Using Zero-Knowledge Proofs to Fight Disinformation

Trisha Datta and Dan Boneh Stanford University

By Alistair Coleman & Shayan Sardarizadeh BBC Monitoring

24 February 2022

By Alistair Cole BBC Monitorin

24 February 202

Fact-checking videos and pictures from Ukraine

Since Russia's attacks on Ukraine began, we have seen several videos and pictures go viral that are actually fake posts.

BBC Monitorin

24 Februarv 202

By Alistair Cole Fact-checking videos and pictures from Ukraine

Since Rus False social media posts are hindering and pictu earthquake relief efforts in Turkey. You can help stop that



Sony Unlocks In-Camera Forgery-Proof Technology

04 Aug, 2022



embedded certified signing key **sk**₁





Sony Unlocks In-Camera Forgery-Proof Technology

04 Aug, 2022



🚺 Adobe 🛛 🖬 🖪 🖸 🛛 intel 📑 Microsoft SONY 🎯 Truepic 🈏

Sony Unlocks In-Camera Forgery-Proof Technology

04 Aug, 2022



Adobe **Arm BBC** intel Hicrosoft SONY



Sony Unlocks In-Camera Forgery-Proof Technology

04 Aug, 2022



Sony Unlocks In-Camera Forgery-Proof Technology





See Rivadeneira.

A Problem: Post-Processing

- Newspapers often process photos before publication
 - At minimum, need to resize (90 MB \rightarrow 8 MB)
 - Allowable operations from the Associated Press: cropping, grayscale, ...

A Problem: Post-Processing

- Newspapers often process photos before publication
 - At minimum, need to resize (90 MB \rightarrow 8 MB)
 - Allowable operations from the Associated Press: cropping, grayscale, ...

Problem: browser cannot verify the C2PA signature of a processed photo

A Problem: Post-Processing

- Newspapers often process photos before publication
 - At minimum, need to resize (90 MB \rightarrow 8 MB)
 - Allowable operations from the Associated Press: cropping, grayscale, ...

Problem: browser cannot verify the C2PA signature of a processed photo















1.

- ZK-SNARK: efficiently verifiable statement about a secret witness
- Complete 2. Sound **Public Statement** 3. Non-interactive Zero-Knowledge 4. 5. Succinct I claim something is true about the Secret Witness secret witness. Here is a proof π . Thanks! I will verify it Prover Verifie















1. Complete/Sound: verifier doesn't need to trust prover



- **1. Complete/Sound:** verifier doesn't need to trust prover
- 2. Non-interactive: interactivity would entail unique proofs for each browser



- **1. Complete/Sound:** verifier doesn't need to trust prover
- 2. Non-interactive: interactivity would entail unique proofs for each browser
- **3. Zero-Knowledge:** useful for ops such as cropping



- **1. Complete/Sound:** verifier doesn't need to trust prover
- 2. Non-interactive: interactivity would entail unique proofs for each browser
- **3. Zero-Knowledge:** useful for ops such as cropping
- **4. Succinct:** web browser can efficiently verify proof

Proofs for Post-Processing Ops

• PhotoProof (Naveh and Tromer, 2016): a few minutes to generate photo editing proofs for 128 x 128 pixel image

Proofs for Post-Processing Ops

- PhotoProof (Naveh and Tromer, 2016): a few minutes to generate photo editing proofs for 128 x 128 pixel image
- New tools enable faster development!



oKrates

Performance for Post-Processing Ops Proofs

For resizing, cropping, grayscale ops on images of about 6000 x 4000 pixels (~30MP) using Circom:

- Proof generation time: <1 second
- Witness generation time: <4 minutes

by newspaper once per image

• Verification time: 2 ms

• Proof size: <1 KB

- by browser







2

I know **Orig** such that:

1. signature is a valid signature on Orig

2. Edited is the result of applying Ops to Orig

3. metadata(*Edited*) = metadata(*Orig*)

Attempt 1

 π (Signature)

```
I know (Orig, hash) such
that:
hash = SHA256(Orig)
```













Poseidon hash of lattice hash [GGH'96, SCMPGLW'08]



• A is random matrix from finite field \mathbb{F}_q



- A is random matrix from finite field \mathbb{F}_q
- \vec{x} is low norm vector in \mathbb{F}_q representing the input



- A is random matrix from finite field \mathbb{F}_q
- \vec{x} is low norm vector in \mathbb{F}_q representing the input



- A is random matrix from finite field \mathbb{F}_q
- \vec{x} is low norm vector in \mathbb{F}_q representing the input

Poseidon hash of lattice hash [GGH'96, SCMPGLW'08]



 Collision-resistant assuming SIS → prover must prove original photo representation is low norm

To prove \vec{x} is low norm, i.e., $\vec{x} \in \{0, 1, ..., R - 1\}^n$:

To prove
$$\vec{x}$$
 is low norm,
i.e., $\vec{x} \in \{0, 1, ..., R - 1\}^n$:

R = 3 \vec{x} 2 0 2

To prove
$$\vec{x}$$
 is low norm,
i.e., $\vec{x} \in \{0, 1, ..., R - 1\}^n$:
• $\vec{y} := [0, 1, ..., R - 1]$

$$R = 3$$
 \vec{x}
 2
 0
 2

 \vec{y}
 0
 1
 2
 3

To prove
$$\vec{x}$$
 is low norm,
i.e., $\vec{x} \in \{0, 1, ..., R - 1\}^n$:
• $\vec{y} := [0, 1, ..., R - 1]$
• $\vec{z} := sort(\vec{x} | | \vec{y})$

 R = 3

 \vec{x} 2
 0
 2

 \vec{y} 0
 1
 2
 3

To prove
$$\vec{x}$$
 is low norm,
i.e., $\vec{x} \in \{0, 1, ..., R - 1\}^n$:
• $\vec{y} := [0, 1, ..., R - 1]$
• $\vec{z} := sort(\vec{x} || \vec{y})$

R = 3 \vec{x} \vec{y} \vec{Z}

To prove
$$\vec{x}$$
 is low norm,
i.e., $\vec{x} \in \{0, 1, ..., R - 1\}^n$:
• $\vec{y} := [0, 1, ..., R - 1]$
• $\vec{z} := sort(\vec{x} || \vec{y})$

R = 3 \vec{x} \vec{y} \vec{Z}

To prove
$$\vec{x}$$
 is low norm,

i.e.,
$$\vec{x} \in \{0, 1, ..., R - 1\}^n$$
:
• $\vec{y} := [0, 1, ..., R - 1]$

•
$$\vec{z} := sort(\vec{x} | | \vec{y})$$

Must show:

• \vec{z} is permutation of \vec{x} and \vec{y}

•
$$\vec{z}[i+1] - \vec{z}[i] \in \{0,1\}$$

R = 3 $\vec{\chi}$ \vec{y} \vec{Z}

Performance Results for Proving Signatures

Time to Generate Proof vs. Input Size



Performance Results for Proving Signatures

Time to Generate Proof vs. Input Size





- Proof systems have greatly improved due to their need in blockchains ⇒ non-blockchain applications benefit
- Proofs about large images (4000 × 6000) can be done in reasonable time
- Applicable to C2PA for image authenticity
- Open problem: ZK proofs for videos?