

Online Child Sexual Exploitation

CS 152 — Trust and Safety

Alex Stamos

Stanford Cyber Policy Center

2026

Content Warning:

This lecture will contain frank discussion of criminal activity that many people may find challenging or upsetting. We will also mention cases that led to victim suicide. **We will not show images of this activity during the lecture.** Our learning objective is to understand how platforms are misused to commit offences and the responses of individuals, law enforcement agencies, governments, and the platforms themselves to minimize offense and victimization.

- How technology is misused to exploit children
- How technology is deployed for good to protect them
- CSAM: detection, reporting, NCMEC's role
- Grooming and sextortion of minors
- **AI-assisted OCSE: the fastest-growing front**

Online Child
Sexual
Exploitation

Alex Stamos

CSE and CSAM

Reporting
Obligation

Responses and
Mitigations

Grooming and
Sextortion of
Minors

AI-Assisted
OCSE

What's Next?

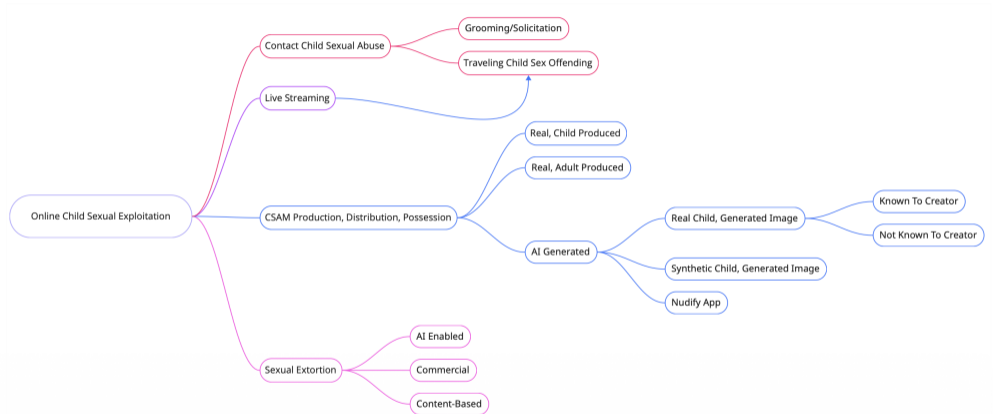
Questions

CSE and CSAM

CSE: *Child Sexual Exploitation*

CSAM: *Child Sexual Abuse Material*

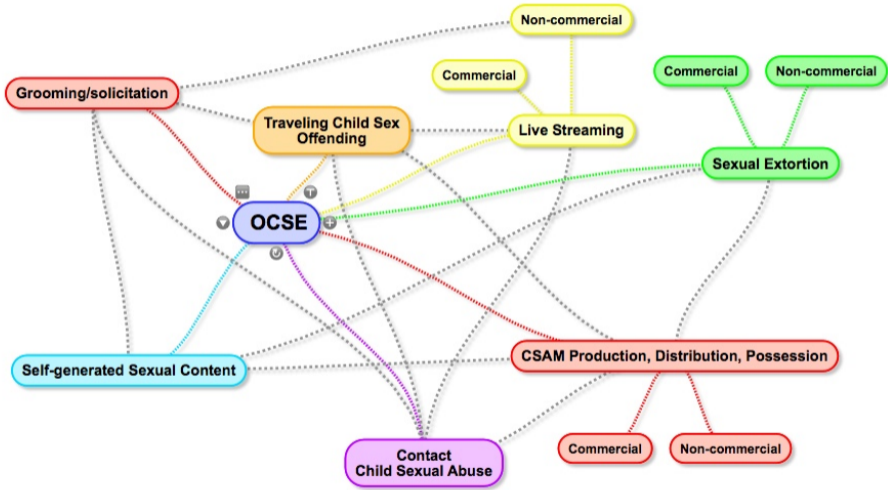
How Does CSE Manifest Online?



Stamos and Grossman (2025) adopted from Baines (2018)

[Stamos and Grossman (2025). adopted from Baines (2018)]

OCSE Categories Are Deeply Interconnected



A grooming relationship may produce self-generated content, which is then used for

Where/How Are Offenses Committed?

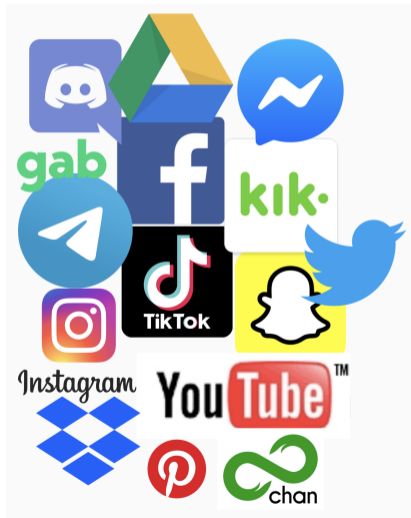
The short answer is: **everywhere.**

CSAM production:

- Abuse recorded offline, then distributed online
- Solicited/coerced online from children
- **Generated entirely by AI from public photos**
(new)

CSAM distribution:

- Social platforms (“severe memes”)
- P2P filesharing - volume offending
- Dark Net and hidden services
- Commercial vs. non-commercial
- “Live streamed” abuse



- Every piece of CSAM can be linked to the exploitation of a child who was coerced or victimized into sharing a photo
- Saving and sharing such imagery extends the worst moment of victims' lives
- Viewing CSAM creates demand for more CSAM, which leads to more exploitation of children

"I wonder if the men I pass in the grocery store have seen them. Because the most intimate parts of me are being viewed by thousands of strangers, and traded around, I feel out of control. They are trading my trauma around like treats at a party, but it is far from innocent. It feels like I am being raped by each and every one of them."

-Vicky (a pseudonym)

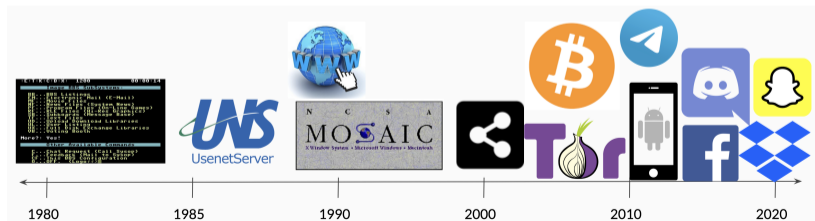
Project Arachnid



A Brief History: Five Overlapping Eras (Chad Steele)

Technology trends toward increased anonymity, storage, space and scale

1. The Networking Era (1987-1996)	2. The Internet Goes Mainstream (1996-2004)	3. Peer-to-Peer Software (2004-2008)	4. Dark Web Technologies (2008-2014)	5. Mobile Consumption (2014-2025)
---	---	--	--	---



What makes CSE/CSAM different from other forms of abuse?

Online Child
Sexual
Exploitation

Alex Stamos

CSE and CSAM

Reporting
Obligation

Responses and
Mitigations

Grooming and
Sextortion of
Minors

AI-Assisted
OCSE

What's Next?

Questions

Everyone agrees it's bad

- This means you can rely upon local law
- Not politically complicated
- Massive incentive to over-enforce

- CSAM Statutes:
 - 18 U.S.C. 2252/2258A/2258B
 - “Child porn statutes” (called CSAM by practitioners).
- Beyond **requirement to remove**, have a **reporting obligation**

For the purposes of this chapter, the term—

① “minor” means any person under the age of eighteen years;

(2)(A) Except as provided in subparagraph (B), “sexually explicit conduct” means actual or simulated—

- ① sexual intercourse, including genital-genital, oral-genital, anal-genital, or oral-anal, whether between persons of the same or opposite sex;
- ② bestiality;
- ③ masturbation;
- ④ sadistic or masochistic abuse; or
- ⑤ lascivious exhibition of the anus, genitals, or pubic area of any person;

How do courts determine if a photo is lascivious exhibition?

The “Dorr factors” are considered by the jury and judge:

- 1 Whether the focal point of the visual depiction is on the child’s genitalia or pubic area.
- 2 Whether the setting of the visual depiction is sexually suggestive, i.e., in a place or pose generally associated with sexual activity.
- 3 Whether the child is depicted in an unnatural pose, or in inappropriate attire, considering the age of the child.
- 4 Whether the child is fully or partially clothed, or nude.
- 5 Whether the visual depiction suggests sexual coyness or a willingness to engage in sexual activity.
- 6 Whether the visual depiction is intended or designed to elicit a sexual response in the viewer.

Do not need all factors, different courts use this in different ways

(a) Duty To Report.—

(1) In general.—

(A) Duty.—In order to reduce the proliferation of online child sexual exploitation and to prevent the online sexual exploitation of children, a provider—

- (i) shall, as soon as reasonably possible after obtaining actual knowledge of any facts or circumstances described in paragraph (2)(A), take the actions described in subparagraph (B); and
- (ii) may, after obtaining actual knowledge of any facts or circumstances described in paragraph (2)(B), take the actions described in subparagraph (B).

(B) Actions described.—The actions described in this subparagraph are—

- (i) providing to the CyberTipline of NCMEC, or any successor to the CyberTipline operated by NCMEC, the mailing address, telephone number, facsimile number, electronic mailing address of, and individual point of contact for, such provider; and
- (ii) making a report of such facts or circumstances to the CyberTipline, or any successor to the CyberTipline operated by NCMEC.

Dimensions on Which CSAM Laws Vary

- **Is there a law?** In 2018, most recent year for which we have data, 16 countries lacked laws specifically about CSAM (Iran, Iraq, Kuwait, Lebanon, Libya, etc - some prohibit all pornography, eg Iran)
- **Is CSAM defined?** 51 countries had laws but did not define CSAM
- **Is tech-facilitated CSAM prohibited?** 25 countries had laws but did not specify prohibition of technology-facilitated CSAM
- **Is possession prohibited?** 38 countries had laws but did not prohibit possession of CSAM (eg Russia)
- **Is fictional CSAM prohibited?** Variation in whether laws prohibit fictional CSAM (legal in Japan, Denmark, elsewhere)
- **Is ISP reporting required?**

Dimensions on Which CSAM Laws Vary

Country	Legislation Specific to CSAM	"Child Sexual Abuse Material" Defined	Technology-Facilitated CSAM Offenses	Simple Possession	ISP Reporting
Hungary	✓	✓	✓ ²⁰¹	✓	✗ ²⁰²
Iceland	✓	✓	✓	✓	✗
India	✓	✓	✓	✓	✓
Indonesia	✓	✗	✓ ²⁰³	✓	✗
Iran	✗	✗	✗	✗	✗
Iraq	✗	✗	✗	✗	✗
Ireland	✓	✓	✓	✓	✗ ²⁰⁴
Israel	✓	✗	✓ ²⁰⁵	✓	✗

CSE and CSAM

Reporting
Obligation

Responses and
Mitigations

Grooming and
Sextortion of
Minors

AI-Assisted
OCSE

What's Next?

Questions

What Does That Mean for Platform Policies?

- Difficulty in determining lascivious means that many public platforms ban all photos of nude children
- Private services like iCloud Backup often do not have specific policies
- Public/private services like Google Photos might _____ all nudity but then apply a higher standard on sharing
- Many platforms only look for “worst” content

CSE and CSAM

Reporting
Obligation

Responses and
Mitigations

Grooming and
Sextortion of
Minors

AI-Assisted
OCSE

What's Next?

Questions

	A	B
1	Prepubescent minor, engaged in a sexual act	Pubescent minor, engaged in a sexual act
2	Prepubescent minor, engaged in lascivious exhibition	Pubescent minor, engaged in lascivious exhibition

CSE and CSAM

Reporting
Obligation

Responses and
Mitigations

Grooming and
Sextortion of
Minors

AI-Assisted
OCSE

What's Next?

Questions

lizard

Online Child
Sexual
Exploitation

Alex Stamos

CSE and CSAM

**Reporting
Obligation**

Responses and
Mitigations

Grooming and
Sextortion of
Minors

AI-Assisted
OCSE

What's Next?

Questions

Reporting Obligation

But what about Section 230 immunity???

Breaking federal criminal law is NOT immunized by Section 230.

CSAM, terrorism, trafficking, etc.

National Center for Missing and Exploited Children (NCMEC)

Online Child
Sexual
Exploitation

Alex Stamos

CSE and CSAM

Reporting
Obligation

Responses and
Mitigations

Grooming and
Sextortion of
Minors

AI-Assisted
OCSE

What's Next?

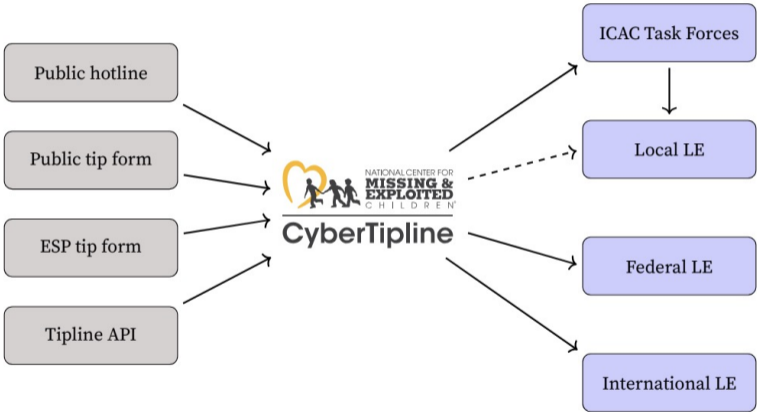
Questions



Is a child being sexually exploited online?

[Report it here.](https://www.cyberTipline.org)

Modern Responses to Known CSAM



NCMEC CyberTipline volume by user location, CY 2022–2024

CyberTipline Reports Relating to U.S. and International Users	CY 2022	CY 2023	CY 2024
U.S.	1,920,963 6%	1,621,937 4.5%	1,276,041 6.2%
International	28,825,411 89.9%	33,220,027 91.7%	17,279,122 84.2%
Unknown Location ³	1,312,655 4.1%	1,368,404 3.8%	1,957,640 9.5%
Total	32,059,029	36,210,368	20,512,803

Vast Majority Are the Trading of Known CSAM

CyberTipline Reports by Reported Incident Type

Child Pornography (possession, manufacture, and distribution)	19,854,300	■
Extraterritorial Child Sexual Abuse and Exploitation	2,911	■
Child Sex Trafficking	26,823	■
Child Sexual Molestation	22,242	■
Misleading Domain Name	14,496	■
Misleading Words or Digital Images on the Internet	12,216	■
Online Enticement of Children for Sexual Acts	546,333	■
Unsolicited Obscene Material Sent to a Child	33,482	■
Total	20,512,803	

Online Child
Sexual
Exploitation

Alex Stamos

CSE and CSAM

Reporting
Obligation

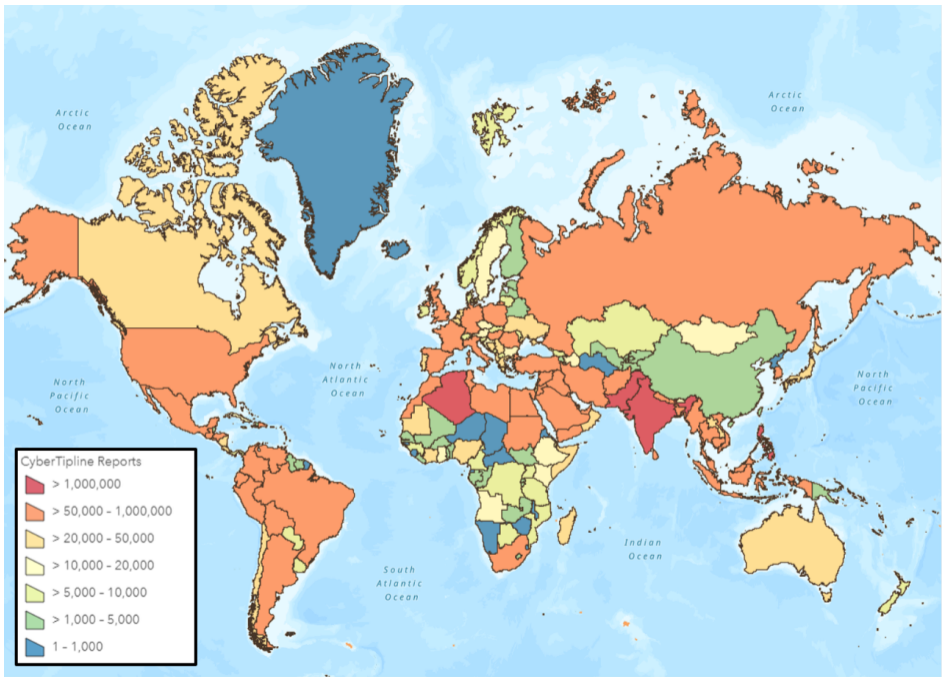
Responses and
Mitigations

Grooming and
Sextortion of
Minors

AI-Assisted
OCSE

What's Next?

Questions



“Malicious” users

- Preferential offenders (motivation: sexual interest in children)
- Commercial offenders (financial motive)
- Situational offenders (opportunistic)

“Nonmalicious” users

- Unintentional offenders (sharing image with statement of abhorrence; meme; vigilante)
- Minor nonexploitative users (eg sexting between teens)
- Situational “risky” offenders (sharing adult pornography, that included an image of a child, maybe a 16 year old who might look like an adult)

“We evaluated 150 accounts that we reported to NCMEC for uploading CSAM in July and August of 2020 and January 2021, and we estimate that more than **75% of these did not exhibit malicious intent** (i.e. did not intend to harm a child), but appeared to share for other reasons, such as outrage or poor humor.”

A Typology of Online CSE Offenders — Indirect

Type	Networking	Security
Browser — stumbles on material accidentally	Nil	Nil
Private fantasy — creates digital images for own use	Nil	Nil
Trawler — actively seeks via open browsers	Low	Nil
Non-secure collector — P2P and chat	High	Nil
Secure collector — encrypted networks, hidden services	High	High

All five are *indirect* offenders — they collect or view but do not directly abuse a child.

A Typology of Online CSE Offenders — Direct

Type	Networking	Security
Groomer — builds online relationship to enable abuse	Varies	Depends on child
Physical abuser — offline abuse; may record privately	Varies	Physical contact only
Producer — records abuse to share or sell	Varies	Varies
Distributor — disseminates material at any level	Varies	Tends to be secure

These categories *directly* contact, abuse, produce, or disseminate.

Online Child
Sexual
Exploitation

Alex Stamos

CSE and CSAM

Reporting
Obligation

Responses and
Mitigations

Grooming and
Sextortion of
Minors

AI-Assisted
OCSE

What's Next?

Questions

Responses and Mitigations

- 1 Reach out to NCMEC early
- 2 Create detailed policies for relevant type of CSAM/CSE
- 3 Enlist lawyers and policy teams who understand this area
- 4 Invest in coordinating bodies
- 5 Publish transparency reports on child harm
- 6 Create resiliency programs for employees working on these issues

A Global Phenomenon Requires a Global Response



[Adapted from Baines (2018)]

Online Child
Sexual
Exploitation

Alex Stamos

CSE and CSAM

Reporting
Obligation

Responses and
Mitigations

Grooming and
Sextortion of
Minors

AI-Assisted
OCSE

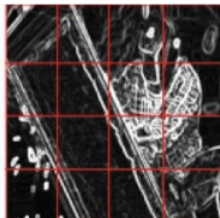
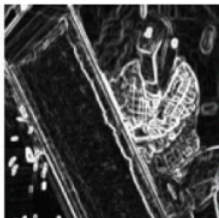
What's Next?

Questions

- 1 Proactively detect and report known CSAM
- 2 Provide reporting mechanisms for users and victims
- 3 Flag potential grooming behaviors using known indicators
- 4 Enforce identity indicators around age
- 5 Restrict discovery of children by unfamiliar adults

596

FARID



CSE and CSAM

Reporting
Obligation

Responses and
Mitigations

Grooming and
Sextortion of
Minors

AI-Assisted
OCSE

What's Next?

Questions

Perceptual Hashes vs. Manipulation

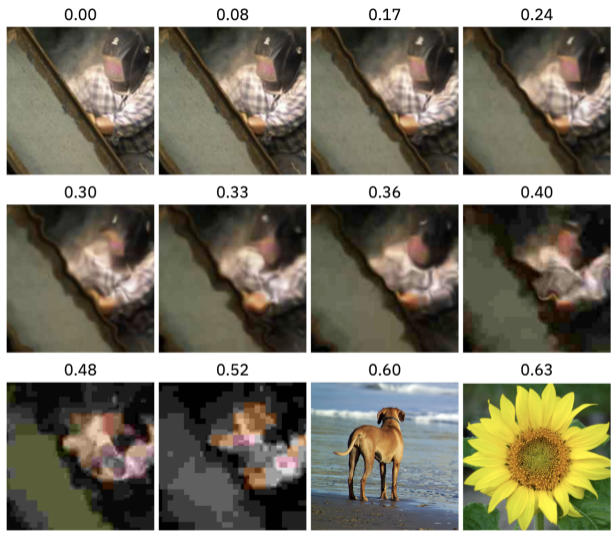


Figure 4: Ten versions of the original welder image (top left) with increasing levels of photometric and geometric distortions. Shown above each image is the numeric hash distance relative to the original (the larger the hash distance, the more dissimilar the

Automated Classification:

- Tools like PhotoDNA are great at identifying known images if they're fed the correct hashes.
- But what about new CSAM generated by children who have been groomed or coerced online?
- Is there more we can do to intervene before material is shared, i.e. in the "chat" phase?

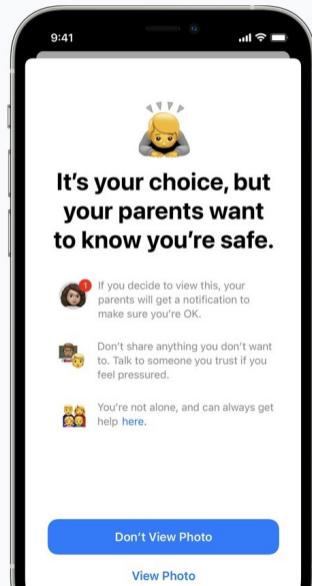
Human review, meanwhile is:

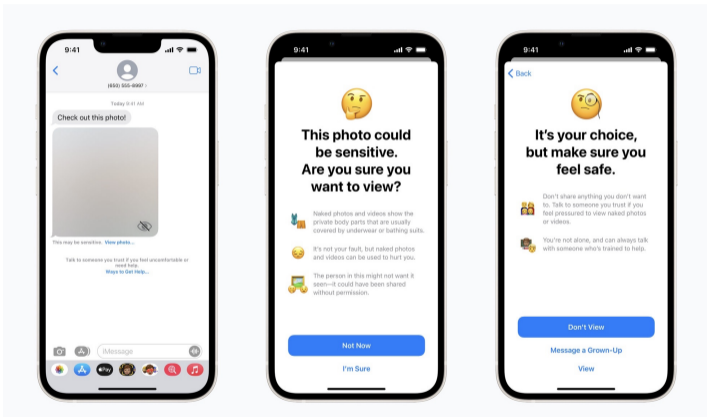
- Fallible and time consuming
- Liable to be challenged by cultural differences (so is technology)
- Emotionally tough, requiring companies to think about their reviewers' resilience and well-being.

Why are there trade offs in designing for user privacy and designing for user safety?

The Apple Controversy

- Apple announces plan to:
 - Alert parents when child sends or receives nude photo on iMessage
 - Blur nude photos children receive
 - Scan iCloud Photos (on device) for known CSAM
 - Report to NCMEC if threshold is reached
- Criticisms:
 - Known CSAM dataset could be expanded to include...political dissent images
 - Apple says it will let outsiders audit the image dataset, but they don't have a reputation for doing this
 - Parental alert could out a child





- Scan profile and group photos
- Examine reported content
- Machine learning classifiers scan unencrypted text (e.g. group description)
- WhatsApp says they ban 300,000 accounts for CSAM activity per month
- Receive tips from law enforcement when WhatsApp links found in dark web forums

Privacy and Safety Sometimes Go Together

“Meta must also limit recommendation engines and discoverability. ‘People You May Know’ has for years been understood as problematic within Facebook for its propensity to make inappropriate suggestions: recommending a therapist’s clients to each other, or potential targets to a possible abuser.”

“One way to prevent such features from leading to unwanted and unsafe social connections would be to make it so that users with no social connection or those surfaced via search or People You May Know can start chats without end-to-end encryption, and then have the option to mutually opt in.”

-David Thiel

An important distinction: not all sex offenders who target children are pedophiles, and not all pedophiles commit sexual offenses.

- Wurtele et al. (2014): Among men, 6% indicated some likelihood of having sex with a child if they were guaranteed they would not be caught or punished, as did 2% of women. Nine percent of males and 3% of females indicated some likelihood of viewing CSAM on the Internet. Overall, nearly 10% of males and 4% of females reported some likelihood of having sex with children or viewing CSAM.
- Self-report perpetration surveys conducted in Australia, Canada, Sweden, the UK and the US show that between 4% (Seto et al., 2014) and 12% of men, and 3% of women in the general population engage with CSAM (Seigfried-Spellar, 2014). [Reported in Wager et al. 2018]

Online Child
Sexual
Exploitation

Alex Stamos

CSE and CSAM

Reporting
Obligation

Responses and
Mitigations

Grooming and
Sextortion of
Minors

AI-Assisted
OCSE

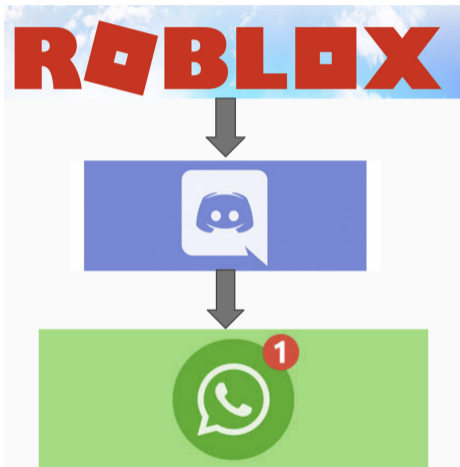
What's Next?

Questions

Grooming and Sextortion of Minors

- **Grooming:** creating a relationship with a child/young person to earn their trust and then exploit and/or abuse them
- **Sextortion:** when an offender threatens to distribute a victim's images/videos of a sexual nature unless they provide money or further sexual content

Leveraging Multiple Technologies and Platforms



- Social and gaming platforms are used to make initial contact.
- Children may then be encouraged to move to less open/more secure messaging apps.
- Children may be groomed to share sexual images/videos they have taken themselves.
- This can rapidly develop into the sextortion blackmail scenario.
- Children may also be groomed for offline meets for sexual activity.

Modern Day: Social Media, Livestream Abuse, and Video Games

Online Child
Sexual
Exploitation

Alex Stamos

CSE and CSAM

Reporting
Obligation

Responses and
Mitigations

Grooming and
Sextortion of
Minors

AI-Assisted
OCSE

What's Next?

Questions



- Took CS152
- Met “Sophia” on Tinder
- Talk for 3 days-set up a skype
- She has no camera-> begins masturbating and encourages him to do the same
- Demanded money or send video to “everyone”



ENGINEERING

Scale of the crisis:

- 2024: FBI received ~55,000 sextortion reports, \$33.5M losses (59% increase from 2023)
- First half 2025: 23,593 reports (vs 13,842 in H1 2024)
- NCMEC aware of **36+ teen suicides** from sextortion victimization

Who is targeted:

- Primary targets: **Boys ages 14–17**
- 30% of victims face demands within 24 hours of contact
- 81% of threats occur exclusively online
- Demands on minors higher than on young adults

The 764 Network:

- FBI now classifies as “tier one” terrorism threat
- Organized networks specifically seeking to drive victims to self-harm
- Often based overseas (Philippines, West Africa, etc.)
- Pose as young girls using fake profile images

Tactics:

- Initial contact on Instagram, Snapchat, gaming platforms
- Move victims to E2E encrypted messaging apps
- Payment demands via Cash App, gift cards



FBI public-awareness campaign targeting minor victims and bystanders.
(<https://tips.fbi.gov>)

- Run by NCMEC: <takeitdown.ncmec.org>
- Uses **client-side hashing** so the image never leaves the child's device
- Companies opt in for proactive blocking; reactive submissions can flow into standard CSAM hash sets
- Free, anonymous, and works for AI-generated images of real children

CLIPS *by* SALE



OnlyFans

Porn **hub**



yubo

Snap Inc.

Threads

REDGIFS

YouTube

AZ NUDE

Online Child
Sexual
Exploitation

Alex Stamos

CSE and CSAM

Reporting
Obligation

Responses and
Mitigations

Grooming and
Sextortion of
Minors

AI-Assisted
OCSE

What's Next?

Questions

AI-Assisted OCSE

NCMEC CyberTipline reports of AI-generated CSAM:

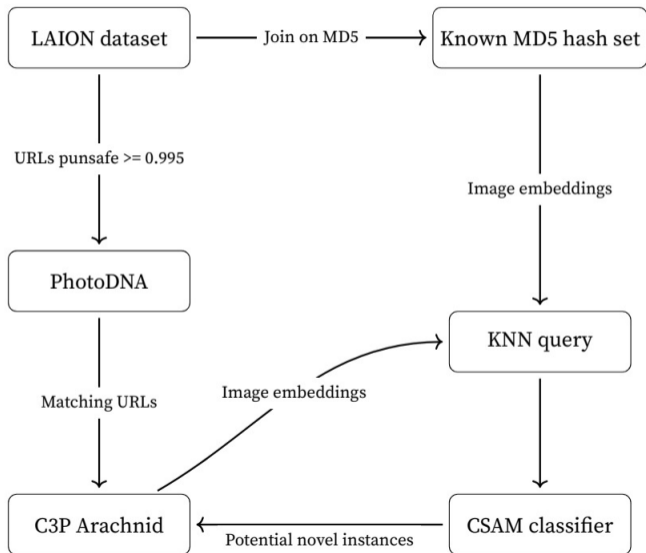
- Full year **2024**: ~67,000 reports
- **First half of 2025 alone**: ~485,000 reports (>**600% jump**)

Internet Watch Foundation:

- AI-CSAM webpages doubled (199 → 426, 2024→2025)
- New: AI-generated *videos* (not just images)
- Surge in depictions of infants 0–2: 5 cases → 92 cases
- Quality is rising: increasingly indistinguishable from real abuse

Thorn: ~1 in 8 teens knows a peer who has been targeted with an AI deepfake nude.

How AI-CSAM Is Made: Training-Data Contamination



A Useful Tell That Is Going Away



- Early diffusion models struggled with hands, text, eyes, jewelry
- Models from 2024+ have largely fixed these
- Cannot rely on visual artifacts as a long-term detection strategy
- Legal status doesn't depend on whether a real child was depicted

2 juvenile males charged with creating AI-generated nude photos of Lancaster Country Day School students



ASHLEY STALNECKER | Staff Writer | Dec 6, 2024

f

x

in

Get unlimited access to breaking news, ancestry archives, our daily E-newspaper, games and more. [Subscribe Today >](#)



- Free and \$5 “nudify” apps; trivially accessible to minors
- Source images: school photos, sports rosters, Instagram
- Cases at high schools and middle schools across the U.S. (Westfield NJ, Beverly Hills, Lancaster PA, ...)
- Victims often only learn after material circulates among peers

- Generative AI removes the need to coerce real intimate content
 - Public photos → fabricated nudes → “we already have them, pay or we send to your school”
 - Voice-cloning lets predators impersonate a romantic interest in real time
 - LLM chatbots scale grooming conversations across thousands of victims simultaneously
- “Enticement is not always necessary for this type of crime.” — NCMEC, 2024 CyberTipline Report*

- **PhotoDNA / hashing:** blind to novel AI-generated images by definition
- **Classifiers:** degrade against synthetic distributions; need constant retraining
- **Provenance / C2PA:** only works if generators sign output — offenders won't
- **Open-weights models:** cannot be recalled once released; Stable Diffusion 1.5 still circulates
- **Trust & Safety teams:** now triaging AI-generated and real material that look identical

Federal:

- *Real child depicted (modified or not):* CSAM statutes (18 U.S.C. §2252)
- *Wholly AI-generated, no real child:* obscenity law (Miller test)

State: 45 states explicitly criminalize AI- or computer-generated CSAM; outliers without specific coverage: AK, CO, MA, OH, VT, DC.

New federal frameworks:

- **TAKE IT DOWN Act** (May 19, 2025): platforms must remove NCII — including AI deepfakes — within **48 hours** of a verified request.
- **ENFORCE Act** (introduced 2025): modernizes federal statutes; stiffens AI-CSAM sentencing.
- **KOSA** (pending): reintroduced May 2025; Senate passed 91–3 in 2024 (with COPPA 2.0); imposes a “duty of care” on platforms for minors.

Thorn / All Tech Is Human “Safety by Design” commitments (signed by OpenAI, Anthropic, Google, Microsoft, Meta, Stability AI, Mistral, ...):

- *Develop*: train on clean data, red-team for CSAM generation, evaluate before release
- *Deploy*: prevent misuse via system prompts, classifiers, watermarking
- *Maintain*: monitor for abuse, remove model versions when needed, share signals industry-wide

Open problem: none of this binds the open-source releases that already have weights in the wild.

- **Origin:** founded ~2020 by Bradley Cadenhead (“Felix”), age 15, in Stephenville, TX (ZIP 76401). Cadenhead now serving 80 years.
- **Ideology:** nihilist, drawing on the *Order of Nine Angles* (O9A) — abuse is the goal, not gratification. Splinters: 676, CVLT, Court, Harm Nation, collectively “the Com.”
- **Tactics:** groom minors (often LGBTQ+ or at-risk) on Discord, Roblox, Telegram, Minecraft. Coerce escalating content — CSAM, “fan signs” carved into skin, self-harm, animal abuse, livestreamed suicide attempts. Doxxing and swatting as enforcement.
- **FBI response:** designated **Tier 1 national-security threat** in 2024 — the first time for a CSE network. ~250–450+ open investigations across all 55 field offices.
- **First terrorism charge:** *U.S. v. Baron Martin* (D. Ariz., Oct 2025) — material support to terrorists alongside running a child-exploitation enterprise.
- **Why T&S struggles:** ideologically motivated rather than gratification-driven; cross-platform; offenders are often minors themselves; abuse artifacts are *trophies* shared as social capital, accelerating virality.

[Sources: FBI Boston open letter; DOJ 764 indictments (2024–2025); GNET, 764:

Online Child
Sexual
Exploitation

Alex Stamos

CSE and CSAM

Reporting
Obligation

Responses and
Mitigations

Grooming and
Sextortion of
Minors

AI-Assisted
OCSE

What's Next?

Questions

What's Next?

- **U.S.** — *Free Speech Coalition v. Paxton* (SCOTUS, Jun 2025) upheld Texas's porn-AV law; ~24 states now have similar laws. State social-media age laws (UT/FL/TX/NY/OH/CA) mostly enjoined. **App-store AV** (TX/UT/LA, Jan 2026) forces Apple's *Declared Age Range* and Google's *Play Age Signals* APIs.
- **UK — Online Safety Act Part 3** (Jul 2025): Ofcom requires “highly effective” age assurance; self-declaration and basic payment checks rejected; fines up to £18M or 10% of global turnover.
- **EU — DSA Art. 28(1)** (Jul 2025): age-verification “mini-wallet” blueprint, ZKP-based, interoperable with the EUDI Wallet.
- **France (ARCOM, Apr 2025) & Italy (AGCOM, Nov 2025)** — “double-anonymity”: verifier doesn't see the site, site doesn't see the verifier.
- **Australia** — world-first **under-16 social-media ban** in force **Dec 10, 2025**. Platforms (not parents) liable; fines up to A\$49.5M.

[Sources: SCOTUS *Paxton*; Ofcom; EU Commission; Australia eSafety; ARCOM/AGCOM]

Age Verification: Approaches & Tradeoffs

Approach	Accuracy	Privacy	Equity / exclusion
Government-ID upload	very high	worst — full PII, doc images	excludes ~10% of US adults without ID
Credit-card check	proxy, not age	financial PII	unbanked excluded; minors have debit cards
Facial age estimation	±1.5–2 yrs; bias risk	better if on-device, no retention	skin-tone and gender bias in non-leading vendors
Digital-ID wallets / mDLs	high — issuer-sourced	selective disclosure of just “≥18”	requires smartphone + credential
Anonymous creds / ZKPs	inherits issuer	best — predicate proof, no identity	ZKP plumbing still maturing
OS / app-store signals	mixed	avoids per-site PII; centralizes power	low-confidence age-down locks out adults
Token-based “double-blind” (FR, IT)	inherits verifier	unaggregable by design	gated on upstream verification

Tensions: privacy vs. accuracy; exclusion of legal adults; gatekeeper concentration (Apple, Google, or states); chilling effect on lawful speech (EFF / Paxton dissent).

Section 230 Cracks: Design Is a Product

Courts are now treating platform **design** — ranking, defaults, friction, age-gating — as a *product* that §230 does not immunize.

- **Anderson v. TikTok** (3d Cir. 2024) — 10-year-old died doing TikTok's algorithmically-recommended "Blackout Challenge." For You is the platform's own speech, so §230 does **not** cover recommendation. First federal appellate narrowing of §230 under *Moody v. NetChoice*.
- **Social Media Adolescent Addiction MDL 3047** (N.D. Cal., 2023) — Meta, TikTok, Snap, and YouTube failed to win blanket §230 dismissal. Claims survive on **design defects** (no parental controls, weak age-gating, no easy deletion) and **failure-to-warn**. First bellwether opened Feb 2026.
- **Garcia v. Character Technologies** (M.D. Fla., May 2025) — after a teen's chatbot-encouraged suicide, the court rejected the First-Amendment "speech" defense. Chatbot output treated as a **product**.

- **42-state AG coalition v. Meta** (N.D. Cal., Oct 2023) — IG and FB knowingly designed to addict minors; under-13 data collection violated COPPA. State consumer-protection theories **survived** §230.
- **New Mexico v. Meta** (NM, Dec 2023) — a decoy-account investigation showed Instagram serving sexual content to minors and connecting predators with kids. Jury awarded **\$375M** — first standalone state-AG verdict on platform child safety.
- **Doe #1 v. MG Freesites** (MindGeek/Aylo, N.D. Ala.) — §230 rejected as a shield for knowing CSAM transmission; class certified Sep 2024. Reinforced by **FTC + Utah v. Aylo** (Sep 2025) for deceiving users on moderation.

- Recommendation systems for minors must be designed *defensively*, not just moderated post-hoc.
- DM/messaging architecture and predator-contact pathways are now litigable design choices.
- Age verification and default privacy for minors are becoming part of a *duty of care*.
- Promoted/featured content is increasingly treated as platform first-party speech.
- **Courts and AGs — not Congress — are now setting de facto child-safety policy.**

Online Child
Sexual
Exploitation

Alex Stamos

CSE and CSAM

Reporting
Obligation

Responses and
Mitigations

Grooming and
Sextortion of
Minors

AI-Assisted
OCSE

What's Next?

Questions

Questions