de l'UE à rester compétitive et à innover²⁷.

La mission considère au contraire que le RIA complète un droit écrit avant que l'IA générative ne se diffuse à une large échelle. Il permet d'améliorer la cohérence du paysage normatif. C'est vrai en particulier de l'articulation avec le RGPD, ainsi que la CNIL l'a rappelé²⁸. C'est vrai aussi des règles applicables au droit d'auteur et aux droits voisins, dont l'effectivité est remise en cause par le manque de transparence notamment.

Bien loin d'entraver la compétitivité des entreprises européennes, le RIA devrait permettre un alignement sur le mieux-disant. Le règlement reprend en effet le **principe d'extraterritorialité** présent dans de nombreux textes européens sur le numérique et s'applique, conformément à l'article 2, éclairé par le considérant 106²⁹, à toutes les entreprises qui veulent opérer sur le marché européen. Or les entreprises souhaitent rarement s'exclure du marché européen. Et le coût élevé de l'entraînement d'un modèle de fondation dissuadera sans doute les fournisseurs d'entraîner un modèle régional « adapté » au seul marché européen.

²⁴ Anu Bradford, [The False Choice Between Digital Regulation and Innovation](#), 6 octobre 2024, Columbia university.

²⁵ Règlement (UE) 2024/1689 établissant des règles harmonisées concernant l'intelligence artificielle et modifiant les règlements (CE) 300/2008, (UE) 167/2013, (UE) 168/2013, (UE) 2018/858, (UE) 2018/1139 et (UE) 2019/2144 et les directives 2014/90/UE, (UE) 2016/797 et (UE) 2020/1828 (règlement sur l'intelligence artificielle), 13 juin 2024.

²⁶ Considérant 1.

²⁷ [Europe needs regulatory certainty on AI](#).

²⁸ CNIL, 12 juillet 2024, [Entrée en vigueur du règlement européen sur l'IA : les premières questions-réponses de la CNIL](#).

²⁹ « Tout fournisseur qui met un modèle d'IA à usage général sur le marché de l'Union devrait se conformer à cette obligation, quelle que soit la juridiction dans laquelle se déroulent les actes pertinents au titre du droit d'auteur qui sous-tendent l'entraînement de ces modèles d'IA à usage général. Cela est nécessaire pour garantir des conditions de concurrence équitables entre les fournisseurs de modèles d'IA à usage général, lorsqu'aucun fournisseur ne devrait pouvoir obtenir un avantage concurrentiel sur le marché de l'Union en appliquant des normes en matière de droit d'auteur moins élevées que celles prévues dans l'Union. »

A l'instar de ce qui s'est produit pour le RGPD notamment, on peut aussi s'attendre à un **rayonnement du RIA** au-delà de l'Union européenne, illustrant à nouveau l'« *effet Bruxelles* »³⁰. Plusieurs projets de loi prennent ainsi d'ores et déjà exemple sur le RIA. Aux Etats-Unis, un projet de loi fédéral [AI Research, Innovation and Accountability Act](#) a été soumis au Sénat. Par ailleurs, bien qu'il ait mis son veto le 29 septembre 2024 au projet de loi [SB 1047](#) de l'Etat de Californie, le gouverneur Gavin Newsom a approuvé huit textes relatifs à l'IA qui s'appliqueront dans cet Etat. L'un de ces textes, la **loi AB 2013** (mise en annexe), prévoit une **obligation de transparence** qui impose de rendre publiques des informations relatives aux données d'entraînement des modèles d'IA générative, avec notamment l'obligation de préciser si les jeux de données incluent des données protégées par le droit d'auteur. On notera aussi avec intérêt la **proposition de jugement final**, divulguée en novembre 2024 par le ministère américain de la Justice (DOJ) dans son procès contre Google, accusé de position monopolistique. Le DOJ demande à la plateforme de fournir aux éditeurs en ligne, aux sites et aux créateurs de contenus « un mécanisme facilement utilisable » afin d'exercer leur droit d'opposition (« *opt-out* ») et empêcher que leurs contenus soient utilisés pour l'entraînement³¹.

Au Canada, c'est l'[Artificial Intelligence and Data Act](#) (AIDA) qui s'inspire de la réglementation européenne de manière assumée dans le document compagnon.

- *Le RIA prévoit une obligation de transparence.*

L'accès aux données étant cruciale pour les fournisseurs de modèles, l'obligation de transparence a été l'une des questions les plus débattues, notamment lors du dernier Trilogue de décembre 2023. Elle s'applique à tous les modèles, y compris ouverts.

Garantir la transparence sur les données utilisées pour entraîner les modèles d'IA à usage général constitue l'une des conditions indispensables à l'émergence de ce marché. Cette transparence est en effet nécessaire pour assurer **une concurrence équitable** entre fournisseurs de modèles d'IA à usage général, comme l'exprime d'ailleurs le RIA³².

De l'efficacité de cette obligation de transparence dépendra aussi l'émergence d'un **marché éthique et compétitif**, respectant la chaîne de valeur et donc rémunérant tous les intrants. Il n'advient jamais qu'un modèle économique puisse se construire de manière pérenne sur la base de l'utilisation gratuite d'objets appartenant à des tiers et dans l'opacité.

Une fois la transparence acquise, le marché pourra s'établir et les modèles de rémunération se préciser³³. L'Union européenne soutient à ce titre l'émergence d'un « marché des licences ».

C'est dans ce contexte que l'article 53 du règlement soumet les fournisseurs de modèles, y compris lorsque ceux-ci sont publiés dans le cadre d'une licence libre et ouverte, à deux obligations :

- Mettre en place une politique visant à se conformer à l'acquis communautaire en matière de droit d'auteur et de droits voisins : c'est le point c du 1 de l'article 53 ;

³⁰ Anu Bradford, [The Brussels Effect](#), Columbia Law School, 2012.

³¹ US District Court for the District of Columbia, [Case No. 1:20-cv-03010-APM](#), nov. 2024, p. 12 : « *Google must provide online Publishers, websites, and content creators an easily useable mechanism to selectively opt-out of having the content of their web pages or domains used in search indexing; used to train or fine-tune AI models, or AI Products; used in retrieval-augmented generation-based tools; or displayed as AI-generated content on its SERP, and such opt-out must be applicable for Google as well as for users of the Search Index. Google must provide for an opt-out specific to itself and each index user on a user-by-user basis and must transmit all opt-outs to index users in a useable format. Google must offer content creators on Google-owned sites (all Google owned or operated properties including YouTube) the same opt-out provided to Publishers, websites, and content creators. Google must not retaliate against any Publisher, website, or content creator who opts-out pursuant to this provision* ».

³² Considérant 106.

³³ C'est l'objet de la mission économique-juridique « rémunération IA », en cours au CSPLA.

- Rédiger « un résumé suffisamment détaillé du contenu utilisé pour entraîner le modèle d'IA à usage général ». Ce résumé doit être mis « à la disposition du public » et conforme à un modèle fourni par le Bureau de l'IA³⁴. C'est le point d du même paragraphe.

Les fournisseurs de modèles d'IA pourront s'appuyer sur des codes de bonnes pratiques, dont la rédaction est encouragée et facilitée par le Bureau de l'IA³⁵.

Ces différents processus sont en cours au niveau européen.

- *Ses modalités de mise en œuvre font l'objet d'échanges nourris avec la Commission et l'AI Office.*

Une consultation multipartite sur des modèles d'IA à usage général fiables au titre de la législation sur l'IA a été conduite par le Bureau de l'IA et s'est terminée le 18 septembre 2024³⁶.

La consultation s'inscrit dans le calendrier de mise en œuvre du règlement³⁷ qui prévoit, notamment, que les codes de bonnes pratiques seront prêts le 2 mai 2025 et que l'obligation de transparence s'appliquera à compter du 2 août 2025³⁸ ou, pour les fournisseurs de modèles mis sur le marché ou en service avant cette date, le 2 août 2027³⁹. Toute nouvelle version d'un modèle devrait être incluse dans la première échéance.

Des groupes de travail ont été formés par la Commission, avec la mission d'élaborer le premier *Code of Practice* pour les IA à usage général⁴⁰.

La Commission a également annoncé le 21 novembre dernier, lors d'une réunion du groupe de travail « droit d'auteur » du code de bonnes pratiques, une première ébauche du modèle de résumé diffusé par le Bureau de l'IA pour janvier 2025.

Parallèlement, la Hongrie, qui préside le Conseil de l'UE jusqu'au 31 décembre 2024, a adressé à chaque Etat membre un questionnaire relatif, notamment, à l'intelligence artificielle. Elle a par ailleurs indiqué qu'elle porterait une attention particulière au sujet et à la préparation de la mise en œuvre du RIA aux niveaux européen et nationaux⁴¹.

c. La mission créée par la ministre de la Culture a pour objet de préciser la portée des dispositions de l'article 53, 1, d, et de proposer un modèle de résumé qui puisse être porté au nom de la France au niveau européen.

Mise en place pour une courte durée (du 15 avril 2024 au 30 novembre 2024), la mission travaille à **droit constant**. Il ne s'agit pas d'imaginer ce qu'aurait pu ou aurait dû être le RIA, et pas davantage de préconiser la mise en place de solutions nouvelles qui nécessiterait une réforme, mais d'**éclairer la portée des dispositions existantes**, dans le champ, limité, tant par son objet (le droit d'auteur et les droits voisins) que par le calendrier qui est le sien.

Ce travail nourrira les positions de la France au niveau européen, ainsi que l'a rappelé Thomas Courbe, directeur général des entreprises et représentant de la France au Comité européen de l'IA, organe consultatif créé par l'article 65 du RIA et chargé d'assurer la mise en œuvre

³⁴ Bureau établi par la décision de la Commission européenne du 24 janvier 2024 (https://eur-lex.europa.eu/legal-content/FR/TXT/PDF/?uri=OJ:C_202401459) publiée au Journal officiel de l'Union européenne du 14 février 2024.

³⁵ Article 56.

³⁶ [Législation sur l'IA : Donnez votre avis sur Trustworthy General-Purpose AI | Bâtir l'avenir numérique de l'Europe](#) ; les réponses pouvaient être soumises via un formulaire.

³⁷ [Calendrier de mise en œuvre de la loi européenne sur l'intelligence artificielle \(artificialintelligenceact.eu\)](#)

³⁸ Article 113, b.

³⁹ Article 111, 3.

⁴⁰ Un premier *draft* a été révélé le 14/11/2024.

⁴¹ [Programme of the Hungarian presidency of the council of the European Union in the second half of 2024](#), page 37.

effective de la législation sur l'IA dans l'ensemble de l'UE, notamment en coordonnant les autorités nationales⁴².

Il importe toutefois d'ajouter que, au vu de la rédaction finale de la disposition, l'obligation de transparence **s'étend bien au-delà des contenus protégés** par le droit d'auteur et les droits voisins, visés au premier chef⁴³. La formulation intègre désormais dans le champ de l'obligation notamment les données à caractère personnel, ce qui suppose de préciser l'articulation du RIA avec le RGPD – question que la CNIL a commencé à expertiser, en organisant, le 11 octobre 2023, une première consultation publique sur la constitution de bases de données d'apprentissage des systèmes d'IA qui s'est traduite par la publication de fiches pratiques⁴⁴. La CNIL a également mis en place, du 10 juin au 1^{er} octobre 2024, une nouvelle consultation publique sur de nouvelles fiches, accompagnées d'un questionnaire consacré à l'encadrement du développement des systèmes d'intelligence artificielle⁴⁵.

La mission considère qu'il s'agit ici d'un enjeu fondamental qui concentrera sans doute l'attention des autorités dans les prochains mois. Le croisement du droit des données à caractère personnel avec la propriété littéraire et artistique est sur ce point très pertinent, précisément avec le développement des IA génératives reprenant la voix ou l'image d'artistes, concentrant alors les atteintes – droit d'auteur, droits voisins, droit des données à caractère personnel, droits de la personnalité.

En outre, l'obligation de transparence devrait aussi conduire à traiter la problématique de la **représentativité des données** (biais d'entraînement pouvant, notamment, générer et amplifier des discriminations) et celle de la **diversité des expressions culturelles** ou du **rayonnement de la culture française et francophone**⁴⁶.

Enfin, l'articulation avec les dispositions du **droit de la concurrence**, et les compétences de l'Autorité de la concurrence (ADLC), constitue un autre champ qui reste à expertiser, notamment les éventuelles sanctions qui pourraient être infligées lorsque le non-respect de l'obligation de transparence porte atteinte au fonctionnement du marché (précisément au titre de l'abus de position dominante) ou encore constitue une pratique commerciale déloyale dont des concurrents pourraient se saisir⁴⁷. De manière générale, toute violation d'une obligation de conformité est susceptible de constituer une atteinte concurrentielle. Les autorités, européenne comme nationales, investissent ce champ. En France, l'ADLC recommande, dans son avis 24-A-05 du 28 juin 2024, de vérifier que le marché des données se construit en assurant un

⁴² Le Monde, [Régulation européenne de l'IA : la bataille se poursuit entre créateurs de contenu et entreprises de la tech](#), 20 juin 2024 : « Il y a un travail à faire pour détailler la mise en pratique des résumés suffisamment détaillés. C'est pour cela que le gouvernement a confié deux missions sur le droit d'auteur au Conseil supérieur de la propriété littéraire et artistique. Ces propositions nourriront nos positions au niveau européen. »

⁴³ Considérant 107 : l'élaboration du résumé s'inscrit dans un objectif de transparence sur les données d'entraînement « y compris le texte et les données protégés par la législation sur le droit d'auteur » (nous soulignons). Si les données d'entraînement se limitait aux textes et données protégés par la législation sur le droit d'auteur, cette précision serait inutile. L'ensemble des données d'entraînement sur lequel porte l'obligation de transparence est donc plus large que les seuls contenus protégés par le droit d'auteur.

⁴⁴ [Consultation publique – fiches pratiques sur la constitution de bases de données pour la conception de systèmes d'IA - Synthèse des contributions \(cnil.fr\), février 2024.](#)

⁴⁵ [Intelligence artificielle : nouvelle consultation publique sur le développement des systèmes d'IA | CNIL, juin 2024.](#)

⁴⁶ V. [Sommet de la Francophonie \(oct. 2024\) : Déclaration de Villers-Cotterêts](#), art. 20.

⁴⁷ V. notamment [Cass. Com., 27 septembre 2023, n° 21 - 21.995](#) jugeant que le fait pour une entreprise de manquer à ses obligations légales de conformité peut être constitutif d'un acte de concurrence déloyale. Et [CJUE \(Grande chambre\), 4 octobre 2024, Lindenapotheke, C. 21/23](#), qui juge que le RGPD ne s'oppose pas à ce que les Etats membres prévoient dans leur droit national la possibilité pour les concurrents de l'auteur présumé d'une violation de ce règlement de se prévaloir de cette violation devant les juridictions civiles en tant que pratique commerciale déloyale.

« *équilibre entre juste rémunération des ayants droit et accès des développeurs de modèles aux données nécessaires pour innover, en prenant en compte la diversité des cas d'usage des données* »⁴⁸.

Ces aspects excèdent le champ de la mission et ne seront donc pas évalués. Le modèle de résumé proposé ne traitera que des données culturelles et n'intégrera pas d'éléments intéressant la représentativité des données.

II. Analyse

1. L'obligation de mettre en place une politique de conformité et celle de mettre à disposition du public un résumé suffisamment détaillé participent d'un même objectif : améliorer la transparence.

a. Le RIA semble considérer qu'il s'agit de deux obligations à traiter en silo.

Telles que présentées dans le RIA, il semblerait que le modèle de résumé (article 53, 1, d) et la politique interne de respect des droits (même disposition, point c) soient deux obligations indépendantes. A traiter les deux conjointement, on pourrait craindre d'être conduit à interpréter des stipulations relatives aux droits d'auteur, et donc à dépasser le cadre du RIA.

En réalité, il n'en est rien. D'ailleurs, dans le [First Draft of the General-Purpose AI Code of Practice](#), publié par le Bureau de l'IA le 14 novembre 2024, les deux sujets sont traités dans le même point (« *Rules related to copyright* »).

b. Pour la mission, les deux obligations sont indissociables.

La mission estime que le RIA crée pour les fournisseurs d'IA une obligation nouvelle et autonome de conformité à la racine (« *compliance by design* ») laquelle n'est effective que si la question de la politique de conformité et celle du résumé suffisamment détaillé sont traitées conjointement.

Ces deux obligations participent en effet d'un **même objectif de transparence**. Cette analyse est soutenue par la formulation même du considérant 107, qui prévoit que les résumés suffisamment détaillés sont mis à la disposition du public « *afin d'accroître la transparence* ». En outre, le considérant 108 mentionne ces deux obligations comme n'en constituant qu'une seule⁴⁹. Ce lien ressort également de la portée matérielle que le texte donne aux deux alinéas c et d du 1 de l'article 53, puisqu'ils sont les seuls, au sein de ce 1, à s'imposer à l'ensemble des fournisseurs d'IA, y compris ceux produisant des modèles libres⁵⁰.

⁴⁸ [Avis 24-A-05 du 28 juin 2024](#), relatif au fonctionnement concurrentiel du secteur de l'intelligence artificielle générative, p. 95.

⁴⁹ « *En ce qui concerne l'obligation imposée aux fournisseurs de modèles d'IA à usage général de mettre en place une politique visant à respecter la législation de l'Union sur le droit d'auteur et de mettre à la disposition du public un résumé du contenu utilisé pour l'entraînement, le Bureau de l'IA devrait vérifier...* », souligné par la mission.

⁵⁰ Ce que le considérant 104 justifie ainsi : « *...étant donné que la publication de modèles d'IA à usage général sous licence libre et ouverte ne révèle pas nécessairement des informations importantes sur le jeu de données utilisé pour l'entraînement ou de réglage fin du modèle et sur la manière dont le respect de la législation sur le droit d'auteur a été assuré, l'exception prévue pour les modèles d'IA à usage général en ce qui concerne les exigences en matière de transparence ne devrait pas concerner l'obligation de produire un résumé du contenu utilisé pour l'entraînement des modèles ni l'obligation de mettre en place une politique visant à respecter la législation de l'Union sur le droit d'auteur...* ».

Certes, contrairement au résumé suffisamment détaillé, il n'est pas indiqué dans le texte du règlement que la politique de conformité doit être mise à la disposition du public. Le règlement prévoit seulement que les fournisseurs de modèles sont tenus de la « *mettre en place* »⁵¹.

Toutefois, la pertinence des informations du résumé suffisamment détaillé s'apprécie nécessairement à la lumière des mesures mises en œuvre par le fournisseur pour se conformer à ses obligations en droit d'auteur. La politique de conformité est en quelque sorte le négatif du résumé détaillé : ce que le second dit en plein, la première le dit nécessairement en creux.

En conséquence, par souci de cohérence et de bonne articulation entre le modèle de résumé et le *Code of practice*, et parce que les obligations sont complémentaires et concernent le même champ d'application ici (la propriété littéraire et artistique), la politique de conformité devrait être mentionnée dans le résumé, **au moins dans ses grandes lignes**. Le *Code of practice* sera sans doute plus complet sur les éléments exigés.

c. Le modèle de résumé doit intégrer des éléments relatifs à la conformité, et notamment au respect de la réserve de droits.

Sans exiger que la politique de conformité soit détaillée dans le résumé, la mission considère que ses principaux éléments devraient y figurer.

Le modèle de résumé devrait notamment inviter les fournisseurs à préciser quels **protocoles** sont reconnus par les moissonneurs de données qu'ils utilisent, soit directement, soit via des tiers, par exemple indiquer si le protocole robots.txt est respecté. Cela étant, ce dernier protocole ne saurait être l'unique système accepté. Comme indiqué, il n'y a pas de raison d'exclure un autre procédé « lisible par machine », d'autant que le protocole robots.txt est jugé inefficace pour certains contenus.

Pour les jeux de données obtenus à titre gratuit ou onéreux auprès d'un tiers, il convient notamment d'indiquer si des mesures sont prises afin de s'assurer que ces données ont été collectées en conformité avec le droit (garantie de l'existence d'autorisation ou de licence).

S'agissant de l'exception de fouille de textes et de données prévue par la directive 2019/790, certains titulaires de droits contestent qu'elle soit applicable à l'entraînement des modèles d'IA⁵². L'exception, conçue avant l'émergence des modèles d'IA générative, viserait exclusivement l'exploitation du contenu sémantique des œuvres, alors que l'exploitation des données par les modèles d'IA ne se limiterait pas aux contenus sémantiques, mais s'étendrait également aux contenus syntaxiques⁵³. On pourrait de même s'interroger sur la conformité de l'exception au test en trois étapes, filtre applicable à toutes les exceptions en vertu de l'article 5, 5 de la directive 2001/29 et par renvoi à l'article 7 de la directive 2019/790.

Cette position subtile implique de préciser le cadre juridique actuel, sans doute à l'issue de plusieurs années de négociation (ou de contentieux jusqu'à la CJUE) qui porteront nécessairement préjudice aux ayants droit et à l'émergence de nouvelles sociétés innovantes.

A l'inverse, dans le sens d'une applicabilité de cette exception, on peut relever que le législateur européen a explicitement renvoyé à l'exception de fouille de textes et de données dans le RIA,

⁵¹ Art. 53, 1, c.

⁵² V. not. la position du European Writers' Council, *EWC second Statement on the AI Act Proposal*, juill. 2023 : https://europeanwriterscouncil.eu/23ewc_on_aiact/ et plus largement : *Joint Statement to Ursula von der Leyen and the new elected European Parliament on the impact of AI on the European creative community*, juill. 2024 : https://europeanwriterscouncil.eu/247js_aiimpact_europeancreativecommunity/.

⁵³ Tim W. Dornis, Sebastian Stober, *Urheberrecht und Training generativer KI-Modelle - technologische und juristische Grundlagen*, 29 août 2024.

à la fois dans les considérants⁵⁴ et dans les dispositions, l'article 53, 1, c, mentionnant que la politique de conformité doit « *notamment* » viser « *à identifier et respecter (...) une réservation de droits exprimées conformément à l'article 4, paragraphe 3, de la directive (UE) 2019/790* ». La Commission européenne, assez logiquement, a également pris position dans le sens d'une application de l'exception, ainsi qu'il ressort de la réponse à une question parlementaire fournie par M. Thierry Breton en mars 2023⁵⁵.

On notera à nouveau ici avec intérêt que dans sa décision Robert Kneschke contre LAION e.v du 27 septembre 2024, la Cour du district de Hambourg applique la clause de réserve de droits de l'article 4 de la directive (UE) 2019/790.

En tout état de cause, la lettre du RIA prévoit explicitement que l'obligation de transparence s'étend à la question du respect de la clause de réserve de droits. Ne pas l'évoquer dans le modèle de résumé reviendrait donc, en l'état du droit, à rogner sur le champ de l'obligation de transparence.

2. La transparence ne consiste pas à laisser les acteurs s'autoréguler et peut aller jusqu'à exiger la production d'une liste des contenus utilisés.

a. Le résumé ne saurait se borner à lister les principales sources de données dans l'attente de la création d'un marché de la donnée.

Certains fournisseurs de modèles considèrent que l'obligation de transparence doit se limiter à lister les principales sources de données utilisées pour l'entraînement. Ils font valoir, d'une part, que le considérant 107 précise que pourraient figurer dans le résumé les « *principaux jeux ou collections de données utilisés pour entraîner le modèle, tels que les archives de données ou bases de données publiques ou privées de grande ampleur* » et, d'autre part, qu'il y a lieu de tenir compte « *de la nécessité de protéger les secrets d'affaires et les informations commerciales confidentielles* ».

Ne pourrait notamment être exigée la liste des URLs des sites moissonnés, dont la révélation porterait atteinte au secret des affaires.

Le secret des affaires

La directive (UE) 2016/943 du Parlement européen et du Conseil du 8 juin 2016 sur la protection des savoir-faire et des informations commerciales non divulgués (secrets d'affaires) contre l'obtention, l'utilisation et la divulgation illicites définit comme « secret d'affaires » des informations qui répondent à toutes les conditions suivantes : « *a) elles sont secrètes en ce sens que, dans leur globalité ou dans la configuration et l'assemblage exacts de leurs éléments, elles ne sont pas généralement connues des personnes appartenant aux milieux qui s'occupent normalement du genre d'informations en question, ou ne leur sont pas aisément accessibles, b) elles ont une valeur commerciale parce qu'elles sont secrètes, c) elles ont fait l'objet, de la part de la personne qui en a le contrôle de façon licite, de dispositions raisonnables, compte tenu des circonstances, destinées à les garder secrètes* » (article 2, 1).

⁵⁴ Considérants 105 et 106.

⁵⁵ https://www.europarl.europa.eu/doceo/document/E-9-2023-000479-ASW_EN.html

Ces trois critères cumulatifs sont repris en droit interne à l'article L. 151-1 du code de commerce.

Suivant cette logique, il suffirait de mentionner les noms des principaux ensembles de données, d'indiquer si des données publiques ont été utilisées (sans préciser lesquelles), de mentionner la nature des données (image, texte, etc.) et d'explicitier les principes guidant le traitement des données.

Ce niveau d'information ne peut cependant permettre de donner au dispositif légal adopté un « **effet utile** ». Il est en effet **insuffisant pour atteindre l'objectif fixé par le législateur** : « *aider les parties ayant des intérêts légitimes, y compris les titulaires de droits d'auteur, à exercer et à faire respecter les droits que leur confère la législation de l'Union* »⁵⁶.

En outre, le considérant 107 mentionne les principaux jeux ou collections de données à titre d'**exemples**, et non de manière limitative.

Ce considérant précise même, au contraire, que si le secret des affaires peut limiter le degré de détail **technique** fourni par le résumé, ce « *résumé devrait être généralement **complet en termes de contenu*** » (souligné par la mission).

Enfin, il faut rappeler que l'invocation du **secret des affaires** a bien entendu des **limites**. En droit interne, l'article L. 151-7 du code de commerce dispose que le secret des affaires ne peut être opposé aux autorités juridictionnelles et administratives agissant, notamment, dans l'exercice de leurs pouvoirs d'enquête, de contrôle, d'autorisation ou de sanction. Et l'on notera avec intérêt que, dans l'affaire Dun & Bradstreet Austria GmbH C-203/22 relative au traitement de données personnelles par une IA, ayant conduit à refuser la conclusion ou la prolongation d'un contrat de téléphonie mobile au motif que la personne ne présentait pas une solvabilité financière suffisante, l'avocat général Jean Richard de la Tour a considéré, le 12 septembre 2024⁵⁷, que le secret des affaires ne pouvait conduire à écarter le droit qu'un individu tire du RGPD de comprendre comment une décision qui l'affecte est prise. Cette position paraît transposable aux droits qu'une personne tient des dispositions de droit d'auteur issues des textes européens. Le secret des affaires ne peut conduire, en vidant de toute substance le résumé suffisamment détaillé, à écarter le droit qu'un titulaire de droits tire du RIA à disposer d'éléments pouvant l'aider « *à exercer et à faire respecter les droits que leur confère la législation de l'Union* »⁵⁸. Enfin, la directive relative au secret des affaires envisage même l'hypothèse d'une règle de l'Union qui exigerait la révélation d'informations au public, y compris des secrets des affaires, pour des motifs d'intérêt public⁵⁹.

Fournir une liste d'URLs, aussi longue soit-elle, n'apparaît donc pas contraire aux provisions du RIA si cela est nécessaire pour atteindre l'objectif visé par le législateur européen.

Le secret des affaires est par ailleurs difficile à soutenir lorsqu'on utilise le Common Crawl.

Enfin, si le respect du secret des affaires s'impose pour la mise à disposition publique du résumé, il n'en est pas de même lorsque des discussions bilatérales sont engagées (des accords de confidentialité sont d'ailleurs fréquemment signés en d'autres domaines ; il suffirait donc de

⁵⁶ Considérant 107.

⁵⁷

<https://curia.europa.eu/juris/document/document.jsf?text=&docid=290022&pageIndex=0&doclang=FR&mode=req&dir=&occ=first&part=1&cid=2434261>

⁵⁸ Considérant 107.

⁵⁹ Article 1^{er}, 2 : « *La présente directive ne porte pas atteinte à / (...) / b) l'application de règle de l'Union (...) exigeant des détenteurs de secrets d'affaires qu'ils révèlent, pour des motifs d'intérêt public, des informations, y compris des secrets d'affaires, au public/ (...)* ».

faire de même) et encore moins lorsque la demande émane d'une autorité administrative ou judiciaire.

b. Le texte n'exclut pas de lister les contenus protégés utilisés pour l'entraînement des modèles.

Le considérant 108 exclut explicitement que la vérification, par le Bureau de l'IA, du respect, par le fournisseur, de ses obligations dans le domaine des droits d'auteur et droits voisins s'effectue « œuvre par œuvre », qu'il s'agisse de la mise en place de la politique visant à respecter la législation de l'Union sur le droit d'auteur ou de la mise à disposition du public du résumé relatif aux données d'entraînement⁶⁰.

Certaines parties prenantes déduisent de cette précision relative au travail de vérification du Bureau de l'IA que la granularité du résumé ne peut descendre au niveau de l'œuvre, ou plutôt du contenu protégé.

Mais il ressort seulement de ces considérants, qui, on l'a vu, mentionnent un résumé « *généralement complet en termes de contenu* », que si le résumé comprend une liste des contenus utilisés, le Bureau de l'IA n'est pas tenu de vérifier que cette liste est exhaustive, et pas davantage que l'usage qui a été fait de ces contenus est licite. Il s'agit d'une sorte de **contrôle de la « recevabilité »**, vérifiant le respect de la formalité, et **non d'un examen substantiel**, sur le fond. Dans cette phase initiale de contrôle du respect de la conformité, une évaluation « œuvre par œuvre » est exclue, comme l'indique le considérant 108. Mais il est possible qu'**un examen sur le fond puisse intervenir dans un second temps, notamment en cas de réclamation**⁶¹.

Cela est cohérent avec le positionnement en surplomb que le règlement confère au Bureau de l'IA, ainsi qu'avec les moyens qui y sont attachés, lesquels ne permettent pas d'envisager ce type de vérification exhaustive en première intention.

c. La portée normative du résumé doit être proportionnée à l'objectif poursuivi : aider les intéressés à faire valoir leurs droits.

L'article 53, 1, d, a, nécessairement, un « effet utile ». L'obligation de transparence n'est ainsi pas une simple obligation formelle dont on pourrait s'affranchir en renseignant un long formulaire administratif sollicitant des informations sans portée. Sa mise en œuvre doit avoir un sens, ainsi que le rappelaient des organisations de créateurs et titulaires de droits dans une lettre du 29 octobre 2024 adressée aux députés européens⁶².

Pour saisir la portée normative de cette obligation, il faut considérer l'objectif poursuivi. Le législateur l'a indiqué explicitement : le résumé suffisamment détaillé a pour objet d'aider les titulaires de droits d'auteur « *à exercer et à faire respecter les droits que leur confère la législation de l'Union* »⁶³.

Il serait vain – et inefficace – de s'attacher au seul terme de « résumé », sauf à ignorer la volonté législative. Celui-ci doit être compris en contemplation des autres exigences – en réseau. Il ne serait pas audible de requérir un résumé qui ne remplirait pas ses objectifs et n'aurait donc aucune portée normative.

⁶⁰ Considérant 108.

⁶¹ V. *infra* point d : « L'effort doit être poursuivi pour donner à la transparence les conséquences attendues, à savoir créer un marché et permettre la rémunération des contenus. »

⁶² [Joint letter of creators and rightholders organisations](#), 29 octobre 2024, Bruxelles.

⁶³ Considérant 107.

Ainsi, une **lecture finaliste et globale de l'article 53, 1, d** donne toute sa portée à une expression qui semble être, de prime abord, un oxymore si l'on se concentre sur les seuls termes « résumé » et « détaillé ». **Le résumé est suffisamment détaillé s'il permet d'atteindre cet objectif.** Autrement dit, le degré de détail s'apprécie au regard de l'objectif, avec une limite, celle du secret des affaires.

Pour atteindre l'objectif, il faut et il suffit de mettre les titulaires de droits en capacité de déterminer si leurs œuvres et objets protégés *ont pu être utilisés*. Peu importe si des milliards de lignes doivent être renseignées. Ce n'est pas techniquement impossible pour des acteurs du numérique habitués à manipuler des données massives et les titulaires de droits (à travers parfois leurs représentants) savent de plus en plus gérer de tels volumes.

Mais il n'est pas nécessaire, pour atteindre cet objectif, de détailler, dans un résumé public, *comment* cette donnée a été utilisée.

En particulier, exiger de rendre publiques des informations relatives au processus de filtrage des données ou au processus de tokenisation, serait contraire au secret des affaires.

L'utilisation du terme « résumé » vise donc cette absence de complétude liée à la non-révélation des seules informations techniques.

On peut d'ailleurs relever que le RIA précise que le résumé porte sur le « contenu » utilisé pour l'entraînement et non sur les « données » utilisées. Le premier terme désigne de façon plus englobante les sources, alors que les « données » visent une manière plus structurée (organisée, filtrée) de représenter l'information. L'emploi du terme « contenu » va donc dans le sens de l'approche qu'il est proposé de retenir.

Il s'agit là d'un premier pas sur le chemin qui permettra à l'intéressé de faire respecter ses droits, mais pas du dernier. Comment, dans ces conditions, franchir les pas suivants ? Car, sauf à vouloir vider cette obligation de toute substance, il faut bien qu'une suite puisse être donnée.

Comment, concrètement, exercer et faire respecter ses droits ? Le RIA prévoit-il une forme de « droit de suite » permettant d'obtenir des informations complémentaires en cas de besoin ?

Ces questions ne sont pas sans lien avec la mission puisque préciser l'utilisation que l'on peut faire du résumé rétroagit sur son contenu.

d. L'effort doit être poursuivi pour donner à la transparence les conséquences attendues, à savoir créer un marché et permettre la rémunération des contenus.

Il est certain qu'aujourd'hui, la situation est insatisfaisante. D'un côté, la mission a eu accès à des éléments tangibles qui montrent que lorsque le titulaire de droits suspecte une utilisation non autorisée d'œuvres protégées et demande des informations complémentaires, le fournisseur de modèle exige qu'il lui précise le nom des contenus (ou identifiants) et les sources à partir desquelles ce contenu aurait été récupéré, ce qui revient à faire peser sur le titulaire de droits une charge de la preuve insurmontable. D'un autre côté, il peut s'avérer complexe pour les fournisseurs de modèle de répondre efficacement à un afflux constant de demandes plus ou moins précises que pourraient leur adresser individuellement des millions de personnes.

Un cadre de discussion doit donc être proposé.

Le RIA esquisse seulement la suite du processus. Mais ce n'est pas surprenant. Ce règlement n'est pas un texte relatif au droit d'auteur. Il cantonne le rôle de contrôle de la Commission à

vérifier le respect des seules obligations qu'il a créées et le Bureau de l'IA n'a pas vocation à repérer et sanctionner de potentielles violations du droit d'auteur⁶⁴.

Une méconnaissance de l'obligation de transparence n'est pas équivalente à une violation du droit d'auteur. Et honorer l'obligation de transparence ne garantit pas que le droit d'auteur a été systématiquement respecté. C'est même parce que l'obligation de transparence est respectée qu'un titulaire de droits est mis en capacité d'identifier une potentielle violation de ses droits. Mais la procédure qui s'engage à partir de là relève plus de l'« *enforcement* » que du droit substantiel.

En résumé, l'obligation de transparence crée **un pont vers le droit d'auteur et les droits voisins**, mais le RIA n'offre pas de procédure spécifique à la matière.

Toutefois, sauf à vouloir judiciaire systématiquement la question, en requérant un droit d'information prévu par la [directive 2004/48/EC du 29 avril 2004](#) relative au respect des droits de propriété intellectuelle⁶⁵ ou, sur le terrain du droit commun, des mesures d'instruction pour établir la preuve (CPC, art. 145), toutes les parties ont intérêt à ce qu'une procédure postérieure permette d'organiser un **échange d'informations dans des conditions satisfaisantes**, notamment au regard du secret des affaires.

Deux voies, non exclusives l'une de l'autre, sont envisageables à **droit constant** (pour rappel, la lettre de mission demande d'expertiser la portée de l'obligation de transparence et de lister les informations nécessaires, sans requérir de modifications du droit applicable), en sollicitant une simple interprétation du RIA.

La première voie est celle de l'**échange direct entre les titulaires de droits (ou leurs représentants) et les fournisseurs d'IA**. De manière évidente, un tel échange entre acteurs de bonne volonté pourrait permettre de régler certaines situations. Dans ce cadre, il est impératif que le modèle de résumé désigne un point de contact unique pour permettre la communication et le traitement de plaintes directes éventuelles. Dans cette phase unissant des professionnels, des précisions sur l'usage effectif des contenus protégés doivent pouvoir être apportées sans opposer le secret des affaires, préservé grâce à la signature d'accords de confidentialité (pratique fréquente lors des négociations en d'autres domaines). Il est en effet essentiel, pour évaluer de manière transparente une atteinte éventuelle, de pouvoir accéder à des informations plus poussées. Il en est de même de la rémunération⁶⁶. Cette première voie n'exige pas de modification du droit positif. Il suffit d'imposer dans le résumé la mention du point de contact.

La deuxième voie est prévue de manière générale par le RIA. Même si le sujet spécifique du droit d'auteur et des droits voisins n'a sans doute été anticipé, les textes n'en sont pas moins applicables. A ce titre, la Commission dispose des pouvoirs de « surveiller et contrôler » le respect des dispositions du chapitre V (dont les dispositions concernées – chapitre relatif aux modèles d'IA à usage général) et l'exécution de ces tâches est confiée au **Bureau de l'IA** (art. 88, 1), lequel dispose pour cela de « tous les pouvoirs d'une autorité de surveillance du marché » (RIA, art. 75, 1) au sens du [règlement \(UE\) 2019/1020](#) du 20 juin 2019. Notamment, en application de l'article 88, 2 du RIA, le Bureau est susceptible d'être à même de demander de la documentation et des informations, comme prévu à l'article 91. En outre, l'article 85 du RIA consacre le droit d'introduire une **réclamation auprès de l'autorité de surveillance du**

⁶⁴ V. en particulier l'article 88 ainsi que les considérants 108 et 161.

⁶⁵ V. l'article 8 relatif au droit d'information, transposé par la loi n° 2007-1544 du 29 octobre 2007. – CPI, art. L. 331-1-2.

⁶⁶ Cette nécessité d'une transparence poussée pour l'évaluation des utilisations et de la rémunération est identique en matière de droits voisins des éditeurs de presse (CPI, art. L. 218-4, al. 3) et c'est bien le manque de transparence qui suscite la multiplication des actions judiciaires qu'on connaît en ce moment.

marché chargée de la vérification de la conformité⁶⁷. Pour rappel, la conformité exigée est double : l'article 53, 1, c prescrit la mise en place de mesures visant à respecter le droit d'auteur et les droits voisins pendant que l'article 53, 1, d implique de déclarer les sources collectées pour l'entraînement du modèle. Disposant de pouvoirs d'enquête et d'exécution larges (art. 14 du règlement 2019/1020), l'autorité serait donc chargée de cette vérification et son statut la mettrait à l'abri d'une contestation sur le terrain du secret des affaires.

Pour autant, si l'autorité se voit attribuer une **capacité juridique d'action** par le RIA, il est utile à ce stade de s'interroger sur sa **capacité opérationnelle** à traiter les réclamations susceptibles d'advenir des titulaires de droits de 27 Etats membres⁶⁸, sauf à imaginer une augmentation substantielle des moyens mis à disposition de ce nouvel organisme. Considérant le nombre de plaintes potentielles et la complexité des sujets, les délais de traitement pourraient ne pas être satisfaisants. Aussi, pour alléger la charge du Bureau, une délégation de cette mission (ou d'une partie de cette mission) à une autorité nationale, selon une procédure qui resterait à préciser, pourrait être imaginée. Cette proposition resterait à expertiser à droit constant.

A tout le moins, que la procédure soit européenne (Bureau de l'IA) ou déléguée à une autorité nationale⁶⁹, en cas d'information jugée insuffisante, inexacte ou incomplète sur la politique interne de conformité comme sur le résumé, le fournisseur devra répondre aux injonctions de l'autorité, notamment en rapportant la preuve que le contenu objet de la réclamation n'a pas été utilisé. En effet, à ce stade (réclamation devant une autorité), on doit considérer que **le contrôle alors effectué par l'autorité n'est plus formel, mais substantiel**, en cohérence avec l'étendue des pouvoirs confiés à une autorité de surveillance du marché.

Cette analyse est favorable tant aux titulaires de droits qu'aux fournisseurs d'IA. Pour les premiers, l'établissement de **la preuve de l'utilisation de leurs contenus** est facilité. Pour les seconds, cette procédure est aussi le moyen de compléter les éléments divulgués **sans risquer de compromettre des informations concurrentielles**. En effet, puisque, comme vu plus haut, le RIA requiert une complétude du résumé s'agissant du contenu et que la communication des règles techniques, notamment de filtrage, pourrait porter atteinte au secret des affaires, cette procédure permet, en révélant les résultats dans un cadre sécurisé, de confirmer ou d'infirmer l'utilisation d'une donnée protégée.

Cette procédure administrative offrira aussi aux acteurs de bonne volonté un cadre souple de **médiation**, destiné à faciliter la résolution d'un différend. Elle ne saurait cependant être obligatoire et se ferait donc **sans préjudice d'un recours judiciaire**.

En cas d'action judiciaire, le juge pourrait par ailleurs **tirer toutes les conséquences des constatations effectuées par l'autorité (violation des dispositions du RIA)**. En effet, si la révélation de la méthode globale est susceptible de porter atteinte au secret des affaires, la confirmation ou l'infirmerie de l'utilisation d'un contenu n'est pas à mettre sur le même plan.

En définitive, par cette procédure de réclamation – telle qu'interprétée –, le RIA permet une sorte de « **droit de suite** », à la main des titulaires de droits comme des fournisseurs d'IA.

En définitive, il s'agit, au stade du résumé public, d'identifier les sources collectées pour l'entraînement – les ingrédients. Mais la recette de cuisine avec les instructions de préparation

⁶⁷ La mission n'ignore pas que l'objet de la réclamation ne pourra porter que sur les mesures du RIA et non sur la démonstration d'une violation du droit d'auteur ou des droits voisins. La question de la territorialité des directives de droit d'auteur reste un sujet délicat à expertiser et à articuler avec l'extra-territorialité du RIA.

⁶⁸ Pour rappel, les réclamations ne peuvent porter que sur une violation de la conformité exigée.

⁶⁹ La compétence territoriale devrait ici dépendre de l'auteur de la réclamation.

(méthodes de filtrage, processus de tokenisation et de vectorisation, *etc.*) relève du secret des affaires et n'a pas à figurer dans un résumé public. Cette recette doit cependant être accessible, en cas de réclamation ou d'action judiciaire, pour rendre au droit son **effectivité**.

La mission a déduit de ces éléments plusieurs lignes directrices pour l'élaboration du modèle de résumé.

III. Lignes directrices pour le modèle de résumé.

1. Le modèle doit être « simple et utile » pour permettre au fournisseur d'IA d'élaborer son résumé.

Comme indiqué, le résumé lui-même doit être **complet en termes de contenu**.

Le modèle, quant à lui, doit rester **simple et utile**⁷⁰. Il doit guider efficacement les fournisseurs de modèles d'IA.

Ces deux directives ne sont pas contradictoires : la complétude n'est pas le symptôme de la complexité mais le signe de l'efficacité. L'utilité renforce l'idée d'une **lecture finaliste** de la disposition.

2. Les principaux éléments de la politique de conformité doivent figurer en amont, puisqu'ils justifient, en aval, la présence ou l'absence de certains éléments.

L'entraînement d'un modèle d'IA ne conduisant pas, logiquement, à liciter l'illicite, des éléments relatifs à la politique de conformité devront être exigés pour les données collectées directement comme pour celles en provenance de tiers, qu'il s'agisse de bases de données type Common Crawl ou d'informations fournies par les utilisateurs (*prompts*).

3. S'agissant, ensuite, des informations relatives aux contenus, le degré de détail attendu est fonction du degré de fiabilité des sources.

L'accès à des contenus libres de droit, c'est-à-dire relevant du domaine public ou dont l'utilisation est expressément autorisée par leur titulaire (« licence libre »), étant, par construction, licite, il n'y a pas à exiger une granularité fine d'information. En revanche, les informations disponibles (comme les identifiants) doivent être mentionnés pour permettre une vérification aux titulaires. La mission met également en garde sur l'absence d'harmonisation à l'échelle mondiale de la durée de protection et sur la nécessité de vérifier la réalité de l'absence de protection et l'étendue de l'autorisation éventuelle.

Il en va de même pour les contenus relevant d'arrangements contractuels. Exiger d'un fournisseur de modèle d'IA qu'il précise avec qui il a signé de tels contrats, voire qu'il livre des informations relatives au contenu de ces contrats, contreviendrait au secret des affaires. Il peut certes paraître paradoxal que des entreprises à la fois indiquent que la liste des contrats conclus relève du secret des affaires et ne peut être rendue publique et, en même temps, révèlent par voie de presse la signature d'accords ou déclarent publiquement (ou par réponse à des courriers de titulaires) qu'ils n'ont pas besoin de signer de tels accords. Mais ces déclarations font elles-mêmes partie de la stratégie de l'entreprise. Bien que, dans une situation idéale, une telle liste

⁷⁰ Considérant 107.

des contrats passés contribuerait à plus de transparence, la mission considère que le respect du secret des affaires fait obstacle à ce qu'elle soit rendue publique. En revanche, il ne paraît pas hors de portée d'exiger des fournisseurs qu'ils précisent s'il existe ou non de tels accords.

S'agissant des autres contenus, notamment les contenus publics qui ne sont pas libres de droit ou les contenus qui ont fait l'objet d'une autorisation d'utilisation (par ex. bases de données mises à disposition), des informations plus détaillées seront nécessaires puisqu'il s'agit des gisements de données les plus susceptibles d'intégrer des contenus piratés. Certes, les informations sur le nom des contenus ou des titulaires de droits ne sont pas disponibles et il ne peut être exigé de mentionner ce qui ne peut l'être. En revanche, les métadonnées associées et les identifiants doivent être intégrés.

Il est aussi indispensable de fournir une liste d'URLs précisant notamment les dates de moissonnage : sans cela, les ayants droit seront dans l'impossibilité d'identifier l'utilisation potentielle de leurs œuvres à des fins d'entraînement, et l'objectif fixé par le législateur ne pourra pas être atteint. Il n'existe pas d'obstacle technique dirimant : les entreprises ayant créé des modèles comprenant plusieurs milliards de paramètres ont la capacité de lister des milliards d'URLs, et les gestionnaires de droits savent désormais de plus en plus manipuler de telles quantités de données. Juridiquement, la circonstance qu'une URL révèle des informations sur le contenu consulté ne suffit pas, puisqu'il s'agit de données non filtrées, à caractériser une atteinte au secret des affaires.

Si un jeu de données contient à la fois des contenus libres de droit et ces « autres contenus », ou si un fournisseur de modèle est dans l'incapacité de faire la part entre les deux, il conviendra d'assimiler l'ensemble des données aux « autres contenus » soumis à l'exigence la plus haute de transparence. Cela ne pourra qu'encourager les entreprises, en amont, à consolider leur politique de conformité afin d'être mis en capacité de distinguer les données libres de droit des autres données.

De manière générale, le fournisseur d'IA doit s'assurer du respect des droits d'auteur et voisins au titre du RIA, en mettant en place une politique interne. Notamment, lorsqu'un fournisseur contracte avec un tiers pour utiliser un jeu de données, il doit s'assurer que le droit d'auteur et les droits voisins ont été respectés (respect de l'*opt out*, existence de licences en amont...). Dans un sens comparable, le bénéfice de l'exception « fouille de textes et de données » implique d'avoir un « accès licite » aux contenus.

Aussi, le fournisseur doit partager dans le résumé la **méthodologie** lui permettant de révéler cette **double injonction de licéité**.

4. Le modèle de résumé doit enfin requérir en amont des informations contextuelles importantes.

Pour les raisons déjà indiquées, le résumé doit indiquer le point de contact chez le fournisseur d'IA émetteur.

Il doit encore mentionner si le modèle est créé *ex nihilo* ou à partir d'un autre modèle. Et s'il réalise une simple mise à jour⁷¹.

Si le contenu des conventions intègre le secret des affaires, comme déjà évoqué, il serait utile d'informer *a minima* de l'existence de tels accords.

⁷¹ Considérant 109.

IV. Modèle de résumé

Point de contact pour toute demande ou procédure.

S'agit-il de la mise à jour d'un résumé après évolution du modèle ? o / n

Si oui, fournir le résumé du modèle initial.

L'entraînement du modèle repose-t-il en tout ou partie sur un modèle d'IA déjà existant ? o / n

Si oui, lien vers le résumé de ce modèle.

Certains des contenus utilisés font-ils l'objet d'accords ? o / n

Si oui, dans quel secteur ?

Si le partenaire est un organisme public, le mentionner.

Contenus libres de droit*					
Moissonnés depuis internet				Récupérés auprès de tiers	
Jeu de données	Méthode de moissonnage	Nom de domaine moissonné	Taille et type de données (image, son, multimodal, etc.)	Jeu de données	Le cas échéant, lien vers le jeu de données
<i>(liste jeux de données)</i>	<i>(description) (1)</i>	<i>(pour chaque jeu de données, liste des domaines moissonnés)</i>	<i>(pour chaque jeu de données, taille et type de données)</i>	<i>(liste des jeux de données)</i>	<i>(pour chaque jeu de données, lien si applicable)</i>

** uniquement pour les jeux de données constitués exclusivement de données libres de droit, c'est-à-dire relevant du domaine public à l'issue de la protection au titre du droit d'auteur et des droits voisins ou contenus dont l'utilisation a été expressément autorisée par l'ayant droit en licence libre. La mission met en garde sur le périmètre de la licence libre qui peut néanmoins interdire ce type d'usage. A défaut ou en cas de doute sur le licence, considérer que le jeu de données relève de la catégorie « Autres contenus ».*

Autres contenus (non libres de droits)											
Moissonnés directement ou par un tiers mandaté depuis internet					Jeux de données récupérés auprès de tiers			Prompts		Données synthétiques générées à partir de données humaines	
Jeu de données	Méthode de moissonnage	Méthodologie pour assurer la conformité avec le droit de l'Union, (notamment clause de réserve de droits, exclusion des sites pirates...)	URLs moissonnées + date de moissonnage	Taille et type de données (image, son, multimodal, etc.)	Jeu de données	Méthodologie pour assurer la conformité avec le droit de l'Union (ex. : clause de réserve de droits...)	Description	Description	Méthodologie pour assurer la conformité avec le droit de l'Union	Jeu de données	Méthodologie pour assurer la conformité avec le droit de l'Union
<i>(liste jeux de données datés)</i>	<i>(description) (1)</i>	<i>(description) (2)</i>	<i>(pour chaque jeu de données, liste des URLs avec date de collecte)</i>	<i>(pour chaque jeu de données, taille et type de données...)</i> (3)	<i>(liste des jeux de données datés)</i>	<i>(description, notamment fondement légal de la collecte de données) (4)</i>	(5)	<i>(préciser si les prompts utilisateurs peuvent être enregistrés et utilisés pour l'entraînement du modèle)</i>	<i>(description) (6)</i>	<i>(liste des jeux de données datés, modèles utilisés)</i>	<i>(description, notamment méthodes permettant d'atténuer le risque de contrefaçon...)</i>

A titre non exhaustif, type de questions qui pourraient être intégrées au modèle de résumé par le Bureau de l'IA :

- (1) Identification des moissonneurs utilisés ? Base légale du moissonnage ? Dates de moissonnage ?
- (2) Fondement légal de la collecte ? Méthode utilisée pour identifier les métadonnées pertinentes ? Méthodes pour éviter la suppression de ces métadonnées au cours de l'entraînement ou la génération de contenu ? Standards et protocoles reconnus ? Méthodes mises en œuvre pour respecter les instructions lisibles par la machine ? Autres méthodes mises en œuvre pour garantir un traitement licite ? Méthodes mises en œuvre pour assurer la mise en

conformité a posteriori si du contenu protégé a été utilisé sans autorisation ? Certification ou contrôle externe mobilisés pour assurer la licéité de l'opération ? Méthodes pour neutraliser les données dont les droits ont été réservés ? Identification des standards reconnus (type ISBN, etc.) ?

- (3) Taille des données d'entraînement ? Taille de chaque type de donnée ?
- (4) Par ex., garantie du fournisseur de données ?
- (5) Identification et taille du jeu de données ? Type de données ? Liste des types d'identifiants uniques (ISBN, DOI, ISAN, etc.) inclus ? Lien vers le jeu de données si applicable ou, à défaut, liste des contenus (par les identifiants disponibles) du jeu de données à date avec suffisamment de détail relatif au contenu pour faciliter aux ayants droit l'exercice de leur droit ? Pour chaque type d'identifiant, pourcentage d'item avec cet identifiant ?
- (6) Méthode pour apprécier si le *prompt* reproduit des données protégées ? Méthode pour neutraliser le cas échéant ces données ?

Annexes

Liste des contributeurs et des personnes auditionnées

Considérant le court délai d'action (et les congés estivaux), la mission a organisé une séance ouverte à tous (y compris les autres ministères ou encore les acteurs de la société civile extérieurs au CSPLA) dès le début des travaux et a sollicité prioritairement des contributions écrites. Des auditions et échanges oraux à partir des contributions ont également eu lieu.

La mission remercie les acteurs qui ont travaillé dans l'urgence, dans une période où les sollicitations étaient nombreuses.

Autorité de la concurrence (ADLC)

Société civile pour l'administration des droits des artistes et musiciens interprètes (ADAMI)

Aday

Alliance de la presse d'information générale (APIG)

Allonia

Autorité de régulation de la communication audiovisuelle et numérique (ARCOM)

Association des services internet communautaires (ASIC)

Bibliothèque nationale de France (BNF)

Botscorner

Canal plus

Center of the Picture Industry (CEPIC)

Centre français d'exploitation du droit de copie (CFC)

Centre national du cinéma et de l'image animée (CNC)

Chambre syndicale de l'édition musicale

Commission nationale de l'informatique et des libertés (CNIL)

CMI France

Eurocinema

European Authors' societies (GESAC)

European Magazine Media Association

European Newspaper Publishers' Association

European Publishers Council

Federation of european publishers (FEP)

Fireflies.ai

Gaumont

Google

GESTE

GFII

Institut national de l'audiovisuel (INA)

International Federation of the Phonographic Industry (IFPI)

La Sofia

Les Voix (association)

Ligue des auteurs professionnels

Linkup

Microsoft

Mistral AI

Motion Picture Association EMEA

News Media Europe

Numeum

Open Future Foundation

Panodyyssey

Ministère de la culture / Ministère de l'économie, des finances et de l'industrie – Pôle d'expertise de la Régulation Numérique

Ministère de l'économie, des finances et de l'industrie – Direction générale des entreprises

Paramount

Photoroom

PopScreen Games

Procirep

RELX

Société civile des Auteurs Réalisateur Producteurs (ARP)

Société des auteurs et compositeurs dramatiques (SACD)

Syndicat des éditeurs de la presse magazine (SEPM)

Trust My content

Société des auteurs, compositeurs et éditeurs de musique (SACEM)

SAIF

Société civile des auteurs multimédia (SCAM)

Société civile des producteurs phonographiques (SCPP)

Société des auteurs dans les arts graphiques et plastiques (ADAGP)

Société des Gens de Lettres (SGDL)

Syndicat des catalogues de films de patrimoine (SCFP)

Syndicat des Editeurs de la Presse Magazine (SEPM)

Syndicat national des auteurs et compositeurs (SNAC)

Syndicat national de l'édition (SNE)

Syndicat national de l'édition phonographique (SNEP)

Syndicat des producteurs indépendants (SPI)

SPEDIDAM

STM

The Walt Disney Company

Union des producteurs phonographiques français indépendants (UPFI)

Union Nationale des Syndicats d'Artistes Musiciens

Vivendi

Lettre de mission

Conseil supérieur de la propriété littéraire et artistique



MINISTÈRE DE LA CULTURE

Liberté
Égalité
Fraternité

Madame Alexandra BENSAMOUN
Professeure des universités

Paris, le 12 avril 2024

OBJET : Mission relative à la mise en œuvre du règlement européen établissant des règles harmonisées sur l'intelligence artificielle

Madame,

Le projet de règlement européen établissant des règles harmonisées sur l'intelligence artificielle (RIA) intègre en son article 53 une obligation pour les fournisseurs de modèles d'IA à usage général de prendre des mesures visant à respecter le droit d'auteur et, en particulier, le cadre posé par la directive du 17 avril 2019 sur le droit d'auteur et les droits voisins dans le marché unique numérique (DAMUN). Parmi ces mesures, les fournisseurs devront élaborer et rendre disponible publiquement un « *résumé suffisamment détaillé* » des données ayant servi à l'entraînement de leur modèle.

Cette transparence sur les sources ayant permis l'entraînement des systèmes d'IA en amont s'avère primordiale pour permettre aux titulaires de droits d'auteur et de droits voisins de vérifier que les conditions d'accès licite et d'utilisation de leurs œuvres et prestations – et notamment leur opposition éventuelle à toute fouille de données (« *opt out* ») – ont été respectées.

La portée de cette obligation de transparence prévoit toutefois un certain nombre de limitations, prévues par le projet de règlement, dont la mise en œuvre mérite d'être précisée. C'est notamment le cas du périmètre des fournisseurs concernés par cette obligation, du niveau de précision des informations à fournir, de l'impact des secrets industriels et commerciaux sur la divulgation des informations ou encore de la forme de la divulgation ainsi imposée.

Afin de faciliter la mise en œuvre de cette obligation de transparence, le projet de règlement confie au Bureau européen de l'intelligence artificielle, créé par une décision de la Commission européenne du 24 janvier 2024, le soin d'élaborer un modèle de résumé simple et efficace des données d'entraînement utilisées par les IA. Pour l'exécution de cette tâche,

le bureau consultera les parties prenantes, parmi lesquelles des experts de la communauté scientifique et de l'éducation, des citoyens, des organisations de la société civile et des partenaires sociaux.

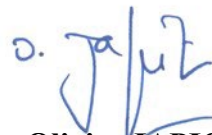
Le récent rapport de la Commission de l'intelligence artificielle, constituée par le Gouvernement en 2023 (*IA : notre ambition pour la France*), préconise, quant à lui, de « *mettre en œuvre et évaluer les obligations de transparence prévues par le règlement européen sur l'IA en encourageant le développement de standards et d'une infrastructure adaptée* ».

La ministre de la Culture souhaite que, dans le prolongement de vos précédents travaux (rapport sur les enjeux juridiques et économiques de l'intelligence artificielle dans les secteurs de la création culturelle, de janvier 2020, et rapport sur les exceptions de fouille de textes et de données – « *text and data mining* » – de décembre 2020), le CSPLA lance une nouvelle mission ayant pour objet, d'une part, d'expertiser la portée de l'obligation de transparence prévue par le règlement européen, compte tenu des interrogations mentionnées plus haut, et, d'autre part, d'établir la liste des informations qui vous paraissent devoir nécessairement être communiquées, selon les secteurs culturels concernés, pour permettre aux auteurs et aux titulaires de droits voisins d'exercer leurs droits.

Je vous confie cette mission, pour laquelle vous serez assistée d'un rapporteur. Vous pourrez également vous appuyer sur les services du ministère de la culture, en particulier le secrétariat général (service des affaires juridiques et internationales). Vous procéderez aux auditions des membres du CSPLA ainsi que des entités et personnalités dont vous jugerez les contributions utiles, en particulier les services du ministère de l'économie, des finances et de la souveraineté industrielle et numérique. Vous pourrez associer à vos travaux le professeur Frédéric Pascal, membre du CSPLA.

Il serait souhaitable que l'état d'avancement de vos travaux puisse être présenté au cours de la prochaine séance plénière du CSPLA qui se tiendra au début de l'été, et que vous puissiez présenter votre rapport au mois de décembre prochain, après avoir échangé avec les membres du CSPLA intéressés.

Je vous remercie d'avoir accepté cette mission et vous prie de croire, Madame, à l'expression de mes sentiments les meilleurs.



Olivier JAPIOT

Président du CSPLA

Règlement (UE) 2024/1689 du Parlement européen et du Conseil du 13 juin 2024 établissant des règles harmonisées concernant l'intelligence artificielle (extraits)

Considérants

(1) L'objectif du présent règlement est d'améliorer le fonctionnement du marché intérieur en établissant un cadre juridique uniforme, en particulier pour le développement, la mise sur le marché, la mise en service et l'utilisation de systèmes d'intelligence artificielle (ci-après dénommés «systèmes d'IA») dans l'Union, dans le respect des valeurs de l'Union, de promouvoir l'adoption de l'intelligence artificielle (IA) axée sur l'humain et digne de confiance tout en garantissant un niveau élevé de protection de la santé, de la sécurité et des droits fondamentaux consacrés dans la Charte des droits fondamentaux de l'Union européenne (ci-après dénommée «Charte»), y compris la démocratie, l'état de droit et la protection de l'environnement, de protéger contre les effets néfastes des systèmes d'IA dans l'Union, et de soutenir l'innovation. Le présent règlement garantit la libre circulation transfrontière des biens et services fondés sur l'IA, empêchant ainsi les États membres d'imposer des restrictions au développement, à la commercialisation et à l'utilisation de systèmes d'IA, sauf autorisation expresse du présent règlement.

(...)

(104) Les fournisseurs de modèles d'IA à usage général qui sont publiés sous licence libre et ouverte et dont les paramètres, y compris les poids, les informations sur l'architecture des modèles et les informations sur l'utilisation des modèles, sont rendus publics devraient faire l'objet d'exceptions en ce qui concerne les exigences en matière de transparence imposées pour les modèles d'IA à usage général, à moins que les modèles ne puissent être considérés comme présentant un risque systémique, auquel cas le fait que les modèles soient transparents et accompagnés d'une licence ouverte ne devrait pas être considéré comme une raison suffisante pour exclure le respect des obligations prévues par le présent règlement. En tout état de cause, étant donné que la publication de modèles d'IA à usage général sous licence libre et ouverte ne révèle pas nécessairement des informations importantes sur le jeu de données utilisé pour l'entraînement ou de réglage fin du modèle et sur la manière dont le respect de la législation sur le droit d'auteur a été assuré, l'exception prévue pour les modèles d'IA à usage général en ce qui concerne les exigences en matière de transparence ne devrait pas concerner l'obligation de produire un résumé du contenu utilisé pour l'entraînement des modèles ni l'obligation de mettre en place une politique visant à respecter la législation de l'Union sur le droit d'auteur, en particulier pour identifier et respecter la réservation de droits au titre de l'article 4, paragraphe 3, de la directive (UE) 2019/790 du Parlement européen et du Conseil [\(40\)](#).

(105) Les modèles d'IA à usage général, en particulier les grands modèles d'IA génératifs, capables de générer du texte, des images et d'autres contenus, présentent des possibilités d'innovation uniques mais aussi des défis pour les artistes, les auteurs et les autres créateurs, et la manière dont leur contenu créatif est créé, distribué, utilisé et consommé. Le développement et l'entraînement de ces modèles requièrent un accès à de grandes quantités de texte, d'images, de vidéos et d'autres données. Les techniques de fouille de textes et de données peuvent être largement utilisées dans ce contexte pour extraire et analyser ces contenus, qui peuvent être protégés par le droit d'auteur et les droits voisins. Toute utilisation d'un contenu protégé par le droit d'auteur nécessite l'autorisation du

titulaire de droits concerné, à moins que des exceptions et limitations pertinentes en matière de droit d'auteur ne s'appliquent. La directive (UE) 2019/790 a introduit des exceptions et des limitations autorisant les reproductions et extractions d'œuvres ou d'autres objets protégés aux fins de la fouille de textes et de données, sous certaines conditions. En vertu de ces règles, les titulaires de droits peuvent choisir de réserver leurs droits sur leurs œuvres ou autres objets protégés afin d'empêcher la fouille de textes et de données, à moins que celle-ci ne soit effectuée à des fins de recherche scientifique. Lorsque les droits d'exclusion ont été expressément réservés de manière appropriée, les fournisseurs de modèles d'IA à usage général doivent obtenir une autorisation des titulaires de droits s'ils souhaitent procéder à une fouille de textes et de données sur ces œuvres.

(106) Les fournisseurs qui mettent des modèles d'IA à usage général sur le marché de l'Union devraient veiller au respect des obligations pertinentes prévues par le présent règlement. À cette fin, les fournisseurs de modèles d'IA à usage général devraient mettre en place une politique visant à respecter la législation de l'Union sur le droit d'auteur et les droits voisins, en particulier pour identifier et respecter la réservation de droits exprimées par les titulaires de droits conformément à l'article 4, paragraphe 3, de la directive (UE) 2019/790. Tout fournisseur qui met un modèle d'IA à usage général sur le marché de l'Union devrait se conformer à cette obligation, quelle que soit la juridiction dans laquelle se déroulent les actes pertinents au titre du droit d'auteur qui sous-tendent l'entraînement de ces modèles d'IA à usage général. Cela est nécessaire pour garantir des conditions de concurrence équitables entre les fournisseurs de modèles d'IA à usage général, lorsqu'aucun fournisseur ne devrait pouvoir obtenir un avantage concurrentiel sur le marché de l'Union en appliquant des normes en matière de droit d'auteur moins élevées que celles prévues dans l'Union.

(107) Afin d'accroître la transparence concernant les données utilisées dans le cadre de l'entraînement préalable et de l'entraînement des modèles d'IA à usage général, y compris le texte et les données protégés par la législation sur le droit d'auteur, il convient que les fournisseurs de ces modèles élaborent et mettent à la disposition du public un résumé suffisamment détaillé du contenu utilisé pour entraîner les modèles d'IA à usage général. Tout en tenant dûment compte de la nécessité de protéger les secrets d'affaires et les informations commerciales confidentielles, ce résumé devrait être généralement complet en termes de contenu plutôt que détaillé sur le plan technique afin d'aider les parties ayant des intérêts légitimes, y compris les titulaires de droits d'auteur, à exercer et à faire respecter les droits que leur confère la législation de l'Union, par exemple en énumérant les principaux jeux ou collections de données utilisés pour entraîner le modèle, tels que les archives de données ou bases de données publiques ou privées de grande ampleur, et en fournissant un texte explicatif sur les autres sources de données utilisées. Il convient que le Bureau de l'IA fournisse un modèle de résumé, qui devrait être simple et utile et permettre au fournisseur de fournir le résumé requis sous forme descriptive.

(108) En ce qui concerne l'obligation imposée aux fournisseurs de modèles d'IA à usage général de mettre en place une politique visant à respecter la législation de l'Union sur le droit d'auteur et de mettre à la disposition du public un résumé du contenu utilisé pour l'entraînement, le Bureau de l'IA devrait vérifier si le fournisseur a rempli cette obligation sans vérifier ou évaluer œuvre par œuvre les données d'entraînement en ce qui concerne le respect du droit d'auteur. Le présent règlement n'affecte pas l'application des règles en matière de droit d'auteur prévues par la législation de l'Union.

(...)

(161) Il est nécessaire de clarifier les responsabilités et les compétences au niveau de l'Union et au niveau national en ce qui concerne les systèmes d'IA qui reposent sur des modèles d'IA à usage général. Afin d'éviter les chevauchements de compétences, lorsqu'un système est fondé sur un modèle d'IA à usage général et que le modèle et le système sont fournis par le même fournisseur, la surveillance devrait avoir lieu au niveau de l'Union par l'intermédiaire du Bureau de l'IA, qui devrait disposer à cette fin des pouvoirs d'une autorité de surveillance du marché au sens du règlement (UE) 2019/1020. Dans tous les autres cas, les autorités nationales de surveillance du marché demeurent chargées de la surveillance des systèmes d'IA. Toutefois, pour les systèmes d'IA à usage général qui peuvent être utilisés directement par les déployeurs pour au moins un usage classé comme étant à haut risque, les autorités de surveillance du marché devraient coopérer avec le Bureau de l'IA pour mener les évaluations de la conformité, et informer le Comité IA et les autres autorités de surveillance du marché en conséquence. En outre, toute autorité de surveillance du marché devrait être en mesure de solliciter l'assistance du Bureau de l'IA lorsqu'elle n'est pas en mesure de conclure une enquête sur un système d'IA à haut risque parce qu'elle ne peut accéder à certaines informations liées au modèle d'IA à usage général sur lequel repose ce système. Dans de tels cas, la procédure relative à l'assistance mutuelle pour les cas transfrontières prévue au chapitre VI du règlement (UE) 2019/1020 devrait s'appliquer mutatis mutandis.

(...)

(156) Afin de garantir un contrôle approprié et efficace du respect des exigences et obligations énoncées par le présent règlement, qui fait partie de la législation d'harmonisation de l'Union, le système de surveillance du marché et de mise en conformité des produits établi par le règlement (UE) 2019/1020 devrait s'appliquer dans son intégralité. Les autorités de surveillance du marché désignées en vertu du présent règlement devraient disposer de tous les pouvoirs d'exécution prévus par le présent règlement et par le règlement (UE) 2019/1020, et elles devraient exercer leurs pouvoirs et s'acquitter de leurs tâches de manière indépendante, impartiale et sans parti pris. Bien que la majorité des systèmes d'IA ne fassent pas l'objet d'exigences et obligations particulières au titre du présent règlement, les autorités de surveillance du marché peuvent prendre des mesures à l'égard de tous les systèmes d'IA lorsqu'ils présentent un risque conformément au présent règlement. En raison de la nature spécifique des institutions, agences et organes de l'Union relevant du champ d'application du présent règlement, il convient de désigner le Contrôleur européen de la protection des données comme autorité compétente pour la surveillance du marché en ce qui les concerne. Cela devrait être sans préjudice de la désignation des autorités nationales compétentes par les États membres. Les activités de surveillance du marché ne devraient pas affecter la capacité des entités surveillées à s'acquitter de leurs tâches de manière indépendante, lorsque cette indépendance constitue une exigence du droit de l'Union.

(157) Le présent règlement est sans préjudice des compétences, des tâches, des pouvoirs et de l'indépendance des autorités ou organismes publics nationaux compétents qui contrôlent l'application du droit de l'Union en matière de protection des droits fondamentaux, y compris les organismes chargés des questions d'égalité et les autorités de protection des données. Lorsque leur mandat l'exige, ces autorités ou organismes publics nationaux devraient également avoir accès à toute documentation créée en vertu du présent règlement. Une procédure de sauvegarde spécifique devrait être mise en place pour garantir une application adéquate et en temps utile opposable aux systèmes d'IA présentant un risque pour la santé, la sécurité et les droits fondamentaux. La procédure

applicable à ces systèmes d'IA présentant un risque devrait être appliquée aux systèmes d'IA à haut risque présentant un risque, aux systèmes interdits qui ont été mis sur le marché, mis en service ou utilisés en violation des interdictions concernant des pratiques définies par le présent règlement, et aux systèmes d'IA qui ont été mis à disposition en violation des exigences de transparence énoncées dans le présent règlement et qui présentent un risque.

(...)

Articles

(...)

SECTION 2

Obligations incombant aux fournisseurs de modèles d'IA à usage général

Article 53

Obligations incombant aux fournisseurs de modèles d'IA à usage général

1. Les fournisseurs de modèles d'IA à usage général:

- a) élaborent et tiennent à jour la documentation technique du modèle, y compris son processus d'entraînement et d'essai et les résultats de son évaluation, qui contient, au minimum, les informations énoncées à l'annexe XI aux fins de la fournir, sur demande, au Bureau de l'IA et aux autorités nationales compétentes;
- b) élaborent, tiennent à jour et mettent à disposition des informations et de la documentation à l'intention des fournisseurs de systèmes d'IA qui envisagent d'intégrer le modèle d'IA à usage général dans leurs systèmes d'IA. Sans préjudice de la nécessité d'observer et de protéger les droits de propriété intellectuelle et les informations confidentielles de nature commerciale ou les secrets d'affaires conformément au droit de l'Union et au droit national, ces informations et cette documentation:
 - i) permettent aux fournisseurs de systèmes d'IA d'avoir une bonne compréhension des capacités et des limites du modèle d'IA à usage général et de se conformer aux obligations qui leur incombent en vertu du présent règlement; et
 - ii) contiennent, au minimum, les éléments énoncés à l'annexe XII;
- c) mettent en place une politique visant à se conformer au droit de l'Union en matière de droit d'auteur et droits voisins, et notamment à identifier et à respecter, y compris au moyen de technologies de pointe, une réservation de droits exprimée conformément à l'article 4, paragraphe 3, de la directive (UE) 2019/790;
- d) élaborent et mettent à la disposition du public un résumé suffisamment détaillé du contenu utilisé pour entraîner le modèle d'IA à usage général, conformément à un modèle fourni par le Bureau de l'IA.

2. Les obligations énoncées au paragraphe 1, points a) et b), ne s'appliquent pas aux fournisseurs de modèles d'IA qui sont publiés dans le cadre d'une licence libre et ouverte permettant de consulter, d'utiliser, de modifier et de distribuer le modèle, et dont les paramètres, y compris les poids, les informations sur l'architecture du modèle et les informations sur

l'utilisation du modèle, sont rendus publics. Cette exception ne s'applique pas aux modèles d'IA à usage général présentant un risque systémique.

3. Les fournisseurs de modèles d'IA à usage général coopèrent, en tant que de besoin, avec la Commission et les autorités nationales compétentes dans l'exercice de leurs compétences et pouvoirs en vertu du présent règlement.

4. Les fournisseurs de modèles d'IA à usage général peuvent s'appuyer sur des codes de bonne pratique au sens de l'article 56 pour démontrer qu'ils respectent les obligations énoncées au paragraphe 1 du présent article, jusqu'à la publication d'une norme harmonisée. Le respect des normes européennes harmonisées confère au fournisseur une présomption de conformité dans la mesure où lesdites normes couvrent ces obligations. Les fournisseurs de modèles d'IA à usage général qui n'adhèrent pas à un code de bonnes pratiques approuvé ou ne respectent pas une norme européenne harmonisée démontrent qu'ils disposent d'autres moyens appropriés de mise en conformité et les soumettent à l'appréciation de la Commission.

5. Afin de faciliter le respect de l'annexe XI, et notamment du point 2, points d) et e), la Commission est habilitée à adopter des actes délégués conformément à l'article 97 pour préciser les méthodes de mesure et de calcul en vue de permettre l'élaboration d'une documentation comparable et vérifiable.

6. La Commission est habilitée à adopter des actes délégués conformément à l'article 97, paragraphe 2, pour modifier les annexes XI et XII à la lumière des évolutions technologiques.

7. Toute information ou documentation obtenue en vertu du présent article, y compris les secrets d'affaires, est traitée conformément aux obligations de confidentialité énoncées à l'article 78.

(...)

CHAPITRE IX

SURVEILLANCE APRÈS COMMERCIALISATION, PARTAGE D'INFORMATIONS ET SURVEILLANCE DU MARCHÉ

(...)

SECTION 4

Voies de recours

Article 85

Droit d'introduire une réclamation auprès d'une autorité de surveillance du marché

Sans préjudice d'autres recours administratifs ou judiciaires, toute personne physique ou morale ayant des motifs de considérer qu'il y a eu violation des dispositions du présent règlement peut déposer des réclamations auprès de l'autorité de surveillance du marché concernée.

Conformément au règlement (UE) 2019/1020, ces réclamations sont prises en compte aux fins de l'exercice des activités de surveillance du marché, et sont traitées conformément aux procédures spécifiques établies en conséquence par les autorités de surveillance du marché.

SECTION 5

Surveillance, enquêtes, contrôle de l'application et contrôle en ce qui concerne les fournisseurs de modèles d'IA à usage général

Article 88

Contrôle de l'exécution des obligations incombant aux fournisseurs de modèles d'IA à usage général

1. La Commission dispose de pouvoirs exclusifs pour surveiller et contrôler le respect du chapitre V, en tenant compte des garanties procédurales prévues à l'article 94. La Commission confie l'exécution de ces tâches au Bureau de l'IA, sans préjudice des pouvoirs d'organisation dont elle dispose ainsi que de la répartition des compétences entre les États membres et l'Union fondée sur les traités.

2. Sans préjudice de l'article 75, paragraphe 3, les autorités de surveillance du marché peuvent demander à la Commission d'exercer les pouvoirs prévus dans la présente section, lorsque cela est nécessaire et proportionné pour contribuer à l'accomplissement des tâches qui leur incombent en vertu du présent règlement.

(...)

Article 91

Pouvoir de demander de la documentation et des informations

1. La Commission peut demander au fournisseur du modèle d'IA à usage général concerné de fournir la documentation établie par le fournisseur conformément aux articles 53 et 55, ou toute information supplémentaire nécessaire pour évaluer la conformité du fournisseur avec le présent règlement.

2. Avant d'envoyer la demande d'informations, le Bureau de l'IA peut entamer un dialogue structuré avec le fournisseur du modèle d'IA à usage général.

3. Sur demande dûment motivée du groupe scientifique, la Commission peut adresser une demande d'informations au fournisseur d'un modèle d'IA à usage général, lorsque l'accès à ces informations est nécessaire et proportionné pour l'accomplissement des tâches du groupe scientifique au titre de l'article 68, paragraphe 2.

4. La demande d'informations mentionne la base juridique et l'objet de la demande, précise quelles informations sont requises, fixe un délai dans lequel les informations doivent être fournies, et indique les amendes prévues à l'article 101 en cas de fourniture d'informations inexactes, incomplètes ou trompeuses.

5. Le fournisseur du modèle d'IA à usage général concerné, ou son représentant, fournit les informations demandées. Dans le cas de personnes morales, d'entreprises ou de sociétés, ou lorsque le fournisseur n'a pas de personnalité juridique, les personnes autorisées à les représenter en vertu de la loi ou de leurs statuts fournissent les informations demandées pour le compte du fournisseur du modèle d'IA à usage général concerné. Les avocats dûment habilités à agir peuvent fournir des informations pour le compte de leurs clients. Les clients demeurent néanmoins pleinement responsables si les informations fournies sont incomplètes, inexacts ou trompeuses.

Loi AB 2013 (Californie): Generative artificial intelligence: training data transparency

THE PEOPLE OF THE STATE OF CALIFORNIA DO ENACT AS FOLLOWS:

SECTION 1.

Title 15.2 (commencing with Section 3110) is added to Part 4 of Division 3 of the Civil Code, to read:

TITLE 15.2. Artificial Intelligence Training Data Transparency

3110.

For purposes of this title, the following definitions shall apply:

- (a) “Artificial intelligence” means an engineered or machine-based system that varies in its level of autonomy and that can, for explicit or implicit objectives, infer from the input it receives how to generate outputs that can influence physical or virtual environments.
- (b) “Developer” means a person, partnership, state or local government agency, or corporation that designs, codes, produces, or substantially modifies an artificial intelligence system or service for use by members of the public. For purposes of this subdivision, “members of the public” does not include an affiliate as defined in subparagraph (A) of paragraph (1) of subdivision (c) of Section 1799.1a, or a hospital’s medical staff member.
- (c) “Generative artificial intelligence” means artificial intelligence that can generate derived synthetic content, such as text, images, video, and audio, that emulates the structure and characteristics of the artificial intelligence’s training data.
- (d) “Substantially modifies” or “substantial modification” means a new version, new release, or other update to a generative artificial intelligence system or service that materially changes its functionality or performance, including the results of retraining or fine tuning.
- (e) “Synthetic data generation” means a process in which seed data are used to create artificial data that have some of the statistical characteristics of the seed data.
- (f) “Train a generative artificial intelligence system or service” includes testing, validating, or fine tuning by the developer of the artificial intelligence system or service.

3111.

On or before January 1, 2026, and before each time thereafter that a generative artificial intelligence system or service, or a substantial modification to a generative artificial intelligence system or service, released on or after January 1, 2022, is made publicly available to Californians for use, regardless of whether the terms of that use include compensation, the developer of the system or service shall post on the developer’s internet website documentation regarding the data used by the developer to train the generative artificial intelligence system or service, including, but not be limited to, all of the following:

- (a) A high-level summary of the datasets used in the development of the generative artificial intelligence system or service, including, but not limited to:
 - (1) The sources or owners of the datasets.
 - (2) A description of how the datasets further the intended purpose of the artificial intelligence system or service.
 - (3) The number of data points included in the datasets, which may be in general ranges, and with estimated figures for dynamic datasets.
 - (4) A description of the types of data points within the datasets. For purposes of this paragraph, the following definitions apply:
 - (A) As applied to datasets that include labels, “types of data points” means the types of labels used.
 - (B) As applied to datasets without labeling, “types of data points” refers to the general characteristics.
 - (5) Whether the datasets include any data protected by copyright, trademark, or patent, or whether the datasets are entirely in the public domain.
 - (6) Whether the datasets were purchased or licensed by the developer.
 - (7) Whether the datasets include personal information, as defined in subdivision (v) of Section 1798.140.
 - (8) Whether the datasets include aggregate consumer information, as defined in subdivision (b) of Section 1798.140.
 - (9) Whether there was any cleaning, processing, or other modification to the datasets by the developer, including the intended purpose of those efforts in relation to the artificial intelligence system or service.

- (10) The time period during which the data in the datasets were collected, including a notice if the data collection is ongoing.
- (11) The dates the datasets were first used during the development of the artificial intelligence system or service.
- (12) Whether the generative artificial intelligence system or service used or continuously uses synthetic data generation in its development. A developer may include a description of the functional need or desired purpose of the synthetic data in relation to the intended purpose of the system or service.
- (b) A developer shall not be required to post documentation regarding the data used to train a generative artificial intelligence system or service for any of the following:
- (1) A generative artificial intelligence system or service whose sole purpose is to help ensure security and integrity. For purposes of this paragraph, “security and integrity” has the same meaning as defined in subdivision (ac) of Section 1798.140, except as applied to any developer or user and not limited to businesses, as defined in subdivision (d) of that section.
 - (2) A generative artificial intelligence system or service whose sole purpose is the operation of aircraft in the national airspace.
 - (3) A generative artificial intelligence system or service developed for national security, military, or defense purposes that is made available only to a federal entity.