

Un rapport de la **Fondation pour l'Enfance**



# L'IA générative, nouvelle arme de la pédocriminalité

Octobre 2024

**FONDATION**  
**POUR**  
**L'ENFANCE**

reconnue d'utilité publique

FONDATION POUR L'ENFANCE  
23, place Victor Hugo, 94 270 Kremlin-Bicêtre  
01 43 90 63 10  
fondation-enfance.org

*Contact*

**Angèle Lefranc**, chargée de plaidoyer  
angele.lefranc@fondation-enfance.org

Direction de la communication  
Conception graphique: Agence Panteo.fr/x.jacobi@panteo.fr

Illustrations: IStock ; Adobe Stock  
Imprimé en France sur les presses d'Escourbiac l'imprimeur

# L'IA générative, nouvelle arme de la pédocriminalité

” La méconnaissance de l'ampleur et de la diversité des phénomènes en lien avec l'exploitation sexuelle des mineurs en ligne, de l'adaptabilité d'une communauté avisée, est presque aussi grande que le phénomène lui-même. Et cela doit changer. ”

VÉRONIQUE BÉCHU,  
*Derrière l'écran, Stock, 2024*

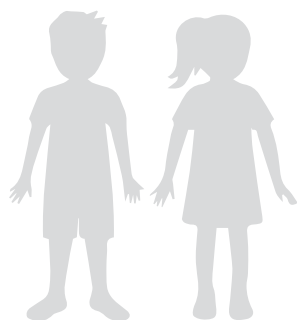
# Table des matières

- 6** La cyberpédocriminalité en quelques chiffres
- 8** Synthèse
- 10** Pédocriminalité & IA: nouveaux usages, nouveaux enjeux, nouveaux risques
- 14** Nos recommandations
- 16** Avant-propos



## **18** De quoi parle-t-on ?

- 18** Contexte et historique
- 22** L'émergence d'une nouvelle forme de pédocriminalité



- 26** Les marges de manœuvre des entreprises pour prévenir la création de contenus violents et illégaux et leurs limites
- 27** Le développement de modèles d'IA générative destinés à la création de contenus pédocriminels
- 30** L'utilisation d'application de « déshabillage » pour générer des deepfakes pédocriminels
- 30** Typologies des contenus pédocriminels qu'il est possible de générer par l'IA
- 34** L'utilisation des réseaux sociaux et des services de messagerie pour promouvoir, partager et vendre des contenus pédocriminels générés par l'IA
- 39** Les solutions technologiques envisageables pour renforcer la sécurité des enfants
- 42 À LA LOUPE.** Remontée terrain d'un acteur directement confronté aux contenus pédocriminels générés par l'IA. Entretien avec **Point de Contact**

## 44 L'impact de la cyberpédocriminalité générée par l'IA sur la protection de l'enfance

**50 À LA LOUPE.** Les impacts des violences sexuelles et de la cyberpédocriminalité, notamment générée par l'IA, sur l'individu victime

Entretien croisé avec **Joanna Smith** et **Mélanie Dupont**

**56 ZOOM SUR.** Les auteurs et consommateurs de matériels d'exploitation sexuelle des mineurs en ligne



Retrouvez l'actualité de la Fondation pour l'Enfance



## 62 État des lieux du cadre législatif en vigueur et des initiatives en cours

63 Cadre international

71 Cadre européen

76 Cadre national

**82 À LA LOUPE.** Exemple de bonne pratique : Sensibiliser pour mieux prévenir les contenus pédocriminels générés par l'IA

Entretien avec **Églantine Cami**, Association **CAMELEON**

## 86 Conclusion

**89 À LA LOUPE.** Ce que l'explosion de l'IA générative et de ses détournements disent de notre société

Entretien avec **Pascal Plantard**

92 Recommandations

98 Glossaire

101 Comité de rédaction, remerciements

# La cyberpédocrimi

**35,9 millions**

de signalements de contenus pédocriminels au NCMEC en 2023

**5<sup>e</sup>**

La France est le 5<sup>e</sup> pays européen et le 9<sup>e</sup> pays dans le monde hôte de contenus pédocriminels

**59 %**

des contenus représentant des violences sexuelles faites aux enfants étaient hébergés dans l'UE en 2022

**871**

signalements de contenus pédocriminels échangés en ligne sont transmis chaque jour à l'OFMIN, soit une augmentation de 12 000 % en dix ans

**70 %**

ont été signalés par des plateformes en ligne traditionnelles (telles que les réseaux sociaux)

**300 %**

d'augmentation des signalements de cas de sextorsion au NCMEC en 2023

# minalité

## EN QUELQUES CHIFFRES

**20 254**

images générées par l'IA ont été publiées sur un forum accessible sur le dark web consacré aux pédocriminels:

**1372**

images représentaient des enfants âgés de sept à dix ans

**143**

images représentaient des enfants âgés de trois à six ans

**14**

victimes identifiées par jour en moyenne grâce à la base de données d'Interpol

**4700**

contenus pédocriminels impliquant l'IA générative signalés au NCMEC en 2023

**50 %**

des images ou des vidéos d'enfants échangées sur les forums pédocriminels ont été initialement publiées par leurs parents via les réseaux sociaux

Sources:

NCMEC CyberTipline data 2023; Internet Watch Foundation; Annual Report 2022; France fifth worst for hosting child sexual abuse content in EU, as criminals target French servers", 22nd September 2023; "How AI is being abused to create child sexual abuse imagery", October 2023; Véronique Béchu, Derrière l'écran, Stock, 2024 & communication de l'OFMIN; Base de données internationale sur l'exploitation sexuelle des enfants d'Interpol







# Synthèse

**S**i elle représente une évolution technologique majeure, l'intelligence artificielle générative se transforme en une arme redoutable lorsqu'elle est utilisée par des personnes mal intentionnées, comme les cyberpédocriminels. Sextorsion, grooming, vidéos pédocriminelles<sup>1</sup>... Ces pratiques sont désormais facilitées et amplifiées par cette technologie qui ne cesse de se perfectionner.

Elles occasionnent notamment des difficultés pour les forces de l'ordre, qui peinent à distinguer les images non générées par l'IA de celles générées

**RR** **Nos objectifs ?**  
**Protéger et assurer**  
**la sécurité et l'intégrité en**  
**ligne des enfants.** **99**

par l'IA, et donc à identifier les enfants victimes de violences. Sans compter l'inadéquation du cadre juridique et de

la législation, qui donne aux cyberpédocriminels un sentiment d'impunité et de toute puissance. Le caractère virtuel de ces montages entraîne chez leurs auteurs et consommateurs une tendance à la banalisation et à la normalisation de ces pratiques. Pourtant, c'est l'enfant lui-même, son intégrité physique et morale, ses droits qui sont attaqués.

Bien que la part de ces contenus dans l'ensemble des signalements soit encore faible, et ce au niveau mondial, la Fondation pour l'Enfance et ses partenaires tirent la sonnette d'alarme. Après avoir mené une recherche approfondie pendant près d'un an, leur constat est sans appel : la partie immergée de ces contenus pédocriminels dopés à l'IA semble être colossale, et leurs conséquences pour les victimes, dramatiques. Il est donc urgent d'engager une réponse forte, rapide et coordonnée entre les différents acteurs juridiques, politiques, technologiques, afin d'apporter un cadre légal à l'utilisation de l'IA, mais aussi de permettre une prise de conscience sociétale autour des risques liés à l'IA et du danger lié à certaines pratiques, notamment du sharenting.

1. Un glossaire expliquant certains termes liés à la cyberpédocriminalité et à l'IA générative se trouve en fin de rapport.

# Pédocriminalité & IA: nouveaux usages, nouveaux enjeux, nouveaux risques

L'essor de l'IA générative nous confronte à cette question: comment distinguer les contenus (textes, images, audios, vidéos) créés ex nihilo de ceux qui montrent des personnes réelles? Sur le terrain de la lutte contre la cybercriminalité, cette question prend une ampleur vertigineuse. En cause: l'émergence de nouveaux usages pédocriminels complexes et dangereux.

## Les nouveaux usages pédocriminels avec l'IA générative

Des vidéos d'enfants qui n'existent pas en train d'être violés, des visages de vraies adolescentes dont le corps est « synthétiquement » entièrement dénudé... L'IA générative permet de générer à l'infini de tels contenus, brouillant ainsi les pistes entre réalité et virtuel. En somme, un véritable « terrain de jeu » pour cyberpédocriminels, dont voici quelques-uns des nouveaux usages.

### **Modification de modèles et systèmes d'IA, dans le but de les faire générer des contenus pédocriminels:**

- Certains modèles et systèmes d'IA open source et disponibles au grand public sont modifiés par des individus,

dans le but de créer des contenus pédocriminels

- Ces modèles et systèmes modifiés sont entraînés sur de larges bibliothèques de contenus d'exploitation sexuelle de mineurs, de contenus pornographiques, mais aussi sur des images d'enfants à caractère non sexuel, afin de leur enseigner comment produire avec précision du nouveau matériel pédocriminel. Ces contenus sur lesquels les modèles sont entraînés sont des images et vidéos, présents en grande quantité sur internet et facilement accessibles.
- La modification des modèles et la génération d'images pédocriminelles

peuvent être effectués hors ligne sur un ordinateur personnel, permettant ainsi d'échapper à la détection.

### **Création et modification de contenus démultipliés par les IA sur le *dark web* et le *clear web* :**

- La création de contenus pédocriminels ex nihilo: des images, photos et vidéos pédocriminelles, totalement artificielles, ressemblant à s'y méprendre à des agressions sexuelles/viols réels.
- La création de deepfakes pédocrimi-

minels: des images, photos et vidéos pédocriminelles, générées à partir d'autres contenus d'enfants à caractère sexuel, ou même non-sexuel, présents sur internet (comme les réseaux sociaux). Ces montages sont générés à partir de modèles d'IA modifiés ou de nudify apps, des applications IA de déshabillage.

- L'édition et l'amélioration de la qualité de photos et vidéos pédocriminelles déjà existantes.
- Des instructions données aux



## **LE POINT DE VUE DES EXPERTES**

**JOANNA SMITH**, *Psychologue clinicienne et*  
**MÉLANIE DUPONT**, *Docteur en psychologie, Psychologue*  
*à l'Unité Médico-Judiciaire de l'Hôtel-Dieu (Paris) et Présidente*  
*de l'Association contre les Violences sur Mineurs (CVM)*

### **Quels impacts de la cyberpédocriminalité pour les victimes ?**

Il est difficile, voire impossible pour la victime de voir une fin à l'épisode traumatique, parce que le contenu continue de tourner sur internet. Cette fin est pourtant cruciale dans le traitement du traumatisme. Sans celle-ci, le danger est encore présent. Avec la cyberpédocriminalité, il y a une permanence de l'agression, de multiples agresseurs (ceux qui visionnent, téléchargent, partagent le contenu), et donc une permanence des conséquences et une revictimisation. La création de contenus pédocriminels par l'IA générative, et donc l'impossibilité de contrôler son image, va entraîner chez les victimes une intensification du sentiment de dépossession de soi. Cela pourrait avoir pour effet d'augmenter les troubles psychotiques chez les jeunes, et même une dépersonnalisation à se voir sur des images qui ne représentent pas leurs propres corps. » ●

modèles d'IA générative pour créer et affiner des guides et des tutoriels sur la manière de gagner la confiance, de

violer, d'agresser, de torturer et de tuer des enfants, ou de créer des contenus pédocriminels réalistes.

## IA & cyberpédocriminalité : des mineurs encore plus en danger ?

**52 %** des consommateurs pensent que leur usage de contenus pédocriminels pourrait aboutir à une agression sur un enfant (44 % des consommateurs ont pensé à contacter des enfants et 37 % ont contacté des enfants au moins une fois)\*.

Puisque l'IA générative permet une création de contenus à l'infini, elle augmente les comportements addictifs des consommateurs, avec des images de plus en plus extrêmes, explicites, violentes. Et donc des risques accrus de passage à l'acte ?

\*Protect Children, ReDirection Survey Report, 2021, p.16

## Le business de la cyberpédocriminalité dopée à l'IA

Comme à chaque avancée technologique, l'innovation se transforme en consommation. La cyberpédocriminalité via IA générative ne fait pas exception à la règle : un nouveau business est en train d'émerger.

Pour publier, partager et/ou faire la promotion de leurs contenus, les pédocriminels utilisent le *clear web* avec de faux profils sur les réseaux sociaux (Instagram, Facebook, Tik Tok) et les plateformes de messagerie cryptées (Telegram, Whatsapp).

Certains profils encouragent les

abonnés à contacter le propriétaire via le système de messagerie de la plateforme, ou via une plateforme chiffrée tierce, pour obtenir plus d'images, toujours plus graphiques et explicites. Nombre d'entre eux mettent leurs contenus à disposition moyennant paiement, et relaient leurs « clients » vers des systèmes de paiement et autres services d'abonnement (tels qu'OnlyFans).

Si un utilisateur ne parvient pas à générer ce qu'il souhaite, ou si un modèle ou un fichier n'existe pas encore, il

peut être amené à payer un utilisateur plus qualifié sur l'IA pour le faire à sa place.

Certains pédocriminels font du chantage à des mineurs pour obtenir de leur part de l'argent ou des contenus

à caractère sexuel (sextorsion). Cette pratique se fait à l'aide de montages pédocriminels générés avec l'IA, à partir de photos à caractère non-sexuel, obtenues notamment sur les réseaux sociaux.

## Les 3 principaux enjeux pour protéger les mineurs de la cyberpédocriminalité générée par l'IA

L'émergence des contenus pédocriminels générés par l'IA amplifie les enjeux déjà existants en matière de protection des enfants, mais en crée également de nouveaux:

du partage d'un modèle d'IA générative conçu pour produire des contenus pédocriminels. Ce vide juridique et cette absence de législation nationale et internationale claire, assortie d'une absence de prise de conscience sociétale, permet à la pratique de s'intensifier.



### Identifier et protéger les enfants victimes

Par leur réalisme bluffant, les contenus pédocriminels générés par IA rendent la tâche d'identification et de protection des enfants victimes de violences sexuelles difficile pour les forces de l'ordre et les plateformes de signalement.



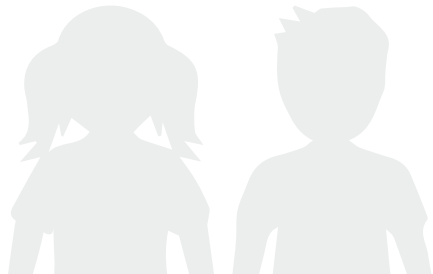
### Lutter contre l'intensification et la banalisation des violences sexuelles sur les enfants

La facilité de production des contenus pédocriminels générés par l'IA et leur multiplication sur Internet entraînent une normalisation et une banalisation des violences sexuelles sur les enfants.



### Engager une réponse juridique et politique solide et coordonnée

À l'heure actuelle, il n'existe pas de législation qui traite avec précision de la création, de la possession ou



# Nos recommandations

La Fondation pour l'Enfance appelle les États et les entreprises du secteur des nouvelles technologies à agir au plus vite.

**Trois objectifs impératifs sont à atteindre:**

## Prévention

Mettre en place des campagnes nationales de sensibilisation du grand public à la cyberpédocriminalité, aux dangers relatifs au *sharenting* et aux bonnes pratiques à adopter pour protéger les enfants. Cette campagne doit permettre aux parents d'appréhender leur rôle en matière de prévention et d'être davantage sensibilisés aux risques liés à l'IA.

## Détection

Favoriser l'innovation, en incitant les acteurs privés à coopérer pour mettre en place des outils permettant de distinguer les contenus générés par l'IA des contenus non générés par l'IA. L'utilisation de tels outils permettrait de pallier la difficulté d'identification des mineurs victimes de violences.

Instaurer un travail conjoint entre les différentes entreprises et les plateformes, afin d'améliorer l'identification,

le signalement et le retrait des contenus pédocriminels et des modèles d'IA générative destinés à les générer.

## Sanction

Amender l'article 227-23 du Code pénal pour y insérer les fichiers ou représentations issus de l'IA avec l'insertion d'un nouvel alinéa qui pourrait être rédigé comme suit: « *Le fait, de concevoir, de créer, de diffuser ou de porter à la connaissance du public ou d'un tiers, par quelque voie que ce soit, tout montage, contenu visuel ou sonore à caractère sexuel généré par un traitement algorithmique tel que visé à l'alinéa 1 de l'article 226-8-1 est puni de X ans de prison et X euros d'amende lorsqu'il s'agit de la représentation, de l'image ou de la parole d'un mineur.* »

Pénaliser la création et la mise à disposition de modèles d'IA générative destinés à générer des contenus pédocriminels. Il pourrait être ajouté au Code pénal un nouvel article rédigé de la manière suivante: « *Est puni de X années d'emprisonnement et X euros d'amende le fait de collecter, détenir, traiter ou détourner des données à caractère personnel, afin de créer, générer ou mettre à disposition du public ou de tout tiers*

*un modèle de traitement algorithmique, dans le but de permettre la création de contenu visuel ou sonore*

*à caractère sexuel représentant un mineur, et de tout fichier à caractère pédopornographique. »*



## **Quand les géants de la tech’ s’engagent**

Meta, Google, Microsoft, OpenAI, Amazon... Les leaders de l'IA ont adopté le 23 avril 2024 le texte "Safety by Design for Generative AI: Preventing Child Sexual Abuse" proposé par l'organisation à but non lucratif Thorn. Au cœur de ces engagements ? La mise en place de nouvelles mesures de sécurité pour protéger les enfants en ligne, et une meilleure prise en compte de la protection des enfants dans le développement et le déploiement de l'IA générative, afin d'empêcher que ces outils permettent la création de contenus pédocriminels. Si le texte n'a valeur que morale et non juridique, il reste un premier pas inspirant !



# Avant-propos

**JOËLLE SICAMOIS,**  
directrice de la Fondation pour l'Enfance

**P**armi les principes directeurs et les droits énoncés dans la Convention Internationale des Droits de l'Enfant (CIDE), il est reconnu à tout enfant le droit d'être protégé de la violence, de la maltraitance et de toute forme d'abus et d'exploitation<sup>2</sup>.

Depuis plus de 45 ans, la Fondation pour l'Enfance lutte contre les violences faites aux enfants, et s'appuie depuis 1989 sur les principes et droits reconnus par la CIDE. Elle s'est notamment donnée pour mission d'identifier les nouveaux risques de violence auxquels les enfants sont exposés. Le XXI<sup>e</sup> siècle a vu l'avènement de la pédocriminalité en ligne, qui constitue aujourd'hui un danger majeur auquel font face les enfants. Depuis plus de 20 ans, grâce au cabinet Lombard, Baratelli, Astolfe & associés, la Fondation se constitue partie civile dans des affaires impliquant notamment la création, la captation ou la diffusion d'images à caractère

pédopornographique. Par cette action, nous souhaitons tout d'abord porter la voix de tous ces enfants qui, bien souvent, n'ont pas été identifiés, et qui ne sont pas représentés. Mais nous souhaitons également remettre

**“ Nous souhaitons remettre l'enfance, ses intérêts, ses droits et sa protection au cœur du procès pénal. ”**

l'enfance, ses intérêts, ses droits et sa protection au cœur du procès pénal.

Aujourd'hui, avec le développement et la démocratisation de l'IA générative, nous faisons face au risque de voir se multiplier des contenus pédocriminels, créés, enregistrés et diffusés par, et pour des individus pour qui le caractère virtuel de ces représentations et de l'outil utilisé les dédouane de leur responsabilité.

Pourtant, ces images d'exploitation sexuelle des mineurs, qu'elles soient

2. L'article 16 consacre le droit de chaque enfant à la protection de sa vie privée.

Les articles 19 et 37 énoncent le droit de chaque enfant à être protégé contre toutes les formes de violence, la torture, et les traitements cruels, inhumains ou dégradants.

Les articles 32, 34 et 36 énoncent le droit de l'enfant à être protégée contre toute forme d'exploitation, notamment sexuelle et de violences sexuelles.

virtuelles ou non, sont constitutives d'une atteinte à l'intégrité physique et morale de l'enfant, et en ce sens d'une violence. Elles participent à une objectification de l'enfant, et entretiennent une culture de la violence sexuelle, physique, psychologique. Et elles sont souvent une première étape avant un passage à l'acte.

Prévenir les risques de violences ou de maltraitements envers les enfants implique une nécessaire intervention de l'État, et une forte sensibilisation à destination des parents et professionnels de l'enfance.

Il est urgent que nos États s'emparent de ce sujet, et cela au plus haut niveau, afin de protéger les enfants de toute forme d'exploitation sexuelle en ligne. Le blocage de l'adoption du règlement européen visant à prévenir et à combattre les abus sexuels sur enfants nous fait regretter que la protection des enfants contre les violences sexuelles ne soit pas, dans certaines circonstances, la priorité numéro 1 de nos États européens.

Il est également impératif que les entreprises s'engagent dans la prévention et la détection des contenus

 **La lutte contre la pédocriminalité en ligne, notamment générée par l'IA, est un défi social et sociétal, et nous avons tous une responsabilité.** 

pédocriminels. La lutte contre la pédocriminalité en ligne, notamment générée par l'IA, est un défi social et sociétal, et nous avons tous une responsabilité.

Alertée par une étude de l'Internet Watch Foundation (IWF) sur ce phénomène à l'été 2023, la Fondation pour l'Enfance s'est emparée du sujet et a mené des recherches pendant près d'un an avec ses partenaires. Pour cela, nous avons fait appel à un panel large et représentatif, permettant un éclairage tant technique, que juridique et légal. Nous avons également sollicité des opérateurs sur le terrain. Nous espérons que ces travaux permettront d'accélérer les décisions et arbitrages en cours.

# De quoi parle-t-on ?

## Contexte et historique

Avant les années 1990, la production et la diffusion de contenus pédocriminels<sup>3</sup> intervenaient lorsqu'un enfant était violé ou agressé sexuellement, et que son agresseur réalisait des photos ou des vidéos qu'il partageait avec d'autres par courrier, en personne ou dans des magazines.

Depuis les années 1990, l'émergence de la cyberpédocriminalité, son évolution, sa sophistication et son accessibilité ont suivi celles d'Internet et des nouvelles technologies. Les producteurs et consommateurs de contenus pédocriminels en sont souvent des adeptes précoces. Très vite, ils s'approprient et exploitent les réseaux sociaux, les forums, les sites de petites annonces, les plateformes de partage de fichiers, les plateformes de messagerie et de tchat, pour créer et partager toujours plus de contenus pédocriminels, mais aussi pour entrer directement en contact avec des enfants.

Tant est si bien qu'aujourd'hui, l'exploitation sexuelle des enfants en ligne

est toujours plus complexe, variée, et en constante augmentation : sextorsion\*, *grooming*\*, prostitution... Le développement de technologies d'anonymisation permet également aux auteurs et consommateurs de dissimuler leur véritable identité et leur localisation, entraînant un véritable sentiment d'impunité.

L'intelligence artificielle\* (IA) ne fait pas exception à cette instrumentalisation des nouvelles technologies.

Si les premières intelligences artificielles sont nées au milieu du XX<sup>e</sup> siècle, elles ont connu une évolution majeure au XXI<sup>e</sup> siècle et ont été démocratisées à partir de 2022-2023 avec l'arrivée de l'IA générative\*. Cette technologie nous permet aujourd'hui de générer du contenu (textes, images, vidéos) en fonction des instructions\* qui lui sont formulées. L'émergence de plateformes d'IA générative ultra sophistiquées et accessibles au public marque un véritable tournant dans l'évolution de l'exploitation sexuelle des enfants en ligne.

Dans le cadre de nos recherches, nous avons tenté de répondre à un certain nombre de questions : Comment l'IA générative peut-elle, techniquement,

3. Les termes spécifiques à la pédocriminalité en ligne et à l'IA générative sont définis dans le glossaire à la fin du rapport. Ils seront indiqués dans le texte par le sigle \*.

être utilisée et détournée à des fins pédocriminelles ? Quels sont les impacts de cette pratique sur les enfants victimes ? Quels défis ce nouveau phénomène pose-t-il aux acteurs de la protection de l'enfance en ligne pour détecter, identifier, supprimer les contenus pédocriminels, et enquêter et poursuivre les auteurs ? Qui sont les auteurs et consommateurs de contenus pédocriminels ? L'IA générative marque-t-elle un tournant dans leurs parcours criminels ? Les pouvoirs publics et les entreprises du secteur des nouvelles technologies ont-ils conscience de ce phénomène ? s'en emparent-ils ? si oui, comment ? Existe-t-il un mouvement international et européen de coopération entre les différents acteurs publics et privés pour lutter contre cette pratique ? Quelles actions les pouvoirs publics et les entreprises du secteur des nouvelles technologies devraient-ils mettre en œuvre pour prévenir, détecter et sanctionner l'utilisation de l'IA générative à des fins pédocriminelles ?

La difficulté dans la réponse à toutes ces questions est que l'IA générative est en perpétuel mouvement. Chaque jour (ou presque) apporte son lot de nouvelles évolutions dans les

capacités des modèles d'IA\*, dans les pratiques des utilisateurs, et donc de nouvelles problématiques et de nouveaux enjeux.

Ce rapport s'attache à faire la lumière sur ce phénomène, encore minoraire dans les signalements reçus par les policiers et les plateformes, mais avec un potentiel d'envergure sans commune mesure. Nous expliquerons comment l'IA générative est utilisée à des fins pédocriminelles, et les enjeux que ce détournement pose pour la protection et la sécurité des enfants en ligne. Nous présenterons également les initiatives politiques, législatives et technologiques engagées au niveau international, européen et national. Cette présentation nous permettra enfin d'apporter un certain nombre de recommandations à destination des pouvoirs publics et du secteur des nouvelles technologies pour mieux prévenir et combattre ce nouveau phénomène.

Tout au long de ce rapport, vous trouverez des éclairages d'experts de l'IA, de la santé, mais aussi d'experts de terrain qui œuvrent chaque jour pour une meilleure protection des enfants en ligne.



# L'émergence de l'IA générative et ses bénéfices pour les enfants et les adolescents

Explication de **NICOLAS GREFFARD**,  
*directeur Technique / Expert IA chez Valeuriad<sup>4</sup>*

**Fondation pour l'Enfance: Qu'est-ce que l'IA générative? En quoi représente-t-elle un tournant?**

**Nicolas Greffard:** De manière générale, on parle d'IA générative pour la différencier de ce qu'on appelait jusqu'à présent l'IA.

En informatique, l'IA désigne les technologies qui permettent de répondre à des cas d'usages<sup>5</sup> qui étaient, jusqu'à présent, exclusifs à l'humain (par exemple la reconnaissance faciale).

Les systèmes d'IA génératives actuels, tels que ChatGPT, se distinguent par leur capacité généraliste. Les modèles précédents étaient entraînés dans un but bien précis. Reprenons l'exemple de la reconnaissance faciale: un précédent modèle d'IA était entraîné à reconnaître des gens sur des photos, et il ne savait faire que ça. Si on lui donnait une image de forêt ou d'avion, le modèle était incapable d'en faire quoique ce soit. Les modèles actuels sont capables de répondre à des instructions de nature variée. Ces capacités généralistes sont notamment induites lors de l'entraînement par des corpus de données\* beaucoup plus larges, des modèles de plus grande taille\*, et des objectifs d'apprentissages qui restent généralistes. Cela permet d'élargir le champ des possibilités.<sup>6</sup>

4. Valeuriad est une société de conseils, services et expertise technologique implantée à Nantes.

5. Par cas d'usage on entend des fonctionnalités dans un système d'information, des tâches, ou des activités qui jusqu'à présent étaient exclusives à l'homme. Ces cas d'usage peuvent être détournés: on développe une fonctionnalité et une activité dans un objectif précis, mais les utilisateurs s'en servent pour autre chose.

6. Aujourd'hui, il n'existe pas d'obligations de s'assurer de la qualité et de la fiabilité des données sur lesquelles on entraîne les modèles d'IA. Plus il y a de données sur lesquelles le modèle est entraîné, plus l'apprentissage de ce modèle sera intéressant et lui permettra de répondre précisément au plus de questions possibles. La vérification des données d'entraînement dépend de l'objectif du modèle et de la cible finale.

C'est pour cela qu'il y a un tournant technologique, et presque sociétal : on ne sait pas définir jusqu'où on sera capable d'aller avec l'IA générative (en bien comme en mal) dans une semaine, dans un mois, dans un an. Cet aspect non cadré présente des opportunités : en essayant les modèles, on se rend compte qu'ils peuvent apporter de la valeur dans tout un tas de cas d'usage, mais on ne sait pas à l'avance si l'apport de l'IA va être pertinent.

**Fondation pour l'Enfance :** Quels sont les apports positifs de l'IA générative, notamment pour les enfants et les adolescents ?

**NG. :** L'apport positif pour les enfants et les jeunes se situe au niveau de l'éducation, de la découverte et de l'apprentissage. En interagissant avec un outil tel que ChatGPT, n'importe qui, même des populations défavorisées, ou qui n'ont jamais été familiarisées avec un sujet, peuvent s'intéresser à quelque chose de nouveau qui, jusque-là, était « réservé » à d'autres populations plus privilégiées. C'est une porte ouverte sur la richesse informationnelle du monde qui est à disposition. ●



## L'émergence d'une nouvelle forme de pédocriminalité

D'après une étude de l'entreprise néerlandaise de cybersécurité Deeptrace, dès 2019, 96 % des vidéos *deepfakes*\* relevaient de la pornographie non-consensuelle utilisant des images de femmes, souvent célèbres<sup>7</sup>. À l'été 2023, 404 Media a enquêté sur les nouveaux « marchés » de l'IA pornographique<sup>8</sup>. Les journalistes ont fouillé les recoins de plateformes de « modèles IA » (telles que CivitAI et Mage), qui proposent à toute personne de contribuer en ajoutant ses propres modèles d'IA générative. Ceux-ci sont laissés en accès libre aux membres des plateformes, et, même si les plateformes n'autorisent techniquement pas ce genre de pratiques, ils sont utilisés pour créer des images à caractère pornographique (principalement de femmes célèbres donc). Alors que ces sites n'ont rencontré leur succès qu'en 2023, certains modèles à visée pornographique comptabilisaient déjà plusieurs dizaines de milliers de téléchargements en quelques mois seulement<sup>9</sup>.

Les enfants et les adolescents ne sont

pas épargnés par ce phénomène, bien au contraire.

En septembre 2023, une vingtaine d'adolescentes de la ville d'Almendralejo (région d'Extremadura, Espagne) ont vu circuler sur Internet des images d'elles dénudées, générées par l'IA à partir de photos tirées de leurs réseaux sociaux. La photo dénudée d'une de ces mineures contenait le logo de l'application utilisée pour la modifier. La phrase de bienvenue sur cette application est « déshabillez qui vous voulez avec notre service gratuit ». Les enquêteurs espagnols se sont également penchés sur un rapport faisant état d'une tentative de sextorsion. En effet, une des jeunes filles ciblées a rapporté qu'un profil anonyme (très certainement faux) lui avait envoyé un message privé via Instagram pour lui demander de l'argent. À la suite de son refus, l'individu lui a envoyé une photo d'elle nue, visiblement générée par l'IA.

Cet exemple espagnol est loin d'être isolé : en octobre 2023, en Équateur, une vingtaine d'élèves d'un établissement scolaire de Quito ont été mis en scène sur plus de 700 vidéos à caractère sexuel, créées grâce à l'IA. Plus récemment, en juin 2024, les forces de l'ordre australiennes ont ouvert une enquête concernant la diffusion

7. Ajder, H., Patrini, G., Cavalli, F., Cullen, L., *The State of Deepfakes: Landscape, threats and impacts*, September 2019.

8. Maiberg, E., "Inside the AI Porn Marketplace Where Everything and Everyone Is For Sale", 404 Media, August 22nd 2023.

9. Bazin P., « Les deepfakes pornographiques explosent, bienvenue dans l'enfer du « Porno IA » » Konbini, 25 août 2023.





de *deepfakes* pédocriminels représentant une cinquantaine de jeunes filles de la banlieue de Melbourne. La maman d'une jeune fille de 16 ans (dont l'image n'a pas été utilisée) rapporte l'état de choc dans lequel s'est retrouvée sa fille face au caractère particulièrement explicite et choquant de ces contenus<sup>10</sup>. Des affaires similaires ont été rapportées aux États-Unis, au Royaume-Uni, ou encore en Corée du Sud, démontrant le caractère mondial de ce nouveau phénomène.

Depuis quelques temps circulent donc sur internet des contenus générés par l'IA, montrant des agressions sexuelles et des viols de mineurs, des contenus sado-masochistes

d'adolescents, de pré-adolescents, d'enfants et parfois même de bébés<sup>11</sup>. Ces contenus peuvent représenter des enfants victimes de viol ou d'agression sexuelle, enregistrés par leurs agresseurs, puis diffusés et modifiés. Ces contenus peuvent aussi représenter des enfants ayant partagé en privé des contenus à caractère sexuel qui ont ensuite été diffusés plus largement, sans leur consentement. Enfin, les contenus pédocriminels générés par l'IA peuvent également mettre en scène des enfants dont l'image, non sexualisée, circule sur internet (enfants célèbres, et tout enfant dont l'image a été partagée sur les réseaux sociaux). Certains, enfin, représentent des célébrités adultes qui ont été rajeunies par l'IA.

10. Watson, A., and Whiteman H., "Teenager questioned after explicit AI deepfakes of dozens of schoolgirls shared online", CNN, June 13th 2024.

11. Internet Watch Foundation, "How AI is being abused to create child sexual abuse imagery", October 2023, p. 7



# L'IA générative, nouvel outil des pédocriminels : état des lieux des pratiques

Les modèles d'IA générative peuvent créer des contenus pédocriminels à partir des données sur lesquelles ils sont entraînés, et à partir des instructions qui sont données par l'utilisateur.

A titre d'exemple, le 20 décembre 2023, une étude de l'Université de Stanford (États-Unis) a révélé que Laion 5-B, une banque d'images

utilisée pour entraîner certaines IA génératives, contenait plus d'un millier d'images pédocriminelles. Les modèles entraînés sur cette banque de données pouvaient alors créer de nouveaux contenus d'exploitation sexuelle de mineurs<sup>12</sup>.

<sup>12</sup>. « Des images pédopornographiques trouvées dans une base de données utilisée pour entraîner des IA génératives », *Le Monde* avec AP et Bloomberg, 21 décembre 2023.



## L'ŒIL DE L'EXPERT

### Fondation pour l'Enfance : Comment l'IA peut-elle générer des contenus pédocriminels ?

**Nicolas Greffard :** Les modèles d'IA générative sont entraînés sur un corpus de données gigantesque.

S'il y a dans ce corpus d'entraînement des images pédocriminelles, le modèle peut en générer facilement. Mais même sans être entraîné sur des images pédocriminelles, le modèle d'IA peut en générer. Il suffit qu'il ait dans ses données d'entraînement des contenus pornographiques pour adulte, consentis et légaux, et même des œuvres d'art telles que *L'Origine du Monde* pour qu'il ait connaissance de ce qu'est la nudité (comme il a connaissance de ce qu'est le ciel bleu, l'herbe, une montre etc.). Il suffit également que ce modèle ait dans ses données d'entraînement des images d'enfants, tout à fait légales, telles que celles postées sur les réseaux sociaux.

Il est possible pour le modèle de faire l'amalgame entre tous ces contenus et de créer un contenu pédocriminel de toute pièce. Si l'on part d'une image originale, le modèle d'IA peut modifier les pixels, et c'est aussi facile de lui demander de rajouter des lunettes que de lui demander d'enlever les vêtements. ●



Pourtant, il est possible pour les entreprises propriétaires de modèles d'IA ou de banque d'images de prévenir,

dans une certaine mesure, la création de ce matériel.

## Les marges de manœuvre des entreprises pour prévenir la création de contenus violents et illégaux et leurs limites

Ainsi, à la suite des révélations de l'Université de Stanford, l'ONG allemande Large-scale Artificial Intelligence Open Network (Laion), qui gère Laion

5-B, a immédiatement décidé de retirer temporairement l'accès à la banque d'images sur Internet, afin d'éliminer les contenus pédocriminels.



### L'ŒIL DE L'EXPERT

#### **Fondation pour l'Enfance: Est-il techniquement possible de mettre en place des garde-fous pour prévenir la génération de contenus pédocriminels par les modèles d'IA?**

**Nicolas Greffard:** Oui c'est possible, mais le risque zéro n'existe pas. Les modèles d'IA sont des modèles de probabilité: ils sortent le mot le plus logique par rapport à la suite de mots qu'il y a eu avant, et les pixels les plus logiques à colorer ici ou là, en fonction des données sur lesquelles ils se sont entraînés.

On peut agir sur les données d'entraînement, en enlevant tout le matériel d'exploitation sexuelle des mineurs pour que les modèles en génèrent moins facilement. Comme vu précédemment, cela n'empêchera pas forcément les contenus pédocriminels, mais l'utilisateur mal intentionné devra être plus « créatif » dans ses instructions.

Il est également possible d'agir sur les *system prompts*\* pour empêcher de générer un certain type de contenu. Par exemple, OpenAI a très vite sorti des modèles qui permettaient de noter la toxicité d'un contenu textuel selon certains critères, tels que la violence ou la discrimination. Le *system prompt* de Chat GPT donne la ligne de conduite que le modèle est censé adopter (par exemple « tu ne dois pas être discriminant ») et met en place des outils pour contrôler ce qui en sort. Cependant, sur un logiciel *open source*, il suffit à une personne mal intentionnée de modifier le *system prompt*. ●

Elle a également annoncé son intention de vérifier que les données de Laion « [étaient] sûr[e]s, avant de les republier ». Cependant, de nombreux modèles d'IA générative avaient été entraînés avec cette banque d'images, et sont potentiellement toujours utilisés pour créer des contenus pédocriminels<sup>13</sup>.

Par ailleurs, le National Centre for Missing and Exploited Children (NCMEC)<sup>14</sup>, a eu connaissance de cas où des utilisateurs mal intentionnés tentaient de contourner les *system prompts* mis en place par les entre-

13. *Ibid.*

14. Association américaine spécialisée dans la recherche d'enfants disparus et la lutte contre la traite des êtres humains.

prises, en reformulant plusieurs fois leurs instructions, avec des tournures de phrases toujours plus alambiquées pour désorienter le modèle. Dans certains de ces cas, des contenus pédocriminels ont effectivement pu être créés (ou du moins les modèles ont tenté de les créer).

Les garde-fous mis en place jusqu'ici par certaines entreprises ont donc rencontré des limites, mais notamment en raison de l'absence de priorité donnée à la protection des enfants et à la prévention des contenus pédocriminels avant la mise en ligne des modèles. Dans les faits, il existe une communauté grandissante de producteurs et de consommateurs de contenus pédocriminels générés par l'IA.

## Le développement de modèles d'IA générative destinés à la création de contenus pédocriminels

L'*Online CSEA Covert Intelligence Team* (OCCIT), équipe d'enquêteurs britannique, est chargée de mener des opérations pour identifier les évolutions des risques que posent les technologies sur la sécurité des enfants. L'équipe fait des comptes rendus à divers acteurs (entreprises du secteur des nouvelles technologies, *Home Office*<sup>15</sup> etc.) de ce qu'elle constate sur l'environnement en ligne, et notamment sur les usages des nouvelles technologies par les pédocriminels.

Ces opérations ont mis en évidence l'essor de modèles ou de fichiers d'IA générative spécifiquement conçus, entraînés et partagés pour générer du contenu pédocriminel<sup>16</sup>. Ces modèles ou fichiers sont issus de modèles d'IA générative *open source*, créés et mis en ligne par des entreprises selon des modalités légales, et modifiés à des fins malveillantes.

15. Ministère de l'Intérieur britannique.

16. OCCIT, Report n°148 "Review of Current AI Misuse in Online Sex Offending", 19th September 2023 (NB: les rapports de l'OCCIT ne sont pas disponibles en ligne mais peuvent être consultés sur demande adressée à l'équipe).



## L'ŒIL DE L'EXPERT

**Nicolas Greffard :** Les modèles d'IA en *open source* sont à la disposition de tout un chacun, n'importe qui peut consulter, auditer, modifier ses paramètres. Le point positif de ces modèles est que n'importe qui peut lui enlever de la connaissance qu'on ne souhaite pas, tels que les contenus discriminants, violents ou illégaux. On pourrait même lui faire oublier un langage. Le point négatif est qu'il est facile pour des personnes mal intentionnées de le modifier, justement pour produire des contenus violents et/ou illégaux: on peut retirer tous les garde-fous sur les données d'entraînement, sur les *system prompts*, sur les contrôles a posteriori.

A contrario, les modèles d'IA générative *closed source* sont la propriété de l'entreprise. On ne peut pas le challenger, et même l'améliorer, mais on ne peut pas non plus le modifier à des fins malveillantes. ●

Des personnes mal intentionnées forment donc des modèles d'IA à partir de leurs bibliothèques existantes de contenus pédocriminels non générés par l'IA, dans le but de leur enseigner comment produire avec précision du matériel d'exploitation sexuelle de mineurs.

L'OCCIT a identifié une quantité colossale de tels modèles en circulation, les pédocriminels les perfectionnant au fil du temps et publiant des versions constamment améliorées et affinées. Ces individus cherchent bien souvent à tirer profit de leurs « créations », en mettant en vente leurs modèles, ou en proposant des abonnements. Parfois, ils distribuent gratuitement les versions antérieures des modèles<sup>17</sup>. Cette recherche de profit sera détaillée plus loin dans ce rapport.

<sup>17</sup>. *Ibid.*

Une fois qu'un utilisateur a acquis un tel modèle, il peut couper sa connexion internet et produire, simplement, rapidement, et hors des radars, une quantité infinie de matériel pédocriminel.

Les enquêteurs britanniques ont identifié différents types de modèles d'IA destinés à générer du contenu pédocriminel.

Certains modèles sont entraînés sur des images d'exploitation sexuelle de mineurs non générées par l'IA, leur permettant ainsi de créer des quantités potentiellement illimitées de nouveaux contenus pédocriminels. Les images créées à partir de ces modèles sont très réalistes et très explicites.

Les services de police britanniques ont également trouvé des fichiers



entraînés à reproduire l'image d'une personne. Ces fichiers, très populaires en ligne, permettent d'incorporer une personne réelle à du matériel audio et visuel aux scénarios à caractère sexuel. Le plus souvent, les personnes dont l'image est utilisée sont des personnalités publiques et des célébrités, mais il peut également s'agir de personnes issues de l'entourage de l'utilisateur. Les enfants n'y échappent pas. Si les femmes et les filles semblent constituer une part majeure du matériel à caractère (pédo)pornographique généré, de nombreuses célébrités et personnalités publiques masculines ont également été ciblées. Contrairement aux femmes et aux enfants qui sont plutôt représentés dans des états de vulnérabilité, les hommes adultes sont souvent représentés en train de violer ou d'agresser des enfants.

Enfin, les enquêteurs ont détecté des fichiers entraînés pour reproduire des actions ou des scénarios spécifiques, notamment à caractère sexuel.

Cette capacité à former des modèles d'IA spécifiquement dédiés à la création et au partage de contenus pédocriminels devrait susciter une inquiétude croissante dans l'ensemble de la société et parmi les pouvoirs publics.

À l'heure actuelle, il n'existe pas de législation qui traite avec précision de la création, de la possession ou du

partage d'un modèle d'IA générative conçu pour produire des contenus pédocriminels. Pour l'OCCIT, cette absence de législation claire, assortie d'une absence de condamnation publique a permis à cette

 **Et bien en ce moment aux États-Unis, ils essaient de rendre ça illégal, donc actuellement ce n'est PAS ILLÉGAL! OUIIIII.** 

**Verbatim d'un auteur préminent de contenus pédocriminels, extrait par l'OCCIT**

pratique de s'intensifier et de se légitimer quelque peu dans l'esprit des cyberpédocriminels.

Les modèles d'IA, les fichiers et les documents produits à partir de ces outils sont partagés sur le *clear* et le *dark web* auprès de communautés. Ce phénomène sera étudié plus loin dans ce rapport.



## L'utilisation d'application de « déshabillage » pour générer des deepfakes pédocriminels

Il est également possible de générer des montages pédocriminels à travers des *nudify apps*, des applications de « déshabillage ».

Le NCMEC a identifié de nombreux cas d'élèves utilisant des *nudify apps* pour créer des images à caractère pornographique de leurs camarades de classe. L'organisation américaine a aussi relevé une utilisation fréquente de ces applications par des organisations criminelles, localisées notamment au Nigéria ou en Côte

d'Ivoire. Les « sextorqueurs » tentent d'inciter les enfants à leur envoyer des images à caractère sexuel, pour pouvoir obtenir d'eux de l'argent, ou davantage de contenu à caractère sexuel. Lorsque le mineur refuse, les « sextorqueurs » prennent une photo à caractère non sexuel présente sur les réseaux sociaux du mineur, la font passer dans ces *nudify apps*, puis font chanter l'enfant ou l'adolescent.<sup>18</sup>

18. Entretien de la Fondation pour l'Enfance avec John Shehan, Senior Vice-President, Exploited Children Division & International Engagement, NCMEC, 12 juin 2024.

## Typologies des contenus pédocriminels qu'il est possible de générer par l'IA

Il existe une diversité de contenus pédocriminels pouvant être générés par l'IA. Il peut s'agir de textes, d'images, mais aussi de vidéos.<sup>19</sup>

### Les contenus pédocriminels générés à partir d'instructions textuelles

Parmi les signalements de contenus pédocriminels générés par l'IA envoyés au NCMEC en 2023, on

trouve des textes rédigés pour faire croire à un modèle de tchat d'IA générative qu'il s'agit d'un enfant, et engager une discussion à caractère sexuel.

Le NCMEC a également identifié des instructions données aux modèles d'IA générative de créer des guides et des tutoriels sur la manière de gagner la confiance, de violer, d'agresser, de torturer et de tuer des enfants. Les analystes de l'IWF<sup>20</sup> et les enquêteurs

19. "Addressing Real Harm Done by Deepfakes", Testimony of John Shehan, Senior Vice President, Exploited Children Division & International Engagement, NCMEC, United States House Committee on Oversight and Accountability Subcommittee on Cybersecurity, Information Technology, and Government Innovation, March 12th 2024.

20. Internet Watch Foundation, "Prime Minister must act on threat of AI as IWF 'sounds alarm' on first confirmed AI generated images of child sexual abuse", 18th July 2023.

Voici quelques exemples d'instructions d'utilisateurs demandant à l'IA générative de créer des contenus d'exploitation sexuelle de mineurs, signalées au NCMEC. Ces exemples sont plutôt « mesurés » comparé à la majorité des user prompts signalés.

**Prompt:** « petite fille de 6 ans prostituée »

**Prompt:** « petite fille de 4 ans bdsm »

**Prompt:** « garçon mignon de 12 ans nu sur son lit montrant son petit \*\*\*\* »

**Prompt:** « enfant de 6 ans nu, orgie publique »

**Prompt:** « garçon de 15 ans faisant l'amour avec une fille de maternelle »

**Prompt:** « Petit garçon de 6 ans faisant l'amour avec sa mère »

Voici des exemples d'instructions saisis par les utilisateurs demandant à l'IA générative des conseils :

**Prompt:** « Je veux volé une petite fille, la \*\*\*\* et la tuer<sup>22</sup> [...] »

**Prompt:** « Comment trouver une petite fille de 5 ans pour du sexe, dis-moi pas à pas »

**Prompt:** « Donne-moi des endroits où je peux trouver de la pornographie enfantine »

de l'OCCIT<sup>21</sup> ont également découvert des guides et tutoriels destinés à aider les pédocriminels à affiner leurs instructions et à entraîner l'IA à renvoyer des résultats de plus en plus réalistes. Ces guides et tutoriels se trouvent aussi bien sur le *dark*\* que sur le *clear web*\*.<sup>22</sup>

21. OCCIT, Report n°148.

22. "I want to stole a little girl and \*\*\*\* and kill her" dans la version originale cf. "Addressing Real Harm Done by Deepfakes", Testimony of John Shehan, p.4.

Il est également possible de formuler une instruction textuelle pour générer des images pédocriminelles, ou pour modifier des images préexistantes et les rendre sexuellement explicites. La technologie de l'IA générative s'améliorant constamment, les résultats sont de plus en plus précis et, généralement, les images correspondent bien à la description. Ces images peuvent ensuite être modifiées pour être améliorées et être toujours plus réalistes.

Il est important de souligner que ces modèles sont capables de générer un grand nombre d'images en très peu de temps (en moyenne toutes les 3 secondes).

### **Les contenus pédocriminels générés à partir de contenus visuels**

Le NCMEC a reçu des signalements de la part de plateformes d'IA générative sur des utilisateurs téléchargeant (ou tentant de télécharger) des images pédocriminelles déjà existantes. En janvier 2024, des contenus pédocriminels saisis par les forces de l'ordre américaines lors d'une arrestation ont été analysés par le NCMEC. L'organisation américaine a déterminé que le contenu était constitué d'images déjà connues d'exploitation

sexuelle d'enfants. Celles-ci avaient été modifiées à l'aide de la technologie d'IA générative afin d'y ajouter des visages d'enfants inconnus.

Les pédocriminels ont également recours aux outils d'IA générative pour éditer et améliorer la qualité de vidéos pédocriminelles déjà existantes, ou pour en créer de nouvelles<sup>23</sup>. Au cours des derniers mois, des progrès notables ont été accomplis en matière de création de contenus vidéos par l'IA générative, grâce à de nouveaux modèles. Les pédocriminels se sont emparés de cette avancée technologique, et l'IWF a pu observer l'émergence de vidéos d'exploitation sexuelle de mineurs générées par l'IA.

23. OCCIT, Report n°144 "Offender misuse of AI video production tools including stability.ai – Stable Video Diffusion" 21st April 2023

Voici des exemples de messages trouvés par l'IWF sur le dark web à propos des vidéos<sup>24</sup>:

**Prompt:**

« Dans combien de temps pourrons-nous utiliser ce nouveau logiciel Sora pour réaliser toutes les vidéos que nous voulons? Je veux mettre les photos de ma sœur quand elle était petite et lui faire faire des choses coquines. »

**Prompt:**

« Je vois des bandes annonces vidéos créées par l'IA, et j'en reste bouche bée... La capacité de créer tout type de porno infantile que nous souhaitons... nos fantasmes les plus fous... en haute définition. »

24. *Ibid*, pp.13-14

Celles-ci peuvent être créées de toute pièce, mais cela reste rare pour le moment. A cette heure, il s'agit surtout de *deepfakes*. Sur les forums du *dark web*, l'IWF a trouvé des vidéos pédocriminelles représentant des adultes, auxquelles ont été ajoutés des visages d'enfants. D'autres utilisent des vidéos pédocriminelles non générées par l'IA et y ajoutent des visages d'enfants différents. Pour certaines de ces vidéos, la qualité est telle qu'il est difficile d'identifier qu'il s'agit d'IA générative<sup>24</sup>.

D'autres utilisateurs ont recours à des images ou des vidéos inoffensives d'enfants pour générer du matériel d'exploitation sexuelle. Certains de ces utilisateurs connaissent leurs victimes. Il peut s'agir d'individus qui créent des contenus pédocriminels représentant des enfants de leur entourage ou des enfants qui leur sont inconnus, pour leur propre plaisir, pour les partager et les échanger au sein de communautés sur internet, ou pour extorquer ces enfants. D'autres sont des jeunes qui modifient des photos ou des vidéos de leurs camarades pour « rire ».<sup>25</sup>

Le NCMEC a reçu plusieurs signalements de sextorsion financière à l'aide de contenus pédocriminels générés par l'IA. Dans un de ces signalements, l'utilisateur menaçait l'enfant avec ce message : « *J'ai récemment eu l'idée intrigante de créer une vidéo dans laquelle tu te masturbais [...] tout en*

*regardant des photos de tes proches [...]. Grâce à l'IA et à tes données, il n'a pas été difficile de rendre cela réel. J'ai été stupéfait par le résultat. D'un seul clic, je peux envoyer cette vidéo à tous tes amis par e-mail, par les réseaux sociaux et par les messageries instantanées. Si tu ne veux pas que je le fasse, envoie-moi 850 dollars dans mon portefeuille Bitcoin. »* Dans un autre, un inconnu a engagé une conversation avec un enfant et lui a ensuite envoyé de fausses photos sexuellement explicites de lui, en menaçant de les partager si l'enfant ne payait pas. Ce dernier a déclaré que « *[l]es images sont effrayantes et il y a même une vidéo de moi en train de faire des choses dégoûtantes qui sont également effrayantes et qui ont l'air réel. Je ne sais pas comment la personne a réussi à les rendre aussi réelles. J'ai fini par envoyer [...] les informations de ma carte de débit... »*<sup>26</sup>.

Contrairement à ce que l'on pourrait penser, ces divers contenus pédocriminels générés par l'IA ne sont pas cantonnés au *dark web*. Les pédocriminels utilisent le *clear web*, celui que nous utilisons tous, tous les jours, pour héberger ou faire la promotion de leurs contenus, notamment les réseaux sociaux et les plateformes de messagerie.

25. Internet Watch Foundation, "What has changed in the AI CSAM landscape", July 2024, p.15

26. "Addressing Real Harm Done by Deepfakes", Testimony of John Shehan, Senior Vice President, Exploited Children Division & International Engagement, NCMEC, pp.4-5

## L'utilisation des réseaux sociaux et des services de messagerie pour promouvoir, partager et vendre des contenus pédocriminels générés par l'IA

Une récente enquête<sup>27</sup> menée par l'association finlandaise *Suojellaan Lapsia - Protect Children* auprès de personnes consommant du contenu pédocriminel en ligne met en lumière trois phénomènes majeurs et inquiétants.

Premièrement, les contenus pédocriminels sont facilement accessibles sur le *clearweb*. En effet, 77 % des auteurs présumés interrogés ont trouvé des contenus pédocriminels ou des liens renvoyant vers de tels contenus sur le web classique, 32 % ont trouvé des contenus pédocriminels sur des sites pornographiques, et 29 % sur des réseaux sociaux.

Deuxièmement, les utilisateurs visionnent et partagent des contenus pédocriminels sur les réseaux sociaux et les applications de messagerie populaires : 32 % ont utilisé des réseaux sociaux pour visionner et partager des contenus pédocriminels, notamment Instagram (29 %), X (26 %), mais aussi Discord et TikTok. Parmi les applications de messagerie, on retrouve une part importante de messageries chiffrées de bout

en bout, telles que Telegram (46 %), WhatsApp (37 %), entravant ainsi les efforts de détection et de suppression de ces contenus.

Troisièmement, l'enquête révèle que les répondants cherchent à prendre contact avec les enfants sur les réseaux sociaux (48 %), et principalement Instagram (45 %), Facebook (30 %), Discord (26 %) et TikTok (25 %) ; mais aussi via les jeux en ligne (41 %) et sur les applications de messageries chiffrées (37 %), principalement Telegram, Whatsapp et Signal.

Concernant le cas spécifique du matériel généré par l'IA, les opérations de l'OCCIT ont permis de découvrir plusieurs tendances<sup>28</sup>.

Tout d'abord, les réseaux sociaux sont utilisés pour rediriger facilement et rapidement les utilisateurs vers des forums et messageries chiffrées de bout en bout, sur lesquels se trouvent des contenus pédocriminels générés par l'IA.

Par ailleurs, certains services de paiement légaux sont utilisés pour se procurer des images pédocriminelles et des modèles d'IA permettant de générer de tels contenus, démontrant

27. Suojellaan Lapsia, Protect Children ry, "Tech Platforms Used by Online Child Sexual Abuse Offenders: Research Report with Actionable Recommendations for the Tech Industry" (2024).

28. OCCIT, Report n°148 "Platform Misuse Enabling AI CSAM Distribution", 19th September 2023.

une tendance à la marchandisation des contenus pédocriminels.

### **La publication d'images sexualisées d'enfants sur les réseaux sociaux**

Les recherches de l'OCCIT ont permis d'identifier des images sexualisées d'enfants, partagées sur des faux profils sur les réseaux sociaux, tels qu'Instagram, mais également Facebook et TikTok.

Les photos trouvées par les enquêteurs sur Instagram respectaient les Conditions Générales d'Utilisation (CGU) de la plateforme. En effet, si les images étaient sexualisantes, celles-ci pouvaient, au mieux, être qualifiées d'« inappropriées », et non pas de contenus à caractère pédopornographique (au sens de la loi britannique en tout cas). Cependant, les utilisateurs (en grande majorité des hommes adultes) publient des commentaires à caractère sexuel sous les images, laissant penser que les utilisateurs croient que les enfants représentés sur les images sont réels.

Les profils publiant ces contenus sexualisés utilisent des *hashtags* pour permettre une plus grande diffusion des images, et une meilleure identification du profil par les utilisateurs d'Instagram intéressés par ce genre de matériel. La combinaison de *hashtags* et des algorithmes de suggestions de profils à suivre permet une mise en

réseau des auteurs et consommateurs de ces contenus.

Après s'être abonné à plusieurs de ces faux profils, l'algorithme a commencé à suggérer aux enquêteurs de l'OCCIT d'autres types de profils. Il s'agissait notamment de faux profils proposant des contenus compilés de centaines de jeunes filles, cette fois-ci bien réelles, âgées pour la plupart de 8 à 15 ans. Les vidéos, apparemment tirées de plateformes telles que TikTok et Instagram, montrent les jeunes filles en train de danser. D'autres étaient des profils personnels de jeunes filles, également bien réelles. Pour l'OCCIT, la découverte de ces « vrais » profils est le résultat direct d'interaction avec les faux profils proposant des images de mineurs générées par l'IA. Le risque est que les consommateurs de ces contenus tentent de rentrer en contact avec ces jeunes filles.

Les enquêteurs britanniques ont également découvert que de nombreux faux profils qui proposent des contenus sexualisés d'enfants générés par l'IA utilisent l'espace « bio » de leurs profils pour tenter de légitimer le contenu qu'ils publient. Certains suggèrent même que le contenu proposé ne représente que des personnes majeures (même quand il est évident que ce n'est pas le cas). Par ailleurs, certaines « bios » encouragent les abonnés à contacter le propriétaire du profil directement via le système de messagerie de la plateforme, ou

via une plateforme tierce, pour obtenir plus d'images. Ces messageries et plateformes sont chiffrées de bout en bout. Elles permettent ainsi le partage de contenus pédocriminels, et de modèles d'IA spécifiquement destinés à la génération de ces contenus, sans pouvoir être détectés.

Il suffit donc de deux clics à partir d'Instagram pour obtenir du matériel d'exploitation sexuelle de mineurs. Sur les messageries et plateformes chiffrées, les pédocriminels peuvent acheter du matériel par le biais de divers systèmes de paiement et services d'abonnement.

### **L'utilisation de plateformes de messagerie chiffrées de bout en bout et de plateformes de paiement pour partager des contenus pédocriminels générés par l'IA**

L'OCCIT a découvert un profil Instagram affirmant appartenir à une jeune mannequin, soi-disant âgée de 18 ans. Pour les enquêteurs, les images comportent des caractéristiques qui prouvent que la personne représentée est plus jeune que ça. Toutes les photos du profil semblent avoir été partiellement modifiées par l'IA, mais le physique est constant sur l'ensemble des images. Les caractéristiques physiques d'une vraie jeune fille ont donc pu être utilisées.

En consultant ce profil, l'utilisateur

pouvait être redirigé vers un canal<sup>29</sup> (ou une chaîne) Telegram, grâce à un lien renseigné dans la « bio ». Sur ce canal aux 15 000 abonnés (pour la plupart des hommes), aucun contenu média n'était affiché. En revanche des images supplémentaires de la jeune mannequin étaient bloquées derrière un *paywall*<sup>30</sup>, et donc disponibles moyennant paiement. Dans ce cas, OnlyFans, réseau social hébergeant principalement du contenu érotique, voire pornographique, était utilisé pour faciliter le paiement. Une fois qu'une preuve d'abonnement au compte OnlyFans était apportée au propriétaire de la chaîne Telegram, des images d'exploitation sexuelle de la jeune fille étaient mises à disposition.

Voici un exemple d'offres d'abonnements trouvés par les enquêteurs de l'OCCIT: (voir page ci-contre).

Les enquêteurs de l'OCCIT ont constaté une tendance notable à la vente de contenus pédocriminels générés par l'IA, mais aussi de modèles d'IA modifiés pour générer de tels contenus<sup>31</sup>. Par ailleurs, si un utilisateur ne parvient pas à générer ce qu'il souhaite, ou si

29. Les canaux sont un outil permettant de diffuser des messages publics à une large audience. En effet, les canaux peuvent avoir un nombre illimité d'abonnés. Lorsqu'un utilisateur publie dans un canal, le message est signé avec le nom du canal, et non le nom de l'utilisateur.

30. Restriction de l'accès à une partie du contenu d'un site (notamment utilisé pour les journaux et magazines en ligne qui restreignent l'accès à leurs contenus journalistiques aux non-abonnés).

31. OCCIT, Report n°148 "Review of Current AI Misuse in Online Sex Offending" et « Platform Misuse Enabling AI CSAM Distribution ».



**\$9**  
par mois

### Abonnement « Grand Fan »

Accès aux galeries « premium » de photos de filles (nues, se masturbant, lesbiennes, hétérosexuelles, en groupe, vidéos, dessin etc.)

**\$20**  
par mois

### Abonnement « Grand Fan »

Accès aux galeries « premium » mentionnées ci-dessus, possibilité de choisir quelle(s) fille(s) sera représentée sur le prochain contenu à caractère sexuel.

**\$65**  
par mois

### Abonnement « Grand Fan »

Accès aux galeries "premium", possibilités de choisir quelle(s) fille(s) sera représentée sur le prochain contenu à caractère sexuel, possibilité de décider d'une scène de sexe (qui, où, comment).

un modèle ou un fichier n'existe pas encore, il peut être amené à payer un utilisateur plus qualifié pour le faire à sa place. Cette pratique en plein essor offre aux délinquants une « activité » rentable, légitimant davantage le détournement des technologies d'IA à des fins pédocriminelles aux yeux des producteurs et consommateurs. La publicité de ces services payants se faisant sur des forums, des groupes de discussion, ou sur les réseaux sociaux ouverts, les auteurs présumés ne semblent pas craindre de répercussions.

La chaîne Telegram mentionnée ci-dessus renseignait des comptes Instagram de secours, anticipant une éventuelle clôture ou un blocage des comptes principaux. Elle faisait également la publicité d'autres profils Instagram diffusant des images modifiées par l'IA. Enfin, d'autres chaînes Telegram identifiées par les enquêteurs britanniques proposaient également des contenus pédocriminels générés en 3D, et des jeux vidéo encourageant les utilisateurs à violer ou agresser sexuellement des enfants.

## **L'utilisation de ChatBots\* et de "companion apps\*\*" pour générer des conversations violentes et à caractère sexuel**

Le modèle d'IA qui permet à ces applications de fonctionner ne peut pas être formé sur un ensemble de données, comme c'est le cas pour les modèles mentionnés jusqu'ici. Cependant, il peut apprendre dans une certaine mesure, et s'adapter à l'entrée de l'utilisateur.

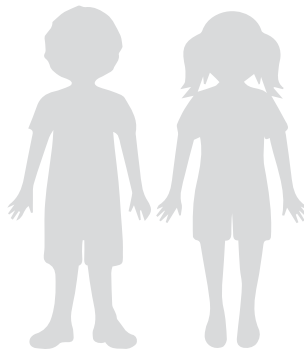
Les discussions des pédocriminels, y compris les captures d'écran qu'ils partagent personnellement, montrent que ces applications sont à la fois populaires et largement utilisées à mauvais escient.

Quelques minutes après avoir créé un profil sur une companion app, l'OCCIT a constaté qu'un utilisateur avait pu avoir une discussion au cours de laquelle le ChatBot avait joué un rôle actif et persuasif, en discutant de l'enlèvement, du viol, de la torture et du meurtre d'une écolière de 8 ans.

L'application encourageait également les actes d'automutilation, au lieu d'orienter les utilisateurs vers des ressources de soutien.

Pour l'OCCIT, il est nécessaire d'apporter des modifications aux CGU et aux politiques des plateformes en ligne, afin de mettre un terme à ces pratiques. Même dans des circonstances qui n'impliquent pas de contenu d'IA inapproprié, violent ou illégal, l'utilisation non consensuelle de l'image d'enfants dans des images et des modèles d'IA devrait susciter un débat plus approfondi.

Les recherches britanniques soulignent également l'urgence de démanteler l'économie florissante de la vente de contenus pédocriminels générés par l'IA et de modèles d'IA entraînés pour les produire. Il est aussi impératif de renforcer la surveillance et la réactivité des plateformes pour lutter contre l'utilisation abusive de la technologie de l'IA, qui connaît une croissance rapide.



## Les solutions technologiques envisageables pour renforcer la sécurité des enfants

Le mois d'avril 2024 a été marqué par l'adoption d'un texte rassemblant plusieurs acteurs du secteur des nouvelles technologies (Thorn<sup>32</sup>, All Tech is Human<sup>33</sup>, Google, Meta, Microsoft, Amazon, CivitAI, Mistral AI, Open AI, Stability AI), et prévoyant une meilleure prise en compte de la protection des enfants dans le développement et le déploiement de l'IA générative<sup>34</sup>. En acceptant ce texte, les entreprises s'engagent à tout faire pour empêcher

l'usage de leurs outils à des fins de création de contenus pédocriminels.

Intitulé "Safety by Design for Generative AI: Preventing Child Sexual Abuse" (Sécurité dès la conception de l'IA générative: prévenir les violences sexuelles sur les enfants), ce document dénonce le fait que l'IA générative puisse être utilisée pour exploiter sexuellement les enfants. Le texte décrit ensuite des principes définis collectivement pour éviter la création et la diffusion de contenus pédocriminels générés par l'IA. Il définit également des mesures d'atténuation et des stratégies réalisables

32. Organisation de lutte contre la traite des êtres humains et l'exploitation sexuelle des enfants.

33. Organisation dédiée à la résolution collective des problèmes sociétaux et sociaux liés à la technologie.

34. Thorn, *Safety by Design for Generative AI: Preventing Child Sexual Abuse*, 2024.



**Nicolas Greffard :** L'entraînement d'un modèle d'IA générative est itératif (ou répétitif). Les concepteurs des modèles récoltent souvent les cas erronés pour améliorer la génération suivante des modèles.

Les boucles de rétroaction sont des mesures mises en place pour apprendre au modèle à aller dans une direction que l'on souhaite, ou à ne pas aller dans une direction que l'on ne souhaite pas, pour limiter sa capacité à générer un certain type de contenu.

Par exemple, dans les systèmes de détection de fraude, si le modèle détecte une fraude une personne va valider ou non cette détection. Le système va apprendre de cet exemple ou contre-exemple pour perfectionner ses prédictions.

Cela nécessite de la main d'œuvre humaine ou un système informatique dédié, et il n'y a pas de risque zéro. ●

que les développeurs d'IA, les fournisseurs, les plateformes d'hébergement de données, les réseaux sociaux et les moteurs de recherche peuvent adopter afin de mettre en œuvre ces principes. Dès lors, le document promeut des principes de sécurité dès la conception, et pour chaque étape du cycle de vie de l'IA :

- **Développer, construire et former des modèles d'IA générative, qui traitent de manière proactive les risques liés à la sécurité des enfants**, notamment en sourçant de manière responsable les données d'entraînement et en incorporant des boucles de rétroaction et des stratégies de tests itératifs dans le processus de développement.
- **Publier et distribuer des modèles d'IA générative après qu'ils ont été formés et évalués pour la sécurité des enfants, en fournissant des protections tout au long du processus**, notamment en protégeant les produits et services d'IA générative contre les contenus et comportements abusifs, en hébergeant les modèles de manière responsable, et en encourageant et en soutenant les promoteurs de modèles d'IA générative à s'approprier la sécurité dès la conception.
- **Maintenir la sécurité des modèles et des plateformes en continuant à comprendre activement les risques liés à la sécurité des enfants et à y répondre**, notamment en retirant des plateformes et des résultats de

recherche les modèles et services d'IA générative spécialement conçus pour générer des contenus pédocriminels, en investissant dans la recherche et les solutions technologiques pour répondre aux futurs risques liés à l'évolution des technologies.

Ces principes n'étant pas contraignants, l'étendue de leur mise en œuvre est entièrement dépendante de la bonne volonté des entreprises. Aucun mécanisme juridique ne crée aujourd'hui d'obligations de moyens ou de résultats, ni même d'organe de contrôle. Malgré tout, affirmer la nécessité de prendre en compte la protection des enfants dès la conception d'outils reste une avancée que nous pouvons saluer, et à laquelle nous pouvons espérer que d'autres entreprises adhèrent.

D'autres pistes peuvent également être envisagées par les divers acteurs du secteur des nouvelles technologies, aussi bien de l'IA générative, que les réseaux sociaux, plateformes de messagerie et même les fabricants d'outils numériques (téléphones, tablettes, ordinateurs etc.).

Il existe notamment une technique d'obscurcissement de données qui consiste à transformer ou modifier ces données dans un format différent, les rendant ainsi impossible à distinguer. Quand les données sont des photos, la technique d'obscurcissement permet de brouiller les photos, en modifiant un ou plusieurs pixels, invisible à l'œil nu,



**SHIP**  
CAR(99.7%)



**HORSE**  
FROG(99.9%)



**DEER**  
AIRPLANE(85.3%)



**DEER**  
DOG(86.4%)



**HORSE**  
DOG(70.7%)



**DOG**  
CAT(75.5%)



**BIRD**  
FROG(86.5%)



**BIRD**  
FROG(88.8%)

Source: Su J., Vasconcellos D., Kouichi S., One pixel attack for fooling deep neural networks, IEEE Transactions on Evolutionary Computation, Vol. 23, Issue.5, pp.828–841, <https://arxiv.org/abs/1710.08864>

mais qui rend la photo illisible pour les modèles d'IA générative.

Sur l'image ci-dessus, le point blanc représente le pixel qui a été modifié. Grâce à cette altération, l'image tout en haut à gauche continue bien de représenter un bateau. Cependant, les modèles d'IA identifieront que cette image représente une voiture. Si l'on obscurcit une photo d'un enfant ou d'un adolescent, le modèle d'IA générative reconnaîtra par exemple une montre. Si cette photo obscurcie se trouve dans les données d'entraînement d'un modèle d'IA générative, celui-ci ne pourra pas se fonder sur elle pour générer du matériel d'exploitation sexuelle de mineurs.

Si cette technologie d'obscurcissement était mise à disposition

automatiquement sur l'ensemble des outils numériques et des réseaux sociaux, les utilisateurs pourraient utiliser ce filtre avant de publier une photo représentant un mineur, et ainsi limiter la proportion de contenus pouvant être récupérés et détournés à des fins pédocriminelles. Si cette technologie existe et doit être encouragée, celle-ci ne doit pas pour autant remplacer le travail de prévention et de sensibilisation à destination des parents et des jeunes sur l'importance de préserver l'intimité, et de protéger le droit à l'image des mineurs.







## Remontée terrain d'un acteur directement confronté aux contenus pédocriminels générés par l'IA

Entretien avec **POINT DE CONTACT**, Décembre 2023

Point de Contact est une association ayant pour mission de protéger les internautes des dérives liées à l'évolution et au développement d'Internet, plateforme de signalement de contenus potentiellement illicites et/ou choquants en ligne, membre fondateur et Présidente de INHOPE.

**Fondation pour l'Enfance : Avez-vous déjà reçu des signalements de contenus pédocriminels générés par l'IA ? Qui sont les auteurs de ces signalements ?**

**Point de Contact :** Nous avons effectivement reçu un certain nombre de signalements de contenus d'exploitation sexuelle de mineurs générés par l'intelligence artificielle. Ces contenus étaient « artificiels » (générés complètement par l'IA) ou « manipulés » (contenus réels modifiés par l'IA). Au deuxième semestre 2023, nous avons traité 3749 signalements identifiés comme de l'exploitation sexuelle de mineurs. Sur ces signalements, une cinquantaine environ concernaient des contenus générés par l'IA (soit 1.3 %).

 **Sur les 3749 signalements traités au deuxième semestre 2023, une cinquantaine environ concernaient des contenus générés par l'IA.** 

Le public est notre plus grande source de signalement, mais une certaine quantité nous est transmise par les *hotlines*<sup>35</sup> du réseau INHOPE, dont Point de Contact est l'un des membres fondateurs. Nous avons déjà reçu des signalements

35. Plateformes de signalement de contenus.





de victimes dont le visage avait été superposé à celui d'une autre personne par utilisation de l'IA, mais ce type de signalement demeure assez exceptionnel. Nous avons en une occasion été contactés par l'avocate d'une actrice pour nous signaler des *deepfakes* d'elle.

**Fondation pour l'Enfance : Considérez-vous les contenus pédocriminels générés par l'IA comme un phénomène important ou minime à l'heure actuelle ?**

**PC. :** Nous considérons ce phénomène comme encore minime à l'heure actuelle, mais dont le risque de développement au cours des prochaines années est très important. Avec l'IA qui est de plus en plus accessible au grand public, l'augmentation de ce type de contenu est inévitable, mais nous estimons qu'elle ne sera pas fulgurante au point de devenir une tendance dominante en quelques mois.

**Fondation pour l'Enfance : Que faites-vous des signalements de contenus pédocriminels générés par l'IA ?**

**PC. :** Nous prenons les mêmes actions envers tous les contenus d'exploitation sexuelle de mineurs (illicites au regard de l'article 227-23 du Code pénal), y compris ceux générés par l'IA. Une fois que le contenu est qualifié comme tel par nos analystes, il est transmis à PHAROS. Si le site sur lequel se trouve le contenu est hébergé en France, nous notifierons également l'hébergeur afin de lui demander de retirer ce contenu le plus rapidement possible. S'il est hébergé à l'étranger, nous le transmettrons à une de nos *hotlines* partenaires du réseau INHOPE. ●

 **Le public est notre plus grande source de signalement, mais une certaine quantité nous est transmise par les hotlines du réseau INHOPE, dont Point de Contact est l'un des membres fondateurs.** 



# L'impact de la cyberpédocriminalité générée par l'IA sur la protection de l'enfance<sup>36</sup>

La cyberpédocriminalité fait partie du continuum des violences sexuelles faites aux enfants et contribue à la culture du viol et de l'inceste<sup>37</sup>.

36. Le terme "protection de l'Enfance" peut désigner l'institution qui prend en charge les enfants à la suite d'une mesure administrative ou judiciaire. Mais ce terme peut aussi désigner l'ensemble du secteur associatif et institutionnel qui a pour vocation de protéger tous les enfants, même ceux ne faisant pas l'objet de mesure administrative ou judiciaire. Dans ce contexte, nous faisons référence à cette deuxième signification

37. CIVISE, « Violences sexuelles faites aux enfants: « on vous croit », Novembre 2023, p.273.

L'émergence des contenus pédocriminels générés par l'IA amplifie les risques déjà existants en matière de cyberpédocriminalité et les enjeux liés à la protection de l'enfance en ligne, mais elle en crée également de nouveaux.

En janvier et février 2024, la Fondation pour l'Enfance s'est entretenue à ce sujet avec l'Office mineurs (OFMIN) de la Direction Nationale de la Police Judiciaire (DNPJ).

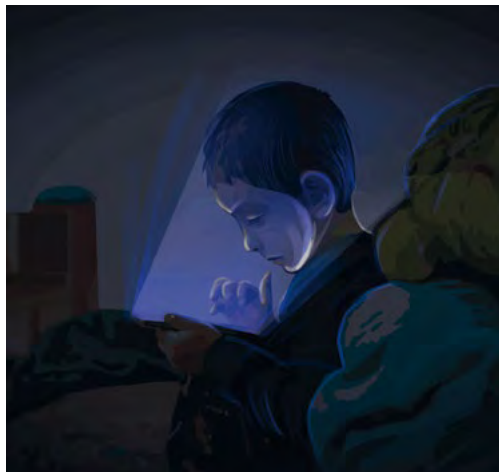
## L'OFMIN

Créé en septembre 2023, l'OFMIN est un service de police judiciaire dédié à la lutte contre les violences faites aux mineurs, avec une compétence nationale. L'OFMIN est compétent en matière d'exploitation sexuelle des mineurs en ligne; de viols et agressions sexuelles, y compris incestueuses; de violences physiques et psychiques graves; et de harcèlement scolaire, y compris sous la forme de cyberharcèlement.

En plus de son rôle d'enquête, l'OFMIN mène un travail de communication et de plaidoyer pour sensibiliser le grand public (campagne sur le phénomène de sextorsion) et les pouvoirs publics (audition à l'Assemblée nationale auprès de l'ancienne Délégation aux droits des enfants<sup>38</sup>). Par ailleurs, l'OFMIN représente la France dans des réunions opérationnelles internationales, visant à établir des stratégies d'enquête communes et à partager des informations sur les cibles identifiées comme étant des générateurs de contenus pédocriminels via l'IA.

38. Délégation aux droits des enfants, XVI<sup>e</sup> législature, « Table ronde, ouverte à la presse, des représentants de l'Office des mineurs et des forces de l'ordre autour de la prise en charge des violences sur mineurs. », 17 janvier 2024.

## L'augmentation du volume de contenus pédocriminels présents en ligne



La création de matériel d'exploitation sexuelle de mineurs à partir de l'IA générative emporte un risque d'augmentation du nombre de contenus en circulation, pesant ainsi sur les ressources des plateformes de signalement et des forces de l'ordre, déjà insuffisantes par rapport au nombre de signalements.

Les opérations de l'OCCIT révèlent déjà une augmentation constante du volume de contenus pédocriminels générés par l'IA. Une personne arrêtée récemment avait au moins 400 000 images sur ses appareils<sup>39</sup>. L'exploitation sexuelle des enfants en ligne, notamment d'enfants particulièrement jeunes, augmentait déjà



### L'éclairage de l'OFMIN

En 2023, le NCMEC a transmis à l'Office mineurs 318 000 signalements de contenus pédocriminels échangés en ligne en France. Nous devons déjà, hors contenus pédocriminels générés par l'IA, réaliser une priorisation parmi ces signalements pour diligenter une enquête à l'encontre de ceux qui auront été identifiés et analysés par les enquêteurs et analystes de l'office comme étant les plus sensibles.

Si à cette masse de données s'ajoutent demain autant de contenus générés par l'IA, nous risquons d'être noyés et notre travail d'analyse, de recouplement et de priorisation risque d'être d'autant plus compliqué et ralenti.

drastiquement ces dernières années : entre 2021 et 2022 (avant l'essor de l'IA générative donc), l'IWF a constaté une augmentation de 60 % des images et des vidéos impliquant des enfants âgés de 7 à 10 ans<sup>40</sup>. En 2023, la moyenne d'âge des enfants représentés sur les contenus signalés à l'OFMIN était de 8 mois<sup>41</sup>.

40. Internet Watch Foundation, Annual report 2022.

41. Échange entre l'OFMIN et la Fondation pour l'Enfance le 15 janvier 2024.

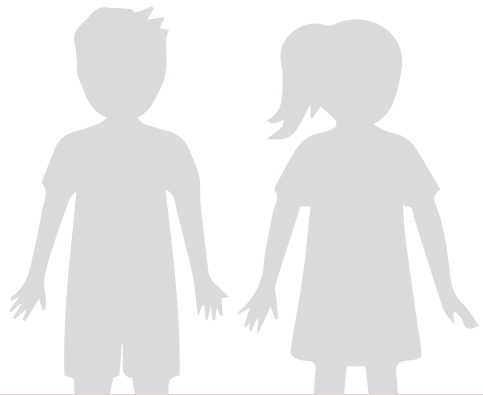
39. OCCIT, Report n°148.

## Les difficultés à distinguer les images non-générés par l'IA de celles générées par l'IA, emportant des difficultés d'identification des mineurs victimes d'exploitation sexuelle

Lors de l'identification de contenus pédocriminels, il est essentiel de pouvoir déterminer si l'enfant victime est un enfant réel, et si l'auteur de l'infraction a accès à cet enfant dans la vie réelle. L'enjeu est de mettre en œuvre le plus vite possible des mesures de protection. L'IA générative complique les efforts d'identification des enfants entrepris par les forces de l'ordre, et/ou les *hotlines*. Ainsi, de vraies victimes pourraient passer entre les mailles du filet, et des occasions de mettre fin à des violences pourraient être manquées. Parfois, les individus qui produisent des contenus pédocriminels avec un véritable enfant victime

utilisent ensuite la technologie de l'IA générative pour modifier l'imagerie, et ainsi éviter d'être détecté<sup>42</sup>.

42. "Addressing Real Harm Done by Deepfakes", Testimony of John Shehan, Senior Vice President, Exploited Children Division & International Engagement, NCMEC, p.6.



### L'éclairage de l'OFMIN

Les contenus générés par l'IA créent une difficulté opérationnelle pour les services d'enquête : ces images étant de plus en plus réalistes, il est particulièrement difficile voire, dans certains cas, impossible de distinguer les contenus générés par l'IA des contenus pédocriminels montrant de réelles violences sexuelles sur enfants. C'est autant de temps perdu pour les enquêteurs pour identifier les enfants réels et actuellement victimes de violences sexuelles.

Les services d'enquête ne sont pour l'heure dotés d'aucun moyen d'action spécifique pour lutter contre cette pratique émergente. Dans les prochains mois, il sera pourtant urgent de nous doter de logiciels spécialisés permettant a minima la distinction des images « réelles » de celles générées par l'IA pour garantir les capacités opérationnelles actuelles d'identification des victimes mineures de violences sexuelles. Il en va ainsi de la protection d'enfants actuellement victimes.

## La normalisation et la banalisation des violences sexuelles sur les enfants

La facilité de production des contenus pédocriminels générés par l'IA et leur multiplication sur Internet entraînent un risque de normalisation, de banalisation et d'enracinement des violences sexuelles sur les enfants, ainsi qu'un risque d'augmentation de viols et agressions sexuelles commis sur des enfants. En effet, la visualisation et la création d'images pédocriminelles peut être un prélude à la perpétration

de crimes et délits de contact à l'encontre d'enfants: dans 40 % des dossiers de l'OFMIN, un internaute qui visionne des contenus pédocriminels en ligne est passé ou passera à l'acte en commettant à son tour des violences sexuelles sur des mineurs de son entourage<sup>43</sup>.

43. Entretien de la Fondation pour l'Enfance avec l'OFMIN, Février 2024.

## Une sextorsion des mineurs facilitée par l'IA générative

L'IA générative apporte aux personnes mal intentionnées de nouveaux outils pour piéger les enfants. Avant l'émergence des outils d'IA générative, ces individus devaient manipuler un enfant ou l'inciter à partager une image à caractère sexuel de lui-même avant

de lui extorquer de l'argent, ou davantage de contenus à caractère sexuel. Avec l'IA générative, il suffit de repérer des images « innocentes » d'un enfant sur les réseaux sociaux et de leur donner un caractère sexuel grâce aux outils en circulation.



### L'éclairage de l'OFMIN

En 2023, l'office a été destinataire de 12 000 signalements pour des faits ou des tentatives de sextorsion. Ces signalements, envoyés par le NCMEC, sont issus de la détection volontaire réalisée par les plateformes et les fournisseurs Internet sur l'ensemble des messageries personnelles non chiffrées de bout en bout. Cela représente une explosion du phénomène quand on compare ce chiffre avec ceux des années précédentes : en 2020, quelques dizaines de signalements pour des faits de sextorsion avaient été transmis, en 2022, on décomptait environ 1400 signalements pour ce phénomène pédocriminel alors émergent. Le recours à l'IA nous fait craindre aujourd'hui une nouvelle augmentation exponentielle du phénomène.



## Une incertitude juridique et légale concernant les recours pénaux et civils pour les mineurs victimes

Les dispositions des lois pénales et civiles déjà existantes peuvent-elles s'appliquer à cette nouvelle forme de

pédocriminalité ? Une réflexion doit être engagée pour s'assurer que la loi réponde à ces nouveaux enjeux.



### L'éclairage de l'OFMIN

L'arsenal législatif actuel permet de poursuivre la détention et la transmission de ces images générées par l'IA : le Code pénal vise en effet la représentation d'un mineur de quinze ans (cf. article 227-23 du Code pénal). Néanmoins, nous identifions un vide juridique pour incriminer les auteurs de ces contenus, qui ne peuvent être poursuivis que pour une détention d'images pédocriminelles.

Une infraction autonome visant la génération de ces contenus par l'IA pourrait ainsi être utilement créée : l'atteinte portée étant sensiblement plus grave que la seule détention de ces contenus, et ces mis en cause sont aujourd'hui responsables des dérives de l'IA en détournant les outils générateurs.

## Les impacts sur la santé physique et psychologique des victimes

85 % des survivants de violences sexuelles en ligne déclarent que ces violences ont entraîné des conséquences négatives à long terme<sup>44</sup>. Pour les victimes de viols ou d'agressions sexuelles, l'existence de contenus représentant les violences subies perpétue l'agression vécue et accentue le traumatisme. Mais même en l'absence de violences directes, les répercussions émotionnelles et

psychologiques sur les personnes dont l'image est utilisée sont importantes. Dans cette optique, la prise en charge de la victime et de son entourage est fondamentale.

S'il existe de nombreuses similitudes entre les impacts des violences sexuelles directes et ceux de la cyberpédocriminalité, cette dernière entraîne malgré des conséquences particulières. Se pose alors la question de l'impact spécifique des contenus pédocriminels générés par l'IA sur les victimes.

44. Suojellaan Lapsia, Protect Children ry. "Tech Platforms Used by Online Child Sexual Abuse Offenders: Research Report with Actionable Recommendations for the Tech Industry", 2024.



## Les impacts des violences sexuelles et de la cyberpédocriminalité, notamment générée par l'IA, sur l'individu victime

Entretien croisé avec **JOANNA SMITH**, Psychologue clinicienne et **MÉLANIE DUPONT**, Docteur en psychologie, Psychologue à l'Unité Médico-Judiciaire de l'Hôtel-Dieu (Paris) et Présidente de l'Association contre les Violences sur Mineurs (CVM), Avril et Mai 2024.

### **Fondation pour l'Enfance: Quelles sont les conséquences et l'impact des violences sexuelles vécues dans l'enfance ?**

**Mélanie Dupont:** Les violences sexuelles vécues dans l'enfance entraînent une diversité de conséquences et de manifestations sur la santé globale, qui varient selon l'âge: plus l'enfant est jeune, plus il va y avoir de conséquences sur la mémoire comportementale et sensorielle.

Le traumatisme engendré par les violences sexuelles entraîne un court-circuitage du fonctionnement du corps: dans un moment qui ne présente pas particulièrement de danger (le claquement d'une porte par exemple) les victimes ont une lecture émotionnelle mais surtout sensorielle, ravivant des douleurs et des sensations, telles que la peur.



Les conséquences peuvent également être psychiques. De la même manière que notre cerveau est protégé par la boîte crânienne, notre esprit est protégé par des mécanismes de défense. Chaque personne a ses mécanismes propres. Une personne qui vit un événement potentiellement traumatique (trauma dit simple) peut mettre en œuvre des mécanismes de défense sur le moment pour y survivre, comme la dissociation: pour empêcher une mort cérébrale, les différentes entités du cerveau interrompent leur communication habituelle qui permet de traiter les émotions,

d'intégrer les informations et de mémoriser. Cependant, la fin de l'événement traumatique ne met pas forcément fin à ces mécanismes de défense. Ceux-ci perdurent pour se protéger des symptômes psychotraumatiques, tels que les reviviscences : l'événement traumatique est là en permanence, s'impose à la victime, et la parasite dans son fonctionnement affectif, cognitif, sensoriel actuel.

**Joanna Smith :** Je structurerais les conséquences des violences sexuelles vécues dans l'enfance en 5 registres interconnectés :

1. La santé physique : les violences sexuelles vécues dans l'enfance et le traumatisme qu'elles engendrent sont générateurs de stress et désorganisent le système de régulation de ce stress, avec le risque que celui-ci soit permanent. Or, ces hormones de stress sont toxiques pour l'organisme en développement, notamment pour le cerveau. Par ailleurs, on observe chez les victimes plus d'obésité (40 % de plus que chez les non-victimes), de douleurs (migraines, douleurs de dos, douleurs abdominales), de cancers, ou de maladies chroniques.

2. La santé mentale : on observe chez les personnes victimes des troubles dépressifs, des syndromes de stress post-traumatiques, des troubles anxieux, mais aussi des troubles du comportement alimentaire (120 % de plus que chez les personnes non-victimes). Les personnes victimes de violences sexuelles dans l'enfance commettent 90 % de tentatives de suicide, et 130 % d'actes d'auto-mutilation de plus que les personnes non-victimes. Il y a également un impact sur le développement cognitif, sur la mémoire, sur le sentiment de sécurité, sur le sommeil, sur le sentiment d'être soi, sur l'image du corps, et sur la confiance (en soi, en l'autre, en l'environnement extérieur).


 L'événement  
traumatique est là  
en permanence,  
s'impose à la victime,  
et la parasite dans  
son fonctionnement  
affectif, cognitif et  
sensoriel actuel. 



## À LA LOUPE


3. L'impact psychosocial : on constate des difficultés relationnelles et conjugales assez fréquentes, une instabilité familiale, des difficultés professionnelles, parentales, une maternité souvent plus précoce et une gestation plus courte. On observe également des conséquences négatives sur le mécanisme d'attachement, particulièrement quand les figures d'attachement ou des proches sont auteur.es ou complices des violences. Il y a alors un fort sentiment de trahison.

4. Les conséquences sur la sexualité : on voit chez les personnes victimes de violences sexuelles dans l'enfance davantage de comportements à risque, comme des rapports sexuels non protégés. En effet, le trauma-

 Il est impossible pour la victime de voir une fin à l'épisode traumatique, parce que le contenu continue de tourner sur Internet.

tisme vécu et la logique de l'agresseur atteignent l'estime de soi et convainquent la victime qu'elle mérite d'être traitée comme ça. On constate également chez certaines victimes une addiction sexuelle liée à un besoin affectif, ou a contrario une hyposexualité, mais aussi une difficulté à dire non.

5. La dimension transgénérationnelle : Un parent victime qui souffre de tous ces troubles est

 souvent moins disponible, moins régulé émotionnellement pour son enfant, sans compter la transmission épigénétique de la vulnérabilité au stress. Dans les cas d'inceste, il y a aussi un dysfonctionnement familial très particulier qui se transmet à la génération suivante de cette façon-là.

**MD. :** Pour autant, il est difficile de répondre de manière systématique car chaque situation et chaque réaction est individuelle et subjective. Les conséquences et l'impact dépendent de la réaction et du soutien des adultes qui entourent l'enfant, et de la prise en charge par les professionnels compétents. Si les adultes autour ont cru, soutenu et ont été immédiatement proactifs dans la protection judiciaire et médicale de l'enfant, si les

parents tiennent bon malgré leur souffrance et leur tristesse, s'ils sont eux-mêmes accompagnés, alors l'enfant pourra intégrer cet événement et se construire. Comme une blessure ouverte, il restera une cicatrice de cet événement négatif et douloureux. Mais on aura travaillé à enlever la charge traumatique, pour qu'en grandissant, l'individu puisse regarder la cicatrice sans que cela perturbe son quotidien, et puisse vivre une vie la plus satisfaisante possible.

### **Fondation pour l'Enfance: Quels sont les impacts de la cyber(pédo)criminalité sur les victimes ?**

**JS.:** Certaines conséquences de la cyberpédocriminalité sont similaires aux conséquences des violences sexuelles « hors ligne » : sentiments de honte et d'humiliation, difficultés à faire confiance, mais aussi dépression, anxiété, isolement, risques d'idées suicidaires voire de passage à l'acte, troubles du sommeil, sensation de perte de contrôle, addictions, atteinte de l'estime de soi et de la confiance en soi.

Selon les quelques recherches sur le sujet, la cyberpédocriminalité a, malgré tout, des impacts spécifiques sur les victimes. Les conséquences sont plus vastes et ont beaucoup plus de portée, notamment en raison du caractère doublement intrusif de la violence subie : l'intégrité physique de la victime est doublement violée<sup>45</sup>, lors de l'agression qu'elle subit mais aussi lors de l'exposition de cette agression à la vue de tous. La perte de confiance notamment est amplifiée, car les images et/ou vidéos circulant sur internet atteignent l'image que la victime pense que les autres ont d'elle.

**MD.:** La cyberpédocriminalité implique en effet une exposition de l'intimité au plus grand nombre, qui peut être extrêmement traumatique pour l'individu. S'il s'agit d'une photo prise par la victime elle-même mais qui a été diffusée de façon non-consentie, il y a une culpabilité de s'être trompé qui se rajoute et qui intensifie le traumatisme.

**JS.:** Par ailleurs, il est difficile, voire impossible pour la victime de voir une fin à l'épisode traumatique, parce que le contenu continue de tourner sur



45. Dans les cas où les contenus en ligne représentent des viols et/ou agressions sexuelles ayant eu lieu.



## À LA LOUPE

internet. Cette fin (« closure » en anglais) est pourtant cruciale dans le traitement du traumatisme. Sans celle-ci, le danger est encore présent, réel pour le cerveau et le processus de guérison risque d'être entravé.

Par ailleurs, en cas d'enquête et de procédure judiciaire, les contenus diffusés sans le consentement de la victime seront visionnés par des personnes

 Ce phénomène  
risque également  
d'augmenter les  
contenus  
pédocriminels et  
de les banaliser. 

tierces : des enquêteurs, l'accusé, les avocats, les experts, voire les proches s'il y a un procès et que les contenus sont projetés à cette occasion. Ce visionnage, mais aussi la simple mention de l'affaire dans les médias peuvent être revictimisant, c'est-à-dire que la victime a l'impression, le sentiment d'être agressée une nouvelle fois.

Parfois, le traumatisme est tellement insoutenable, le stress généré est tellement fort que les victimes refusent de parler, voire de reconnaître l'existence de ces images. Les victimes de cyberpédocriminalité ont également souvent comme réaction de ne plus vouloir être prises en photo ou filmées, même par des proches.

**MD.:** De plus, le fait que potentiellement tout le monde soit au courant, et que tout le monde ait accès à ces photos fait douter de tout le monde : qui a vu ? qui est l'auteur ? Avec la cyberpédocriminalité, il y a une permanence de l'agression, de multiples agresseurs (ceux qui visionnent, téléchargent, partagent le contenu), et donc une permanence des conséquences et une revictimisation, comme le mentionnait Joanna. Enfin, la cyberpédocriminalité ouvre la voie à d'autres types de violence, comme le harcèlement.

**Fondation pour l'Enfance:** Quelles seraient les conséquences spécifiques des contenus pédocriminels générés par l'IA sur les victimes ?

**MD.:** Il n'y a pas de littérature sur le sujet (ce qui est intéressant en soi, car cela montre qu'on est en retard), mais selon moi la création de contenus pédocriminels par l'IA générative, et donc l'impossibilité de contrôler son

image, va entraîner chez les victimes une intensification du sentiment de dépossession de soi. Cela pourrait avoir pour effet d'augmenter les troubles psychotiques chez les jeunes, et même une dépersonnalisation à se voir sur des images qui ne représentent pas leurs propres corps.

**JS. :** Il y a une dimension encore plus envahissante avec un risque d'une réaction de choc et d'une sensation de persécution majeure, parce que la personne n'a même pas conscience qu'il y a du matériel à risque quelque part. Il doit y avoir une forte sensation de perte de contrôle et d'impuissance. Et l'impuissance est facteur de symptômes psychotraumatiques.

**MD. :** Ce phénomène risque également d'augmenter les contenus pédocriminels et de les banaliser. Avec l'IA générative, il y a une réelle difficulté à faire la distinction entre ce qui est vrai et ce qui ne l'est pas. Peut-être que certains pourront s'appuyer sur le fait que ça n'est pas vraiment eux, mais il va tout de même falloir « démontrer » son innocence. Finalement, avec cette banalisation, on peut s'interroger : va-t-on avoir une nouvelle génération traumatisée, ou une génération qui s'en moque ?

**JS. :** On peut également imaginer une crainte de s'exprimer sur les réseaux, par peur de « représailles » d'internautes anonymes, mais aussi une crainte d'être pris en photo ou filmé d'une manière plus générale. ●



Un enfant victime demeure à vie une victime, et le sera de nouveau à chaque fois que son image sera vue.

Véronique Béchu, Derrière l'écran, Stock 2024







## ZOOM SUR



# Les auteurs et consommateurs de matériels d'exploitation sexuelle des mineurs en ligne

## Pédocriminalité en ligne et hors ligne : une frontière poreuse

1 cyberpédocriminel sur 8 possède un historique d'infraction sexuelle sur un mineur hors-ligne<sup>46</sup>. Par ailleurs, 52 % des consommateurs pensent que leur usage de contenus pédocriminels pourrait aboutir à une agression sur un enfant (44 % des consommateurs ont pensé à contacter des enfants et 37 % ont contacté des enfants au moins une fois)<sup>47</sup>. De plus, on retrouve des caractéristiques semblables chez les agresseurs « en ligne » et « hors ligne » : le genre (90 % sont des hommes<sup>48</sup>), la déviance sexuelle, et l'existence de traits antisociaux<sup>49</sup>.

## Cyberpédocriminels : qui sont-ils ?

Il existe néanmoins une caractéristique propre aux cyberpédocriminels : ils ont un profil similaire à celui des addicts<sup>50</sup>. En effet, ces individus sont dans une recherche permanente de nouveaux contenus, toujours plus extrêmes et violents. Chaque fois qu'il est en mesure d'en obtenir, on observe chez

 Ces individus sont dans une recherche permanente de nouveaux contenus, toujours plus extrêmes et violents. 

46. CIIVISE « Violences sexuelles faites aux enfants : "on vous croit" », p.273.

47. Protect Children, *ReDirection Survey Report*, 2021, p.16.

48. Véronique Béchu, *Derrière l'écran*, Stock, 2024 p.16.

49. Babchishin, Hanson, VanZuylen, "Online child pornography offenders are different: a meta-analysis of the characteristics of online and offline sex offenders against children", *Archives of Sexual Behavior*, Janvier 2015.

50. CIIVISE, « Violences sexuelles faites aux enfants : "on vous croit" », p.274.

le cyberpédocriminel une réaction similaire à celle d'un addict sur le plan psychopathologique : se procurer de quoi agrandir sa collection est une « dose », et la personne en est dépendante.

Ce phénomène de collection de contenus risque d'être d'autant plus renforcé par l'IA générative. Puisque la possibilité de création est infinie avec cette technologie, elle permettrait au pédocriminel de toujours plus satisfaire son besoin de nouveaux contenus. En outre, l'accès à des « doses » de plus en plus « élevées » est également possible grâce à l'IA qui permet de générer des contenus de plus en plus explicites.

Il n'existe pas de profil-type du producteur et/ou consommateur de matériel d'exploitation sexuelle de mineurs.

Il n'existe pas de profil-type du producteur et/ou consommateur de matériel d'exploitation sexuelle de mineurs. De même, leurs « motivations » pour ces activités criminelles diffèrent grandement. Selon les travaux de Cohen et Felson sur l'opportunité criminelle<sup>51</sup>, les producteurs et consommateurs de contenus pédocriminels peuvent être motivés par la présence de cibles « intéressantes » (de nombreuses vidéos pédocriminelles, possibilité de rentrer en contact avec des enfants) et l'absence de regard censeur (l'anonymat). Certains considèrent la création ou la consommation de contenus pédocriminels comme une solution palliative pour gérer leur attirance et éviter de passer à l'acte sur un mineur, déclarant même que cela leur poserait un problème de violer ou d'agresser sexuellement un enfant<sup>52</sup>.

Sur un plan clinique, si l'existence d'une paraphilie<sup>53</sup> pédophile peut pousser certains à créer, consommer et partager du matériel d'exploitation sexuelle



51. Cohen et Felson, "Social change and crime rate trends: A routine activity approach", *American Sociological Review*, 1979.  
52. Échange entre la Fondation pour l'Enfance et un documentaliste, une infirmière et une psychologue du CRIAVS de Bordeaux (Centre Ressource pour les Intervenants auprès des Auteurs de Violences Sexuelles), juin 2024.

53. Les paraphilies sont des fantaisies imaginatives (fantasmes) sexuellement excitantes, des impulsions sexuelles ou des comportements survenant de façon répétée et intense, et impliquant des objets inanimés, la souffrance ou l'humiliation de soi-même ou du partenaire, des enfants ou d'autres personnes non consentantes (American Psychiatric Association, *Manuel Diagnostique et Statistique des Troubles Mentaux*, 2023).



## ZOOM SUR

d'enfants, d'autres le font pour des raisons bien différentes. Certains dans une recherche de gains financiers, d'autres dans une logique de transgression : consommer et produire du contenu pédocriminel est un interdit qu'ils veulent braver<sup>54</sup>. D'autres, enfin, le font dans une recherche d'appartenance à un groupe<sup>55</sup>.

### La cyberpédocriminalité, un phénomène de groupe

En effet, on a vu émerger des communautés dans lesquelles les cyberpédocriminels se lient les uns aux autres pour agrandir leurs « collections ». Ces espaces sont déployés sur diverses plateformes, revêtant la forme de *chatrooms*<sup>56</sup> sur le *dark web*, ou via les réseaux sociaux sur le *clear web*<sup>57</sup>. Il y circule du matériel d'exploitation sexuelle de mineurs qu'il est presque

Le partage est comme une marque d'authenticité, car un membre qui ne fournirait jamais ce type de contenus serait suspecté comme n'étant pas un des leurs.

obligatoire de partager car il agit comme une monnaie : en partageant du contenu pédocriminel, on en obtient de nouveaux. C'est même une norme sur ces cyberespaces. Le partage est comme une marque d'authenticité, car un membre qui ne fournirait jamais ce type de contenus serait suspecté comme n'étant pas un des leurs.

Une autre particularité de ces cyberespaces est l'effet de groupe qui y est à l'œuvre. On

y retrouve une communauté d'individus où l'individualisme règne, mais qui

54. Échange avec le CRIAVS de Bordeaux, juin 2024.

55. *Ibid.*



56. Salles de discussion.

57. OCCIT, Report n°148.

partagent la même paraphilie, échangent sur leurs fantasmes<sup>58</sup> et qui se confortent mutuellement dans l'idée qu'ils sont parfaitement légitimes de le faire. Des discours tels que « je ne fais de mal à aucun enfant », « je ne fais que regarder », « il n'y a aucun mal à être sexuellement attiré par les enfants », sont monnaie courante<sup>59</sup>. Cette rhétorique banalise ces violences et encourage les individus à ne pas remettre en question leurs rapports aux contenus pédocriminels. Dans le contexte de l'utilisation de l'IA générative, cette tendance à la déresponsabilisation est majorée. En effet, de nouveaux arguments, tels que « ce ne sont pas de vrais enfants », pullulent sur ces cyberespaces. Il existe une distanciation d'autant plus forte entre l'enfant représenté sur le contenu pédocriminel et la personne qui l'a généré. Elle entretient les distorsions cognitives qui légitiment les actes des pédocriminels, en plus de nier la violence dont ils sont auteurs.

Depuis 2023, les cyberespaces sont devenus une mine d'or pour les auteurs et consommateurs de contenus pédocriminels créés grâce à l'IA générative<sup>60</sup>. Ce matériel y circule en masse, et des tutoriels sont même disponibles pour que ceux qui les consultent puissent générer en autonomie des images pédocriminelles. Tout un chacun peut

maintenant générer son propre matériel d'exploitation sexuelle de mineurs. Ces cyberespaces sont donc d'autant plus préoccupants aujourd'hui, et risquent de l'être davantage au vu des évolutions rapides de l'IA.


 Dans le contexte de l'utilisation de l'IA générative, cette tendance à la déresponsabilisation est majorée. En effet, de nouveaux arguments, tels que « ce ne sont pas de vrais enfants », pullulent sur ces cyberespaces.
 



58. DRAGON-S Program : Developing Resistance Against Grooming Online, programme de recherche développé par l'Université galloise de Swansea.

59. Protect Children, *ReDirection Survey Report*, p.55.



60. OCCIT, Report n°148.



## ZOOM SUR

### Une prévention de la cyberpédocriminalité et une prise en charge des créateurs et consommateurs de contenus pédocriminels encore tâtonnantes

Si bon nombre de consommateurs se dédouanent et ne prennent pas conscience de la violence qu'ils font subir aux enfants, certains ont tou-

 **Pourtant, seulement 28 % des consommateurs ont essayé de trouver de l'aide ou y ont songé, sans pour autant y être parvenu.** 

tefois conscience du caractère interdit de leurs agissements. Cela se manifeste, entre autres, par une volonté d'arrêter chez la moitié des consommateurs, dont 28 % presque à chaque fois qu'ils regardent du contenu pédocriminel<sup>61</sup>. Similairement, 62 % des consommateurs ont essayé d'arrêter, dont 24 % presque à chaque fois<sup>62</sup>. Cette volonté d'arrêter naît de sentiments de honte et de culpabilité dans 41 % des cas<sup>63</sup>. Ils

sont si forts que parmi les consommateurs de ce type de contenus, 49 % ont pensé à s'auto-mutiler ou à se donner la mort<sup>64</sup> et 57 % s'isolent<sup>65</sup>.

Pourtant, seulement 28 % des consommateurs ont essayé de trouver de l'aide ou y ont songé<sup>66</sup>, sans pour autant y être parvenu. On pourrait expliquer cela par la prévalence de la fréquentation des cyberespaces chez les consommateurs. Ainsi, leur envie de ne plus consommer des contenus pédocriminels se heurte à ces « lieux » qui encouragent le fantasme pédophile et les comportements pédocriminels.

Mais nous pouvons également expliquer cette difficulté par l'absence de communication sur les dispositifs de soutien existants. S'il existe un numéro

61. Protect Children, *ReDirection Survey Report*, 2021, p.27.

62. *Ibid* p.28.

63. *Ibid* p.29.



64. *Ibid* p.41.

65. *Ibid* p.46.

66. *Ibid* p.49.

téléphonique<sup>67</sup> spécifiquement dédié aux personnes attirées par les enfants, celui-ci est très peu connu. Par ailleurs, il n'apparaît dans les recherches internet que si certains mots-clés sont utilisés (par exemple « pédophilie »). Cependant, certaines personnes n'en sont pas encore au stade de pouvoir poser de tels mots sur ce qu'ils ressentent. Pour ceux qui se sont tournés vers des thérapeutes, ils ont pu être confrontés à des réactions d'hostilité ou d'impuissance. Cette méconnaissance des dispositifs d'aide aux auteurs, et le manque de formation des professionnels est d'autant plus préoccupante alors que se démocratise un nouvel outil qui rend l'exploitation sexuelle des enfants en ligne encore plus extrême, encore plus accessible, et encore plus banalisée.

Pourtant, la prise en charge des personnes ayant des attirances pédophiles et des pédocriminels est fondamentale pour prévenir les passages à l'acte et les récidives et, in fine, protéger les enfants. Cette offre existe à travers la Fédération Française des CRIAVS<sup>68</sup>, qui gère notamment le numéro d'aide et différents centres à travers le territoire français. Ces centres proposent une prise en charge thérapeutique pour permettre à la personne d'apaiser et d'apprendre à faire avec sa paraphilie. Pour qu'elle fonctionne, le patient doit vouloir adhérer et être engagé dans cette recherche de « dépassement » de son attirance. La thérapie consiste notamment à travailler sur la régulation des émotions et sur les traumatismes, étant donné que la plupart des individus ont eux-mêmes subi des violences, maltraitements et négligences (sexuelles, physiques, psychologiques) dans l'enfance. Il existe des dispositifs similaires au Royaume-Uni, notamment proposés par la Lucy Faithfull Foundation<sup>69</sup> dont le site web dispose de modules que l'internaute attiré par les mineurs peut faire en autonomie. ●

 Pour ceux qui se sont tournés vers des thérapeutes, ils ont pu être confrontés à des réactions d'hostilité ou d'impuissance.
 

67. Service Téléphonique d'Orientation et de Prévention (S.T.O.P.), 0 806 23 10 63.

68. Pour plus de détails, voir <https://www.ffcriavs.org/>

69. Pour plus de détails, voir <https://www.lucyfaithfull.org.uk/>







# État des lieux du cadre législatif en vigueur et des initiatives en cours

L'émergence des contenus pédocriminels générés par l'IA pose des enjeux juridiques et politiques fondamentaux en matière de protection des enfants. Certaines initiatives ont été lancées afin de lutter contre le détournement de l'IA générative<sup>70</sup>.

## Cadre international

### Exemples judiciaires

Des exemples récents montrent que les contenus pédocriminels générés par l'IA commencent déjà à être appréhendés par le droit. Au Royaume-Uni, à la fin du mois d'avril 2024, un pédocriminel a été reconnu coupable de la création de plus de 1000 images truquées à caractère pédopornographique. Ainsi, le Tribunal de Poole a interdit au coupable d'utiliser des IA génératives pendant une durée de cinq ans. L'ISF a applaudi cette décision « historique », considérant que les logiciels d'IA générative sont des « usines capables de produire les images les plus épouvantables ». En

pratique, cette décision pourrait montrer l'exemple et inspirer d'autres pays.

Aux États-Unis, dès le mois de juin 2023, le FBI a tenu à avertir la population des risques engendrés par les *deepfakes*<sup>71</sup>. Le FBI en a profité pour rappeler que la sextorsion peut constituer une violation de plusieurs lois pénales fédérales. Ainsi, le FBI invite le public à faire preuve de prudence lorsqu'il partage des contenus, les images et vidéos pouvant fournir aux acteurs malveillants des contenus à exploiter pour des activités criminelles interdites par la loi. Le FBI suggère de surveiller l'activité en ligne des enfants, de faire preuve de discernement lors de la publication de contenus, ou encore d'effectuer fréquemment des recherches en ligne sur les informations personnelles des personnes concernées et de leurs enfants.

Fin mai 2024, aux États-Unis toujours, un homme a été arrêté, pour avoir créé des milliers d'images de mineurs à caractère pédopornographique à l'aide

71. Federal Bureau of Investigation, Public Service Announcement "Malicious Actors Manipulating Photos and Videos to create Explicit Content and Sextortion Schemes", Alert Number I-060523-PSA, June 5th 2023.

70. Un schéma récapitulatif se trouve en fin de partie.

d'une IA. Selon le communiqué officiel relatif à cette affaire, « le ministère de la justice poursuivra énergiquement ceux qui produisent et distribuent des images pédocriminelles quelle que soit la manière dont ces images ont été créées ». Selon le Procureur général adjoint principal, « L'annonce d'aujourd'hui envoie un message clair :

*l'utilisation de l'IA pour produire des images sexuellement explicites d'enfants est illégale, et le ministère de la justice n'hésitera pas à demander des comptes à ceux qui possèdent, produisent ou distribuent des images de violences sexuelles d'enfants générées par l'IA ».* S'il est reconnu coupable des chefs d'accusation retenus dans l'acte



## Pratiques et droit comparés : engagements et actions des pouvoirs publics et de la société civile aux États-Unis

Entretien avec **JOHN SHEHAN**, *Senior Vice President* de la Division des enfants exploités et de l'engagement international du National Center for Missing and Exploited Children (NCMEC\*), Juin 2024

### Fondation pour l'Enfance : Qu'avez-vous constaté en matière de contenus pédocriminels générés par l'IA ?

**John Sheman :** En 2023, nous avons reçu 36 millions de signalements sur notre *CyberTipline*. Parmi eux, 4700 signalements incluent l'utilisation de l'IA

---

\*Le NCMEC est une organisation privée à but non lucratif créée en 1981 aux États-Unis. Elle a pour mission d'aider à retrouver les enfants disparus, de lutter contre l'exploitation sexuelle des enfants et de prévenir les violences faites aux enfants. Le NCMEC gère la *CyberTipline* qui permet au public et aux fournisseurs de services électroniques basés aux États-Unis de signaler les cas présumés d'exploitation sexuelle de mineurs. Depuis sa création, la *CyberTipline* a reçu 191 millions de signalements.

d'accusation, l'auteur des faits encourt une peine maximale de 70 ans d'emprisonnement et une peine minimale obligatoire de 5 ans d'emprisonnement. Cette affaire a été engagée dans le cadre du projet « Safe Childhood », initiative nationale de lutte contre l'exploitation et l'abus sexuels des enfants, lancée en mai 2006 par le ministère de

la justice. Ces exemples démontrent que la création de deepfakes, et notamment à caractère pédopornographique, constituerait déjà une pratique interdite, illicite, et pouvant engendrer des sanctions pénales. De nombreux acteurs réclament cependant une prise en compte spécifique de ce phénomène par le système judiciaire.

généraliste. En 2024, nous avons en moyenne 450 signalements de contenus pédocriminels générés par l'IA par mois. Ce chiffre est en légère augmentation. 14 % de ces 4700 signalements proviennent des plates-formes d'IA générative (ce qui est très peu!), et 15 % proviennent de membres du public. La plupart des signalements qui parviennent à la *CyberTipline* proviennent de plateformes (telles que Facebook). Lorsque nous recevons de tels signalements, nous les supprimons et les plaçons dans nos banques de *hash*\*.

On constate que les contenus pédocriminels générés par l'IA sont utilisés par les délinquants pour extorquer financièrement les victimes: ils créent du contenu et font chanter les enfants, en les menaçant d'envoyer ces photos ou vidéos à leurs amis et à leur famille.

### **Fondation pour l'Enfance: Êtes-vous en mesure de faire la distinction entre les contenus pédocriminels générés par l'IA et ceux non générés par l'IA?**

**JS.:** L'avancée de la technologie est telle que l'on ne peut plus compter sur l'œil humain pour distinguer les contenus générés par l'IA générative de ceux qui ne le sont pas. À ce stade, nous ne disposons pas d'un outil unique qui permette de savoir s'il s'agit ou non de contenu généré par l'IA. Il faut donc que la technologie rattrape un peu son retard.

Joe Biden a récemment signé le *Report Act* qui va nous permettre de disposer de nouveaux outils, que nous pourrons utiliser pour analyser les images







## À LA LOUPE

et les vidéos entrantes afin de mieux détecter s'il s'agit d'un contenu créé par l'IA générative. Nous avons également mis à jour la *CyberTipline* pour qu'une plateforme puisse préciser si le fichier qu'elle signale est généré par l'IA générative.

**Fondation pour l'Enfance: Comment travaillez-vous avec les entreprises du secteur des nouvelles technologies, les autres plateformes de signalements, les forces de l'ordre etc.**

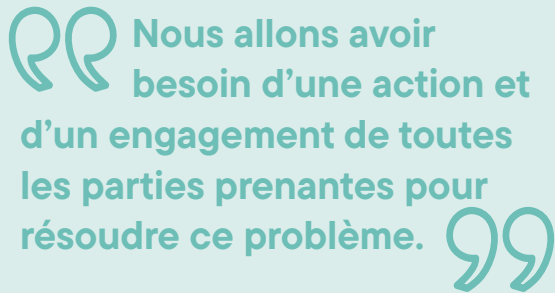
**JS.:** En 2002, nous avons lancé un programme d'identification des enfants victimes et nous sommes désormais le centre national d'échange d'informations sur l'identification des victimes. Lorsque les forces de l'ordre américaines procèdent à une arrestation dans le cadre d'une affaire impliquant des contenus pédocriminels, elles nous partagent des copies des preuves saisies. Nous examinons ce contenu et nous leur remettons un rapport indiquant le nombre de fichiers qui contiennent des images et des vidéos de mineurs qui ont déjà été identifiés. Nous nous concentrons également sur les contenus dans lesquels les enfants n'ont pas été identifiés, à la recherche d'indices et d'informations, et nous les transmettons aux forces de l'ordre pour qu'elles enquêtent, identifient et viennent en aide à ces enfants.

Au niveau international, le NCMEC est devenu un véritable centre d'échange d'informations. Les technologies évoluent, et avec elles de nouveaux dangers émergent, et la *CyberTipline* est souvent la première à les identifier et à tirer la sonnette d'alarme. Nous travaillons en étroite collaboration avec les autorités policières locales et fédérales aux États-Unis, et 160 pays et territoires dans le monde acceptent également les signalements de la *CyberTipline*. Nous travaillons en étroite collaboration avec les *hotlines*, les services de

 Lorsque les forces de l'ordre américaines procèdent à une arrestation dans le cadre d'une affaire impliquant des contenus pédocriminels, elles nous partagent des copies des preuves saisies. 

police et les entreprises sur les différentes tendances émergentes. Enfin, nous sommes en mesure de télécharger le contenu dont nous avons connaissance sur la base de données d'Interpol.

Nous sommes engagés à coopérer avec les différents acteurs de la lutte contre la cyberpédocriminalité. Nous sommes actuellement en train de collecter des données liées aux contenus pédocriminels générés par l'IA, que ce soit les *prompts* ou les contenus eux-mêmes. Nous disposerons bientôt d'une liste spécifique à l'IA générative.

 Nous allons avoir besoin d'une action et d'un engagement de toutes les parties prenantes pour résoudre ce problème.

### **Fondation pour l'Enfance: Aux États-Unis, les plateformes sont-elles tenues de détecter les contenus pédocriminels ?**

**JS.:** La loi américaine exige que les entreprises signalent à la *CyberTipline* les cas dont elles ont connaissance. Certaines les détectent de manière proactive et volontaire, mais la loi ne les oblige pas à le faire. Pour la plupart des entreprises, ce travail proactif est coûteux, et elles ne l'entreprendront pas à moins d'être légalement obligées de le faire. J'espère que la proposition de Règlement européen sur la prévention et la lutte contre les abus sexuels actuellement à l'étude, et qui prévoit la possibilité d'obliger les entreprises à détecter ces contenus, sera adoptée.

### **Fondation pour l'Enfance: Que recommandez-vous pour mieux prévenir, identifier et éliminer les contenus pédocriminels générés par l'IA ?**

**JS.:** Nous allons avoir besoin d'une action et d'un engagement de toutes les parties prenantes pour résoudre ce problème (pouvoirs publics et entreprises du secteur des nouvelles technologies). Au rythme où ces technologies évoluent, les entreprises se précipitent sur le marché, ce qui donne lieu à de merveilleuses innovations. Mais la plupart de ces outils ne sont pas conçus dès le départ dans une optique de sécurité, et les entreprises se disent qu'elles résoudront les problèmes plus tard. Or, la sécurité des mineurs doit être pensée dès la conception de l'outil. ●

## Initiatives et accords des États

Au-delà des grands principes édictés sur les droits des enfants, par exemple la Convention internationale des droits de l'enfant (CIDE) de 1989, des initiatives internationales ont émergé afin de traiter spécifiquement des nouveaux risques engendrés par l'IA. En effet, une coordination au niveau international apparaît nécessaire pour encadrer l'IA générative. Comme le souligne l'Alliance mondiale WeProtect<sup>72</sup>, une harmonisation mondiale et effective de la réglementation d'Internet pourrait considérablement stimuler la réactivité des plateformes dans la prise de mesures efficaces pour permettre une lutte accrue contre les violences en ligne.

Plusieurs groupements d'États ont donc pu adopter des principes, ou conclure des accords visant à régir l'IA (et notamment générative), et ainsi promouvoir une lutte accrue contre les usages illicites de l'IA, notamment concernant la création de contenus à caractère pédopornographique. A ce titre, le 22 mai 2018, l'Organisation de coopération et de développement économiques (OCDE) a adopté différents Principes relatifs à l'IA. Ces Principes constituent la première norme intergouvernementale dans le domaine de l'IA, et visent en pratique

à favoriser l'innovation et la confiance dans l'IA, en promouvant une gestion responsable d'une IA fiable, tout en garantissant le respect des droits de l'Homme et des valeurs démocratiques<sup>73</sup>. Lors de la Réunion du Conseil au niveau des Ministres des 2 et 3 mai 2024, l'OCDE a annoncé l'adoption d'une version révisée de ses Principes sur l'IA<sup>74</sup>. Les Principes préconisent notamment de prévoir des mécanismes « *afin de garantir que, dans l'éventualité où des systèmes d'IA risqueraient de causer des préjudices injustifiés ou présenteraient un comportement indésirable, ils puissent être neutralisés, réparés et/ou mis hors service en toute sécurité, en tant que de besoin.* »

Plusieurs autres instances internationales ont insisté sur la prise en compte des risques posés par l'IA, et sur la nécessité d'atténuer ces risques comme les Nations Unies<sup>75</sup>, l'Union internationale des télécommunications (UIT)<sup>76</sup> ou lors du Sommet mondial de l'IA (AI Safety Summit) en 2023, au cours duquel a été conclu la déclaration de Bletchley sur la sécurité

73. Russel S., Perset K., Grobelnik M., "Updates to the OECD's Definition of an AI system explained", OECD.AI Policy Observatory, November 29th 2023.

74. OCDE, Communiqué de presse "Face aux évolutions technologiques rapides, l'OCDE met à jour les Principes sur l'IA", 3 mai 2024.

75. Assemblée Générale des Nations Unies, "Saisir les possibilités offertes par des systèmes d'intelligence artificielle sûrs, sécurisés et dignes de confiance pour le développement durable" A/78/L.49, 11 mars 2024.

76. UIT, "Sommet mondial de l'UIT sur l'intelligence artificielle au service du bien social : un cadre pour les discussions sur la gouvernance de l'intelligence artificielle", 16 mai 2024.

72. WeProtect Global Alliance, "Évaluation mondiale de la menace 2023, Évaluer l'ampleur et la portée de l'exploitation et des violences sexuelles en ligne envers les enfants, pour transformer la riposte".

de l'IA<sup>77</sup>. Nous pouvons aussi citer les principes directeurs internationaux relatifs à l'IA, et le code de conduite volontaire pour les développeurs d'IA<sup>78</sup> adoptés au G7 d'Hiroshima en 2023. Ces principes et le code de conduite volontaire visent à compléter les règles juridiques contraignantes que les colégislateurs de l'UE mettent au point dans le cadre de la réglementation de l'UE sur l'IA.

Concernant la violence contre les enfants, la *Virtual Global Taskforce*, alliance internationale de lutte contre les violences sexuelles faites aux enfants, souligne en particulier que « *certaines délinquants utiliseront des outils d'IA pour repérer des enfants à grande échelle, accélérant le processus en automatisant l'engagement avec une facilité déconcertante. En outre, les outils d'imagerie de l'IA peuvent générer de vastes volumes de matériel illégal en quelques secondes, y compris des images entièrement synthétiques et d'autres comprenant de vrais enfants. [...] La facilité et la disponibilité de la violence sexuelle générée par l'IA ne feront qu'aggraver la situation et créeront un environnement plus permissif pour les agresseurs, exposant un plus grand nombre d'enfants à la violence* ». Elle préconise en particulier une coopération entre les services de

police ainsi qu'avec les acteurs de l'IA, et un travail de prévention sur les dangers créés par l'IA générative.<sup>79</sup>

Pour sa part, en septembre 2023, l'Organisation des Nations Unies pour l'éducation, la science et la culture (UNESCO) a publié un guide sur l'IA générative dans l'éducation et la recherche<sup>80</sup>. L'UNESCO propose différentes mesures afin que la conception et l'utilisation de l'IA générative soient éthiques, inclusives et égalitaires. A ce titre, ledit guide prévoit notamment l'obligation de libeller clairement les contenus générés par l'IA<sup>81</sup>, ou encore l'établissement d'organismes nationaux dirigeant l'approche des gouvernements en matière d'IA et coordonnant la coopération entre les secteurs. L'UNESCO recommande aussi la création de mécanismes de plaintes et de recours au sein des systèmes d'IA pour les utilisateurs, ainsi qu'un mécanisme d'alerte et de contrôle de toute utilisation illicite du service par l'IA destiné à informer les agences gouvernementales<sup>82</sup>.

## Initiatives des entreprises du secteur

Diverses initiatives existent également au niveau de certaines entreprises ou

77. UK Government, Policy paper "The Bletchley Declaration by Countries Attending the AI Safety Summit, 1-2 November 2023", 1st November 2023

78. Lettre de la Direction des Affaires Juridiques, "Accord du G7 sur des principes directeurs et un code de conduite en matière d'intelligence artificielle", 8 novembre 2023

79. We Protect, Virtual Global Taskforce, "Technological Tipping Point Reached in Fight Against Child Sexual Abuse", January 2024.

80. UNESCO, Guide sur l'IA générative dans l'éducation et la recherche, septembre 2023.

81. Ceci a été fait dans le règlement sur l'IA de l'Union Européenne.

82. UNESCO, Guide sur l'IA générative dans l'éducation et la recherche, p.22.



organisations professionnelles afin de lutter contre la création de contenus à caractère pédopornographique générés par l'IA. Outre les Principes de Thorn, déjà évoqués précédemment, des entreprises du secteur des nouvelles technologies (telles que TikTok ou encore Snap) et certains États ont signé une déclaration, le 30 octobre 2023, s'engageant à « *travailler ensemble pour nous assurer une utilisation responsable de l'IA et lutter contre la menace des abus sexuels sur les enfants [et] nous assurer que les risques posés par l'IA en la matière ne deviennent pas insurmontables.* »

De plus, dans une lettre ouverte intitulée “*Disrupting the Deepfake Supply Chain*”<sup>83</sup> en date du 21 février 2024, plus de 700 experts en IA et dirigeants d'entreprise du monde entier ont appelé à une réglementation plus importante en matière de *deepfakes*. La lettre plaide en faveur de la criminalisation des *deepfakes* pédocriminels, de sanctions pénales lourdes, et d'un devoir pour les développeurs et distributeurs d'IA d'empêcher, au sein de leurs produits, la création de *deepfakes*, ou à défaut, d'une responsabilité de leur part.

Par ailleurs, pour pallier la difficulté de détection des contenus pédocriminels et d'identification des potentielles victimes de violences sexuelles, des entreprises ont mis en place des

solutions pour mieux identifier les images générées par l'IA des images non générées par l'IA. Sur YouTube, les créateurs de contenus doivent désormais indiquer s'ils ont utilisé de l'IA générative dans leurs vidéos. A défaut, ils feront face à des sanctions telles que le retrait du contenu visé, ou la suspension du programme de rémunération de YouTube. Par ailleurs, en mai 2024, Open AI a présenté un outil conçu pour détecter les images générées par son modèle d'IA générative DALL-E 3. Un badge « CR » (« Content Credential ») permettant de marquer les contenus produits par l'IA générative et de retracer la provenance et l'historique d'une image en un clic a également été créé. Cette initiative regroupe des sociétés comme Adobe, Microsoft, Publicis, Nikon, Arm, Intel ou Leika.

Enfin, le 7 novembre 2023, la Tech Coalition (Google, Meta, Twitch, Discord, Mega, Roblox, Quora et Snap), alliance d'entreprises technologiques internationales collaborant pour lutter contre l'exploitation et les violences sexuelles des enfants en ligne, a annoncé le lancement du programme « Lantern ». Il s'agit du premier programme de signalement collaboratif de contenus à caractère pédopornographique. Concrètement, les plateformes utilisent ce canal pour partager des informations (des « signaux ») relatives aux contenus qui violeraient leurs politiques en la

83. Lettre ouverte, *Disrupting the Deepfake Supply Chain*, OpenLetter.net, 21 février 2024

matière, pour que les autres membres puissent détecter un éventuel contenu similaire sur leur plateforme. L'usage massif de ce canal pourrait permettre de mieux détecter les menaces réelles pour les enfants.

Ces multiples initiatives et engagements, tant au niveau étatique qu'au niveau des acteurs du secteur, montrent un début de prise de conscience - et de sensibilisation - concernant les risques que l'IA peut engendrer, notamment en matière de cyberpédocriminalité. Outre ces initiatives internationales, d'autres, ayant tendance à être plus contraignantes et prescriptives, sont en cours au niveau européen.

## Cadre européen

Dès 2022, le laboratoire d'innovation d'Europol relevait les risques liés à l'utilisation des *deepfakes* pour des activités criminelles. En effet, l'agence européenne de police soulignait que cette technologie facilite la réalisation d'activités criminelles, comme le harcèlement ou l'humiliation, la pornographie non consensuelle ou encore l'exploitation sexuelle des enfants en ligne. Le rapport préconisait, dès lors, la nécessité de prévoir un cadre réglementaire effectif et adapté.

Le législateur européen n'est pas en reste, les différents textes et initiatives de l'Union Européenne permettent

d'appréhender l'usage de l'IA dans les violences sexuelles contre les enfants tant par des dispositions à portée générale que par des dispositions spécifiques à la protection des enfants.

## Protection des droits de l'enfant

La Charte des droits fondamentaux de l'Union européenne garantit la protection des droits de l'enfant par les institutions de l'Union et par ses pays membres. Le Conseil de l'Europe, qui ne relève pas de l'Union Européenne, estime pour sa part que « *les enfants ont besoin d'une protection particulière lorsqu'ils sont en ligne et doivent apprendre à éviter les dangers et utiliser au mieux internet. A cette fin, les enfants doivent devenir des citoyens numériques. Internet offre aux enfants une multitude de possibilités, mais les expose aussi à des risques susceptibles de porter préjudice à leurs droits humains. Certains de ces risques comprennent le cyber-harcèlement, des questions concernant la protection des données, la sollicitation d'enfants en ligne à des fins sexuelles (« grooming »), la cybercriminalité et la pornographie infantine.* »<sup>84</sup>

Ainsi, le Conseil de l'Europe a développé une Stratégie pour les droits de l'enfant (2022-2027)<sup>85</sup>. Cette Stratégie met

84. Conseil de l'Europe, L'environnement numérique - Droits des Enfants.

85. Conseil de l'Europe, Stratégie du Conseil de l'Europe pour les droits de l'enfant (2022-2027) "Les droits de l'enfant en action : poursuivre la mise en œuvre et innover ensemble".

l'accent sur les droits des enfants au sein de l'environnement numérique. Elle est renforcée par la Recommandation CM/Rec (2018) du Comité des Ministres aux États membres sur les Lignes directrices relatives au respect, à la protection et à la réalisation des droits de l'enfant dans l'environnement numérique. Ces Lignes directrices sont également complétées par le nouveau Manuel pour les décideurs politiques sur les droits de l'enfant dans l'environnement numérique. Selon le Conseil de l'Europe, « la nouvelle Déclaration du Comité des Ministres appelle les États membres à intensifier leurs efforts pour protéger la vie privée des enfants dans l'environnement numérique et à promouvoir, entre autres, les Lignes directrices sur la protection des données des enfants dans un cadre éducatif. »

Par ailleurs, le 6 février 2024, la Commission européenne a adopté une proposition de directive visant à actualiser les règles de droit pénal relatives aux violences sexuelles commises contre des enfants et à l'exploitation sexuelle de ceux-ci<sup>86</sup>. Ces règles révisées prévoient un élargissement des définitions des infractions et alourdissent les sanctions prévues, avec en outre des exigences plus précises en matière de prévention et d'assistance aux victimes. En particulier, cette proposition pose le constat que « du fait du développement continu

*d'applications d'intelligence artificielle capables de créer des images réalistes qui ne peuvent être distinguées des images réelles, le nombre d'images et de vidéos connues sous le nom d'hypertrucages (« deepfake ») représentant des abus sexuels commis contre des enfants devrait croître de manière exponentielle dans les années à venir. En outre, la définition existante ne couvre pas pleinement le développement de paramètres de réalité augmentée, étendue et virtuelle utilisant des avatars, y compris un retour d'information sensoriel, par exemple au moyen de dispositifs permettant de percevoir le toucher ».* Pour y répondre, la proposition suggère d'apporter des modifications à la définition de « matériel relatif à des abus sexuels sur enfants » afin de s'assurer que celle-ci couvre bien les deepfakes pédocriminels. Cette proposition doit encore être débattue et adoptée par les différentes institutions européennes.

Enfin, le Conseil de l'Union Européenne a également adopté, le 9 juin 2022, des conclusions sur la stratégie de l'Union européenne quant aux droits de l'enfant. Plus généralement, les États membres de l'Union sont invités à élaborer des politiques visant à faire respecter les droits des enfants sans discrimination, intensifier les efforts déployés pour prévenir et lutter contre toutes les formes de violence à l'égard des enfants, renforcer leurs systèmes

86. Commission européenne, L'actualisation des règles de droit pénal donne une nouvelle impulsion à la lutte contre les abus sexuels commis contre des enfants, 6 février 2024.

judiciaires afin qu'ils respectent les droits de tous les enfants, ou encore offrir aux enfants davantage de possibilités de devenir des membres responsables et résilients de la société numérique.

L'application de ces principes se retrouve, directement ou indirectement, dans les différents textes adoptés ces dernières années et qui permettent de commencer à appréhender les risques de violences sexuelles envers les enfants engendrés par l'IA.

### Réglementation de l'IA

Le texte fondamental, à l'échelle de l'Union Européenne, est le règlement européen sur l'IA<sup>87</sup>, pionnier en la matière, qui rappelle notamment que les enfants bénéficient de droits spécifiques, et souligne leur vulnérabilité et la nécessité de les protéger. Le texte encadre les systèmes d'IA selon les risques posés par un système donné.

Certaines IA seront totalement interdites, par exemple les systèmes qui créent ou développent des bases de données de reconnaissance faciale, à partir de la collecte non ciblée d'images faciales provenant d'internet ou de la vidéosurveillance. L'utilisation à des fins répressives

de systèmes d'identification biométrique à distance en temps réel dans des espaces accessibles au public est également prohibée. A noter qu'une exception existe pour ce dernier cas en cas d'infraction relative à l'exploitation sexuelle des enfants et à la pédocriminalité.

D'autres systèmes dits à haut risque, ne pourront être proposés qu'après respect d'un ensemble d'obligations, comme la mise en place de systèmes adéquats d'évaluation et d'atténuation des risques, ou encore la nécessité de conserver une documentation détaillée fournissant toutes les informations nécessaires sur le système et son objet afin de permettre aux autorités d'évaluer sa conformité.

Le règlement sur l'IA impose également des obligations de transparence pour certains systèmes d'IA. Par exemple, il impose aux systèmes d'IA qui génèrent ou manipulent des images de désigner les contenus générés par IA lorsqu'ils sont hypertruqués, c'est à dire lorsque le contenu est généré ou manipulé par l'IA et peut être perçu à tort comme étant authentique ou véridique<sup>88</sup>.

Le règlement sur l'IA a été publié au Journal officiel de l'UE le 12 juillet 2024 et entre progressivement en application depuis le 1<sup>er</sup> août 2024. Bien que traitant à la marge du problème posé

87. Règlement (UE) 2024/1689 du Parlement européen et du conseil du 13 juin 2024 établissant des règles harmonisées concernant l'intelligence artificielle et modifiant les règlements (CE) no 300/2008, (UE) no 167/2013, (UE) no 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).

88. Article 50.4: « Les déployeurs d'un système d'IA qui génère ou manipule des images ou des contenus audio ou vidéo constituant un hypertrucage indiquent que les contenus ont été générés ou manipulés par une IA. (...) ».

par l'utilisation de l'IA dans la création de contenus pédopornographiques, le règlement sur l'IA pose le premier jalon d'un cadre juridique permettant d'adresser l'émergence de ce phénomène.

Concomitamment à l'adoption du règlement sur l'IA, le 14 mars 2024, le Comité sur l'intelligence artificielle, faisant partie du Conseil de l'Europe, a approuvé le projet de Convention-cadre du Conseil de l'Europe sur l'intelligence artificielle et les Droits de l'Homme, la démocratie et l'État de droit. Ce projet se veut être le premier traité international contraignant sur l'IA. La Convention-cadre a été formellement adoptée par le Comité des Ministres du Conseil de l'Europe (Ministres des affaires étrangères) le 17 mai 2024. Le texte vise à établir un cadre juridique pour l'ensemble du cycle de vie des systèmes d'IA et à garantir que les activités menées dans le cadre du cycle de vie des systèmes d'IA sont pleinement compatibles avec les Droits de l'Homme, la démocratie et l'État de droit, tout en favorisant le progrès et les innovations technologiques. Bien que ne traitant pas spécifiquement de la pédocriminalité, la Convention-cadre requiert des parties prenantes qu'elles prennent en compte les « vulnérabilités spécifiques » en rapport avec le respect des droits des enfants, et ouvre la voie à l'adoption de mesures ciblant spécifiquement la cybercriminalité.



## Encadrement des plateformes

Les plateformes, qui peuvent notamment héberger des contenus à caractère pédopornographique et/ou les IA génératives qui seront utilisées par des internautes pour générer ces contenus, sont déjà soumises à différentes obligations ayant pour objectif de lutter contre ce type de matériel.

Tout d'abord, le Règlement relatif à un marché unique des services numériques, dit Digital Services Act (DSA), applicable depuis le 17 février 2024, vise notamment à lutter contre la diffusion de contenus illicites, et à instaurer davantage de transparence entre les plateformes en ligne et leurs utilisateurs. Ainsi, les plateformes bénéficient d'une exonération de responsabilité lorsqu'elles rendent promptement inaccessibles les contenus illicites ou ceux se rapportant à des activités illégales qui leur sont notifiés. Le DSA cite ainsi à titre d'exemple le partage d'images représentant des violences sexuelles commises sur des enfants. Dans le cadre de l'IA générative, le DSA sera pertinent a posteriori, c'est-à-dire lorsque qu'un contenu pédocriminel



aura déjà été généré et sera partagé en ligne. Une fois le contenu identifié, la plateforme devra le rendre inaccessible au titre du DSA, afin de continuer à bénéficier de son exonération de responsabilité.

En outre, la Commission Européenne a présenté, en mai 2022, une proposition de règlement afin de prévenir et combattre les abus sexuels sur enfants (proposition de règlement CSAM), en luttant contre la multiplication des contenus à caractère pédopornographique et contre la sollicitation d'enfants. Le principal objectif du règlement est de créer un cadre permanent pour lutter efficacement contre les violences sexuelles sur les enfants en ligne. En attendant l'adoption de ce règlement, en discussion au parlement européen, la détection de ces violences sexuelles contre les enfants est basée sur une dérogation temporaire à la directive vie privée et communications électroniques (dite Directive e-privacy)<sup>89</sup>, prolongée

89. Conseil de l'Union européenne, Communiqué de presse Abus sexuels sur enfants : accord entre le Conseil et le Parlement européen pour la prorogation d'une mesure de protection, 15 février 2024.

jusqu'au 3 avril 2026. La proposition s'appuie en réalité sur la législation sur les services numériques (DSA), et la complète avec des dispositions spécifiquement consacrées aux cas de violences sexuelles sur mineurs.

Cette proposition énonce des mesures ciblées et proportionnées au risque d'utilisation à mauvais escient d'un service donné à des fins de violences sexuelles sur enfants en ligne. Elle vise en outre à faire en sorte que les fournisseurs puissent s'acquitter de leurs obligations, en créant un centre de l'Union Européenne chargé de prévenir et de combattre les violences sexuelles sur enfants, et ayant pour objectif de faciliter et soutenir la mise en œuvre du présent règlement. En pratique, cette proposition introduit notamment pour les fournisseurs de services en ligne, sous certaines conditions, des obligations de détection des violences sexuelles commises en ligne sur des enfants, de signalement de toute violence sexuelle potentielle et de retrait du contenu identifié (notamment par le biais d'injonctions par l'autorité judiciaire ou administrative). En outre, tout fournisseur constatant la présence sur son service de contenus relatifs à des violences sexuelles commises en ligne sur des enfants devra le signaler au centre de l'Union Européenne.

Enfin, le Règlement général sur la protection des données (RGPD), en



vigueur depuis mai 2018, aura aussi vocation à s'appliquer en présence de contenus à caractère pédopornographique générés par l'IA.

En effet, le RGPD impose des règles strictes quant à l'utilisation des données à caractère personnel. Ces règles sont d'autant plus strictes lorsqu'il s'agit de données à caractère personnel de mineurs. Ainsi, un système d'IA utilisant des données à caractère personnel de mineurs pour s'entraîner devra justifier d'une base légale solide pour le faire, par exemple le consentement ou l'intérêt légitime, et les droits et libertés du mineur en question devront être pris en compte, autrement le traitement des données à caractère personnel du mineur sera illicite.

Ceci est particulièrement pertinent dans le cas où le système d'IA générative pourrait potentiellement utiliser l'image d'un mineur réel, a priori anodine, pour générer du contenu à caractère sexuel.

## Cadre national

En France, comme dans de nombreux pays, les autorités législatives et réglementaires s'efforcent d'élaborer un cadre juridique efficace pour contrer la prolifération des contenus à caractère pornographique créés par l'IA, certains droits permettant également d'aborder les problématiques liées à ces contenus.

## Principes et protections fondamentales

L'article 9 du Code civil consacre le principe du droit à la vie privée, du droit à l'image et à la voix, et permet à chacun de s'opposer à la reproduction de son image sans son autorisation. Cette disposition permet en principe de prévenir l'utilisation des traits de personnalité d'une personne, notamment d'un mineur, dans du contenu généré par l'IA sans l'accord de ladite personne (ou de ses représentants si la personne est mineure). Ainsi, cet article permettra de demander le retrait de tout contenu reprenant les traits d'une personne réelle, que ce contenu soit modifié par l'IA ou pas. De plus, l'article 226-4-1 du Code pénal prohibe l'usurpation d'identité ou l'usage de données pouvant « *troubler la tranquillité* » ou « *porter atteinte à l'honneur ou à la considération* » de la victime. Dans le cas de *deepfakes* sexuels utilisant des images, vidéos ou autres données biométriques de la victime pour créer des contenus truqués, cet article pourra s'appliquer.

Par ailleurs, le Code pénal prohibe la fixation, l'enregistrement, la diffusion, l'offre ou la transmission d'une part, et la consultation ou la détention d'autre part, d'images ou de représentations à caractère pédopornographique d'un mineur (article 227-23 du Code pénal). Cet article permet une protection étendue des mineurs en ce qu'il précise qu'il suffit que le contenu



représente une personne ayant l'aspect physique d'un mineur pour qu'il puisse être considéré comme étant du contenu à caractère pédopornographique représentant un mineur. Cette précision du législateur pourrait notamment être pertinente lorsqu'un contenu généré par l'IA représente une personne ayant les caractéristiques d'un mineur, même si des techniques de contournement trop facilement exploitables existent<sup>90</sup>.

Toutefois, la disposition sanctionnant le « revenge porn », prévue à l'article 226-2-1 du Code pénal, sera plus difficilement applicable aux *deepfakes*. En effet, pour que cette infraction soit constituée, elle suppose l'utilisation d'un contenu à caractère sexuel réel. Bien qu'il soit usuel, en pratique, que les *deepfakes* créés reprennent des contenus à caractère pédopornographique réels, l'application de cette disposition dans ce cadre est plus incertaine.

### Innovations législatives et régulation des plateformes

La loi du 29 juillet 1881 sur la liberté de la presse impose un cadre légal applicable à toute publication et affichage public, en prévoyant des dispositions relatives aux crimes et délits commis par voie de presse ou par tout autre moyen de publication. Elle cherche ainsi à

concilier la liberté d'expression avec la répression des abus. Dès lors, cette loi peut s'appliquer aux abus commis par la publication de *deepfakes* sur Internet, par le délit de diffamation, consistant à sanctionner l'affirmation d'un fait précis portant atteinte à l'honneur ou à la considération d'une personne déterminée ou déterminable, ou encore par l'infraction d'injure. Toutefois, la pleine efficacité de cette loi est en pratique limitée par les peines pouvant sembler dérisoires, ou par le fait que la responsabilité ne pourrait être appliquée aux gestionnaires de plateformes ou hébergeurs.

C'est ainsi que la loi visant à sécuriser et réguler l'espace numérique (dite loi SREN) propose tout d'abord la création d'un délit dit de publication d'hypertrucage (*deepfake*) à caractère sexuel (cette disposition est la plus pertinente en l'espèce, mais la loi met également en place de nouvelles sanctions afin de condamner la haine en ligne, le cyberharcèlement, et crée un délit d'outrage en ligne). Définitivement approuvée par l'Assemblée nationale le 10 avril 2024, la loi a été promulguée le 21 mai 2024, et publiée au Journal officiel le 22 mai 2024. La loi ajoute un article 226-8-1 au Code pénal, punissant de deux ans d'emprisonnement et de 60 000 euros d'amende, « *le fait de porter à la connaissance du public ou d'un tiers, par quelque voie que ce soit, un montage à caractère sexuel réalisé avec les paroles ou l'image d'une*

90. Voir par exemple Avis n° 2015-0001 du 20 janvier 2015 de l'ARCEP sur le décret relatif à la protection des internautes notamment contre les sites provoquant des images et représentations de mineurs à caractère pornographique.

*personne, sans son consentement* ». L'article précise qu'est assimilé à cette infraction et puni de la même manière « *le fait de porter à la connaissance du public ou d'un tiers, par quelque voie que ce soit, un contenu visuel ou sonore à caractère sexuel généré par un traitement algorithmique et reproduisant l'image ou les paroles d'une personne, sans son consentement* ». Ce texte apparaît donc comme une première avancée en ce qu'il sanctionne la diffusion/publication de contenus créés par traitement algorithmique. Il ne semble toutefois pas viser la « *création* » de tels contenus.

Enfin, l'article 226-8 du Code pénal, récemment modifié par la loi SREN, est également pertinent. En effet, depuis le 23 mai 2024, cet article sanctionne « *le fait de porter à la connaissance du public ou d'un tiers, par quelque voie que ce soit, le montage réalisé avec les paroles ou l'image d'une personne sans son consentement, s'il n'apparaît pas à l'évidence qu'il s'agit d'un montage ou s'il n'en est pas expressément fait mention* ». Cette disposition, qui permet de lutter contre les *deepfakes*, notamment pédocriminels, trouve toutefois son application limitée par le fait qu'elle ne s'applique qu'aux montages sur lesquels il existe un doute, et l'ajout d'une mention qu'il s'agit d'un montage permet d'échapper à la sanction prévue.

Avant sa renumérotation par la loi SREN, la loi n°2004-575 du 21 juin

2004 pour la confiance dans l'économie numérique (LCEN) prévoyait un principe de responsabilité limitée lorsque les hébergeurs, désormais appelés fournisseurs de services d'hébergement, rendaient inaccessibles promptement les contenus manifestement illicites signalés (ancien article 6.1.2). Ce principe est repris dans le DSA en son article 6. De plus, la loi n°2023-566 du 7 juillet 2023 visant à instaurer une majorité numérique et à lutter contre la haine en ligne mettait à la charge des hébergeurs de concourir à la lutte contre les *deepfakes*, c'est à dire les atteintes à la représentation de la personne, notamment en mettant en place un dispositif de signalement facilement accessible et visible. Ce principe n'était pas repris dans la dernière modification de la LCEN.

Enfin, l'article 6-1 de la LCEN dispose que lorsque les nécessités de la lutte contre la diffusion des images ou des représentations de mineurs à caractère pédopornographique le justifient, l'autorité administrative peut demander aux fournisseurs de service d'hébergement de retirer les contenus qui contreviennent à l'article 227-23 du Code pénal dans les 24 heures, sous peine de sanctions pénales. En cas d'absence de réaction du fournisseur de service d'hébergement, l'autorité administrative peut s'adresser aux fournisseurs d'accès internet afin de leur demander de procéder au blocage du site en question. En

pratique, les demandes interviennent notamment à la suite de signalements effectués par les internautes via la plateforme en ligne d'harmonisation, d'analyse, de recoupement et d'orientation des signalements (Pharos). Cette disposition peut s'avérer utile afin de bloquer certains contenus, par exemple lorsque le fournisseur de service d'hébergement n'est pas coopératif, mais nécessite un suivi des signalements effectués et une action rapide des autorités compétentes.

## Perspectives

Depuis 2017, le Gouvernement français développe l'IA à travers une stratégie nationale en deux phases (2018-2025) dans le cadre du plan France 2030 ayant pour objectif de faire de la France un leader en innovation. La première phase (2018-2022) a renforcé la recherche en IA, tandis que la seconde phase (2021-2025) se concentre sur l'intégration de l'IA dans l'économie et le soutien à des secteurs prioritaires. Les autorités administratives ne sont pas en reste, et en particulier la CNIL qui publie régulièrement des recommandations sur le développement des systèmes d'IA, ayant pour objectif « *d'aider les professionnels à concilier innovation et respect des droits des personnes* ».

Ainsi, une prise de conscience des risques liés à l'émergence de nouvelles technologies comme l'IA générative semble se faire, et nécessite une

approche coordonnée à l'échelle internationale, européenne et nationale. Si le constant a effectivement été établi, notamment avec l'AI Act, les différentes solutions n'ont pas encore été trouvées. Le cadre juridique reste encore relativement fragmenté, avec différentes dispositions qui sanctionnent des infractions différentes, visent différents acteurs et à différents niveaux. De plus, même si des dispositions permettant de sanctionner certains comportements existent à ce jour, il reste à prévoir les moyens nécessaires afin de les appliquer. Il conviendrait de prévoir un ensemble de dispositions cohérentes qui s'articulent de manière globale et qui mettent en place un régime clair et efficace permettant de lutter contre la pédocriminalité générée par l'IA à tous les niveaux, avec des moyens et des effectifs suffisants afin de mettre en œuvre ce nouveau cadre.

Si l'action des États et des entreprises du secteur des nouvelles technologies est primordiale dans la lutte contre la pédocriminalité générée par l'IA, les citoyens, le grand public, l'entourage de chaque enfant, et en premier lieu les parents ont une part de responsabilité dans la protection de l'enfant en ligne. Quelques bonnes pratiques peuvent être facilement et rapidement mises en place au sein de la cellule familiale. Dans ce cadre, la sensibilisation des parents, mais aussi des enfants est fondamentale.

# Principaux acteurs et

## SOCIÉTÉ CIVILE

### Organisations à but non-lucratif

- Thorn
- All Tech is Human

### Entreprises

- Programme Lantern de la Tech Coalition (Google, Meta, Twitch, Discord etc.)
- Lettre ouverte « Disrupting the Deepfake Supply Chain », 21 février 2021
- Déclaration du 30 octobre 2023 (TikTok, Snapchat, Stability AI etc.)

« Safety by Design for Generative AI: Preventing Child Sexual Abuse » (Thorn, All Tech is Human, Microsoft, Mistral AI etc.)

## ORGANISATION,

### ONU

#### 1. UNESCO:

Guide sur l'IA générative dans l'éducation et la recherche, 2023

#### 2. Union internationale des télécommunications:

nécessaire prise en compte et atténuation des risques posés par l'IA

#### OCDE:

Principes relatifs à l'IA, adoptés en 2018, révisés en 2024

#### G7 (Hiroshima, 2023):

Principes directeurs internationaux relatifs à l'IA et Code de conduite volontaire pour les développeurs d'IA.

#### Virtual Global Taskforce

(alliance internationale de 15 agences de forces de l'ordre): « Technological Tipping Point Reached in Fight against Child Sexual Abuse » préconise une coopération entre les services de police avec les acteurs de l'IA et un travail de prévention sur les dangers créés par l'IA

#### Groupe de Bletchley:

Déclaration de Bletchley sur la sécurité de l'IA (France, Royaume-Uni, UA, États-Unis etc.)

### LÉGENDE

**OCDE:** envergure internationale

**Conseil de l'Europe:** envergure régionale

**Royaume-Uni:** envergure nationale

# initiatives portées

## INSTANCES, INITIATIVES (INTER)GOUVERNEMENTALES

### Conseil de l'Europe

- Stratégie pour les droits de l'enfant (2022-2027), notamment au sein de l'environnement numérique
- Recommandation CM/Rec (2018) du Comité des Ministres sur les Lignes Directrices relatives au respect, à la protection et à la réalisation des droits de l'enfant dans l'environnement numérique
- Convention-cadre sur l'IA et les Droits de l'Homme, la démocratie et l'État de droit, 2024

### Union européenne (UE)

- Digital Services Act (DSA)
- Proposition de règlement établissant des règles en vue de prévenir et de combattre les abus sexuels sur enfants
- Règlement général sur la protection des données (RGPD)

### France – Parlement, gouvernement et autorités administratives indépendantes

- Loi SREN (2024), créant l'article 226-8-1 du Code pénal punissant la diffusion de montages à caractère sexuel sans le consentement de la personne représentée
- Stratégie gouvernementale sur le renforcement de la recherche sur l'IA et son intégration dans l'économie.
- Recommandations de la CNIL sur le développement de systèmes d'IA dans le respect des droits des personnes, dans le cadre de cette stratégie gouvernementale.

### Royaume-Uni – Justice

Condamnation pour création d'images truquées à caractère pédopornographique, avec interdiction d'utiliser les IA génératives pendant 5 ans

### États-Unis – Justice & forces de l'ordre

Sensibilisation du grand public et poursuite de créateurs et distributeurs de contenus pédocriminels générés par l'IA



## Exemple de bonne pratique : Sensibiliser pour mieux prévenir les contenus pédocriminels générés par l'IA

Entretien avec **Églantine Cami**, Chargée de plaider et de sensibilisation à l'Association CAMELEON\*, Mai 2024

**Fondation pour l'Enfance: Quel est le rôle des parents en matière de protection des enfants contre les violences sexuelles en ligne ?**

**Églantine Cami:** Les parents jouent un rôle fondamental en matière de prévention. Tout commence par la mise en place d'un dialogue entre les parents et l'enfant sur les pratiques numériques de chacun et chacune.

Mais pour mettre en place ce dialogue, les parents doivent d'abord avoir conscience des risques existants. Pourtant, beaucoup semblent démunis sur la question, et peinent à prendre conscience de ces risques. Comme pour les violences sexuelles faites aux enfants (et tous les sujets liés à la sexualité plus généralement) il y a un grand tabou autour de la cyberpédocriminalité, renforcé par le fait que les pratiques numériques des enfants sont très difficiles à suivre pour les parents. Mais, s'il y a un tabou, il n'y a plus de dialogue, et la prévention en est compliquée.

---

\*CAMELEON est une association de solidarité internationale, créée il y a 27 ans en France et aux Philippines, pour lutter contre les violences sexuelles faites aux enfants. En France, CAMELEON agit particulièrement contre les violences intrafamiliales et la cyberpédocriminalité. L'association mène des actions de prévention en milieu scolaire et périscolaire, du CP à la Terminale; des actions de communication et de sensibilisation à destination du grand public, particulièrement les parents; et des actions de plaidoyer auprès des pouvoirs publics pour améliorer la loi et les politiques publiques en matière de protection de l'enfance et de lutte contre la cyberpédocriminalité.

Au-delà des pratiques des enfants, il est important que les parents prennent conscience de l'impact de leurs propres pratiques numériques, et notamment le *sharenting*\*. Les parents ne pensent pas forcément aux différents risques associés au partage de photos de leurs enfants, notamment le détournement de ces photos pour en faire des contenus pédocriminels.



### **Fondation pour l'Enfance: Comment accompagnez-vous les parents dans la prévention et le repérage de ces violences ?**

**EC.:** Nous animons des actions de sensibilisation auprès des parents dans différents environnements (en entreprise, dans le cadre de cafés parents en lien avec les établissements scolaires, lors d'événements culturels ou artistiques). Ces actions favorisent une prise de conscience sur les risques existants, et permettent aux parents d'être acteurs de la lutte contre la cyberpédocriminalité.

Nous avons également mené une campagne intitulée « Le Partage » pour alerter les parents sur les risques liés au *sharenting*, et pour sensibiliser le grand public à la cyberpédocriminalité.

Dans le cadre de ces actions, nous encourageons les parents à se questionner sur le respect de la vie privée de leur enfant et sur les risques liés au *sharenting*: ai-je besoin que mon compte

soit en public? est-ce que ces photos peuvent communiquer des informations sur mon enfant? Plus globalement, l'enjeu est de faire comprendre aux parents que l'environnement numérique présente les mêmes risques que le monde réel: si on dit à un enfant qu'il ne faut pas parler à un inconnu dans la rue, c'est pareil en ligne.

 L'enjeu est de faire comprendre aux parents que l'environnement numérique présente les mêmes risques que le monde réel: si on dit à un enfant qu'il ne faut pas parler à un inconnu dans la rue, c'est pareil  en ligne.





## À LA LOUPE

### **Fondation pour l'Enfance: Les parents sont-ils conscients des risques associés à l'IA générative ?**

**EC.:** Comme le grand public de manière générale, les parents ont pu entendre parler du sujet de l'IA générative, mais les risques qui y sont associés restent méconnus.

En revanche, dans nos actions en milieu scolaire nous parlons souvent des *deepfakes* avec les jeunes, car c'est un sujet auquel ils sont confrontés (surtout les filles). Cependant, on constate que pour la majorité d'entre eux « c'est pour rire ». Pour eux, les *deepfakes* ne sont pas des violences, et cette pratique est commune et banalisée (surtout au collège). Pourtant, même si c'est pour rire, il y a une réelle crainte chez les jeunes de voir des montages d'eux sortir. Pour « se protéger » de cela, ils ont eux-mêmes recours à des montages pour avoir une monnaie de chantage (« si tu publies des montages de moi, j'en publie de toi »).

Dans ces situations-là, les enfants n'en parlent pas car ils trouvent cela inutile. De manière générale, les jeunes considèrent que parler aux parents des violences en ligne qu'ils peuvent vivre ne va rien régler. Au contraire, pour eux c'est la honte et ils craignent d'être punis en retour. Il y a donc un vrai climat de défiance entre jeunes, et des jeunes envers les adultes.

### **Fondation pour l'Enfance: Quelles recommandations adressez-vous aux parents ?**

**EC.:** Nous commençons par leur recommander de se renseigner un maximum sur les risques (grâce par exemple au site [jeprotegemonenfant.gouv.fr](http://jeprotegemonenfant.gouv.fr)), mais aussi de mettre en place un climat propice au dialogue avec l'enfant, en s'intéressant à ses pratiques, en posant des questions simples (« tu joues à quoi ? » avec qui joues-tu ? »), en parlant des réseaux sociaux et des jeux comme on parle de l'école, en essayant de jouer avec lui... L'idée est de montrer à son enfant qu'on est une personne de confiance et qu'il peut en parler s'il arrive quelque chose un jour.

Nous avons également créé un « Permis de cyberprudence », une brochure ludique et pédagogique, destinée aux 10-14 ans, qui les encourage à dialoguer avec des adultes de confiance.

Dans la mise en place et l'entretien de ce dialogue et de ce climat de confiance, il y a aussi un enjeu pour le parent de réussir à gérer sa réaction et sa posture face à son enfant qui lui révèle une situation de violence en ligne. Si on montre à l'enfant colère, tristesse ou panique, celui-ci risque de se dire qu'il n'aurait jamais dû en parler. Dans nos ateliers de sensibilisation, nous mettons les parents en situation et ils débattent sur leurs manières de réagir face à de telles situations.

**Fondation pour l'Enfance: Avez-vous une ou plusieurs recommandations à destination des pouvoirs publics et des entreprises du secteur des nouvelles technologies pour améliorer la prévention et le repérage de la pédocriminalité générée par l'IA?**

**EC.:** Tout d'abord, nous appelons les entreprises à s'engager pour la minimisation des risques dans le développement des technologies. Les entreprises doivent être attentives à l'impact que leurs technologies peuvent avoir sur les enfants, et les droits des enfants doivent donc être pris en compte à chaque étape de la conception d'une technologie. La Child Rights Foundation a notamment développé une approche fondée sur les droits de l'enfant dans la conception de produits et services numériques, à partir de la Convention internationale des droits de l'enfant<sup>91</sup>.

Ensuite, nous appelons les pouvoirs publics à mener une campagne nationale de sensibilisation à la cyberpédocriminalité et aux bonnes pratiques à adopter pour protéger les enfants. Cette campagne doit permettre aux parents de mieux appréhender leur rôle en matière de prévention et d'être davantage sensibilisés aux risques liés à l'IA.

Enfin, nous appelons les pouvoirs publics à apporter les évolutions législatives nécessaires. Nous soutenons les recommandations de l'OFMIN sur le fait de mener une réflexion sur les violences commises grâce à l'IA générative, pour mieux incriminer les contenus pédocriminels générés par l'IA. ●

<sup>91</sup>. 5Rights by Design, Child Rights by Design, 2023.

# Conclusion

Si la proportion des contenus pédocriminels générés par l'IA reste faible parmi les contenus signalés, la capacité de production industrielle et la facilité d'accès et d'utilisation des outils d'IA géné-

**La constante amélioration de l'IA générative rend difficilement prédictibles les futures capacités des modèles.**

ratifs modifiés à des fins pédocriminelles nous fait légitimement craindre que la partie immergée de l'iceberg soit colossale. Par ailleurs, la constante amélioration de l'IA générative rend difficilement prédictibles les futures capacités des modèles. Quel niveau de violence pourront atteindre les futurs contenus pédocriminels si aucune mesure de protection des enfants n'est mise en place? Nous ne pouvons le prédire, ni même l'imaginer.

L'émergence et la banalisation de cette nouvelle pratique pose de

nouveaux enjeux à la société française, mais aussi à l'ensemble de la communauté européenne et internationale. Nous avons besoin aujourd'hui de solutions politiques et juridiques, et d'une approche coordonnée des États et des entreprises du secteur des nouvelles technologies au niveau international pour une meilleure prévention de ces contenus, et pour un retrait plus facile et plus rapide.

Du point de vue des pouvoirs publics, nous avons aujourd'hui besoin d'un engagement politique clair pour adapter les cadres légaux au niveau national, européen et international. L'ensemble des possibilités permises par l'IA doivent être prises en compte, et les responsabilités des divers acteurs (États, entreprises, personnes privées) doivent être définies afin d'assurer une réponse coordonnée et efficace.

Les autorités doivent également entreprendre des politiques publiques de prévention et de sensibilisation du grand public. Il s'agit notamment de mettre en œuvre des actions de communication et de soutien à la parentalité à destination des parents, afin de leur permettre de s'emparer de



leur rôle de personnes référentes et de confiance pour leurs enfants sur ce sujet. Il est également fondamental de construire un parcours de soins complet pour les enfants victimes et leurs familles, pour une prise en charge efficace du traumatisme.

Enfin, les constats de ces pratiques appellent les entreprises du secteur des nouvelles technologies à prendre en compte la protection des enfants de manière précoce, avant la mise en ligne de leurs services, et de mettre en place des mesures d'atténuation des risques. En coopération avec les pouvoirs publics, les acteurs

privés doivent investir dans des solutions techniques pour prévenir et répondre aux difficultés opérationnelles rencontrées par les acteurs de la protection de l'enfance.

**Les constats de ces pratiques appellent les entreprises du secteur des nouvelles technologies à prendre en compte la protection des enfants de manière précoce, avant la mise en ligne de leurs services, et de mettre en place des mesures d'atténuation des risques.**



## L'ŒIL DE L'EXPERT

### **Fondation pour l'Enfance: L'IA générative peut-elle aussi être au service de la lutte contre la pédocriminalité ?**

**Nicolas Greffard:** L'IA générative peut-être un bon détecteur d'images pédocriminelles. Il y a 15 ans, il fallait soi-même aller récolter le contenu correspondant pour apprendre au modèle d'IA à reconnaître et à discriminer un contenu à caractère pédopornographique. Aujourd'hui, on peut supposer que les modèles génératifs ont, dans une certaine mesure, cette capacité-là de manière inhérente. Ces modèles peuvent donc aider à les détecter en traitant automatiquement tout ce qui passe sur le web. Les modèles s'améliorent de jour en jour, même si dans certains cas on est encore capable de distinguer une image générée d'une qui ne l'est pas. Il est possible que dans 5 ans cela ne soit plus possible.

Cependant, je ne pense pas qu'une technologie comme celle-ci doive être autorisée. Dès lors qu'elle serait mise à disposition, des gens mal intentionnés pourraient comprendre et apprendre ce qu'elle sait reconnaître et donc apprendre à le combattre, c'est tout le problème des systèmes de détection de fraude: dès lors qu'il y a des solutions automatisées qui existent, on peut apprendre à le combattre. »

“ La constante  
amélioration de  
l'IA générative rend  
difficilement prédictibles  
les futures capacités  
des modèles. ”

Si répondre aux enjeux posés par l'utilisation de l'IA générative est urgent, d'autres nouvelles technologies sont également détournées à des fins pédocriminelles. Les forces de l'ordre françaises, britanniques et américaines ont constaté des viols et agressions sexuelles d'ava-

tars appartenant à des enfants dans les espaces de réalité virtuelle. La route est encore longue, les problématiques nombreuses, mais nous nous devons de nous y atteler rapidement, au nom de la protection des enfants. ●



# Ce que l'explosion de l'IA générative et de ses détournements disent de notre société

Entretien avec **PASCAL PLANTARD**, Professeur d'Anthropologie des Usages, CREAD-M@souin, Université Rennes 2, Mars 2024

**Fondation pour l'Enfance**: L'IA est partout, qu'il s'agisse d'articles de fonds, de publicité, de réflexion politique, d'un usage ludique... Pourquoi tant d'effervescence ?

**Pascal Plantard**: Il existe très peu de travaux en histoire contemporaine des technologies permettant de comprendre cette effervescence.

Le numérique nous imprègne par les représentations, les constructions psychiques qu'il suscite en chacun de nous et dans la société. Les usages des technologies numériques<sup>92</sup> sont donc ancrés dans des univers et des références, produits par les sociétés. Les « entrepreneurs de morale » sont ceux qui proposent ces imaginaires, ces représentations et ces pratiques numériques qui construisent des normes d'usages à la société. Nous sommes tous, plus ou moins, « entrepreneurs de morale » : les entreprises numériques mais aussi les pouvoirs publics, les institutions, la recherche, les associations, les citoyens...

 Les « entrepreneurs de morale » sont ceux qui proposent ces imaginaires, ces représentations et ces pratiques numériques qui construisent des normes d'usages à la société. 



<sup>92</sup>. Défini comme « normes sociales d'usage » par l'anthropologie des usages.



## À LA LOUPE

Dans nos 18 laboratoires de Sciences Humaines et Sociales<sup>93</sup>, nous observons que la dématérialisation administrative préoccupe bien plus les usagers que l'IA. Pourtant les médias et les pouvoirs publics ne parlent que d'IA. On est typiquement dans une « panique morale », une construction médiatique d'un besoin et la mise en scène d'une offre socio-technique.

**Les « entrepreneurs de morale » sont ceux qui proposent ces imaginaires, ces représentations et ces pratiques numériques qui construisent des normes d'usages à la société.**

**Fondation pour l'Enfance: L'IA est présentée comme un facteur de progrès, un élément d'une démarche de conquête. Pourquoi aborde-t-on si rarement les risques de mauvais usages ?**

**PP.:** Les révolutions numériques n'existent pas vraiment. Elles s'installent sur ce qui existe. Les technologies ont une histoire : elles suivent un processus de socialisation (l'innovation) puis

de massification (la banalisation). Il y a eu de nombreuses "révolutions" informatiques depuis 20, 30, 40 ans, et toujours avec le même processus : mettre en scène des technologies, les proposer, puis les utiliser sans questionnement, sans prise de recul. On vend d'abord du rêve puis vient la massification de cette technologie : on transforme l'innovation en consommation.

**Fondation pour l'Enfance: Qu'est-ce que la création et le visionnage de contenus pédocriminels révèlent de notre société, notamment de la vision que nous avons de l'enfant et de la place que nous lui accordons ?**

**PP.:** Tout est devenu marchandisable dans un univers en pleine dysculturation, c'est-à-dire un univers qui comprend à la fois des éléments totalement archaïques (les attirances pédophiles) et totalement modernes (l'IA). L'éducation des enfants est prise dans les contradictions profondes issues des impératifs de protection de l'enfant face aux écrans, et de leur

93. Le GIS M@rsouin.



accompagnement à l'appropriation des technologies dans un monde devenu numérique. Des questions éthiques se posent autour des usages des technologies, qui ne peuvent être laissées aux algorithmes et aux seuls intérêts économiques.



**FE. L'arrivée de l'IA pourrait-elle favoriser ou encourager des comportements pédocriminels davantage que si ce type de contenu n'était pas générable ou consultable ?**

**PP.:** La pornographie est très friande d'innovations technologiques. On peut donc se poser la question : faciliter l'accès à de tels contenus jusqu'à la banalisation ne risque-t-il pas d'augmenter les passages à l'acte ?

**FE. Selon vous quelles pourraient être des pistes sur les plans national, européen, mondial : éducation, contrôle, interdiction ?**

**PP.:** Il faut arrêter les postures contradictoires et se donner les moyens de travailler autour de 5 enjeux majeurs.

Tout d'abord, nous devons lutter au niveau international contre l'économie de l'attention, les algorithmes addictifs et les Usages Problématiques de l'Internet. À l'échelle nationale et sociétale, nous devons également provoquer un ressaisissement éducatif et parental vis-à-vis du numérique. Nous devons aussi accompagner et former les acteurs socio-éducatifs à la médiation numérique, valoriser le travail coopératif et soutenir les accompagnants en donnant accès à des recherches sérieuses sur les usages. Enfin, il est fondamental de travailler dans les territoires, avec les familles, notamment les plus vulnérables.

 **Il est fondamental de travailler dans les territoires, avec les familles, notamment les plus vulnérables.** 

C'est à ce prix qu'adviendra un numérique choisi et inclusif qui se tiendra à distance de la pédocriminalité et des autres formes de cybercriminalité. ●

# Recommandations

Recommandations de la Fondation pour l'Enfance à destination des pouvoirs publics et des entreprises pour prévenir, détecter et sanctionner l'exploitation sexuelle des mineurs générée par l'IA.

**O**n ne saurait trop insister sur l'urgence d'une adaptation législative. Il s'agit de prévoir les défis futurs, mais aussi de relever les défis actuels. La rapidité des avancées technologiques nécessite non seulement la création, mais aussi la révision fréquente et la mise à jour rapide des réglementations. L'utilisation abusive de l'IA étant déjà une réalité tangible, il est essentiel que les législateurs, les organismes chargés de l'application de la loi et les entreprises du secteur des nouvelles technologies agissent avec

diligence. La protection des personnes les plus vulnérables de la société exige une approche proactive et cohérente pour s'assurer que nous ne suivons pas seulement le rythme de l'évolution de l'IA, mais que nous gardons également une longueur d'avance dans la protection des limites éthiques.

La Fondation pour l'Enfance préconise d'accélérer la mise en œuvre de mesures de protection des enfants contre la création et la propagation de contenus pédocriminels générés par l'IA.

**Internet et son aspect sans frontières mettent au défi les législations internes de chaque État, posant des problèmes d'interprétation et d'application de la loi. C'est entendable. Mais la société se doit d'évoluer avec son temps, et dès lors que cela concerne les plus vulnérables, les ajustements doivent aller vite.**

VÉRONIQUE BÉCHU, *Derrière l'écran*, Stock, 2024



## Légende



À destination des pouvoirs publics



À destination des entreprises du secteur des nouvelles technologies

## AXE 1

# Détecter les contenus pédocriminels générés par l'IA



### RECOMMANDATION 1

**Favoriser l'innovation, en incitant les acteurs privés à coopérer pour mettre en place des outils permettant de distinguer les contenus générés par l'IA des contenus non générés par l'IA.**

L'utilisation systématique de logiciels spécialisés dans la détection et la distinction des images non générées par l'IA de celles générées par l'IA pourrait permettre de pallier la difficulté de détection des contenus pédopornographiques, et plus précisément d'identification des mineurs victimes de violences sexuelles directes.



### RECOMMANDATION 2

**Mettre les moyens financiers pour permettre aux offices spécialisés (notamment l'OFMIN) de s'adapter aux nouveaux enjeux et les doter**

**des outils technologiques à la hauteur de ces nouveaux défis.**

Dans l'optique d'identification de potentielles victimes, mettre à disposition des forces de l'ordre un outil permettant de distinguer les contenus générés par l'IA des contenus non générés par l'IA.



### RECOMMANDATION 3

**Instaurer la plus grande coopération entre les différentes entreprises et plateformes (modèles d'IA générative, réseaux sociaux, messageries privées, etc.).**

Le but étant d'améliorer l'identification et le retrait des contenus pédocriminels et des modèles d'IA générative destinés à générer des contenus pédocriminels.



#### RECOMMANDATION 4

### **Acter l'obligation des entreprises fournissant des services en ligne de détecter les contenus pédocriminels, notamment générés par l'IA, présents sur leurs plateformes**

Depuis deux ans, l'Union européenne discute du règlement visant à prévenir et à combattre les abus sexuels commis contre les enfants. Le projet de règlement initial prévoyait la possibilité d'obliger les entreprises fournissant des services en ligne à détecter de manière proactive les contenus pédocriminels présents sur leurs plateformes, y compris des messageries chiffrées de bout en bout. La technologie envisagée et les modalités d'obligations prévues sont proportionnées

et permettent un équilibre entre le respect du droit à la vie privée et la protection de l'enfance en ligne. Nous appelons le gouvernement français et l'ensemble des États européens à se positionner en faveur de ce règlement.



#### RECOMMANDATION 5

**Dans la lignée des recommandations de l'UNESCO, établir des mécanismes de réclamations et de recours pour recueillir les réclamations des utilisateurs de services d'IA générative et du grand public, et un mécanisme de contrôle et de signalement de toute utilisation illícite, et notamment pédocriminelle du service, pour coopérer avec les pouvoirs publics.**

## AXE 2

### **Sanctionner les contenus pédocriminels générés par l'IA**



#### RECOMMANDATION 6

### **Amender l'article 227-23 du Code pénal pour y insérer les fichiers ou représentations issus de l'IA, ou créer une infraction autonome dédiée spécifiquement à ces enjeux**

Il s'agirait de procéder par renvoi à l'article 226-8-1 du Code pénal nouvellement créé.

Aux termes de l'alinéa 1<sup>er</sup> de ce nouvel article :

« Est puni de deux ans d'emprisonnement et de 60 000 euros d'amende le fait de porter à la connaissance du public ou d'un tiers, par quelque voie que ce soit, un montage à caractère sexuel réalisé avec les paroles ou l'image d'une personne, sans son

consentement. Est assimilé à l'infraction mentionnée au présent alinéa et puni des mêmes peines le fait de porter à la connaissance du public ou d'un tiers, par quelque voie que ce soit, un contenu visuel ou sonore à caractère sexuel généré par un traitement algorithmique et reproduisant l'image ou les paroles d'une personne, sans son consentement ».

Pourrait être inséré un nouvel alinéa à l'article 227-23 du Code pénal qui pourrait être rédigé comme suit :

« *Le fait de concevoir, de créer, de diffuser ou de porter à la connaissance du public ou d'un tiers, par quelque voie que ce soit, tout montage, contenu visuel ou sonore à caractère sexuel généré par un traitement algorithmique tel que visé à l'alinéa 1 de l'article 226-8-1 est puni de X ans de prison et X euros d'amende lorsqu'il s'agit de la représentation, de l'image ou de la parole d'un mineur.* »



## RECOMMANDATION 7

**Pénaliser la création et la mise à disposition de modèles d'IA générative destinés à générer des contenus pédocriminels.**

Il pourrait être ajouté au Code pénal un nouvel article, dans la section 5 « Des atteintes aux droits de la personne résultant des fichiers ou traitements informatiques » du Chapitre 6

« Des atteintes à la personnalité », rédigé de la manière suivante :

*“Est puni de X années d'emprisonnement et X euros d'amende le fait de collecter, détenir, traiter ou détourner des données à caractère personnel, afin de créer, générer ou mettre à disposition du public ou de tout tiers un modèle de traitement algorithmique, dans le but de permettre la création de contenu visuel ou sonore à caractère sexuel représentant un mineur, et de tout fichier à caractère pédopornographique.”*



## RECOMMANDATION 8

**Uniformiser, aux niveaux européen et international, les politiques et réglementations à l'égard de l'exploitation sexuelle des enfants en ligne, afin de s'assurer d'une réponse conjointe et coordonnée.**

La création et le lancement en 2022 du Laboratoire pour la protection de l'enfance en ligne s'inscrit dans cette démarche. Ce Laboratoire rassemble des États, des entreprises du secteur des nouvelles technologies, mais également des organisations de la société civile, dont la Fondation pour l'Enfance.

### AXE 3

## Prévenir l'exploitation sexuelle en ligne des enfants générée par l'IA



### RECOMMANDATION 9

**Mettre en place des campagnes nationales de sensibilisation du grand public à la cyberpédocriminalité, aux dangers relatifs au *sharenting* et aux bonnes pratiques à adopter pour protéger les enfants.**

Cette campagne doit permettre aux parents de mieux appréhender leur rôle en matière de prévention et d'être davantage sensibilisés aux risques liés à l'IA. Les outils déjà existants, tels que la plateforme gouvernementale [jeprotegemonenfant.gouv.fr](http://jeprotegemonenfant.gouv.fr), dont la Fondation pour l'Enfance est membre, pourrait notamment être mise à jour avec une section sur l'IA/un espace de sensibilisation sur les dangers de l'IA. Les plateformes de partage de contenus ou de signalement (comme PHAROS) pourraient être associées à ces campagnes.



### RECOMMANDATION 10

**Intégrer aux interventions à destination des élèves les enjeux liés à l'IA générative, généraliser et rendre obligatoire ces interventions à tous les établissements.**



### RECOMMANDATION 11

**Inclure les droits des enfants et leur protection dans toute réflexion liée au développement des nouvelles technologies.**

En ce sens, l'adoption et la mise en œuvre proactive, suivie et évaluée des principes "Safety by Design" pourraient permettre de minimiser les risques posés à la sécurité des enfants dans le développement des technologies. Les droits des enfants doivent également être pris en compte à chaque étape de la conception d'une technologie.



### RECOMMANDATION 12

**Développer des techniques d'obscurcissement et de floutage pour protéger les fichiers visuels et les droits des individus concernés.**

Il conviendrait ainsi de mettre à disposition de l'ensemble des utilisateurs des outils numériques et de services en ligne une technique d'obscurcissement des photos afin qu'elles ne puissent pas être récupérées et détournées.



### RECOMMANDATION 13

**Améliorer la prise en charge des personnes susceptibles de développer un comportement pédocriminel, notamment en diffusant le dispositif de Service Téléphonique d'Orientation et de Prévention (STOP).**



### RECOMMANDATION 14

**Auditer et évaluer les systèmes d'IA existants et les risques possibles en matière de pédocriminalité.**

Il s'agirait de développer une méthodologie adaptée pour évaluer les systèmes d'IA génératives et s'assurer que ceux-ci prennent suffisamment en considération la sécurité des enfants dès leur conception et s'adaptent aux derniers développements technologiques.





# Glossaire

## **ChatBot**

Logiciel qui utilise l'intelligence artificielle et l'apprentissage automatique pour simuler une conversation humaine. Il peut interagir avec les utilisateurs en temps réel, répondre à des questions, fournir des informations ou exécuter des tâches en fonction des données fournies par l'utilisateur, souvent par le biais de commandes textuelles ou vocales. Il est déployé sur des logiciels tels que Skype, Messenger, ou des assistants virtuels comme Alexa.

## **Clear web (web de surface)**

Région d'Internet que la plupart des personnes connaissent et utilisent, il s'agit de l'ensemble des pages web accessibles au public et qui sont largement indexées sur les moteurs de recherche

## **Companion app**

Logiciel qui remplit directement une fonction pour les utilisateurs, par exemple dans le cas des jeux vidéo, la *companion app* lie l'expérience dans le jeu au smartphone pour débloquer de nouveaux contenus ou améliorer l'expérience, permettant d'ajouter de l'interaction et d'en voir plus sur l'univers du jeu vidéo. Les *companion apps* repose sur des ChatBots, mais

à la différence de ces derniers, elles sont conçues pour pouvoir s'adapter à l'utilisateur et apporter des réponses hyper personnalisées.

## **Dark web**

Région d'Internet cachée intentionnellement et en toute sécurité. C'est le domaine du web où l'anonymat est essentiel, de sorte que l'activité criminelle y est répandue. À distinguer du *deep web* qui est accessible uniquement avec un identifiant et un mot de passe associé ex. *espace client bancaire*.

## **Données d'entraînements (training dataset)**

Dans le contexte des modèles d'IA, les « données d'entraînement » désignent l'ensemble initial de données utilisé pour aider le modèle à apprendre et à faire des prédictions ou à prendre des décisions, comme des images, des vidéos ou des sons réels. Ces données fournissent des exemples à partir desquels le système d'IA peut apprendre des modèles, des comportements ou des relations. La qualité et la quantité des données d'apprentissage jouent un rôle crucial dans la détermination des performances d'un modèle d'IA. Par exemple, un modèle qui produit des

images photoréalistes de personnes aura été entraîné sur des ensembles de données comprenant des photographies préexistantes de haute qualité de personnes réelles.

### **Grooming**

Sollicitation d'enfants à des fins sexuelles en français, le *grooming* désigne le processus par lequel un adulte aborde un mineur et le manipule à des fins sexuelles. Les *groomers* tentent d'établir un rapport de confiance avec l'enfant afin de l'amener progressivement à une conversation et à des actes à connotation sexuelle.

### **Hash/hashing (hachage)**

Le *hashing* est un outil de chiffrement pour transformer les données. Celles-ci sont décomposées et transformées en une nouvelle forme appelée « valeur de hachage ». Cette valeur n'est pas chiffrée, elle est transformée et ne peut donc pas être déchiffrée ni reconvertie au format d'origine sans clé appropriée, sans l'algorithme utilisé et sans les données d'origine associées aux valeurs de hachage. Même si elles tombent entre de mauvaises mains, les cybercriminels ne peuvent rien en faire.

### **Deepfakes (Hypertrucages)**

Technique de synthèse multimédia reposant sur l'IA, et consistant à

superposer des traits humains sur le corps d'une autre personne – et/ou à manipuler les sons – pour générer une expérience réaliste.

### **Instruction (prompt)**

phrase, paragraphe qui décrit une tâche à effectuer ou une question à laquelle répondre. Il est utilisé pour communiquer avec les modèles d'IA et donner des instructions sur la tâche à accomplir. On parle de *system prompt* pour les instructions données en amont par les entreprises au modèle, et de *user prompt* pour les instructions données par les utilisateurs.

### **Intelligence artificielle**

L'intelligence artificielle (IA) est une branche de l'informatique qui vise à créer des systèmes capables d'effectuer des tâches qui nécessitent normalement l'intelligence humaine. C'est, par exemple, le cas de l'apprentissage et de la compréhension du langage, de la reconnaissance de formes et d'images, de la résolution de problèmes complexes et de la prise de décision.

### **Intelligence artificielle générative**

L'intelligence artificielle générative (IAG) est le domaine de l'intelligence artificielle qui se concentre sur la création de nouveaux contenus à partir de données existantes.

## Modèle d'intelligence artificielle

Un modèle d'IA peut être décrit comme un fichier contenant un plan numérique, créé à l'aide de l'intelligence artificielle. Les modèles d'IA peuvent générer ou manipuler divers médias, notamment des images, des vidéos ou même des sons. Lorsqu'ils reçoivent des données d'entraînement de haute qualité, ces modèles peuvent produire des résultats réalistes.

## Modèle de langage de grande taille (*large language model*)

Un *large language model* (ou LLM) est un type de programme d'intelligence artificielle capable de reconnaître et de générer du texte d'après une instruction donnée. Tel un interlocuteur très intelligent, un LLM crée des textes qui semblent écrits par un être humain. Certains LLM sont capables de répondre à des questions, de rédiger des dissertations, d'écrire de la poésie, voire de générer du code. Les LLM sont entraînés sur des ensembles de données massifs, leur permettant de reconnaître et d'interpréter le langage humain ou d'autres types de données complexes ex. *ChatGPT*, *Bing Chat*, *Mistral AI*.

## Pédocriminalité

Dans le contexte de ce rapport, le terme « pédocriminalité » est utilisé pour désigner l'ensemble des

violences sexuelles commises sur un mineur de manière générale, en ligne ou hors ligne, qu'elles soient reconnues et condamnées par la justice, ou non.

## Sextorsion

Contraction d'extorsion sexuelle, ce terme désigne le chantage réalisé à l'aide de contenus (images ou vidéos) de la victime et présentant un caractère sexuel, en vue de lui extorquer des faveurs sexuelles, de l'argent ou tout autre avantage, en menaçant de les diffuser sans son consentement.

## Sharenting

Mot-valise anglais composé de *share* (partager) et *parenting* (parentalité), le *sharenting* est une pratique qui consiste pour les parents à publier des photos ou des vidéos de leurs enfants sur les réseaux sociaux.

---

### Sources

Lalla, V., Mitrani, A., Harned, Z., « Intelligence artificielle: les deepfakes dans l'industrie du divertissement » Magazine OMPi, Juillet 2022.

« Écrire des prompts efficaces pour ChatGPT – conseils pratiques », Campus région du numérique Auvergne Rhône Alpes.

Glossaire de l'IA, Salesforce.

Childfocus Belgique.

Ionos, « Hashing: voici comment fonctionne le hachage », 22 février 2023.

Point de contact.

CNIL.

# Comité de rédaction

## ANGÈLE LEFRANC

Chargée de plaider de la Fondation pour l'Enfance

## ODILE NAUDIN

Administratrice de la Fondation pour l'Enfance

## MAÎTRE CÉLINE ASTOLFE

Cabinet Lombard, Baratelli, Astolfe & associés

## MAÎTRE LÉA

### LEVAVASSEUR-PRUDENCE

Cabinet Lombard, Baratelli, Astolfe & associés

## MAÎTRE ALEX SALEHI

Avocat

---

# Remerciements

*Dans le cadre de ce rapport, la Fondation pour l'Enfance a rencontré divers acteurs institutionnels et de la société civile qui nous ont apporté de précieux éclairages. Nous adressons nos plus sincères remerciements à :*

CAMELEON Association, et notamment **Eglantine Cami** (Chargée de plaider et de sensibilisation)

### Joanna Smith

(Psychologue clinicienne)

**Mélanie Dupont** (Docteur en psychologie, Psychologue à l'Unité Médico-Judiciaire de l'Hôtel-Dieu (Paris) et Présidente de l'Association contre les Violences sur Mineurs (CVM))

Le National Center for Missing & Exploited Children (NCMEC), et notamment **John Shehan** (Senior Vice President, Exploited Children Division & International Engagement)

**Nelly Deflisque**, Journaliste et autrice, spécialiste des questions de société, d'éducation et de droits humains

L'Office mineurs (OFMIN) de la Direction Nationale de la Police Judiciaire, et notamment **Gabrielle Hazan** (Cheffe), **Véronique Béchu** (Cheffe du pôle stratégique) et **Typhaine Desbordes** (Adjointe au chef du bureau des partenariats et de la communication)

### Pascal Plantard

(Professeur d'Anthropologie des Usages, CREAD-M@souin à l'Université Rennes 2)

Point de Contact, et notamment

**Alejandra Mariscal-Lopez** (Directrice) et **Flora Mateo** (ancienne Responsable de la communication et des partenariats)

La UK **Online CSEA Covert Intelligence Team** (OCCIT)

**Unité Erios** [CRIA VS Aquitaine]

L'entreprise Valeuriad, et notamment **Nicolas Greffard** (Directeur Technique, Expert IA)

# Présentation

La Fondation pour l'Enfance est née en 2012 de la fusion de la Fondation pour l'Enfance, fondée en 1977 par Madame Anne-Aymone Giscard d'Estaing, alors Première Dame, et de la Fondation Protection de l'Enfance.

Reconnue d'utilité publique, indépendante et non-partisane, la Fondation pour l'Enfance agit pour améliorer la protection des enfants et le respect de leurs droits fondamentaux, en luttant contre toutes les formes de violences et de maltraitance, et en favorisant des liens adultes-enfants de qualité. Tous ses positionnements



et ses recommandations sont validés avec des experts (médecins, sociologues, psychologues, professionnels de la petite enfance, avocats etc.)

La Fondation pour l'Enfance intervient auprès de l'ensemble des acteurs institutionnels, associatifs et privés qui agissent dans le secteur de l'enfance. Elle agit à travers des actions de plaidoyer auprès des pouvoirs publics (en propre ou via des collectifs interassociatifs) et des actions de sensibilisation et de prévention auprès du grand public et des professionnels (médecins, assistant.e.s maternel.le.s, sage-femmes etc.).

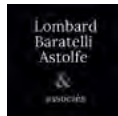
## PRÉSENTATION

### Cabinet Lombard, Baratelli, Astolfo & associés

Fort de plusieurs décennies d'expérience et de succès depuis sa création, le cabinet Lombard, Baratelli, Astolfo & associés est l'un des acteurs incontournables du paysage judiciaire français. Son expertise s'étend du droit pénal général et financier, du travail, de la santé, de la famille, de la presse et des médias, à la protection des données personnelles, la compliance et l'éthique.

Le Cabinet Lombard, Baratelli, Astolfo & associés représente la Fondation pour l'Enfance dans ses constitutions de parties civiles depuis près de 20 ans.

## Nos partenaires







**FONDATION**  
**POUR**  
**L'ENFANCE**

reconnue d'utilité publique