# VerITAS: Verifying Image Transformations at Scale

Trisha Datta, Binyi Chen, Dan Boneh

Stanford University

# These look like prizewinning photos. They're AI fakes.

Artificially generated images of real-world news events proliferate on stock image sites, blurring truth and fiction

By Will Oremus and Pranshu Verma
November 23, 2023 at 6:00 a.m. EST
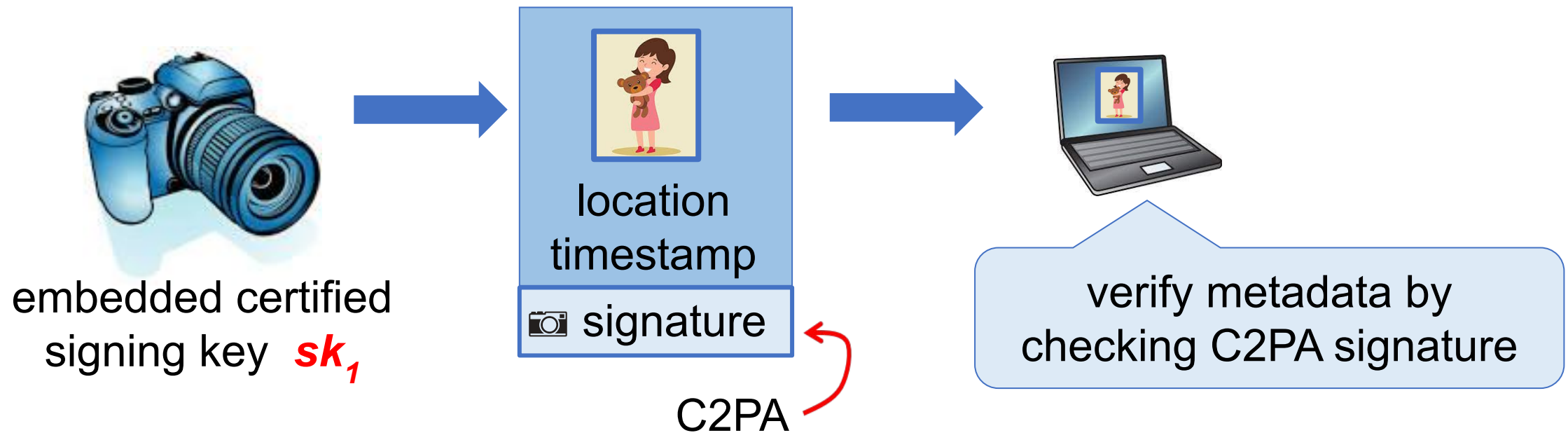
AI-GENERATED FAKE PHOTO

AI-GENERATED FAKE PHOTO

AI-GENERATED FAKE PHOTO

# C2PA: A Content Provenance Standard

**Nikon, Canon, Sony eye tamper-resistant digital signatures to combat deepfakes**
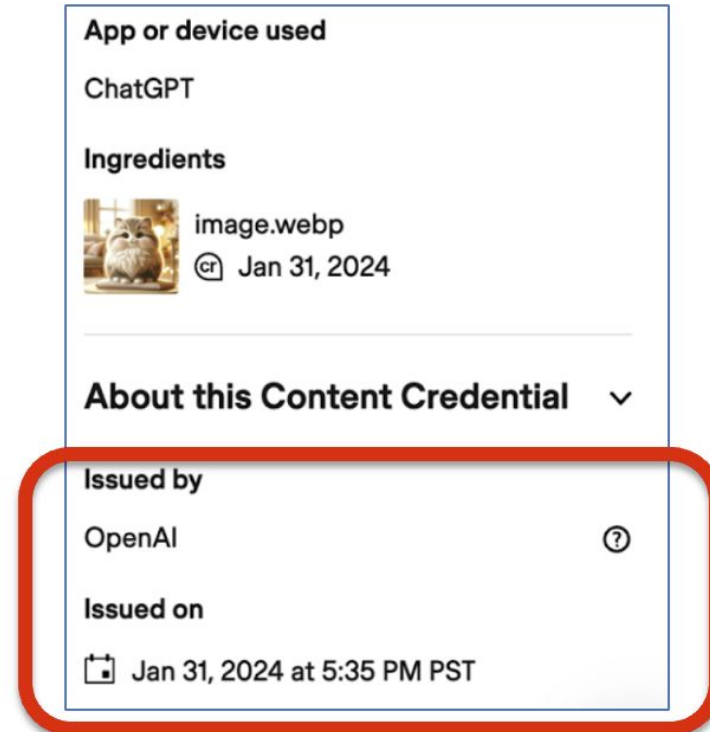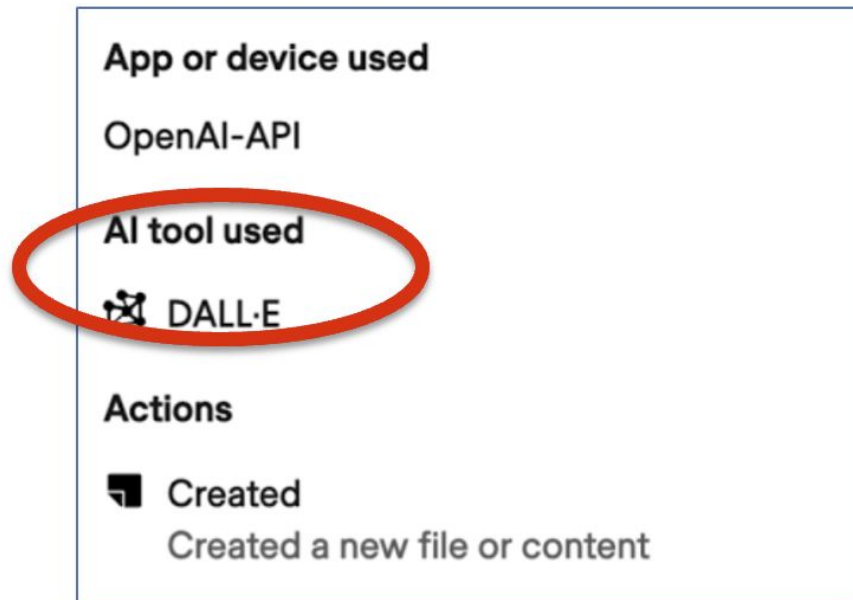
**Partnership for greater trust in digital photography: Leica and Content Authenticity Initiative**



embedded certified signing key  $sk_1$

location timestamp

📷 signature

C2PA

verify metadata by checking C2PA signature

Adobe    BBC    intel    Microsoft    PUBLICIS GROUPE    SONY    Truepic

# Not just news organizations…

# Is this a cat-and-mouse game?

The picture-of-picture attack:

C2PA camera

C2PA Image

Now every v

Many other challenges

(1) Key extraction and revocation (PKI)

(2) Privacy → group signatures

(3) GPS spoofing

- Is this

- Can attacker defeat the filter?

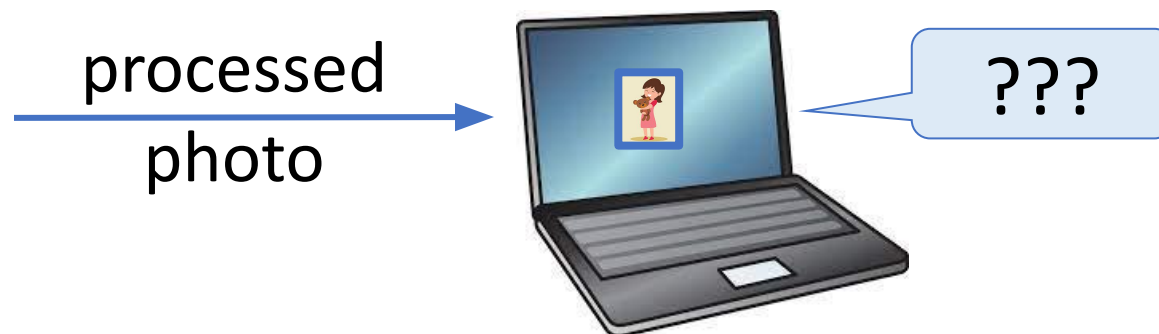[see Chimera, RWC, to appear "Harms Modeling the C2PA," 2022 ]

# A Problem: Post-Processing

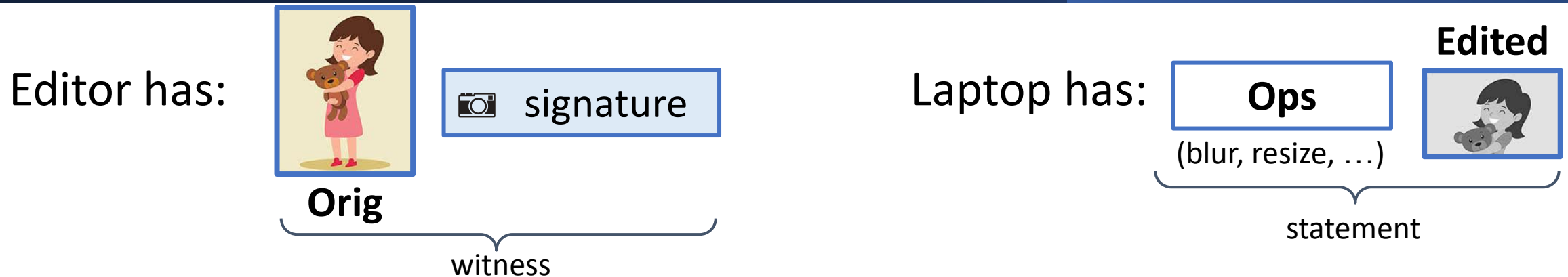Newspapers often process photos before publication

- At minimum, need to resize (90 MB $\rightarrow$ 8 MB)
- Allowable operations from the *Associated Press:*
  
  cropping, grayscale, blurring…

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

processed photo →

???

The solution proposed by C2PA is not ideal … can we do better?

# A Cryptographic Solution: zk-SNARKs!

Editor has: **Orig** signature

**Edited**

Laptop has: **Ops** (blur, resize, …)

witness

statement

Editor creates a proof $\pi$ that:

I know (***Orig, sig***) such that:
  1. ***sig*** is a valid signature on ***Orig***
  2. ***Edited*** is the result of applying ***Ops*** to ***Orig***
  3. metadata(***Edited***) = metadata(***Orig***)

Do we need ZK ZK?

$\Rightarrow$ Laptop verifies $\pi$ and shows metadata to user

# How to prove?

**π**

I know (**Orig**, **sig**) such that:

   1. **sig** is a valid signature on **Orig**

   2. **Edited** is the result of applying **Ops** to **Orig**    **?**

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

# Verifying Edits in a SNARK Prover

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

- New tools enable faster development and bigger statements!
  - Plonky2 library  ("Plonk PIOP" + FRI PCS)
    - Write arithmetic circuit $C_{edit}$   s.t.   $C_{edit}(\textbf{\textit{Orig}}) = \textbf{\textit{Edited}}$

Our work: proof for a 6000 x 4000 image using Plonky2

- resize, crop, grayscale, blur $\rightarrow$ proof gen. time $\leq$ 4 minutes

- Proof size:  $\approx$100 KB ($\ll$image size),  verification time: 0.7 sec
  (can shrink proof with recursion)

# Verifying Signatures in a SNARK Prover

I know (**Orig**, **sig**) such that:

1. **sig** is a valid signature on **Orig**

2. **Edited** is the result of applying **Ops** to **Orig**

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

**?**

**✓**

**Problem**: the SNARK proof must check that a pair (**Orig**, **h**) satisfies **h** = hash(**Orig**)...but **Orig** is 90MB!

# Verifying Signatures in a SNARK Prover

We propose two methods:

### Attempt 1

| Verify Hash |
|---|
| I know (***Orig***, ***hash***) such that: $hash = $ SHA256(***Orig***) |

**Too slow for 90 MB!**

### Attempt 2

| Verify Hash |
|---|
| I know (***Orig***, ***hash***) such that: $hash = $ Poseidon(***Orig***) |

**SNARK-friendly hash ...but still too slow for 90MB!**

| **(1) Verify Lattice Hash** |
|---|
| I know (***Orig***, ***hash***) such that: $hash = $ Poseidon(LatticeHash(***Orig***)) |

**good for camera and prover**

| **(2) Verify a PCS** |
|---|
| I know (***Orig***, ***hash***) such that: $hash = $ **PCS**(***Orig***) |

**great for prover!**

# Verifying Signatures in a SNARK Prover

# How is Method 2 secure?

PLONK proof $\pi$ proves that $C_{edit}(\textbf{\textit{Orig'}}) = \textbf{\textit{Edited}}$

where **Orig'** is provided as witness data
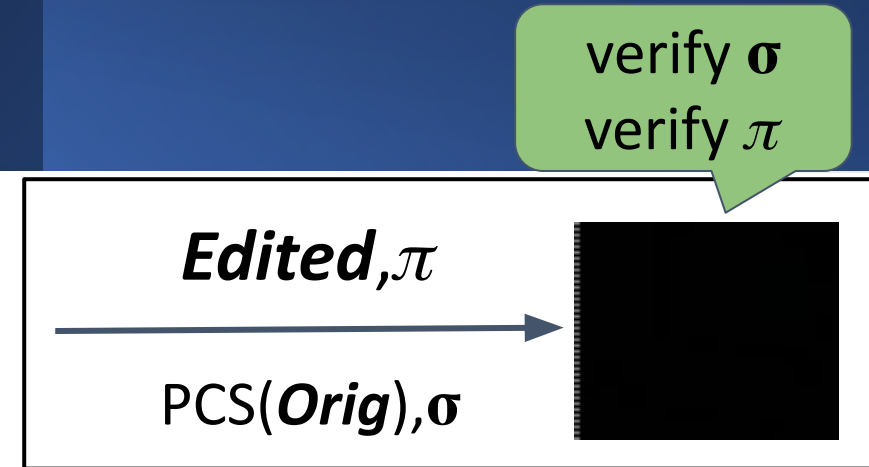
Problem: what if **Orig ≠ Orig'** ??

    → Then edited image is unrelated to camera's image

Solution: non-black-box use of PLONK!

# How is Method 2 secure?

verify $\boldsymbol{\sigma}$
verify $\pi$

PLONK proof $\pi$ proves that $C_{edit}(\boldsymbol{Orig'}) = \boldsymbol{Edited}$
where $\boldsymbol{Orig'}$ is provided as witness data

$\boldsymbol{Edited},\pi$

PCS($\boldsymbol{Orig}$),$\boldsymbol{\sigma}$

Partial explanation of how to produce PLONK proof $\pi$:

1. Encode $C_{edit}$ execution
tableaux as a polynomial

$\boldsymbol{T}(x)$

Some $\boldsymbol{T}(x)$ evaluations
encode the witness $\boldsymbol{Orig'}$

2. Generate proof:

a. Compute $\boldsymbol{com_T}$ = PCS($\boldsymbol{T}(x)$)
b. Prove gates evals in tableaux are correct
c. Prove circuit wiring in tableaux is correct

This works even if tableaux is
committed via multiple polynomials!

# Tradeoffs of Signing PCS vs. Lattice Hash

Signing PCS(**_Orig_**) takes signature verification out of the SNARK circuit

verify $\sigma$
verify $\pi$

**_Edited_**, $\pi$

PCS(**_Orig_**), $\sigma$

… but computing a PCS commitment is
        not feasible on a commercial camera

○ Suitable for a cloud AI image generator

○ Can be offloaded to an untrusted server

# Verifying Signatures in a SNARK Prover

π

I know (**Orig**, **sig**) such that:

   1. **sig** is a valid signature on **Orig** ✓

   2. **Edited** is the result of applying **Ops** to **Orig** ✓

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

# Conclusions

Succinct proofs have become practical and easy to use

- An amazing success of theory of CS

- Development driven by blockchain
  but **many** non-blockchain applications

C2PA: a playground for many cryptographic techniques
- Many challenges to explore…