

IA : Informer les personnes concernées

07 février 2025

Les organismes qui traitent des données personnelles pour développer des modèles ou des systèmes d'IA doivent informer les personnes concernées. La CNIL précise les obligations en la matière.

Assurer la transparence des traitements

Le principe de transparence oblige les organismes qui traitent des données personnelles à informer les personnes concernées afin qu'elles comprennent les usages qui seront faits de leurs données (pourquoi, comment, de quelle manière) et soient en mesure d'exercer leurs droits (droits d'opposition, d'accès, de rectification, etc.).

Ce principe s'applique à tout traitement de données personnelles, que les données soient :

- **directement recueillies auprès des personnes concernées** (aussi appelées *first party data*) : par exemple, dans le cadre d'un contrat de prestation avec des acteurs volontaires pour constituer des données d'entraînement, dans le cadre de la fourniture d'un service, dans le cadre d'une relation entre un citoyen et une administration, etc. ;
- **ou indirectement collectées** (aussi appelées *third party data*) : par exemple, lorsque les données sont collectées sur Internet via le téléchargement de fichiers, le recours à des outils de moissonnage de données (web scraping) ou l'utilisation d'interfaces de programmation applicatives (API) mises à disposition de réutilisateurs par les plateformes en ligne, obtention d'informations auprès de partenaires institutionnels ou commerciaux comme des courtiers de données (data brokers), réutilisation d'une base de données déjà constituée, etc. Cela inclut également les données générées par le responsable du traitement lui-même (CJUE, arrêt du 28 novembre 2024, affaire [C-169/23](#)).

À retenir : Lorsque l'organisme, responsable du traitement, n'a pas directement collecté les données personnelles auprès des personnes concernées, il peut être dispensé de l'obligation de les informer individuellement si cette information est impossible en pratique ou exigerait des efforts disproportionnés. Une information générale (par exemple, sur son site internet) devra toutefois être fournie, et contenir l'ensemble des éléments prévus par l'article 14 du RGPD et détaillés ci-dessous.

Quand fournir l'information ?

Lorsque le responsable du traitement collecte lui-même les données d'apprentissage directement auprès des personnes concernées, il doit informer les personnes au moment de cette collecte. Lorsqu'il a l'intention d'effectuer un traitement ultérieur pour une finalité autre que celle pour laquelle les données ont été collectées (sous réserve de la compatibilité de cette réutilisation, voir la fiche pratique dédiée), il doit en informer les personnes concernées (article 14.4 du RGPD).

Par exemple : une plateforme numérique souhaite réutiliser les données de ses utilisateurs pour entraîner un modèle d'IA. Elle doit les informer préalablement.

En cas de collecte indirecte, l'organisme doit informer les personnes concernées dès que possible, et au plus tard lors de la première prise de contact avec les intéressés ou lors de la première communication des données à un autre destinataire le cas échéant. Dans tous les cas, l'organisme doit informer les personnes concernées dans un délai ne dépassant pas un mois après la date à laquelle il a récupéré leurs données.

À titre de bonne pratique, lorsque les données présentent une sensibilité particulière pour les personnes, la CNIL invite les organismes à respecter un délai raisonnable entre le moment où les personnes sont informées que leurs données sont contenues dans une base de données d'apprentissage et l'entraînement d'un modèle sur cette base (par lui-même ou suite à la diffusion du jeu de données). Cette bonne pratique permettra aux personnes concernées de pouvoir exercer leurs droits pendant ce délai compte tenu des difficultés techniques à exercer ces droits sur le modèle lui-même et des risques que cela engendre (en particulier en fonction de la nature des données mémorisées).

Comment fournir l'information ?

Garantir l'accessibilité de l'information

Les personnes concernées ne doivent pas rencontrer de difficultés dans l'accès à l'information comme dans sa compréhension.

Les mentions d'information doivent être distinguées des autres informations sans lien avec la protection des données (CGU, mentions légales, etc.). À cet égard, alors que les mentions d'information publiées sur les sites web des responsables de traitement peuvent porter sur de nombreux traitements et concerner différentes catégories de personnes (par exemple les utilisateurs du site web en question, les personnes concernées par la phase de développement des systèmes d'IA, les personnes concernées par leur déploiement, etc.), il est recommandé de bien clarifier quelle partie des mentions d'information s'appliquent à quelles catégories de personnes (par exemple en distinguant clairement l'information portant sur les traitements de développement de l'information concernant les autres traitements).

Concrètement, il existe plusieurs moyens pour la fournir :

- en cas de fourniture d'une information individuelle, elle peut figurer sur le formulaire en ligne utilisé par le diffuseur pour collecter des données, être mentionnée dans les courriels ou courriers adressés par un réutilisateur des données lors de son premier contact avec les personnes concernées ou encore être délivrée via un message vocal pré-enregistré, etc.
- en cas de fourniture d'une information générale (c'est-à-dire dans les cas détaillés ci-dessous), elle peut

par exemple prendre la forme de mentions d'information publiées sur un site web librement accessible ou sur un panneau d'affichage.

Garantir l'intelligibilité de l'information

Le RGPD prévoit que l'information doit être fournie de façon concise, transparente, compréhensible et aisément accessible, en des termes clairs et simples. La complexité des systèmes d'intelligence artificielle ne doit pas empêcher la bonne compréhension de l'information par les personnes concernées.

À cet égard, il est recommandé que les responsables du traitement définissent clairement les principales conséquences du traitement : autrement dit, quel sera réellement l'effet du traitement spécifique.

L'information pourrait ainsi détailler, par exemple au moyen de schémas, la manière dont les données sont utilisées lors de l'apprentissage, le fonctionnement du système d'IA développé, ainsi que la distinction qui doit être faite entre la base de données d'apprentissage, le modèle d'IA et les sorties du modèle.

Point d'attention : s'il est envisageable de fournir ces informations au sein de modèles de documentation existants (comme les cartes de données, de modèles ou de systèmes d'IA), elles doivent être facilement accessibles, claires et compréhensibles pour les personnes concernées. Cela implique qu'elles ressortent clairement de ces documentations.

Par exemple : la description de la constitution d'un jeu de données dans une publication scientifique utilisant un vocabulaire technique difficilement compréhensible par les personnes concernées ne serait pas suffisant pour remplir les exigences d'intelligibilité de l'information.

Pour atteindre ces objectifs, la CNIL recommande de mettre en place **une information en plusieurs niveaux, priorisant les informations essentielles** (identité du responsable du traitement, finalités et droits des personnes) **au premier niveau** mais offrant une information complète par ailleurs.

S'agissant des traitements portant sur des données de mineurs, l'information devrait faire l'objet d'une attention particulière pour être suffisamment compréhensible.

- En savoir plus sur [la conception des mentions d'informations](#).

Les dérogations à une information individuelle

En principe, le contenu de l'information précédemment invoqué doit être porté à la connaissance des personnes concernées de manière individuelle, c'est-à-dire directement (par exemple sur un formulaire de collecte de données, de création de compte, par courriel, etc.)

Le RGPD prévoit plusieurs dérogations à l'obligation d'informer individuellement les personnes (par exemple quand un texte de droit européen ou national permet de l'exclure en vertu de l'article 23). Les développements ci-dessous se concentrent sur les dérogations les plus pertinentes en matière de développement d'IA, mais ne sont pas pour autant exhaustives (voir l'article 14.5 du RGPD).

Situation n° 1 : La personne concernée a déjà obtenu les informations sur les traitements à des fins de développement (14.5.a du RGPD)

Lorsque les personnes concernées ont déjà été informées de toutes les caractéristiques du traitement, en particulier de la finalité et de l'identité du responsable du traitement d'apprentissage, **une nouvelle information n'est pas nécessaire.**

À retenir : lorsque les données sont collectées auprès d'un tiers, le responsable du traitement devra s'assurer que l'intégralité des informations sur son propre traitement ont déjà été fournies aux personnes concernées.

À titre de bonne pratique, la CNIL encourage les réutilisateurs de données à s'appuyer sur le diffuseur de données pour informer les personnes, en particulier lorsque ce dernier est encore en contact avec les personnes concernées.

Par exemple le fournisseur d'un service d'éducation en ligne pourrait informer ses clients du fait que leurs données seront traitées par un tiers nommément désigné afin de développer un système d'IA à destination d'enseignants en renvoyant vers les mentions d'information de ce dernier. Si le fournisseur a renvoyé vers toutes les informations sur le traitement en question, le réutilisateur n'aura plus à le faire.

Par exemple, l'éditeur d'un jeu de données qu'il publie sur une plateforme d'échange de données d'entraînement pourrait utilement centraliser les mentions d'information des réutilisateurs sur la page de téléchargement du jeu de données en question.

À l'inverse, si le fournisseur du jeu de données en question a correctement informé les personnes sur son traitement de mise à disposition du jeu de données, mais n'a pas fourni l'ensemble des éléments d'information sur les traitements des réutilisateurs, ces derniers devront informer les personnes concernées par leurs propres moyens.

À noter que dans ce cas, les réutilisateurs du jeu de données se trouveront le plus souvent dans la situation n°2 (qui permet de se contenter de fournir une information générale).

Situation n° 2 : L'information exigerait des efforts disproportionnés (Article 14.5.b du RGPD)

Le responsable du traitement peut alors se contenter de rendre les informations publiquement disponibles.

Cet argument est souvent invoqué par les organismes qui ne sont pas ou plus en lien avec les personnes dont ils traitent les données (par exemple, en cas de réutilisation d'une base de données constituée par un tiers). En effet, dans ce cas, ils ne disposent généralement pas de leurs données de contact.

Une analyse au cas par cas est à réaliser, tenant compte du contexte spécifique de chaque traitement.

L'organisme doit évaluer et documenter le caractère disproportionné en mettant en balance, d'un côté, l'atteinte portée à la vie privée des personnes dont les données sont traitées et, de l'autre, les efforts qu'impliqueraient une communication individuelle des informations aux personnes concernées.

- Pour évaluer **l'ampleur des efforts** à fournir, il y a lieu de prendre en compte l'absence de moyens de contact des personnes concernées, ou l'ancienneté des données de contacts conservées (à l'exactitude incertaine, par exemple des coordonnées de plus de 10 ans), ou encore le nombre de personnes concernées et le coût de la communication.

Par exemple : le responsable du traitement qui entendrait réutiliser les données de ses clients et dispose encore de leur adresse électronique devrait toujours s'en servir pour les informer de manière individuelle.

À l'inverse, le responsable du traitement souhaitant collecter des données indirectement identifiantes n'aura généralement pas à rechercher l'identité réelle ou les coordonnées des personnes afin de les informer directement (une information générale sur son site internet étant alors suffisante).

- Pour évaluer **l'atteinte portée à la vie privée** des personnes concernées et l'intrusivité du traitement, il convient de tenir compte des risques liés au traitement (nature plus ou moins directement identifiante des données, sensibilité des données, etc.) et des garanties éventuelles mises en place (telles que la pseudonymisation, la réalisation d'[une analyse d'impact relative à la protection des données](#) (AIPD), la réduction de la période de conservation ou encore la mise en œuvre de diverses mesures techniques et organisationnelles de sécurité, voir la fiche sur [l'intérêt légitime](#) pour une liste plus détaillée).

Par exemple, en fonction du risque résultant de la nature des données et du contexte de leur publication, le réutilisateur d'un jeu de données publiquement accessible en ligne pourra se prévaloir des mesures prises par le responsable du traitement initial pour informer les personnes concernées sur la possibilité d'une réexploitation par des tiers à des fins d'apprentissage. Le réutilisateur pourra alors se contenter de fournir une information générale (sur son site internet).

Cas particulier de la collecte de données accessibles en ligne

En cas de collecte licite de **données indirectement identifiantes publiées en ligne, une information individuelle sera le plus souvent disproportionnée** dès lors que trouver des moyens de contacter les personnes suppose de rechercher ou collecter des données supplémentaires ou plus identifiantes comme l'identité réelle des personnes.

- Par exemple : il est possible de recourir à une information générale en cas de moissonnage ou de

réutilisation d'un jeu de données d'apprentissage publié en source ouverte de manière licite et ne contenant que des données indirectement identifiantes (telles que des publications ou commentaires dont le contenu est susceptible de permettre l'identification de son auteur).

Il en ira le plus souvent de même pour la collecte de données publiées avec un pseudonyme si ce dernier n'est pas collecté ou conservé par le responsable du traitement.

En cas de collecte de données directement identifiantes, il conviendra de mener une analyse au cas par cas pour apprécier s'il est nécessaire de chercher à informer les personnes de manière individuelle à travers un moyen de contact (par exemple en recherchant leurs coordonnées ou en utilisant un système de messagerie mis en œuvre sur le site internet en question).

Enfin, bien que le volume des données ne saurait présumer à lui seul du caractère disproportionné d'une information individuelle, il en ira le plus souvent ainsi pour la collecte licite de données ne présentant pas de risques pour les personnes concernées, à partir d'un **grand nombre de sites web à des fins de développement d'un grand modèle de langage**, et dont les personnes concernées ne peuvent ignorer qu'elles sont publiquement accessibles, tels que des encyclopédies en sources ouvertes

À retenir : cette dérogation s'appliquera plus facilement aux organismes constituant des bases de données d'apprentissage de systèmes d'IA à des fins de recherche scientifique.

Par exemple : la fourniture d'une information générale s'avérera suffisante pour l'utilisation d'un jeu de données constituées à partir de photographies de profil librement accessibles dans le cadre du développement d'un algorithme de détection d'hypertrucage (*deepfake*) à des fins de recherche scientifique.

Les mesures appropriées pouvant être prises par l'organisme en plus d'une information générale

Au-delà de la fourniture d'une information générale en rendant les informations publiquement disponibles (par le biais, par exemple, de la publication des informations sur le site internet de l'organisme), d'autres mesures appropriées peuvent être prises par l'organisme dans ce cas telles que :

- la réalisation d'une AIPD, y compris lorsque cela n'est pas imposé par l'article 35 du RGPD ;
- l'application de techniques de pseudonymisation des données ;
- la réduction du nombre de données collectées et de la période de conservation ;
- la mise en œuvre de mesures techniques et organisationnelles pour renforcer le niveau de sécurité.

Quelles informations fournir ?

En cas d'information individuelle

Lorsque le responsable du traitement délivre une information individuelle - soit parce qu'il

collecte les données auprès des personnes (article 13 du RGPD) soit parce qu'il collecte les données auprès de tiers mais qu'il dispose d'un moyen de contact et que cela ne représente pas un effort disproportionné (cf. ci-dessus) - **il est généralement requis de fournir l'ensemble des informations qui suivent**, prévues aux articles 13 et 14 du RGPD.

L'organisme qui constitue ou utilise une base de données d'apprentissage pour développer un système d'IA à partir de données personnelles doit informer les personnes concernées sur les éléments suivants, peu importe que les données aient été collectées de manière directe ou indirecte :

- **son identité** et ses coordonnées (tels que son adresse électronique, son adresse postale ou encore son numéro de téléphone) ainsi que les moyens de contacter son délégué à la protection des données ;
- **la finalité et la base légale** du traitement avec, le cas échéant, des précisions sur l'intérêt légitime poursuivi si le traitement se fonde sur celui-ci ;
- **les destinataires** ou *a minima* les catégories de destinataires des données, avec, le cas échéant, des précisions sur les transferts envisagés de ces données vers un pays tiers à l'Union européenne ;
- **la durée de conservation** des données (ou, à défaut, les critères permettant de la déterminer) ;
- **les droits des personnes** concernées (les droits d'accès, de rectification, d'effacement, à la limitation, le droit à la portabilité, le droit d'opposition ou de retirer son consentement à tout moment) ;
- **le droit d'introduire une réclamation** auprès de la CNIL.

À retenir : Si les informations sur la durée de conservation et l'exercice des droits n'ont pas à être systématiquement fournies pour tous les traitements, elles seront quasi-systématiquement requises s'agissant de la constitution et l'utilisation de jeux de données d'apprentissage. En effet, elles sont nécessaires pour garantir un traitement équitable et transparent à l'égard des personnes concernées.

En cas de collecte indirecte, les organismes doivent fournir, en complément :

- **Les catégories de données** personnelles (par exemple : identité, coordonnées, images, publications sur les réseaux sociaux, etc.) ;
- **Lorsque cela est nécessaire pour garantir un traitement équitable et transparent, une indication précise sur la ou les sources** des données (en indiquant notamment s'il s'agit ou non de sources accessibles au public).

Une telle information doit permettre aux personnes de pouvoir anticiper si elles sont concernées par le traitement et faciliter un éventuel exercice de leurs droits sur le traitement source.

En cas de publication d'une notice d'information générale

Lorsque l'information individuelle n'est pas possible, parce que le responsable de traitement ne peut identifier les personnes présentes dans la base, ne dispose pas de données de contact des personnes concernées ou

encore que les contacter individuellement requerrait des efforts disproportionnés, des mesures appropriées doivent être prises. Il est nécessaire de publier une notice d'information, par exemple sur un site internet, qui comporte, si possible, les informations qui auraient été fournies en cas d'information individuelle.

Par ailleurs, cette information doit comprendre, le cas échéant, **le fait que le responsable du traitement ne sera pas en mesure d'identifier les personnes, y compris pour répondre à leurs demandes d'exercice de droit** (conformément à l'article 11 du RGPD). Dans ce cas, la CNIL recommande, si c'est possible, d'indiquer aux personnes souhaitant exercer leurs droits quelles informations complémentaires elles peuvent fournir pour permettre leur identification.

L'information sur les sources présente des difficultés particulières. Deux cas sont à distinguer :

- **Lorsque le responsable du traitement a utilisé un nombre limité de sources** pour constituer sa base de données d'entraînement, il est généralement requis qu'il fournisse des indications précises sur ces sources, sauf exception dont il pourrait justifier.
- **Lorsque de nombreuses sources sont utilisées**, par exemple un grand nombre de sources accessibles au public en ligne, une information globale, indiquant par exemple des catégories de sources, voire les noms de quelques sources principales ou typiques, est généralement suffisante (le considérant 61 du RGPD prévoit bien que « lorsque l'origine des données à caractère personnel n'a pas pu être communiquée à la personne concernée parce que plusieurs sources ont été utilisées, des informations générales devraient être fournies »).

En cas de réutilisation d'un jeu de données ou d'un modèle d'IA soumis au RGPD

Outre l'indication de la source des données utilisées, la CNIL recommande, *a minima* pour les jeux de données qui présentent le plus de risques pour les personnes, de fournir les moyens de contacter le responsable du traitement auprès duquel il a été récupéré. Une bonne pratique consiste à renvoyer directement vers le site web du responsable du traitement d'origine, et à accompagner l'information d'une explication synthétique et claire sur les conditions de collecte et d'annotation des données.

Exemples de mentions :

Ces données d'entraînement sont issues d'un jeu de données publiquement accessible (lien hypertexte vers la publication) contenant 70 000 images et constitué à partir de photographies publiées en ligne en licence ouverte sur le réseau social _____ entre 2020 et 2021.

Ces données sont issues d'un jeu de données fourni par la société _____, courtier de données (dont les coordonnées sont _____). Cette base de données est constituée de 3 000 images initialement collectées auprès d'acteurs volontaires jouant différentes expressions du visage, annotées pour retranscrire leurs émotions.

Dans le cadre du développement de ce système d'IA, nous réutilisons un grand modèle de

langage développé par la société _____ ayant mémorisé des données personnelles. Pour en savoir plus, nous vous renvoyons à sa politique de confidentialité accessible à l'adresse suivante : _____.

En cas de moissonnage (ou *webscraping*) sur des sites web ou de réutilisation de données moissonnées

Si le moissonnage concerne quelques sites, la CNIL recommande une information précise sur les sources utilisées. Lorsque les sources sont très nombreuses, elle recommande de fournir les catégories de sites sources concernés, *a minima* ceux **qui présentent le plus de risques pour les personnes**. Cette recommandation s'applique aux moissonneurs de données, mais également aux responsables de traitements qui réutilisent des jeux de données constitués à partir de données moissonnées.

Exemples de mentions suffisamment précises :

Pour cela, nous avons collecté des données librement accessibles par moissonnage sur les plateformes suivantes : _____. Ces données consistent des publications rendues manifestement publiques par leurs auteurs sur le sujet _____. Les commentaires liés aux publications n'ont pas été collectés. Les images et pseudonymes ont été traités lors de la collecte mais n'ont pas été conservés.

Les données ont été collectées par moissonnage de sites spécialisés dans le domaine étudié et peuvent contenir des données personnelles telles que le nom de l'auteur, ou d'une personne citée dans un article de blog librement accessible. Les images contenues dans l'article ne sont pas collectées. _.

Lorsque cela est possible, fournir une information sur les noms de domaine et URL des pages web moissonnées ainsi que la date ou période de collecte est une bonne pratique, mais n'est pas requise au titre des obligations d'information.

En cas de développement d'un modèle d'IA à usage général au sens du règlement sur l'IA

En parallèle de l'obligation d'information prévue par le RGPD, l'article 53 du règlement sur l'IA prévoit que les fournisseurs de modèles d'IA à usage général élaborent et mettent à la disposition du public un résumé suffisamment détaillé du contenu utilisé pour l'entraînement de ces modèles, conformément aux indications fournies par le Bureau de l'IA (Commission européenne). Le considérant 107 de ce règlement précise que 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 à exercer et à faire respecter leurs droits.

Il s'agit, par exemple, d'énumérer 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.

Ce résumé pourra également, en principe, servir d'information générale sur les sources de données, pour ces types de modèles d'IA, pour l'application du RGPD sous réserve de la publication définitive de ce résumé par le Bureau de l'IA..

Pour entraîner notre grand modèle de langage décrit précédemment [voir notamment la fiche sur la précision de la définition de la finalité du traitement], nous réutilisons des données publiquement accessibles en ligne sur les catégories de site suivants [à remplir] que nous avons collectées le [date] au sein du jeu de données _____.

Exemples de catégories de sites sources :

- sites institutionnels
- encyclopédie en source ouverte
- plateformes d'échanges de papiers de recherche
- sites de presse nationale [ou internationale]
- site de médias audio-visuels nationaux ou internationaux
- forums de discussions spécialisés [en décrivant ces spécialités, par ex : en matière de développement informatique, d'éducation, de santé, d'évènements culturels, d'évènements sportifs, etc.] ou généralistes, etc.
- plateformes de partage d'images, de contenu audio, ou audiovisuel (musiques, photographies, enregistrements audios, vidéos)
- réseaux sociaux professionnels
- sites de vente en ligne

Concernant le cas particulier des modèles d'IA dont le traitement est soumis au RGPD

Un certain nombre de modèles d'IA sont considérés comme anonymes : le RGPD ne s'applique pas à eux en tant que tel, y compris l'obligation d'information. A l'inverse, l'entraînement d'un modèle d'IA peut parfois conduire à ce que celui-ci « mémorise » une partie des données d'apprentissage (voir l'[avis du CEPD 28/2024](#) sur certains aspects de la protection des données liés au traitement des données à caractère personnel dans le contexte des modèles d'IA). Lorsqu'il s'avère que le modèle est en lui-même soumis au RGPD, il convient d'informer les personnes sur les données mémorisées.

Le fournisseur du modèle ou du système d'IA devrait alors préciser les éléments d'information indiqués plus haut (finalité, responsable de traitement, destinataires etc.).

À titre de bonne pratique, il est également recommandé au fournisseur de préciser :

- la nature du risque lié à l'extraction des données à partir du modèle, comme le risque de régurgitation de données dans le cas de l'IA générative ;
- les mesures prises afin de limiter ces risques, et les mécanismes de recours existants dans le cas où ces

risques se manifesteraient, comme la possibilité de signaler à l'organisme une occurrence de régurgitation.

[< Fiche précédente : La base légale de l'intérêt légitime : fiche focus sur les mesures à prendre en cas de collecte des données par moissonnage \(*web scraping*\)](#)

[Sommaire](#)

[Fiche suivante : Respecter et faciliter l'exercice des droits des personnes concernées >](#)