# Welcome

# The Future of Trust & Safety:

## Child Safety Investigations

**Sergio Zaragoza**

Child Safety Intelligence
Analyst & Consultant

**Alexandra Koptyaeva**

Trust and Safety
Representative at Falkor

**Lior Mordechai**

Customer Success &
Marketing Director at Falkor

# Agenda

# The Importance of Investigating Child Safety Issues

# The Importance of Investigating Child Safety Issues

- Affects platform reputation
- Impacts user retention
  - Negative user experience
- Soon to be legally regulated globally
  - Risk assessment
  - Policy Management and Community Guidelines update
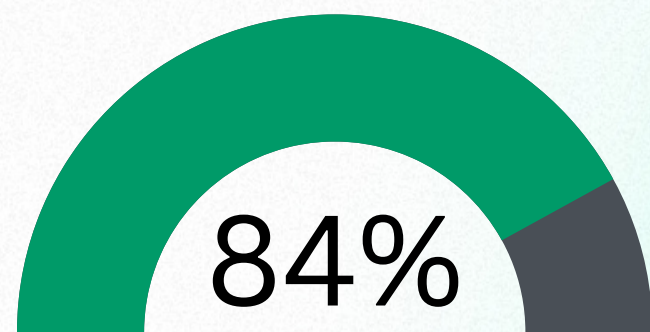  - Compliance: know the laws, your obligations, and consequences

Child Safety is not only about **CSAM.**

There's also **Grooming**, **Sextortion**, **Self-Generated Content**, **Deepfakes**, **Non-Consensual Indecent Imagery**, **Child Sex Trafficking**, **Child Labor Trafficking**, etc.
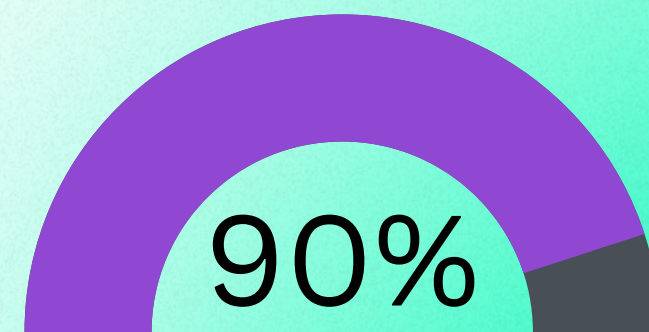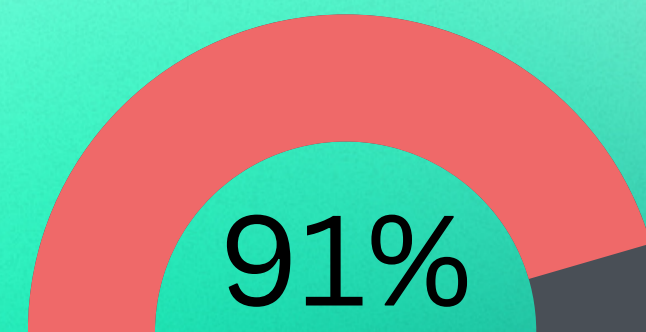
# 2022 Statistics

- 84% of CSAM reports were of never-before-seen material, depicting new child victims of sexual abuse.
- 9 in 10 of the victims depicted are aged 3 to 13 years old and 1% are in the infant category (0-2 years), with children as young as 4 months old victimized.
- 91% of reported CSAM involved female victims.
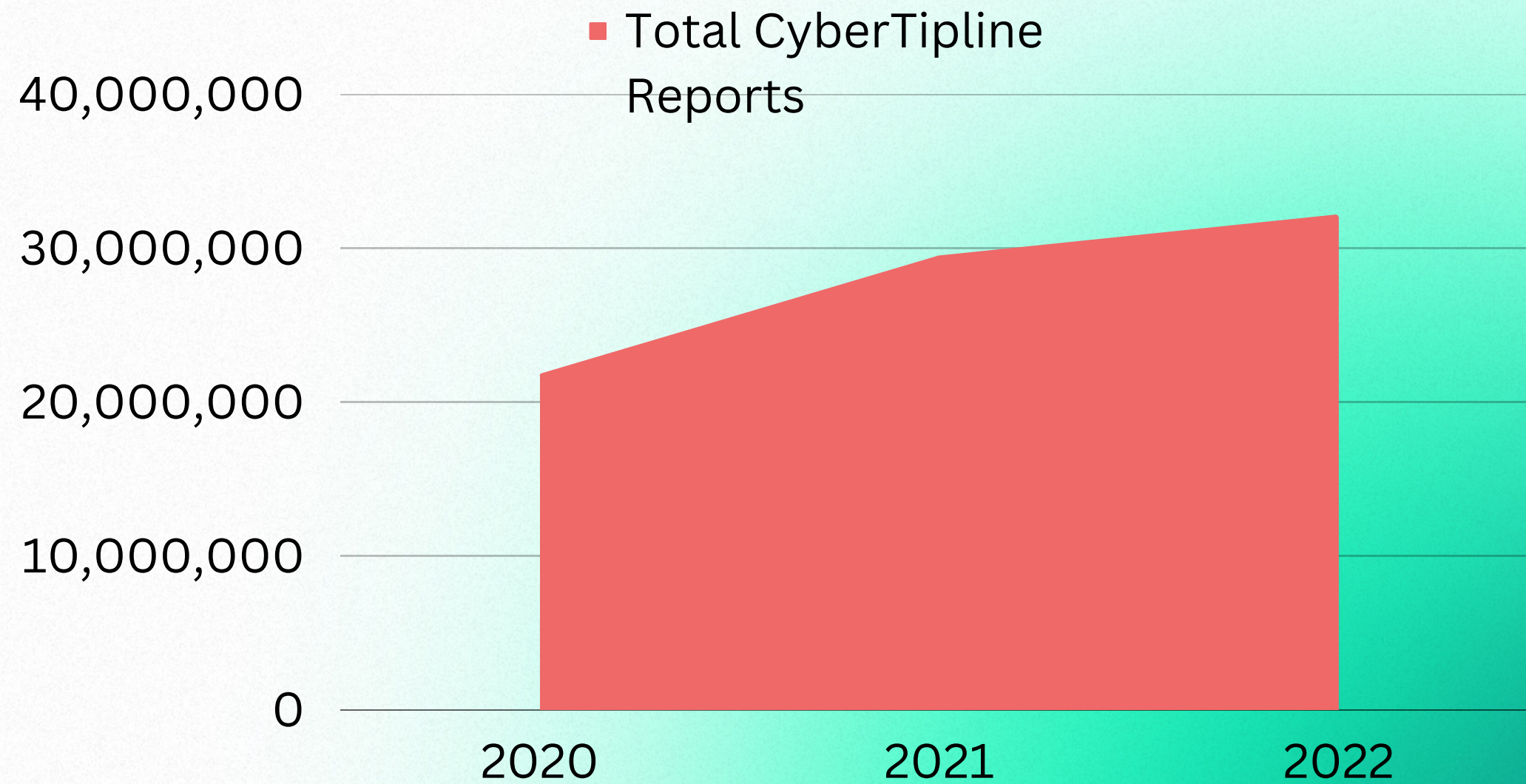
**84%**
Never-before-seen material

**90%**
Victims aged 3 to 13

**91%**
Female victims

Source: INHOPE Annual Report, 2022

# 2022 Statistics

Increase in CyberTipline Reports in the US:



■ Total CyberTipline Reports

| | 2020 | 2021 | 2022 |
|---|---|---|---|

40,000,000

30,000,000

20,000,000

10,000,000

0

Source: NCMEC Annual Report, 2022

# CyberTipline Reports Made Available to Law Enforcement

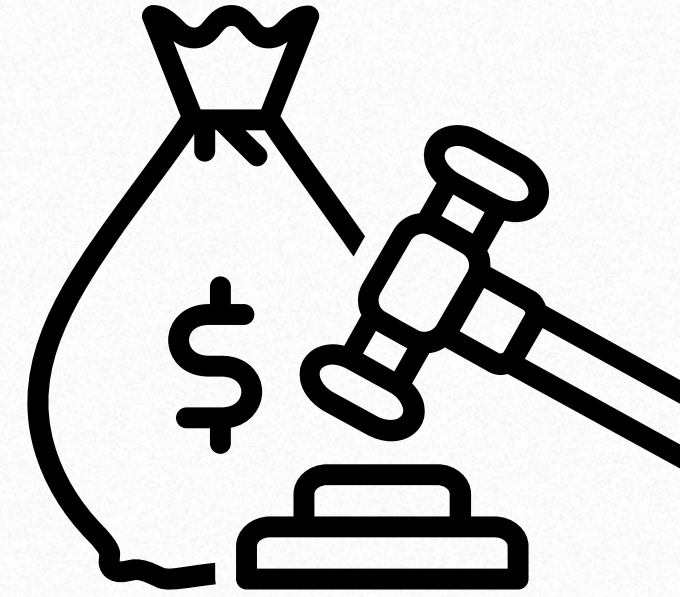| Law Enforcement Recipients of CyberTipline Reports | 2022 |
|---|---|
| Internet Crimes Against Children (ICACs) Units | Actionable: 491,655<br>Informational: 400,715<br>Total: 892,370 |
| Local Law Enforcement | Actionable: 1,462<br>Informational: 3<br>Total: 1,465 |
| Federal Law Enforcement | Actionable: 1,356,988<br>Informational: 997,475<br>Total: 2,354,463 |
| International Law Enforcement | Actionable: 13,995,567<br>Informational: 15,010,413<br>Total: 29,005,980 |

# 2022 Trends

| Generative AI | CAP Sites | Live Streaming |
|---|---|---|
| <ul><li>Creating and sharing (fake) CSAM imagery (e.g., AI-generated body + real victim's face).</li><li>Engaging in Grooming conversations.</li><li>Sextortion, Deepfakes.</li><li>Training LLMs with AI-generated CSAM images so they can have their own image data sets.</li></ul>Recent FBI warning about GenAI | <ul><li>Hotlines have observed the rise of Child Abuse Pyramid Sites.<ul><li>New type of commercial site.</li></ul></li><li>Uses a particular form of "invitation" system to access CSAM content.</li><li>Users are encouraged to share their personal link to invite others to the site.</li><li>The more invitations they share, the more points they accumulate to access CSAM content.</li></ul> | <ul><li>Individuals pay to watch live abuse of a child via a video streaming service</li><li>Incredibly difficult to detect due to real-time nature and lack of digital evidence left behind after the crime</li></ul> |

# New Regulations

# More and More Regulations

1. Online Safety Bill (UK)
2. Online Safety and Media Regulation Act (Ireland)
3. Digital Services Act (Europe)
4. US state-specific regulations (e.g., A.B. 587 in California, HB 20 in Texas, S.B. 7072 in Florida)

# US Regulations: AB-2273 CA Act

California Age-Appropriate Design Code Act

- Signed on **Sep 15, 2022.**
- Takes effect on **July 1, 2024.**
- New legal obligations on companies with respect to online products and services that are "likely to be accessed by children" under the age of 18.
- These platforms must comply with this Act and demonstrate that they are taking steps to protect children.

Violators may be subject to a penalty of **up to $2,500** per affected child for each negligent violation, and **up to $7,500** per affected child for each intentional violation.

# US Regulations: 18 USC 2258A - NCMEC & CSAM

- Duty to Report:
  - To reduce and prevent the online sexual exploitation of children, a provider shall report as reasonably as possible after obtaining actual knowledge and take action.
  - They must inform CyberTipline.

- Failure to Report:
  - Knowingly and willfully failing to make a report — up to **$150,000** (initial knowing and willful failure);
  - Second or subsequent knowing and willful failure to make a report — up to **$300,000**.

# Common Challenges Across T&S Teams

# The Current State of T&S Investigations

Basic workflow:

1. Content is detected by AI & human Content Moderators
2. Content is removed according to guidelines
3. Offending Profile/Account is banned and reported to LE

But! The end problem is not solved.
Bad actors are still on the platform, and are presumably active on other platforms.

# Methodology Tips

# Methodology Tips - Theory

Improve the team's OSINT capabilities & general investigation skills

Keep up with current and emerging trends

Maintain an internal knowledge base of platform trends and user behavior

Attend courses that are relevant to the field

# Methodology Tips - Practice

Increase your proactiveness in detecting threats (stress testing, red teaming)

Dive deep into the data (rather than focusing on one user/bad actor)

Collaborate with other teams and companies (AI, CM, specialized tools)

Increase transparency & data accessibility to relevant stakeholders

# The Proactive Approach

A proactive approach is defined by **leading investigations on the platform**, so we can detect new threats and bad actors ahead of time.
We can follow **leads** that will help us to understand bad actors' **behaviour** with anticipation and act more efficiently independently of AI moderation.

# The Importance of a Proactive Approach

- Depending on the company, it can be an internal process or outsourced.
- Provides the opportunity to detect new trends and threat actors if appropriate tools and resources are available.
- By following different trails, we are able to detect future harms and behaviors before they occur.
- Platform-independent investigations can give a broader perspective on where threats may originate.
- An out-of-the-box approach to gain a deeper understanding of today's landscape.

# Technology for the Modern Analyst

# Deeper Investigations with Falkor

- Analyze available data on problematic users
- Handle all cases in one, shared knowledgebase
- Get clearer insights into user behavior:
  - What content has been shared
  - When, how, and where
  - Possible connected accounts
- Escalating to Law Enforcement:
  - Provide more precise and concise information
  - Reduce their workload so more escalations can be addressed
  - Help speed up the investigations to help more victims

# Supercharge Data Analysis



- Upload any data
- Visualize it in tables, maps, link analysis graphs, & timelines
- Filter it by time, location, or other attributes
- Analyze patterns
- Expand networks

# Streamline Collaboration

- Solve cases together
- Get insights from your teammates
- Share suspicious entities
- Assign tasks and deadlines
- Automatically export reports for escalation

Streamline Collaboration

# Manage Teams



- Get a bird's eye view of your team's work
- Ensure efficiency and employee wellness
- Manage permissions and roles
- Audit activity and improve processes

As a Trust & Safety professional, you do not need to investigate *everything* on the platform in-depth.
But using your data correctly can help you identify the root cause of issues, eliminate them for good, and prevent future problems.

# Resources

- INHOPE Annual Report, 2022

- National Center for Missing and Exploited Children (NCMEC) Transparency Report, 2022

- Sociobits, 'Dark Side of AI Art: Deepfake, Child Pornography and Sextortion', March 2023

- The Guardian, 'AI tools could be used by predators to 'automate child grooming', May 2023

# Thank you!

Sergio Zaragoza

smartinezzaragoza@gmail.com

Alexandra Koptyaeva

alexandra@falkor.ai

Lior Mordechai

lior@falkor.ai