DESIGN YOUR WEAPONS IN THE FIGHT AGAINST DISINFORMATION

15 SEPTEMBER 2024

Associate sponsors: AMD × hp
together we advance_

The Problem

# Is it true?
# Is it real?

# Fakes can be believed

The Problem

"The S&P 500 briefly dropped… as social media accounts… repeated the false claims."

AP, May 23 2023

Reports of an explosion near the Pentagon in Washington DC

# Fakes can be believed

The Problem

WORLD

**Deepfake targets Ukraine's first lady Olena Zelenksa with false claim she bought Bugatti**

CBS News, July 2024

Participants could only verbally identify deepfakes 37% of the time

Source: University of Sydney

'Deepfakes' of Michael Mosley and Hilary Jones being used to promote scams on social media

Sky News, July 2024

34 million AI-generated images are made each day Source: Everypixel

**AI and deepfakes blur reality in India elections**

BBC, May 2024

**Deepfake clips of Gareth Southgate swearing after England match go viral**

The Guardian, July 2024

# The technology is getting better and easier

The Problem

**News needs to separate fact from fiction**

How To Make AI Images Of Yourself (Free)
126K views • 2 weeks ago

Matt Wolfe ✓

Here's how to train your face into the new AI models. Here's the link for $10 free at Replicate: ...

How To Create a Fake YouTube Studio Background Using FREE AI Tools
22K views • 8 months ago

KC Sounds

In today's tutorial, we're diving into the world of AI and creating a jaw-dropping fake background using. Whether you're a ...

10 moments  What You Need | Choose Video | Remove Video Background | Where to find backgrounds | AI source...

IBC2024 **Design Your Weapons in the Fight Against Disinformation** #ACCELERATORS2024
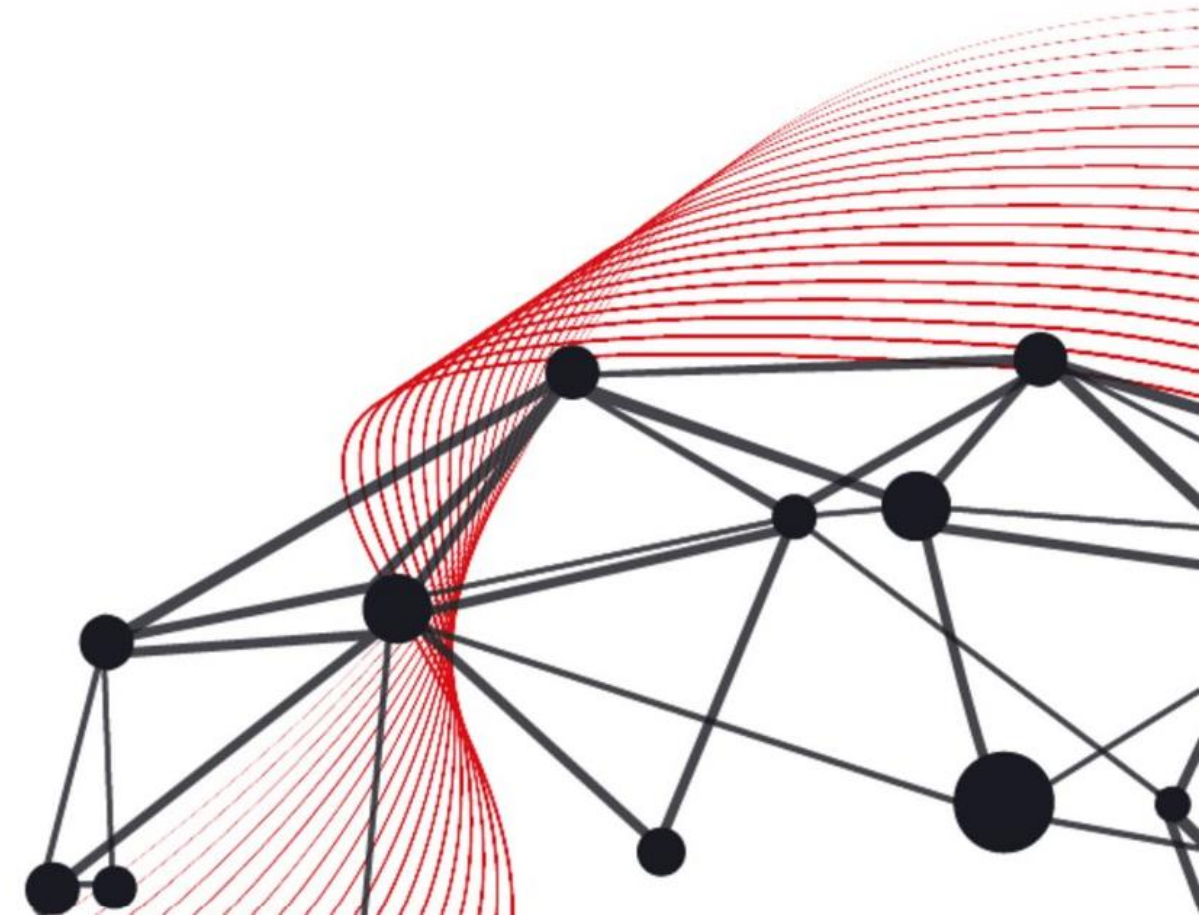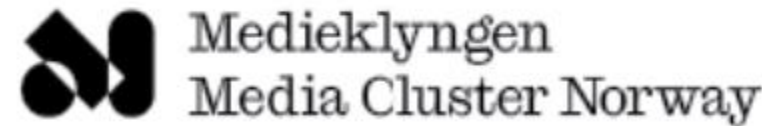
**#ACCELERATORS2024**

**IBC2024**

# Design Your Weapons in the Fight Against Disinformation

**Provenance**

**Detection**

**Collaboration**

# What is provenance, and why is it becoming so significant?
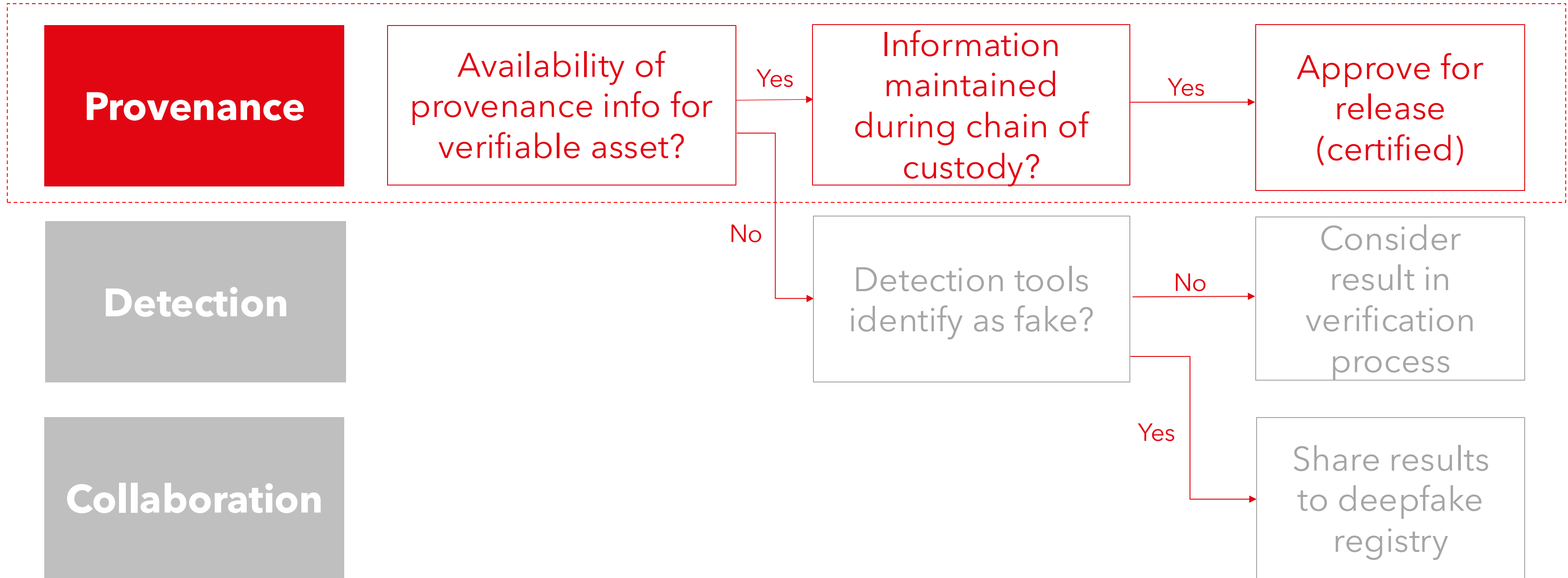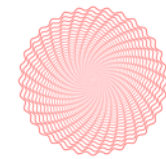
**Digital Media Provenance**

**Metadata**

**Watermarking**

**Fingerprinting**

| Signed | Unsigned |
| --- | --- |

| Visible | Invisible |
| --- | --- |

| Cryptographic Hashing | Perceptual Hashing |
| --- | --- |

# Provenance has impact throughout the media lifecycle

Timeframe:  *Seconds to days*  *May continue for years*



User-Generated Content (UGC)

Processing

Upload to social media channels and/or private circulation

Trusted 3rd parties

Content verification

Capture / origination

Own sources

Newsroom analysis

Media Asset Management

Processing (edits, captions)

Upload to output destination

Transport to users (primary distribution channel)

User sharing (secondary channels)

Archival

IBC2024

# Who does provenance benefit?

**...and how?**

**Content Creators**

e.g. trusted photojournalists, UGC

- Proving authenticity of their work
- Tracking content usage (copyright)

**Journalists**

e.g. newsrooms, editors, archivists

- Rapid & accurate verification by newsrooms
- Safeguarding information on origins (archival content)
- Tracking and legitimising edits and captions

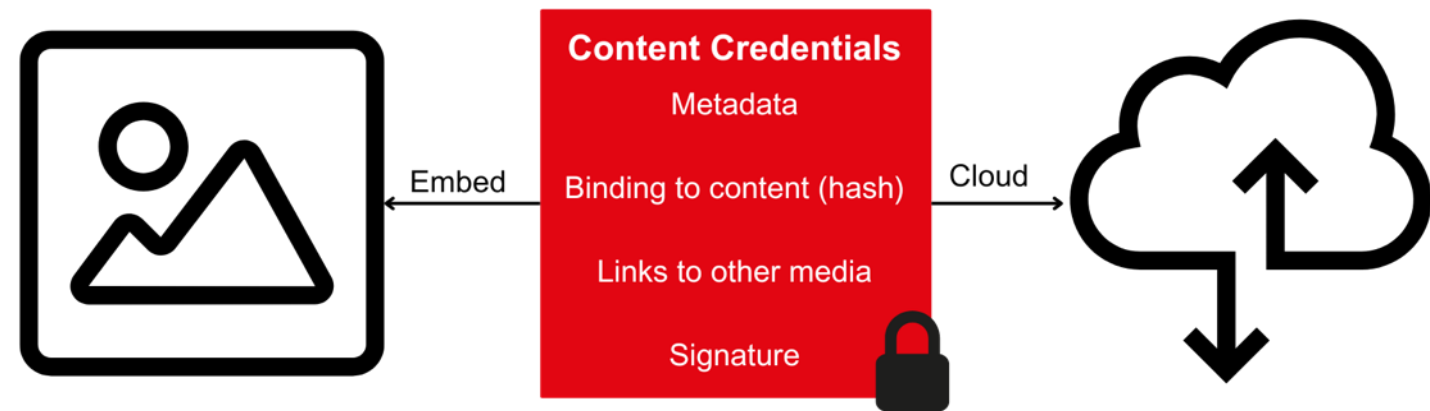**Viewers**

e.g. audiences, copyright owners

- Audience understanding of origin point
- Tracking content ownership & authorship

**IBC2024**

# Standards and solutions are complementary

## Provenance Standards (e.g. C2PA)

- Protocols and guidelines that define how provenance data should be captured
- C2PA is a widely adopted standard that uses cryptographic hashing to securely bind metadata to the content, ensuring authenticity and integrity.



**Content Credentials**
Metadata
Embed | Binding to content (hash) | Cloud
Links to other media
Signature

## Provenance Solutions

- The technical systems and tools that actually implement provenance tracking, ranging from proprietary systems to open-source frameworks
- A range of tools and solutions exist in the market, tailored to specific needs and workflows, offering diverse approaches to ensuring the authenticity and integrity of digital content.

ELUV.IO       TRANSMIXR

OpenOrigins       HUMAN DIGITAL

VIDENTIFIER       NAGRA KUDELSKI

**IBC2024**

# Solutions ecosystem consists of a range of contributing technologies

**ELUV.IO** — **Next-generation content management and distribution solution**, utilizing blockchain technology and cryptographic hashing to allow users to control access and ensure monetization for their content.

**TRANSMIXR** — **Consortium developing the foundation for the future of visual content**, including virtual and augmented reality, with focus on the future of newsroom content creation and distribution.

**OpenOrigins** — **End-to-end provenance solution**, providing tools for authenticated media capture, archival media protection, and verification of online content. Secures provenance information in a tamper-proof blockchain to provide a decentralised trust infrastructure.

**HUMAN DIGITAL** — **Interoperable talent identifier,** enabling ID resolution and provenance verification of notable real (human) public figures, their connected digital replicas, and fictional characters in the media supply chain. A DOI Foundation registration agency (ISO 26324.2022). Member: C2PA, IPTC.

**VIDENTIFIER** — **Automated visual matching and linking of content**, using advanced visual fingerprinting for accuracy and speed. Enables users to verify content matching trusted sources and to detect changesmatching made.

**NAGRA KUDELSKI** — **Content provenance and authenticity solutions** based on forensic watermarking and fingerprinting, enabling users to distribute, protect and monetize their content.
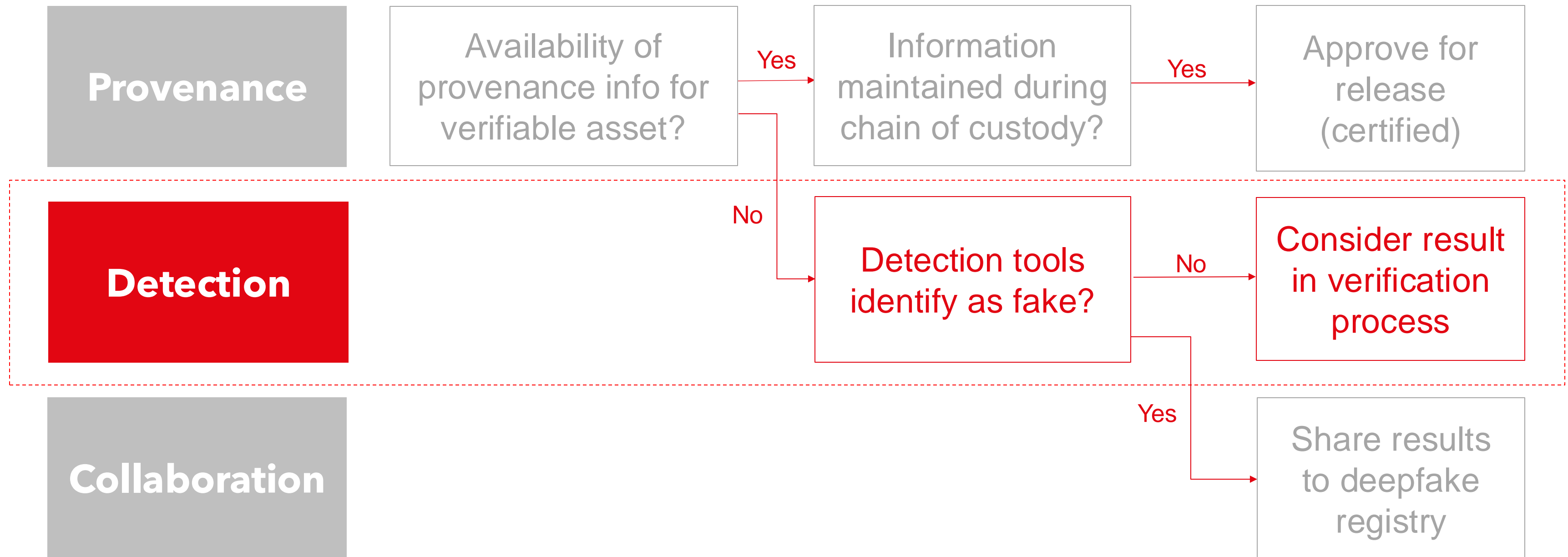
**IBC2024**

# Detection may be leveraged when provenance data is unavailable

- **Detection** is the leveraging of visual analysis and machine learning to **automate detection of manipulated visuals, narrative or wholly synthetically-generated content**.

- Newsrooms and end users can benefit from the ability to detect Generative AI across various media formats, especially when provenance information can't be accessed.

**Challenges in detection today:**

- Generative AI is developing at a faster rate than detection modules.
- Multi-layered editing of AI-generated content can impact the efficacy of detection modules.
- No single solution available to address all issues, with a risk of false results and inaccuracies.

A combinations of detection tools can inform various elements of the verification process, acting as an aid rather than an alternative to human-led content validation.

IBC2024

## Unverified Source

### Scan against existing authenticated database

- Identify duplicates of authenticated original content and assess potential IP implications

- Identify manipulated content based off authenticated original media

### Scan against collaborative database of problematic content

- Identify common disinformation topics or bad actors

- Collaboration or Machine led curation can improve reaction speed and help with allocation of resources

### Visual analysis to identify manipulated or wholly AI-generated content

- Newsrooms are currently mostly reliant on manual checks and contextual verification

- Automated solutions exist but can struggle to keep up with the advances in Generative-AI

# Secure capture of new media

**EXAMPLE: USER-GENERATED CONTENT**

- **Metadata & unique fingerprint** captured

- Hardware verification  ensures that the **capture device is trusted**

- **LiDAR & sensor data** captured along with the image

- Stored on a **tamper-proof blockchain**



3D depth map data of images

# Newsroom analysis

**EXAMPLE: USER-GENERATED CONTENT**

- Photo editor can look up proof information and understand:
  - Image has **verifiable origins**
  - Taken on a **trusted device**
  - Taken at **correct time & location**

Everyday photographer shares photo with journalist

Photo editor looks up origin information

# OpenOrigins

# Transport to viewer

**EXAMPLE: USER-GENERATED CONTENT**

- Chrome plugin displays **green boxes around known verified content**

- Comparison tools allows viewers to check if **an image they see has been modified**



Image similarity

89%

Quad tree file:

Image:  00000-2111185431.png

Threshold:              2

Depth:                  5

Total image pixels:     1687500

Total matched pixels:   1503960

Perceptual algorithm:   PDQ

Highlight image ⬤ Quad tree image

Comparsion tool for edited images

Mock-up of a disinformation site with verified and non-verified content

# Archive Anchoring

**SAFEGUARDING ARCHIVES WITH IMMUTABLE PROVENANCE, PRESERVING OUR WORLD'S SHARED HISTORY**

- Deep database analytics to ensure the media is **validated and authentic.**

- Scalable architecture allows for the **securing of billions of media items** that already exist in archives.

- Without proper authentication, **real and fake become indistinguishable**



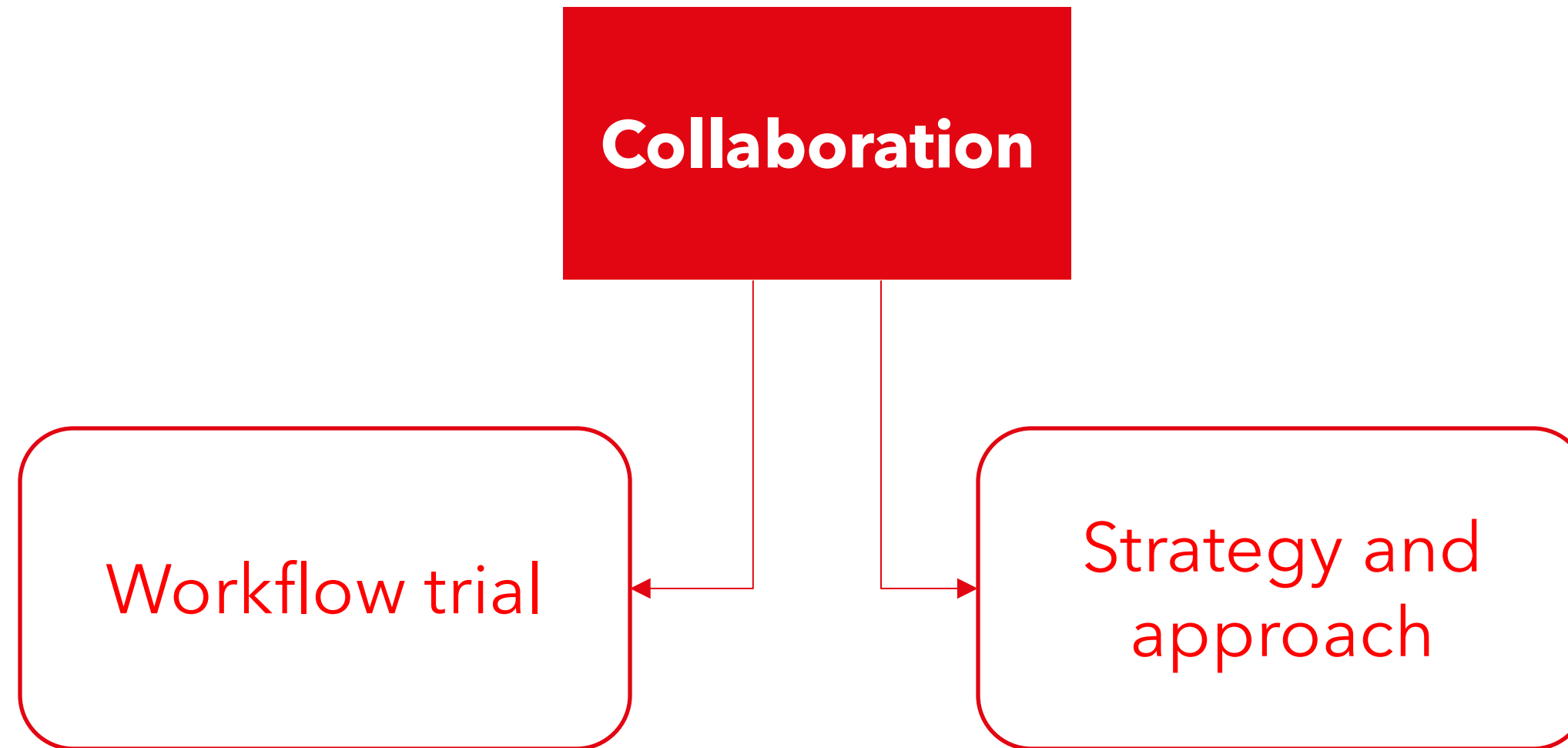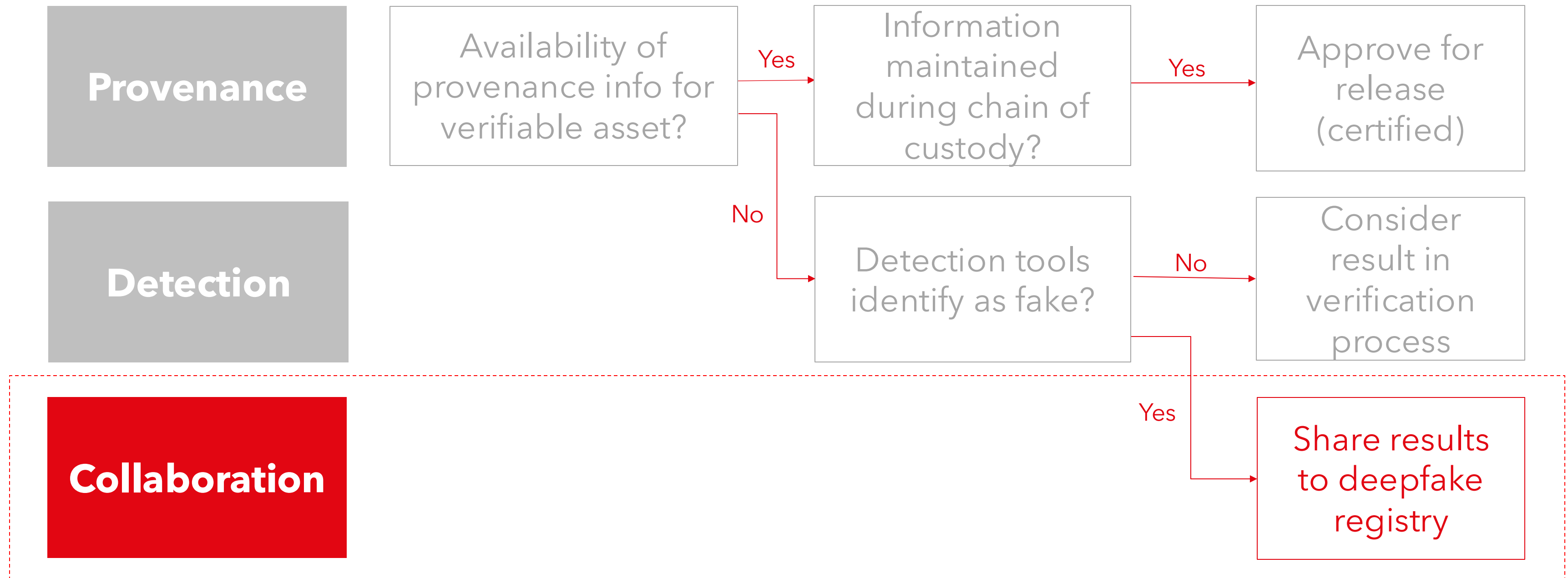Media organisation's view of secured media proofs for archival content

What does Collaboration look like?

Collaboration

Workflow trial

Strategy and approach

# Collaboration – Workflow Trial

Are newsrooms prepared to share knowledge that benefits everyone to produce faster and more informed reporting?

- 1 month

- 7 news organisations

- Shared information on suspect content

Tech support: Google

System development:

# Collaboration – Workflow Trial

# Collaboration – Workflow Trial

## Two major stories during trial period



US Politics



UK Riots

and yet…

| | Content shared |
|---|---|
| Week 1 | 10 stories |
| Week 2 | 5 stories |
| Week 3 | 3 stories |
| Week 4 | 6 stories |

#ACCELERATORS2024

**IBC2024**

# Collaboration – Workflow Trial

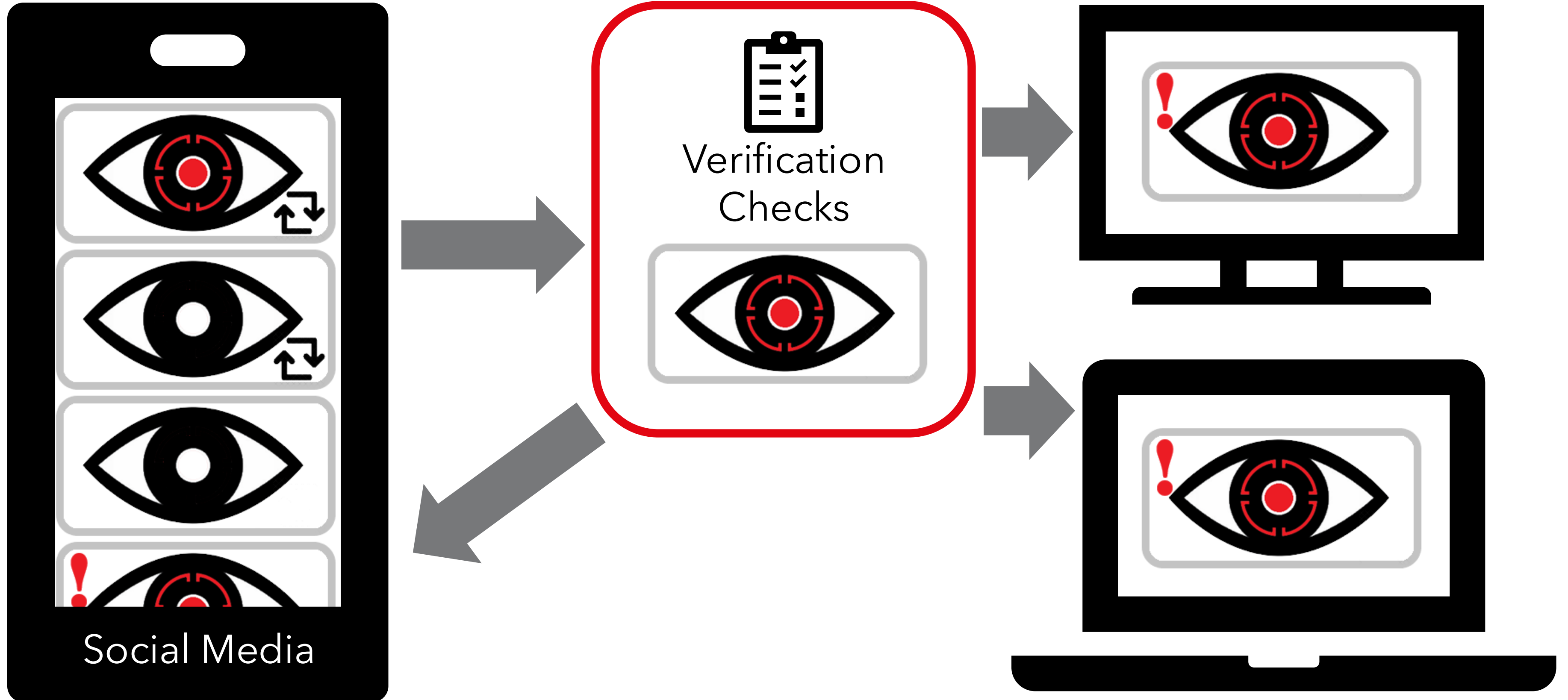|  | Content shared |
|---|---|
| Week 1 | 10 stories |
| Week 2 | 5 stories |
| Week 3 | 3 stories |
| Week 4 | 6 stories |

- **Fake content and disinformation spreads on social media** where it can gain huge traction without being debunked or taken down.

- **News organisations ignore most fake content.** There is too much to debunk it all.

- **Journalists operate with caution.** Fake and genuine content is checked, taking time and allowing disinformation to spread through unverified means.

- **That means mis- and dis-information can spread without authoritative challenge,** causing confusion and affecting narratives.

# Collaboration – Approach & Strategy

Verification
Checks

Social Media

# CALL TO ARMS!

- We all need to work together closely to make sure everyone has access to information they can trust

- We all need to develop a deeper understanding of the products and systems available to fight disinformation and fake content

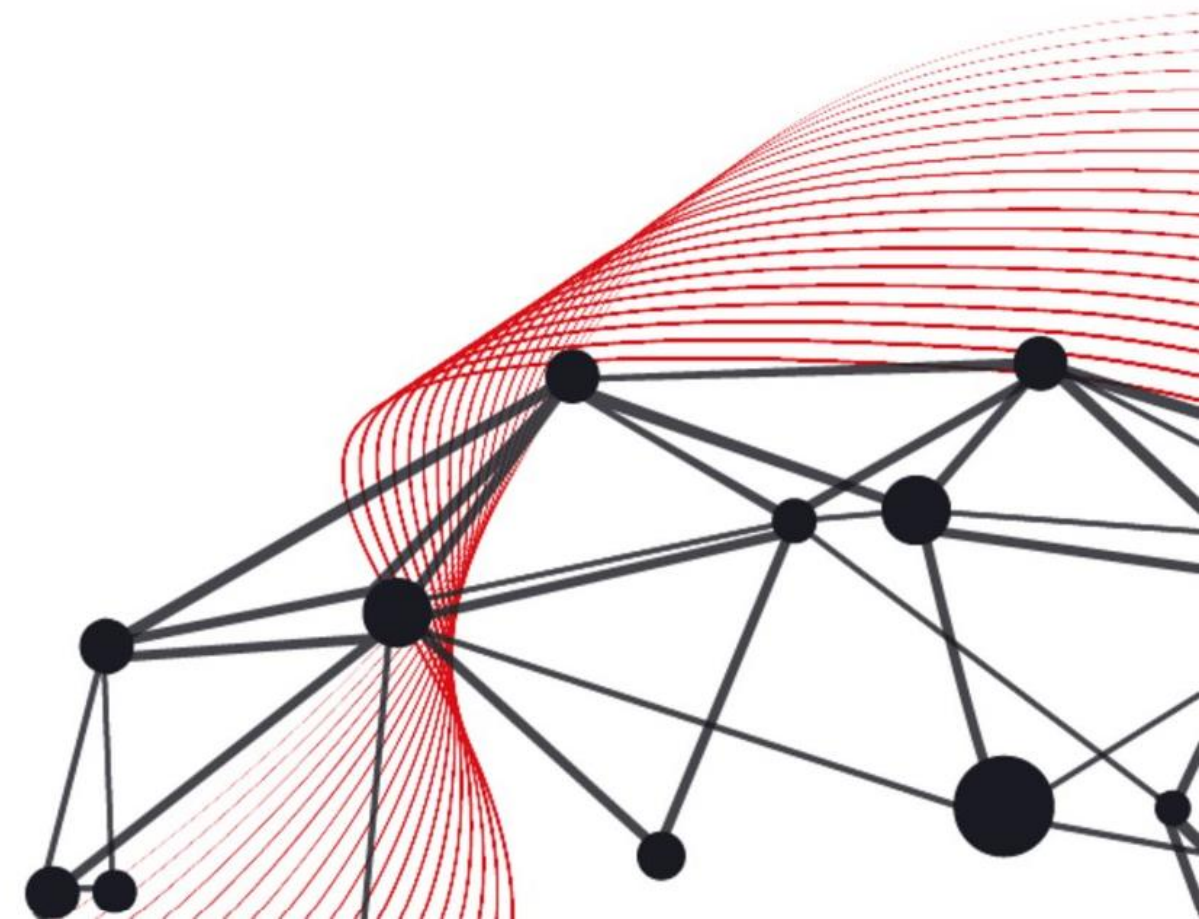- AI generated fake content is getting increasingly hard to detect and easier to produce.

# CALL TO ARMS!

**Join us** on this journey

accelerators@ibc.org