

# Using Zero-Knowledge Proofs to Fight Disinformation

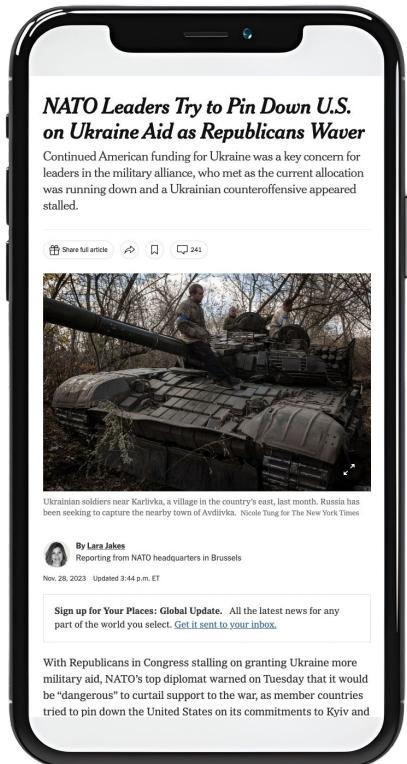
Trisha Datta and Dan Boneh  
Stanford University

# Image Provenance Verification

- Want to verify when and where photos were taken

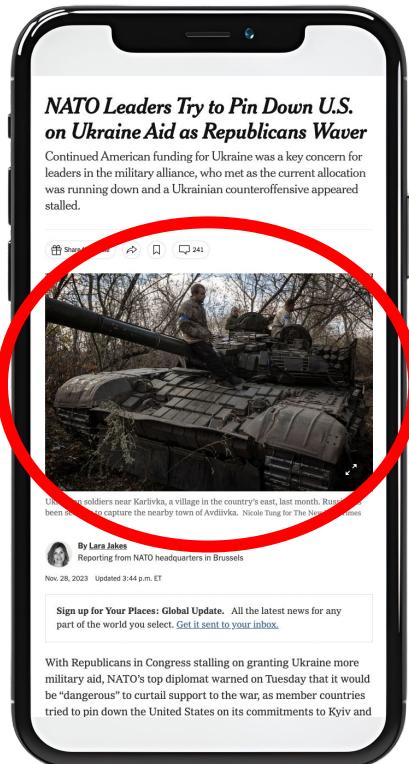
# Image Provenance Verification

- Want to verify when and where photos were taken
- Important for news articles



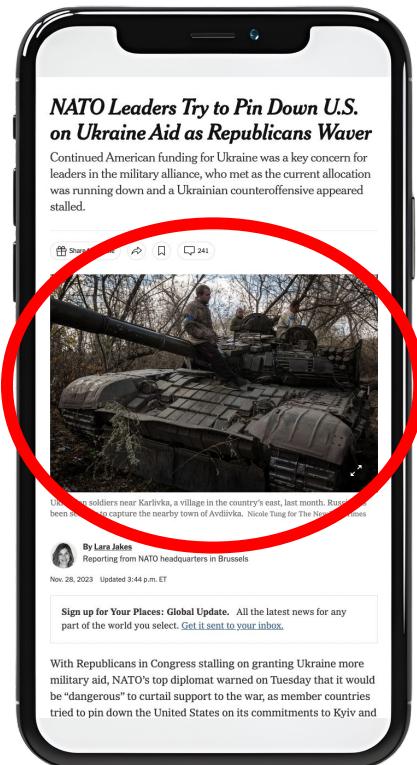
# Image Provenance Verification

- Want to verify when and where photos were taken
- Important for news articles



# Image Provenance Verification

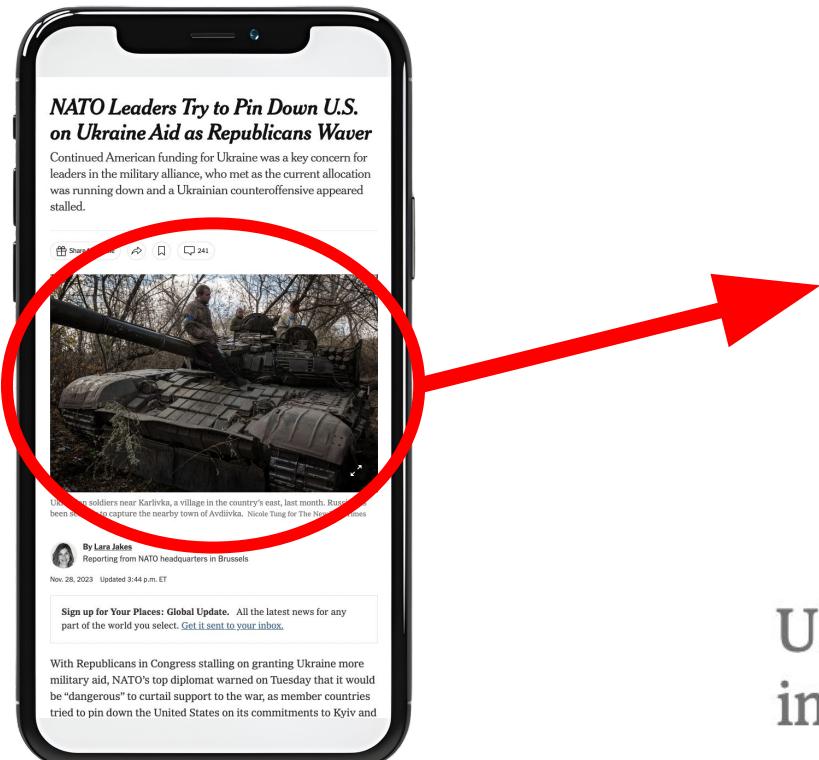
- Want to verify when and where photos were taken
- Important for news articles



Ukrainian soldiers near Karlivka, a village in the country's east, last month.

# Image Provenance Verification

- Want to verify when and where photos were taken
- Important for news articles



Ukrainian soldiers near Karlivka, a village in the country's east, last month.

# Ukraine conflict: Many misleading images have been shared online

By Alistair Coleman & Shayan Sardarizadeh  
BBC Monitoring

24 February 2022

# Ukraine conflict: Many misleading images have been shared online

By Alistair Cole  
BBC Monitoring

24 February 2022

## 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.

# Ukraine conflict: Many misleading images have been shared online

By Alistair Cole  
BBC Monitoring

24 February 2022

## Fact-checking videos and pictures from Ukraine

Since Russia invaded Ukraine, many misleading images and pictures have been shared online.

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

# Ukraine conflict: Many misleading images have been shared online

By Alistair Cole  
BBC Monitoring

24 February 2024

## Fact-checking videos and pictures from Ukraine

Since Russia and pictures

False social media posts are hindering earthquake relief efforts in Turkey. You can help stop them.

THE NEWS / UKRAINE

**BBC Breakfast uses old footage of Russian parade rehearsal to show invasion of Ukraine**

# C2PA: A Content Provenance Standard

## Sony Unlocks In-Camera Forgery-Proof Technology

04 Aug, 2022

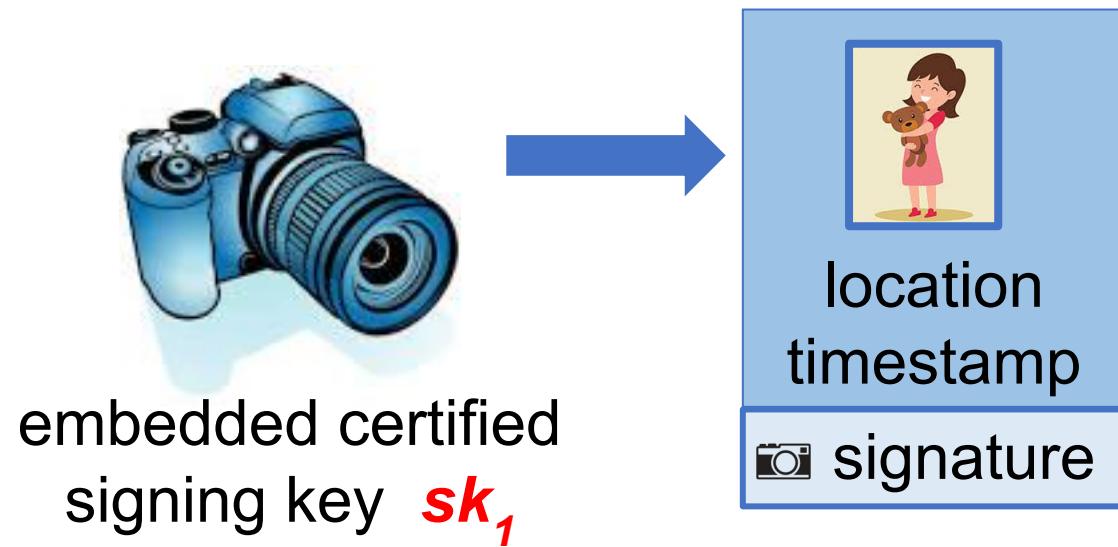


embedded certified  
signing key ***sk<sub>1</sub>***,

# C2PA: A Content Provenance Standard

## Sony Unlocks In-Camera Forgery-Proof Technology

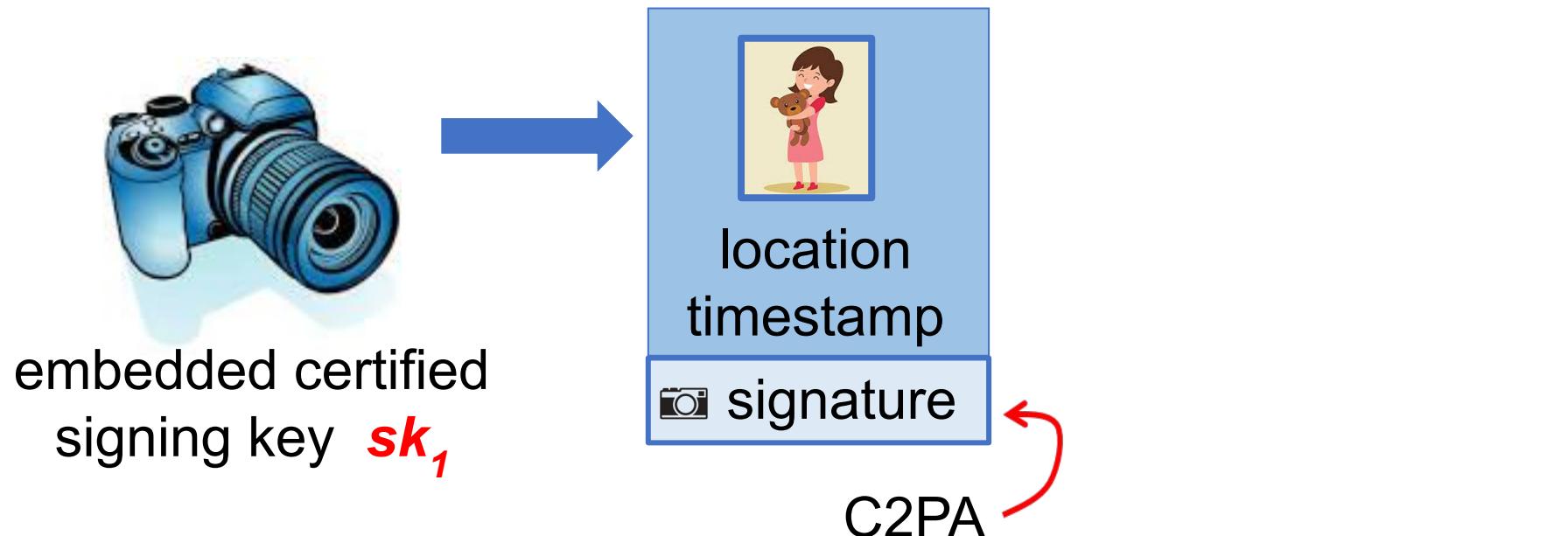
04 Aug, 2022



# C2PA: A Content Provenance Standard

## Sony Unlocks In-Camera Forgery-Proof Technology

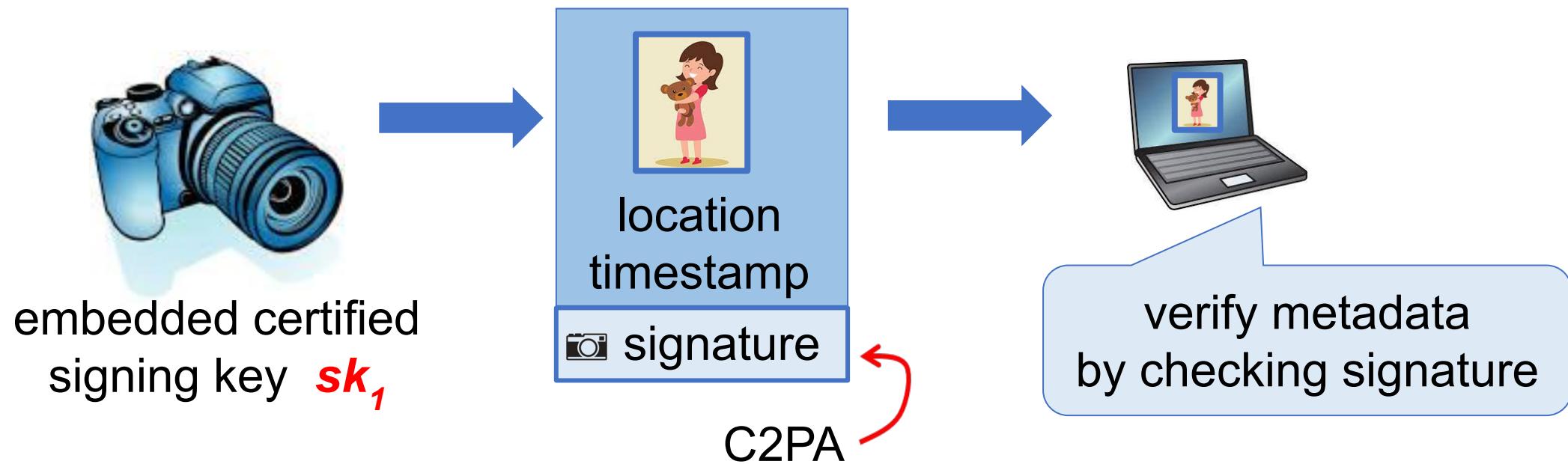
04 Aug, 2022



# C2PA: A Content Provenance Standard

## Sony Unlocks In-Camera Forgery-Proof Technology

04 Aug, 2022

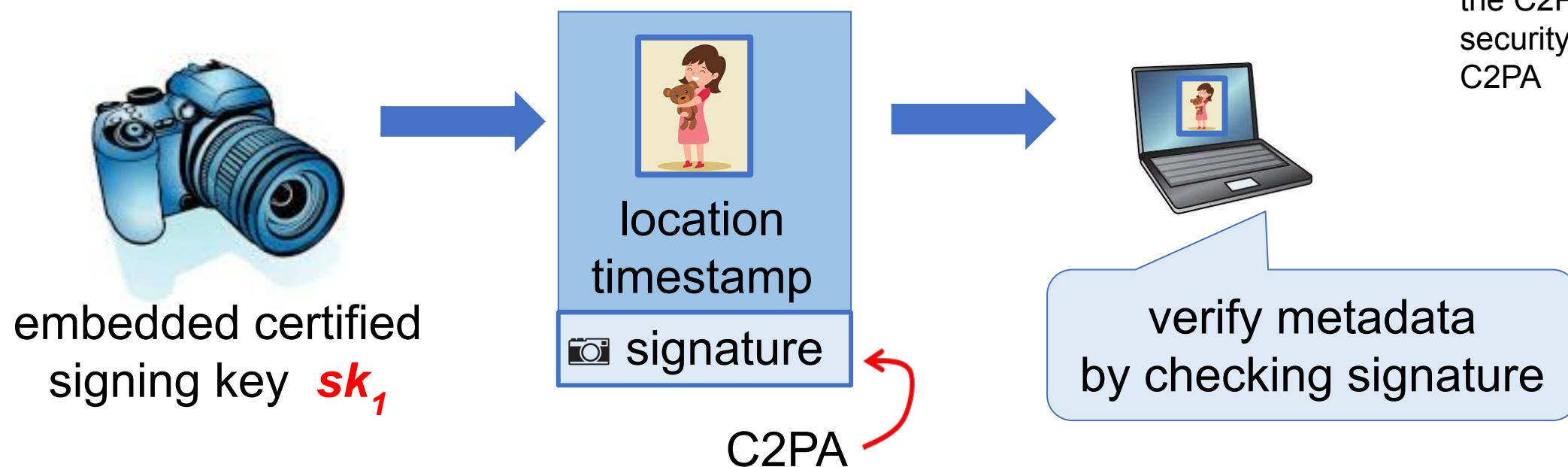


# C2PA: A Content Provenance Standard

## Sony Unlocks In-Camera Forgery-Proof Technology

04 Aug, 2022

See Rivadeneira,  
“Harms Modelling in  
the C2PA,” 2022 for  
security analysis of  
C2PA



# A Problem: Post-Processing

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

# A Problem: Post-Processing

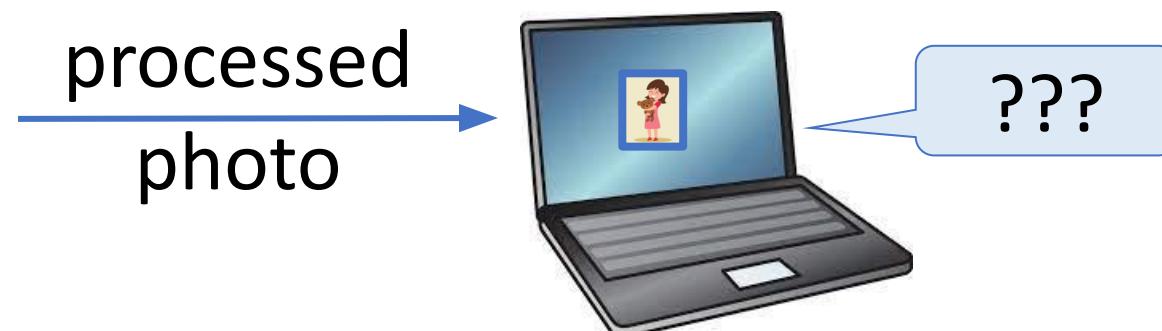
- Newspapers often process photos before publication
  - At minimum, need to resize (90 MB → 8 MB)
  - Allowable operations from the *Associated Press*: cropping, grayscale, ...

**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 → 8 MB)
  - Allowable operations from the *Associated Press*: cropping, grayscale, ...

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

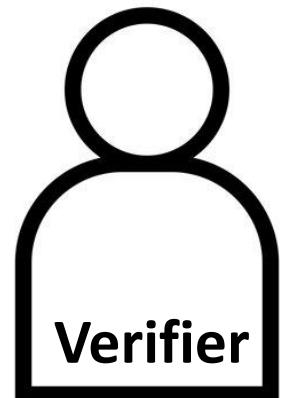
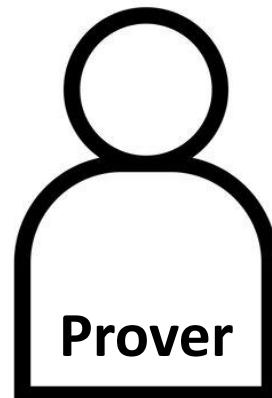


# Proposed Solution: ZK-SNARKs!

- ZK-SNARK: efficiently verifiable statement about a secret witness

# Proposed Solution: ZK-SNARKs!

- ZK-SNARK: efficiently verifiable statement about a secret witness



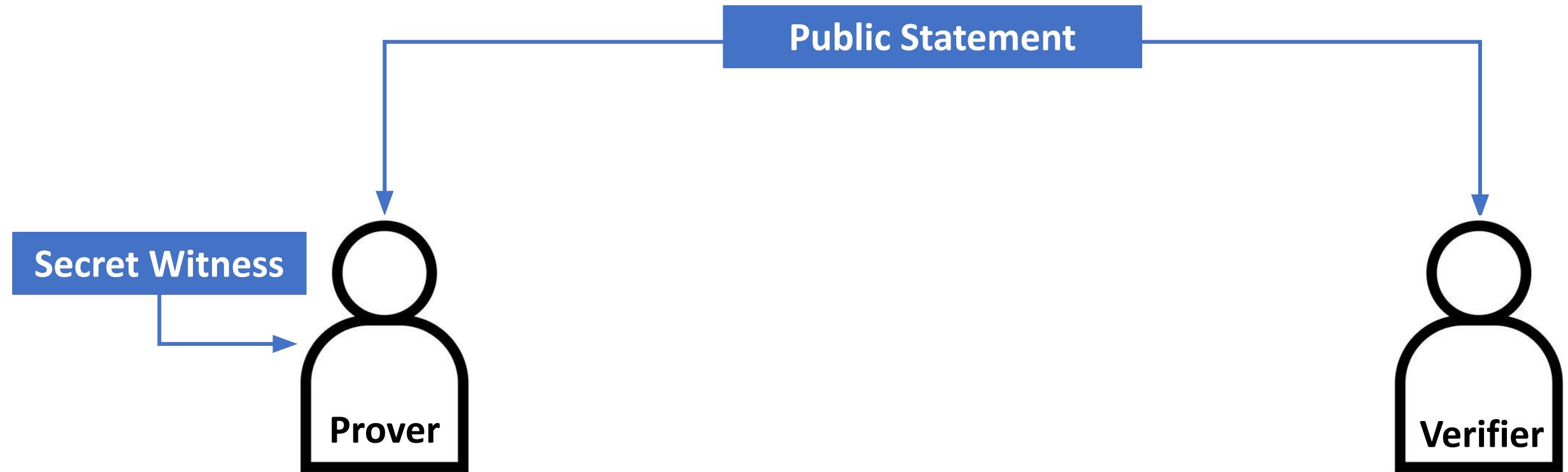
# Proposed Solution: ZK-SNARKs!

- ZK-SNARK: efficiently verifiable statement about a secret witness



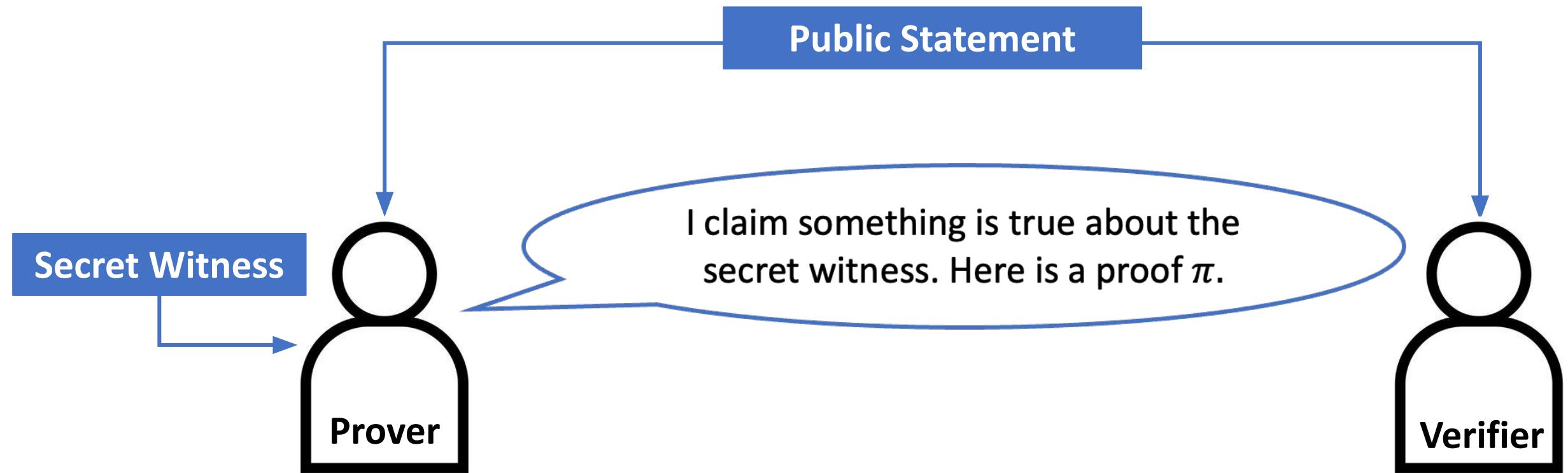
# Proposed Solution: ZK-SNARKs!

- ZK-SNARK: efficiently verifiable statement about a secret witness



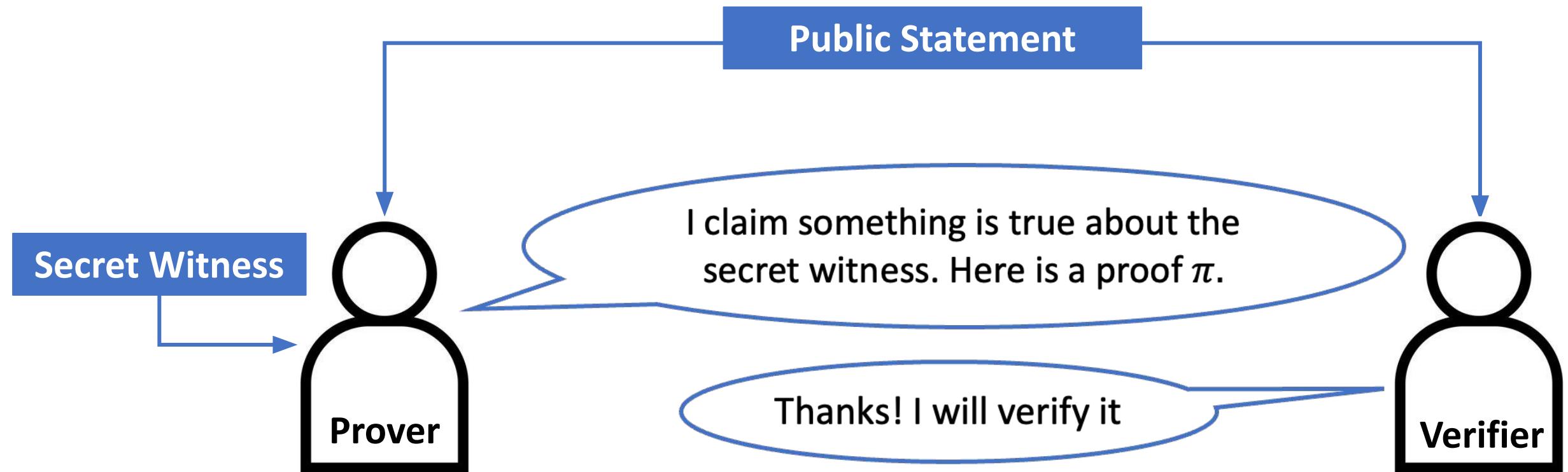
# Proposed Solution: ZK-SNARKs!

- ZK-SNARK: efficiently verifiable statement about a secret witness



# Proposed Solution: ZK-SNARKs!

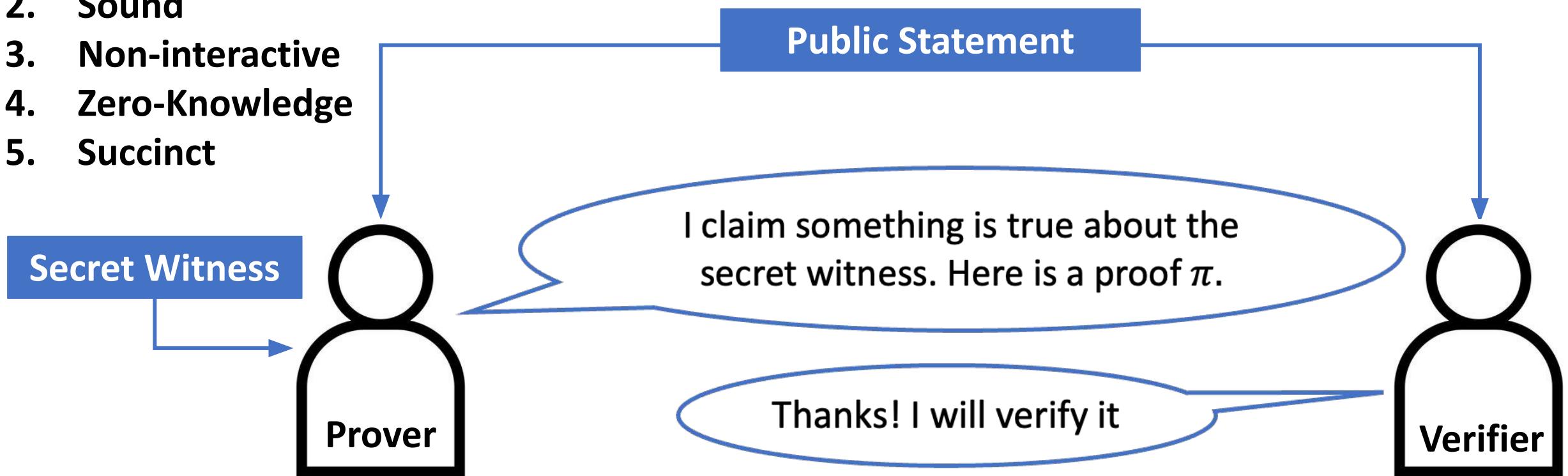
- ZK-SNARK: efficiently verifiable statement about a secret witness



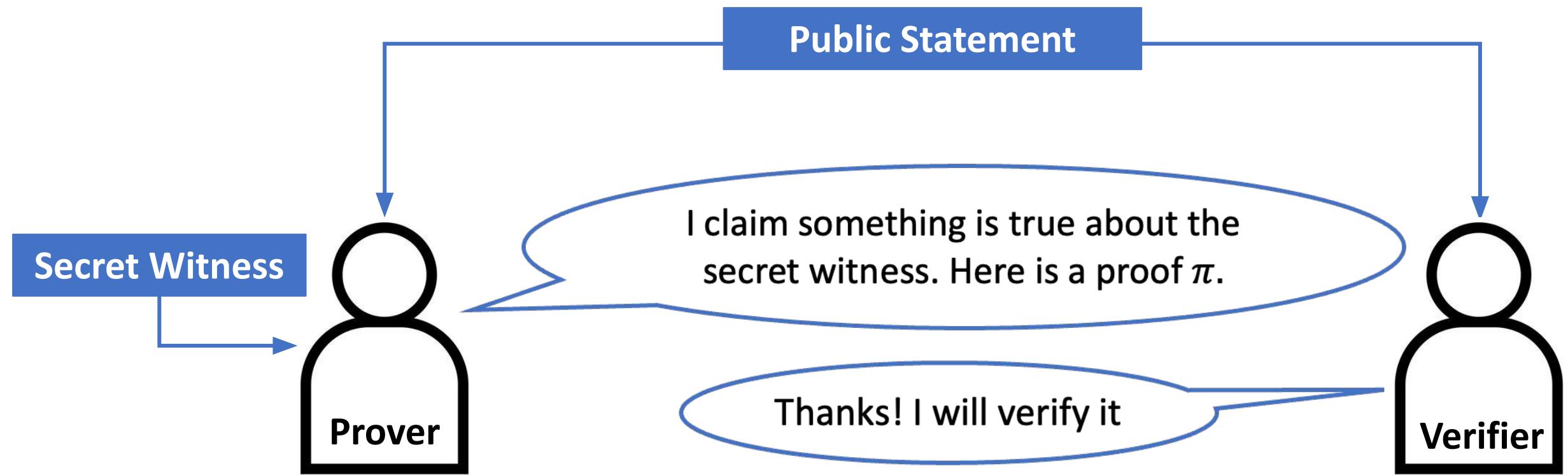
# Proposed Solution: ZK-SNARKs!

- ZK-SNARK: efficiently verifiable statement about a secret witness

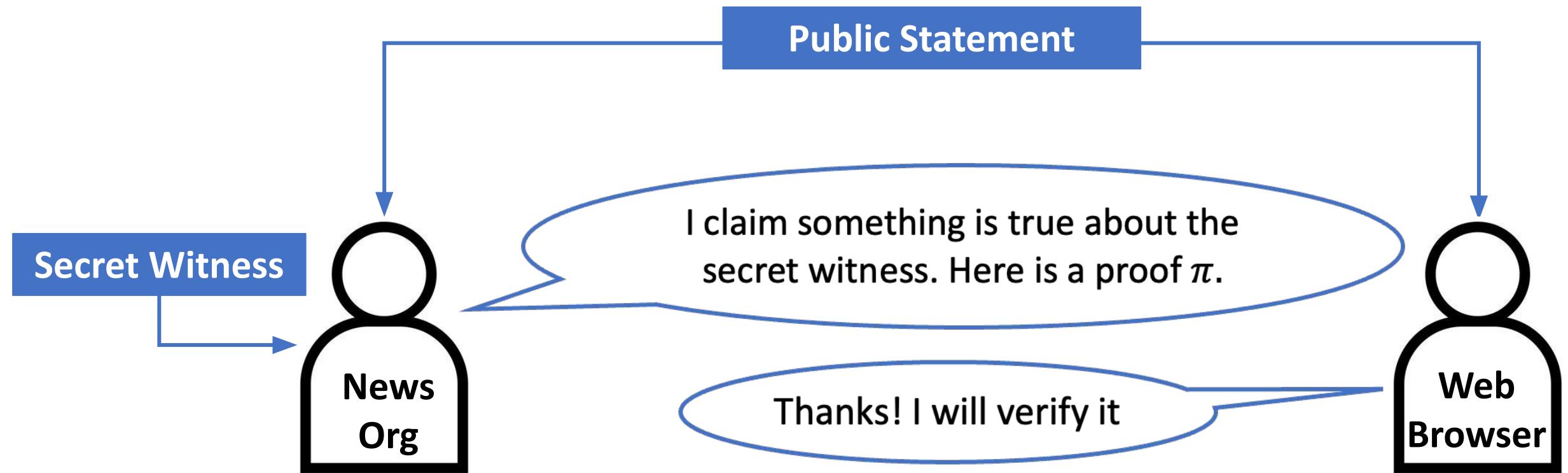
1. Complete
2. Sound
3. Non-interactive
4. Zero-Knowledge
5. Succinct



# Proposed Solution: ZK-SNARKs!



# Proposed Solution: ZK-SNARKs!



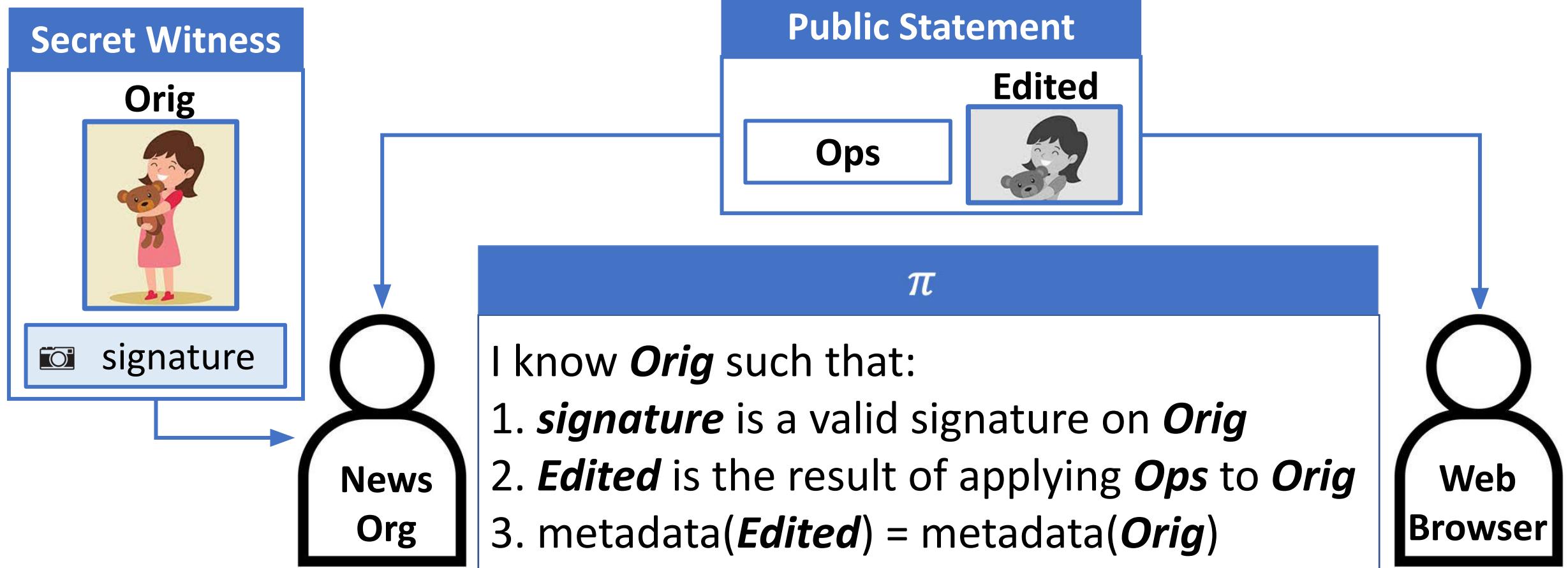
# Proposed Solution: ZK-SNARKs!



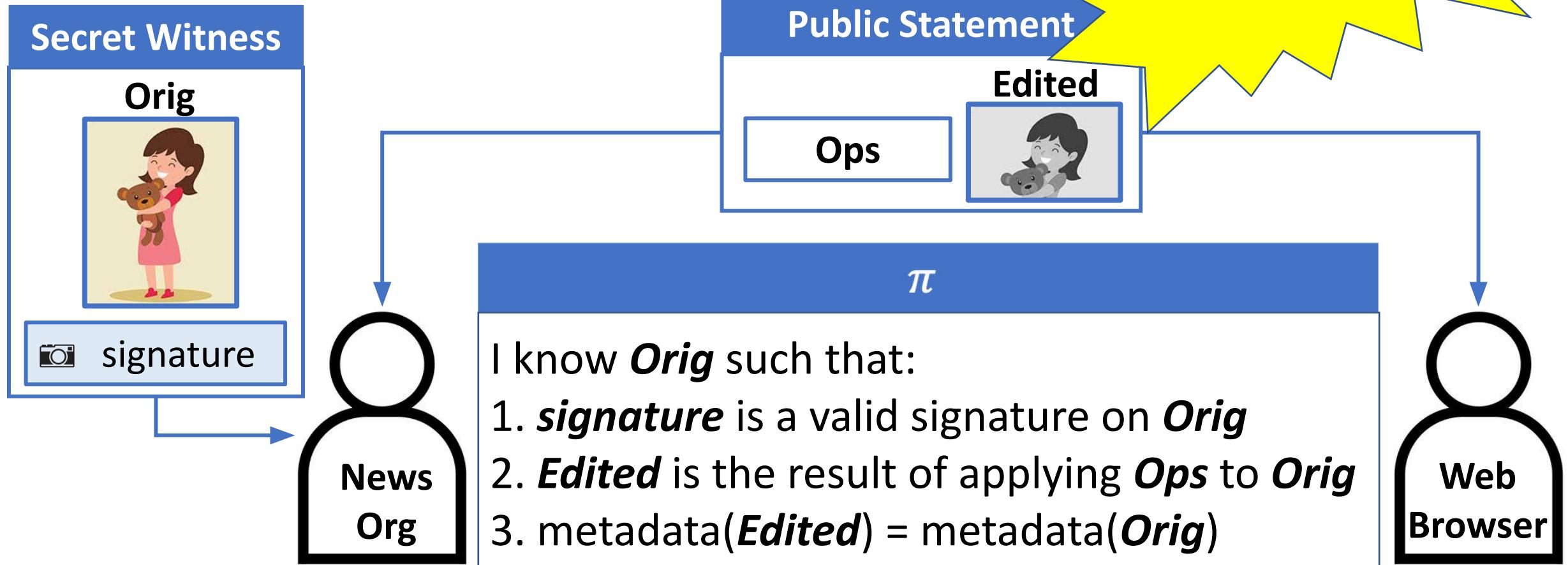
# Proposed Solution: ZK-SNARKs!



# Proposed Solution: ZK-SNARKs!

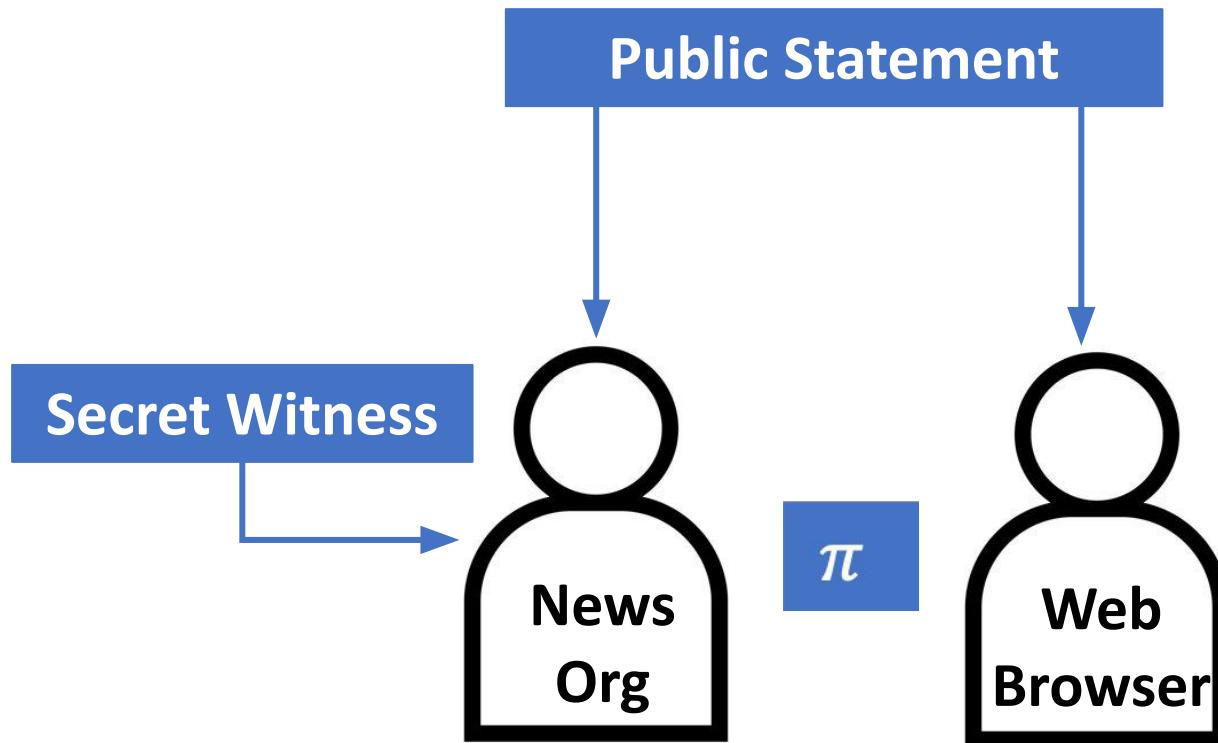


# Proposed Solution: ZK-SNARKs!

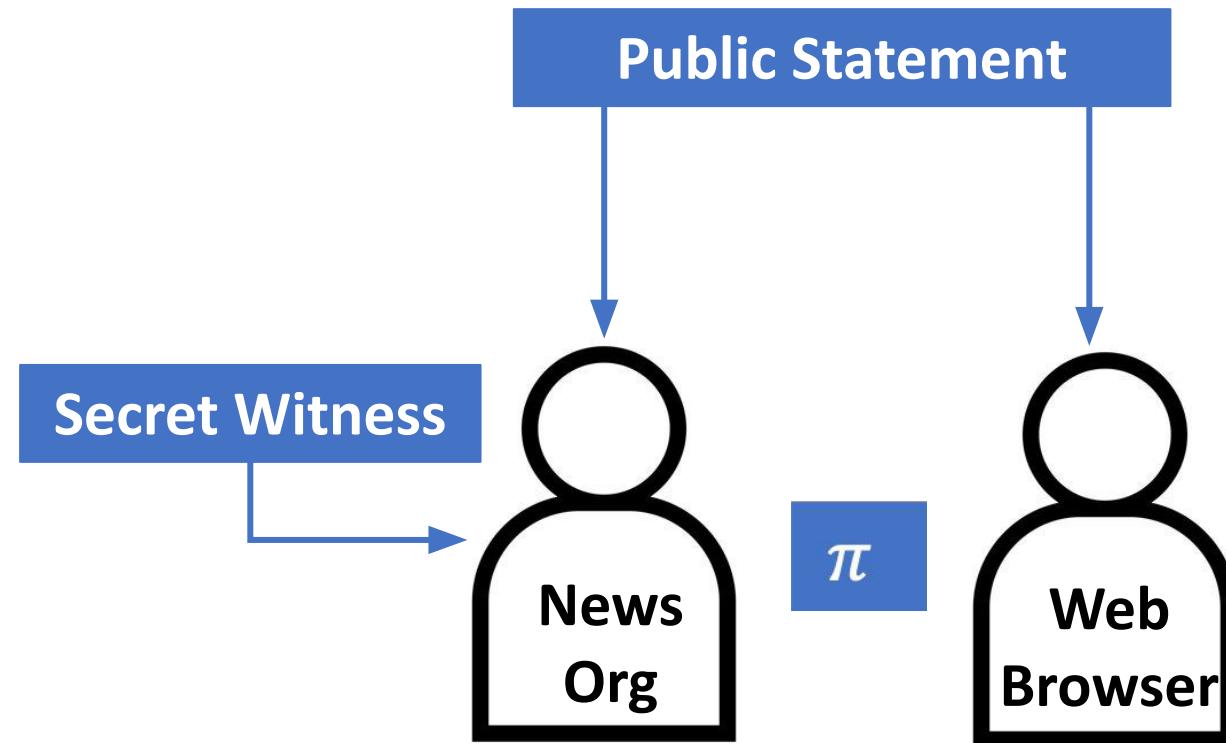


# Proposed Solution: ZK-SNARKs!

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

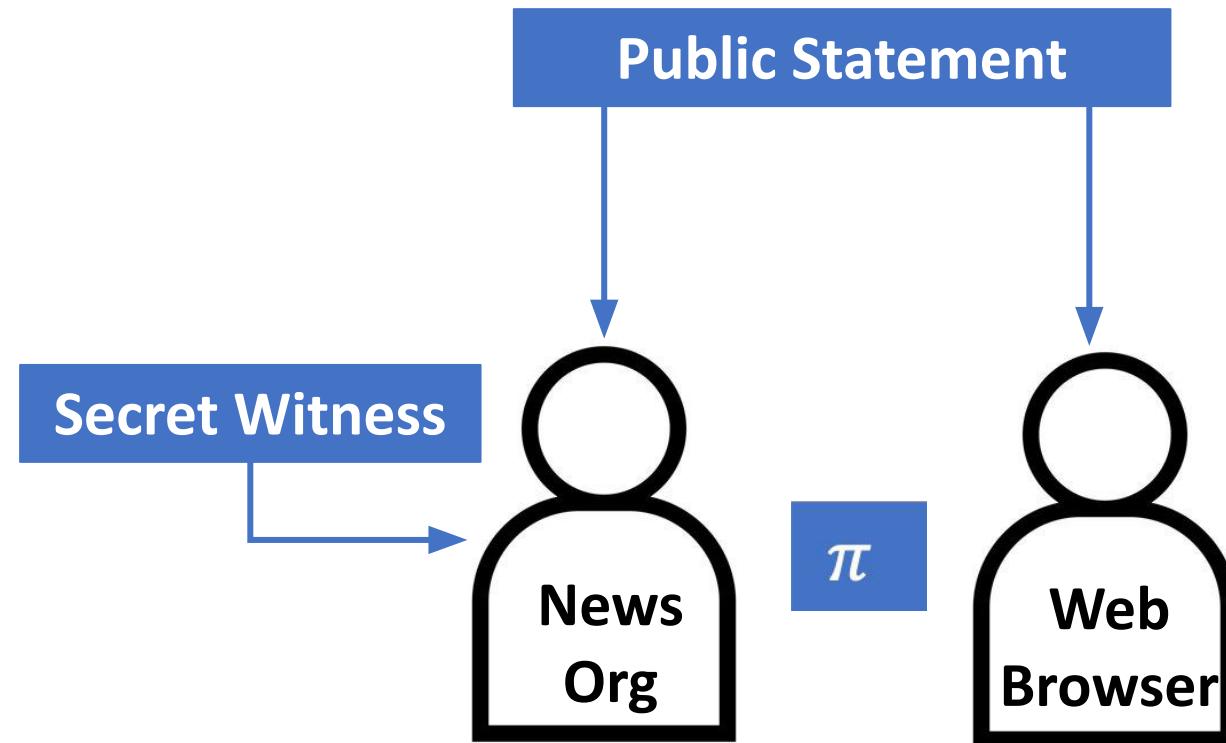


# Proposed Solution: ZK-SNARKs!



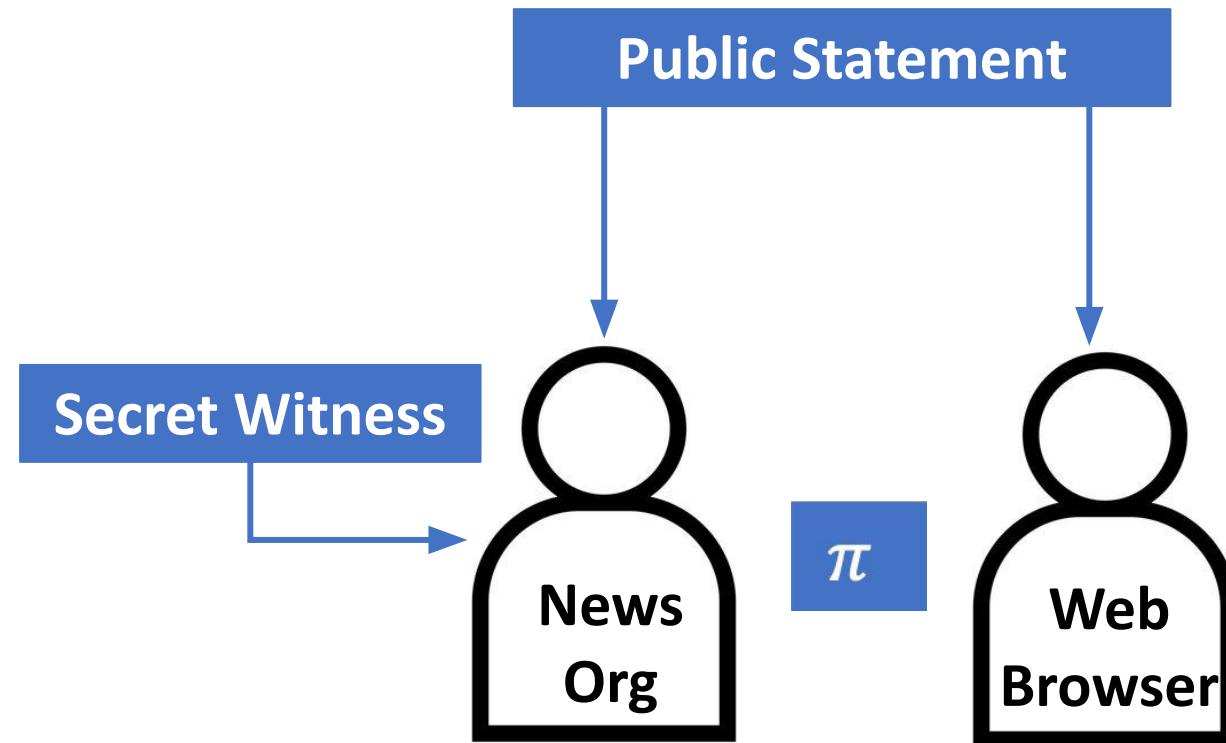
1. **Complete/Sound:** verifier doesn't need to trust prover
2. **Non-interactive:** interactivity would entail unique proofs for each browser

# Proposed Solution: ZK-SNARKs!



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

# Proposed Solution: ZK-SNARKs!



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

# Verifying Signatures in a SNARK Prover

$\pi$

I know  $Orig$  such that:

1. ***signature*** is a valid signature on  $Orig$
2. ***Edited*** is the result of applying  $Ops$  to  $Orig$
3.  $\text{metadata}(\text{Edited}) = \text{metadata}(Orig)$

# Verifying Signatures in a SNARK Prover

$\pi$

I know *Orig* such that:

1. ***signature*** is a valid signature on *Orig*
2. ***Edited*** is the result of applying ***Ops*** to *Orig*
3.  $\text{metadata}(\text{Edited}) = \text{metadata}(\text{Orig})$

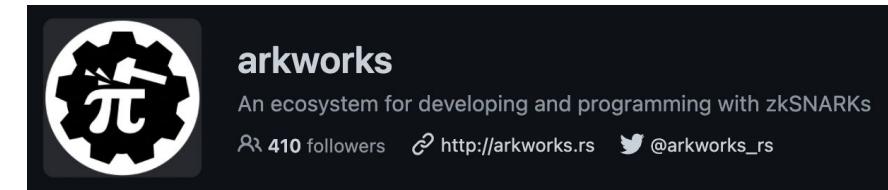
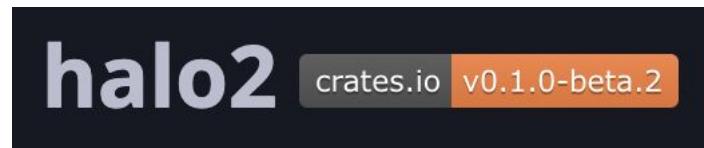
?

# 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!



The Noir Programming Language

# 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

# Verifying Signatures in a SNARK Prover

$\pi$

I know  $Orig$  such that:

1. ***signature*** is a valid signature on  $Orig$
2. ***Edited*** is the result of applying ***Ops*** to  $Orig$
3.  $\text{metadata}(\text{Edited}) = \text{metadata}(Orig)$



# Verifying Signatures in a SNARK Prover

$\pi$

I know  $Orig$  such that:

1. ***signature*** is a valid signature on  $Orig$
2. ***Edited*** is the result of applying ***Ops*** to  $Orig$
3.  $\text{metadata}(\text{Edited}) = \text{metadata}(Orig)$

?



# Verifying Signatures in a SNARK Prover

## Attempt 1

$\pi$  (Signature)

I know  $(\text{Orig}, \text{hash})$  such  
that:

$$\text{hash} = \text{SHA256}(\text{Orig})$$

# Verifying Signatures in a SNARK Prover

Attempt 1

$\pi$  (Signature)

I know  $(\text{Orig}, \text{hash})$  such that:

$\text{hash} = \text{SHA256}(\text{Orig})$



Too slow for  
30 MP!

# Verifying Signatures in a SNARK Prover

Attempt 1

$\pi$  (Signature)

I know  $(\text{Orig}, \text{hash})$  such that:

$$\text{hash} = \text{SHA256}(\text{Orig})$$



Too slow for  
30 MP!

Attempt 2

$\pi$  (Signature)

I know  $(\text{Orig}, \text{hash})$  such that:

$$\text{hash} = \text{Poseidon}(\text{Orig})$$

# Verifying Signatures in a SNARK Prover

Attempt 1

$\pi$  (Signature)

I know  $(\text{Orig}, \text{hash})$  such that:

$$\text{hash} = \text{SHA256}(\text{Orig})$$



Too slow for  
30 MP!

Attempt 2

$\pi$  (Signature)

I know  $(\text{Orig}, \text{hash})$  such that:

$$\text{hash} = \text{Poseidon}(\text{Orig})$$



SNARK-friendly hash

# Verifying Signatures in a SNARK Prover

Attempt 1

$\pi$  (Signature)

I know  $(\text{Orig}, \text{hash})$  such that:

$$\text{hash} = \text{SHA256}(\text{Orig})$$



Too slow for  
30 MP!

Attempt 2

$\pi$  (Signature)

I know  $(\text{Orig}, \text{hash})$  such that:

$$\text{hash} = \text{Poseidon}(\text{Orig})$$



SNARK-friendly hash  
...but still too slow for  
30MP!

# Verifying Signatures in a SNARK Prover

Attempt 1

$\pi$  (Signature)

I know  $(\text{Orig}, \text{hash})$  such that:

$$\text{hash} = \text{SHA256}(\text{Orig})$$



Too slow for  
30 MP!

Attempt 2

$\pi$  (Signature)

I know  $(\text{Orig}, \text{hash})$  such that:

$$\text{hash} = \text{Poseidon}(\text{Orig})$$



SNARK-friendly hash  
...but still too slow for  
30MP!

Attempt 3

$\pi$  (Signature)

I know  $(\text{Orig}, \text{hash})$  such that:  
 $\text{hash} = \text{Poseidon}(\text{LatticeHash}(\text{Orig}))$

# Verifying Signatures in a SNARK Prover

Attempt 1

$\pi$  (Signature)

I know  $(\text{Orig}, \text{hash})$  such that:

$$\text{hash} = \text{SHA256}(\text{Orig})$$



Too slow for  
30 MP!

Attempt 2

$\pi$  (Signature)

I know  $(\text{Orig}, \text{hash})$  such that:

$$\text{hash} = \text{Poseidon}(\text{Orig})$$



SNARK-friendly hash  
...but still too slow for  
30MP!

Attempt 3

$\pi$  (Signature)

I know  $(\text{Orig}, \text{hash})$  such that:  
 $\text{hash} = \text{Poseidon}(\text{LatticeHash}(\text{Orig}))$



# Verifying Signatures in a SNARK Prover

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

# Verifying Signatures in a SNARK Prover

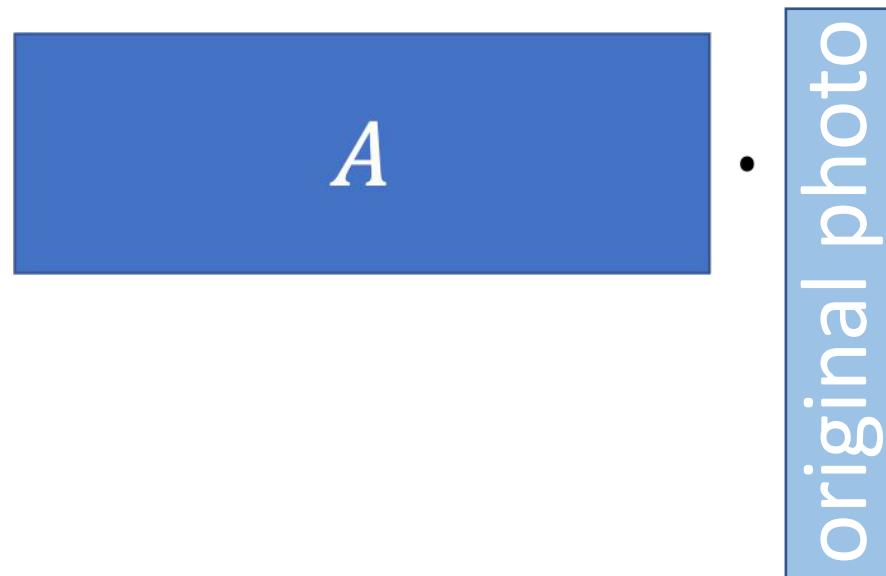
Poseidon hash of lattice hash [GGH'96 , SCMPGLW'08]

$$A$$

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

# Verifying Signatures in a SNARK Prover

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



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

# Verifying Signatures in a SNARK Prover

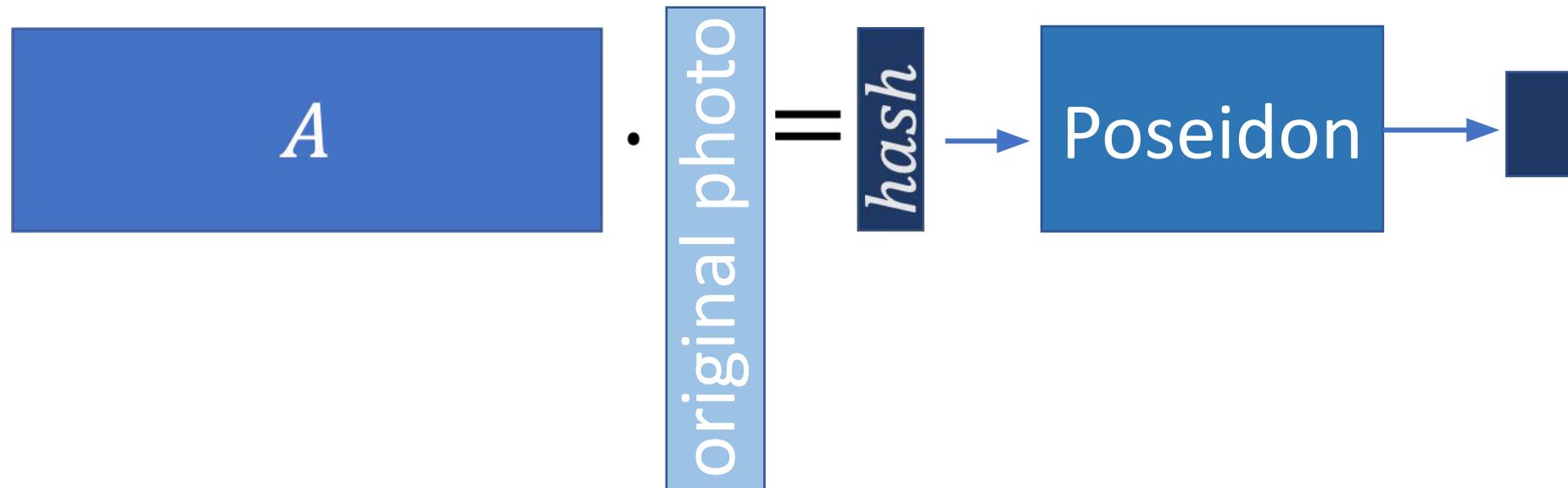
Poseidon hash of lattice hash [GGH'96 , SCMPGLW'08]

$$A \cdot \text{Original photo} = \text{hash}$$

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

# Verifying Signatures in a SNARK Prover

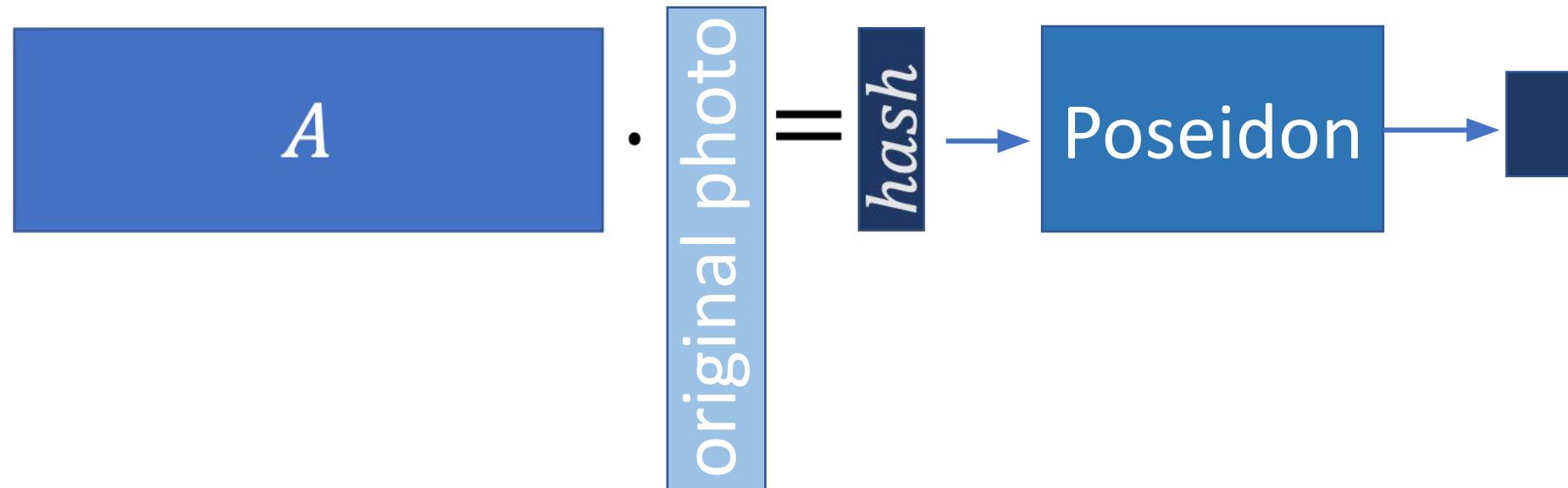
Poseidon hash of lattice hash [GGH'96 , SCMPGLW'08]



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

# Verifying Signatures in a SNARK Prover

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



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

# Verifying Signatures in a SNARK Prover

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

# Verifying Signatures in a SNARK Prover

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

$$R = 4$$

$\vec{x}$	2	0	2
-----------	---	---	---

# Verifying Signatures in a SNARK Prover

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

i.e.,  $\vec{x} \in \{0, 1, \dots, R - 1\}^n$ :

- $\vec{y} := [0, 1, \dots, R - 1]$

$$R = 4$$

$\vec{x}$	2	0	2
-----------	---	---	---

$\vec{y}$	0	1	2	3
-----------	---	---	---	---

# Verifying Signatures in a SNARK Prover

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

i.e.,  $\vec{x} \in \{0, 1, \dots, R - 1\}^n$ :

- $\vec{y} := [0, 1, \dots, R - 1]$
- $\vec{z} := \text{sort}(\vec{x} || \vec{y})$

$$R = 4$$

$\vec{x}$	2	0	2
-----------	---	---	---

$\vec{y}$	0	1	2	3
-----------	---	---	---	---

# Verifying Signatures in a SNARK Prover

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

i.e.,  $\vec{x} \in \{0, 1, \dots, R - 1\}^n$ :

- $\vec{y} := [0, 1, \dots, R - 1]$
- $\vec{z} := \text{sort}(\vec{x} || \vec{y})$

$$R = 4$$

$\vec{x}$	2	0	2
-----------	---	---	---

$\vec{y}$	0	1	2	3
-----------	---	---	---	---

$\vec{z}$	2	0	2	0	1	2	3
-----------	---	---	---	---	---	---	---

# Verifying Signatures in a SNARK Prover

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

i.e.,  $\vec{x} \in \{0, 1, \dots, R - 1\}^n$ :

- $\vec{y} := [0, 1, \dots, R - 1]$
- $\vec{z} := \text{sort}(\vec{x} || \vec{y})$

$$R = 4$$

$\vec{x}$	2	0	2
-----------	---	---	---

$\vec{y}$	0	1	2	3
-----------	---	---	---	---

$\vec{z}$	0	0	1	2	2	2	3
-----------	---	---	---	---	---	---	---

# Verifying Signatures in a SNARK Prover

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

i.e.,  $\vec{x} \in \{0, 1, \dots, R - 1\}^n$ :

- $\vec{y} := [0, 1, \dots, R - 1]$
- $\vec{z} := \text{sort}(\vec{x} \parallel \vec{y})$

Must show:

- $\vec{z}$  is permutation of  $\vec{x}$  and  $\vec{y}$
- $\vec{z}[i + 1] - \vec{z}[i] \in \{0, 1\}$

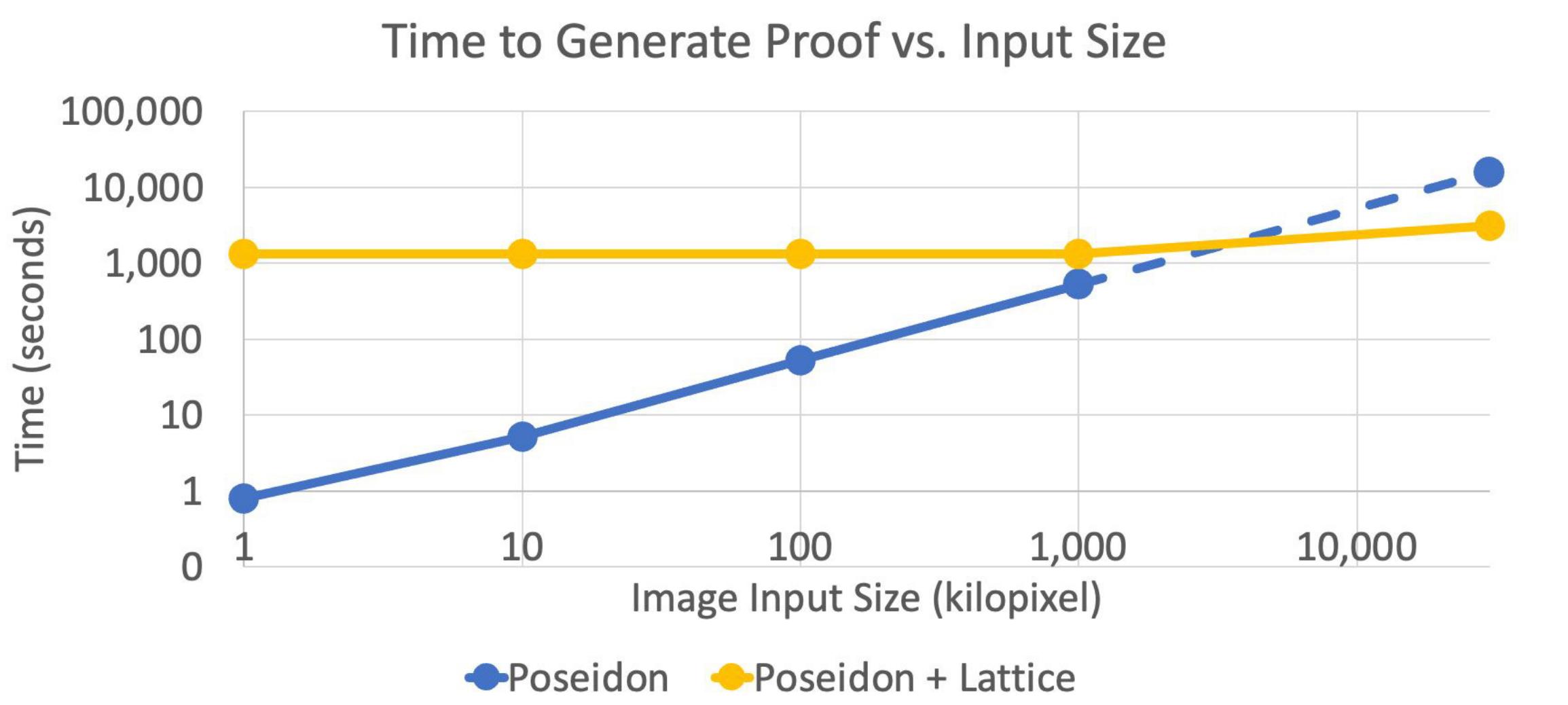
$$R = 4$$

$\vec{x}$	2	0	2
-----------	---	---	---

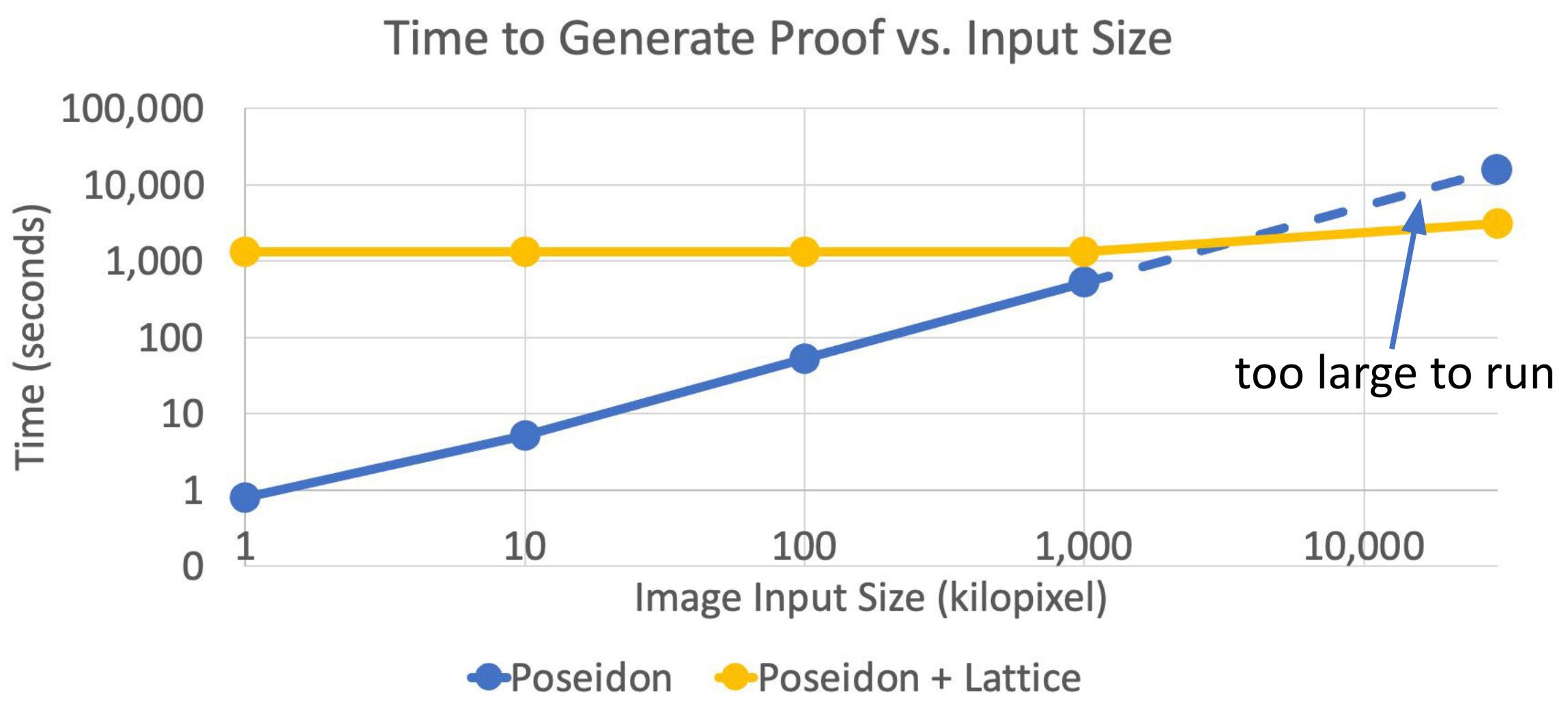
$\vec{y}$	0	1	2	3
-----------	---	---	---	---

$\vec{z}$	0	0	1	2	2	2	3
-----------	---	---	---	---	---	---	---

# Performance Results for Proving Signatures



# Performance Results for Proving Signatures



# Verifying Signatures in a SNARK Prover

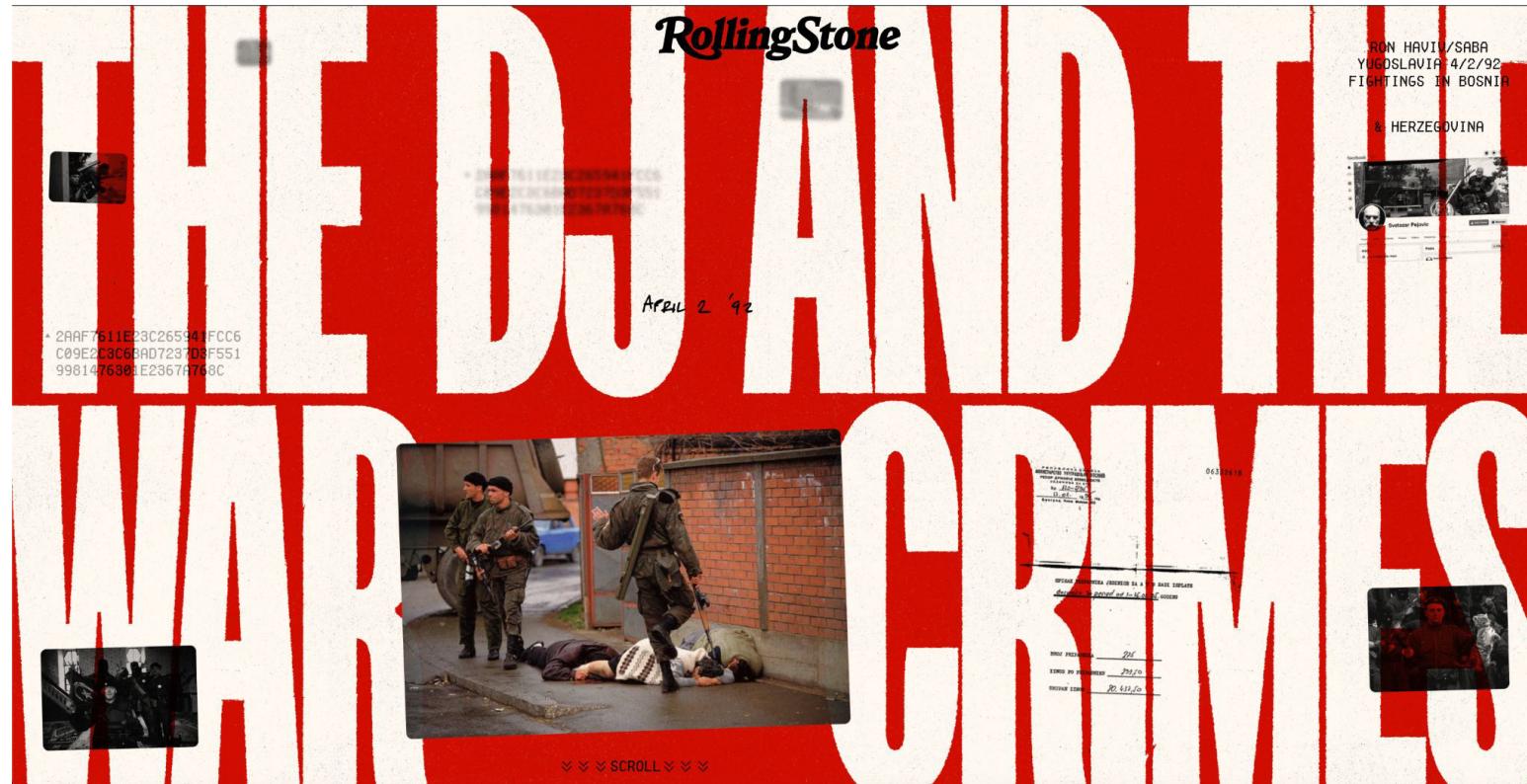
$\pi$

I know  $Orig$  such that:

1. ***signature*** is a valid signature on  $Orig$
2. ***Edited*** is the result of applying ***Ops*** to  $Orig$
3.  $\text{metadata}(\text{Edited}) = \text{metadata}(Orig)$

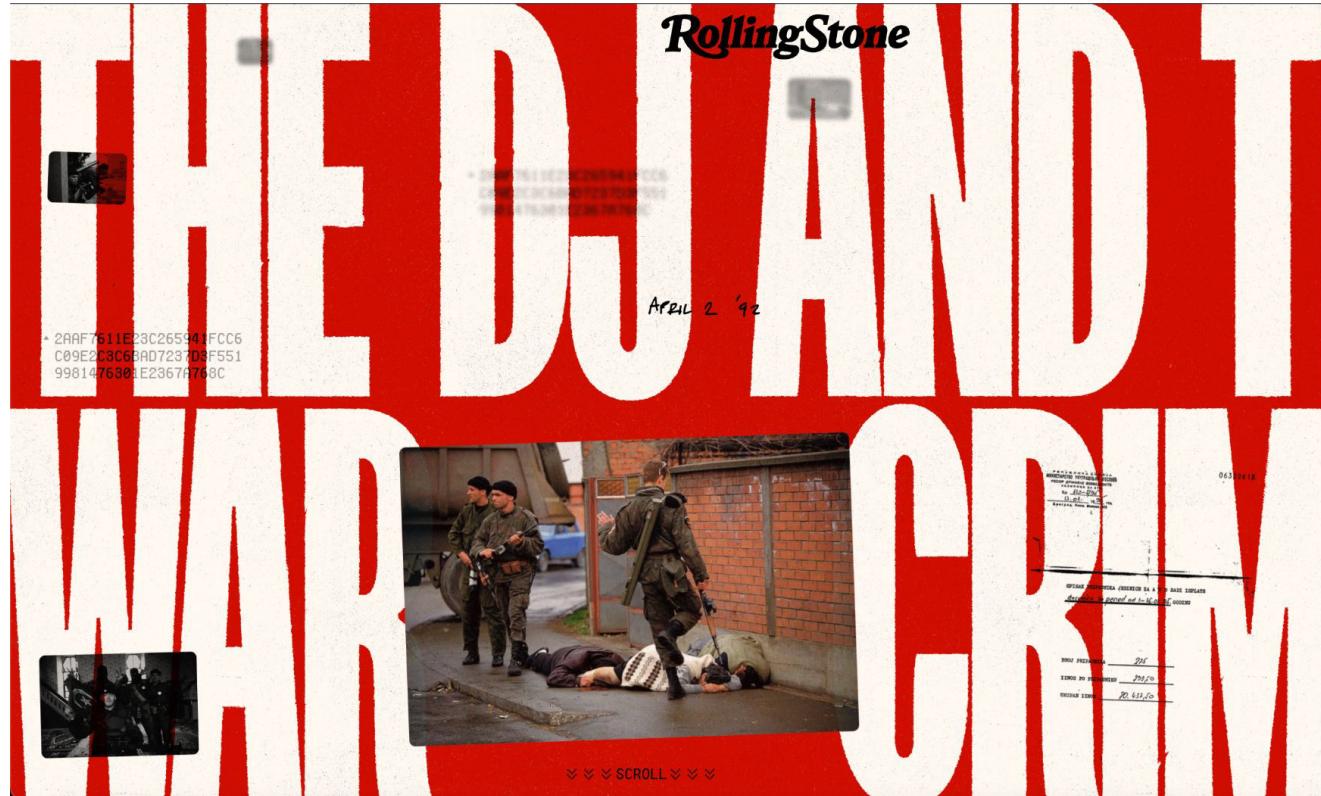


# Real World Example: Redaction Proofs



\*\*\*Joint work with Starling Lab (<https://rb.gy/vcc3wu>)

# Real World Example: Redaction Proofs



\*\*\*Joint work with Starling Lab (<https://rb.gy/vcc3wu>)

# Real World Example: Redaction Proofs

## Index of /ipfs/bafybeigvcncjm5yb4wbttnkpe56ludhaboovh75u5366in6uwfcq534ue 162 kB

 C049-0893_coords.txt	bafk...qila	42 B
 C049-0893_hash.txt	bafk...rymi	10 kB
 C049-0893_proof.txt	bafk...hove	548 B
 C049-0893_red.png	bafk...3tyy	150 kB
 README.txt	bafk...y7ha	217 B

\*\*\*Joint work with Starling Lab (<https://rb.gy/vcc3wu>)

# Conclusions

- Proof systems have greatly improved due to their need in blockchains  
⇒ non-blockchain applications benefit
- Proofs about large images ( $4000 \times 6000$ ) can be done in reasonable time
- Applicable to C2PA for image authenticity
  - If keys extracted, all bets are off ⇒ could rely on hardware enclaves
- Open problem: ZK proofs for videos?