

# Towards a White-Box Secure Fiat-Shamir Transformation

Gal Arnon

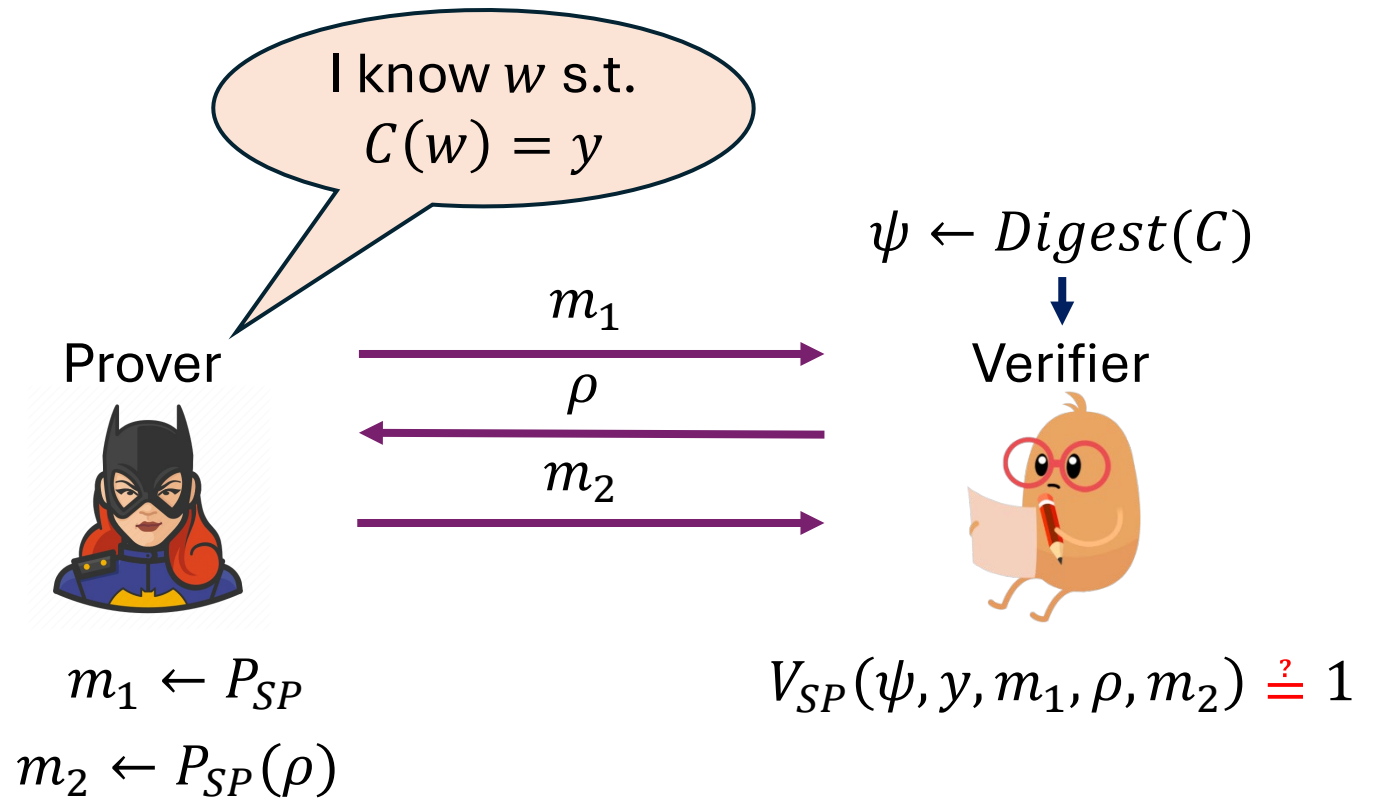


Eylon Yogev



# Sigma Protocol

- 3-message protocol
- Public-coin
- Pre-processing



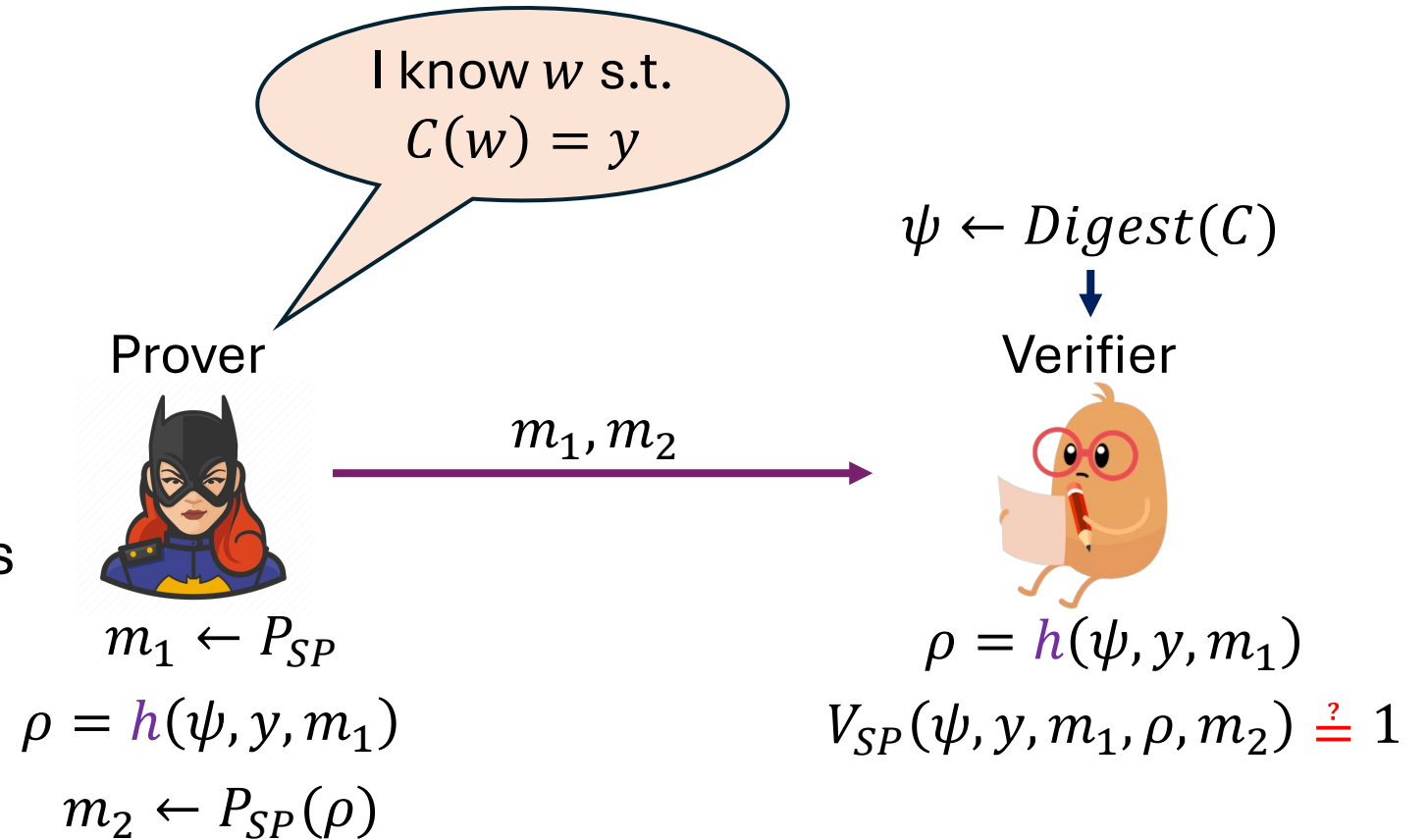
Soundness:

- **Statistical:** against **unbounded** provers (a.k.a. **proof**)
- **Computational:** against **bounded** provers (a.k.a. **argument**)

# Fiat-Shamir (FS)

**Fiat-Shamir transformation:**  
interactive  $\rightarrow$  non-interactive

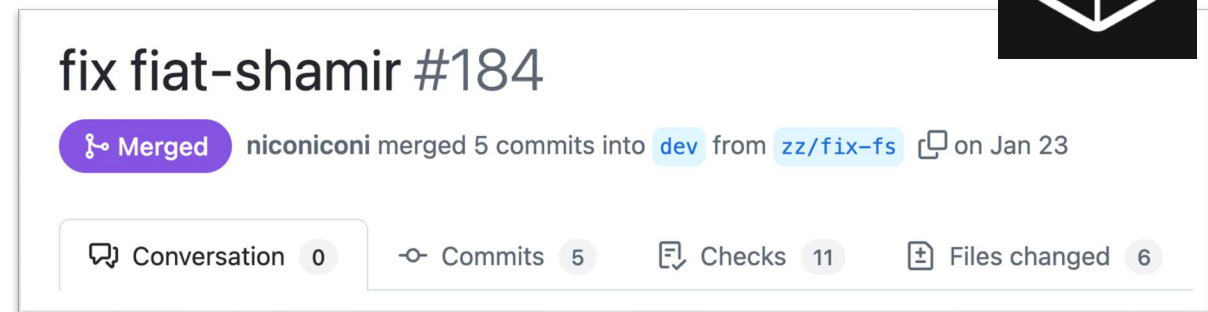
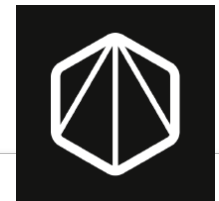
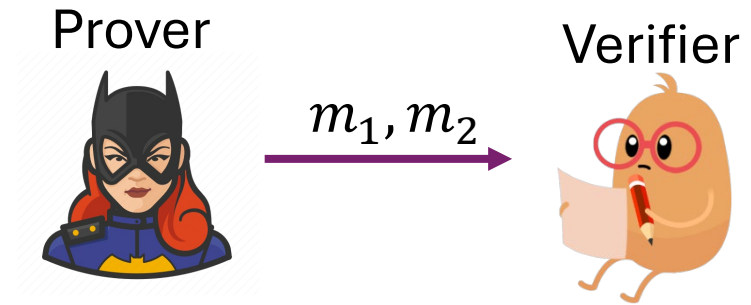
**Idea:** replace verifier randomness  
with hash function  $h$



FS is widely deployed real-world crypto systems, protecting billions of dollars:  
signature schemes, blockchains, ...

# Security of FS

- Pointcheval and Stern [PS96]: secure in the ROM
  - For both **proofs** and **arguments**
- **Proofs**: exist hash functions for which FS is secure  
[CCR16, KRR17, CCRR18, HL18, CCHLRRW19, PS19, BKM20, JJ21, HLR21, CJJ21, HJKS22, KLV23,...]
- **Arguments**: line of attacks using “white-box” techniques (“diagonalization”)
  - **Interactive** protocols that become **insecure with FS** for **any concrete hash**
- Examples of attacks:
  - [Bar01, GK03]: contrived identification schemes
  - [BBHMR19]: contrived CRH for Kilian’s protocol
  - [KRS25]: direct attack on **natural** variant of the [GKR15] protocol
- No attacks on Schnorr’s protocol



# Our Results

# Our Results

A new transformation (XFS) aimed to mitigate white-box attacks

- Focus on **practicality**: negligible overhead to prover and verifier
- Circumvents recent attack on GKR
- **Evidence for security**: **prove secure** in a relativized model where **FS is insecure**

Attacks we don't defend against:

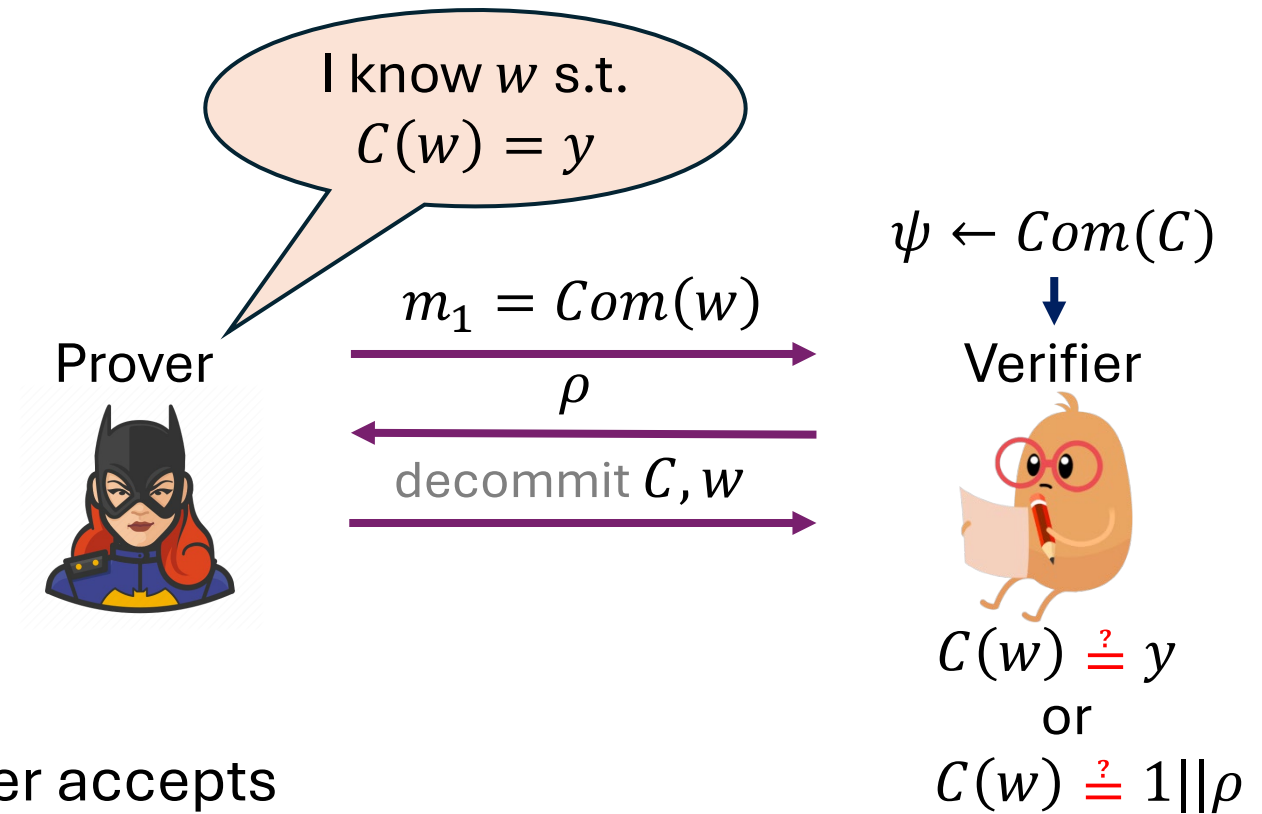
- [Bar01, GK03] for arbitrary poly-time verifiers
- [BBHMR19] CRH attack where to compute CRH need to verify a SNARK

But these attacks are **contrived**

# A Toy Protocol

# Toy Protocol

Let  $Com$  be a succinct commitment scheme



**Completeness:** if  $C(w) = y$  verifier accepts

**Soundness (computational):** follows from the commitment scheme

- the circuit  $C$  along with  $w$  cannot predict  $\rho$



# FS for Toy Protocol

Fiat-Shamir hash function  $h$

$C(w)$ :

1. Parse  $w = \psi || y$
2.  $m_1 \leftarrow Com(w)$
3. Output  $1 || h(\psi, y, m_1)$

Set  $y = 0^m$ ,  $\psi = Com(C)$ , and  $w = \psi || y$

Prover



I know  $w$  s.t.  
 $C(w) = y$

$m_1 = Com(w)$   
decommit  $C, w$

Output starts  
with 1

$\psi \leftarrow Com(C)$

Verifier



$\rho = h(\psi, y, m_1)$

$C(w) \stackrel{?}{=} y$   
or

$1 || \rho = 1 || h(\psi, y, m_1) = C(w) \stackrel{?}{=} 1 || \rho$

**Insecure!**

For any  $h$ :

a prover strategy such that for all  $w$ ,  $C(w) \neq y$  but **verifier accepts**

# Attack on FS

**Main problem:**  $C$  computes “verifier next message”  $\rho = h(\psi, y, m_1)$

**Naïve solution:** make  $h$  “more complex” than  $C$

## Drawbacks:

- **Slow** verifier (computes  $h$ , more complex than  $C$ )
- **Non compatible** with recursion
- Security **unclear**

We propose an alternative solution  
using strong **proof of work**

**Intuitively:** make next message function  
more complex than  $C$ , but easy to verify

# The XFS Transformation

But first, a strong proof of work

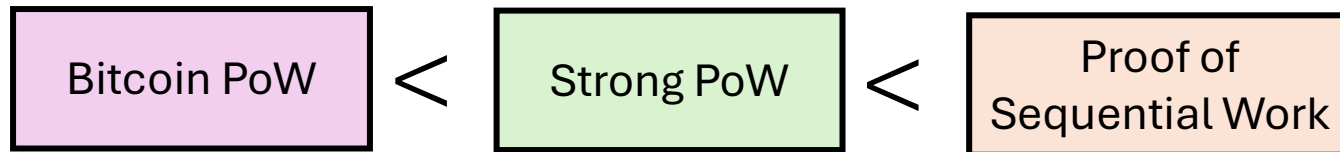
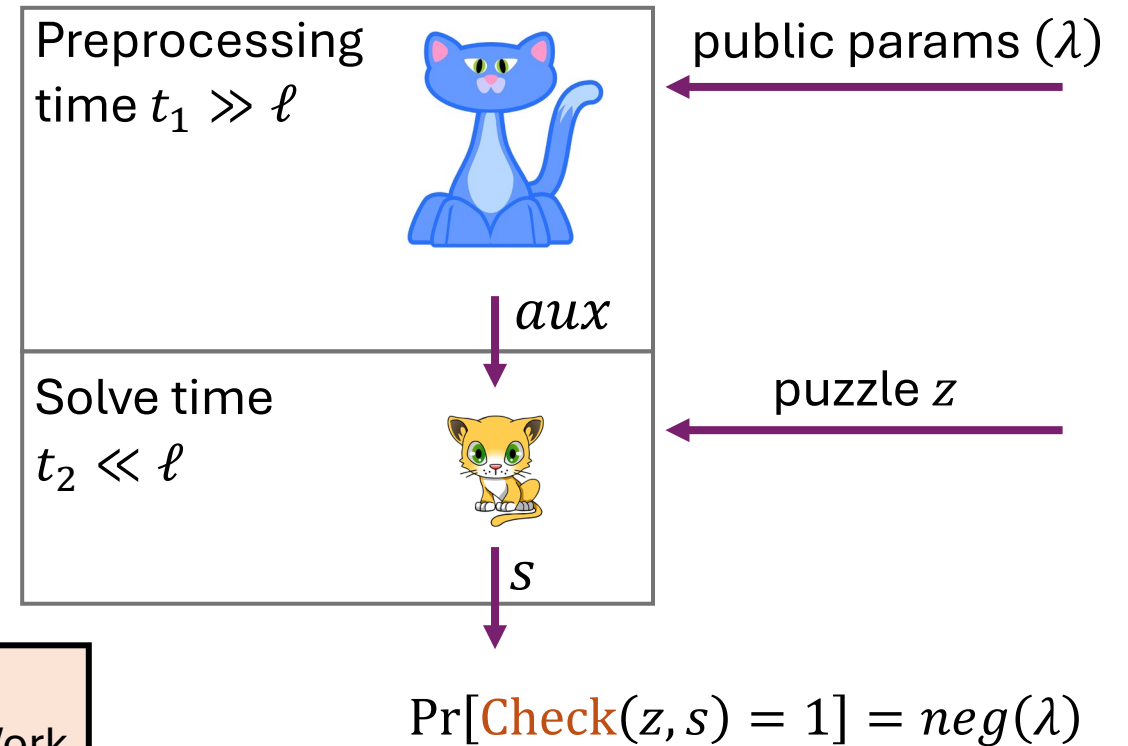
# Strong Proof of Work

Proof of work with hardness  $\ell$  ( $\ll 2^\lambda$ ):

- **Solve**( $z$ ) solves puzzle  $z$  in time  $\ell$
- **Check**( $z, s$ ) verifies a solution  $s$  to puzzle  $z$

**Security:**

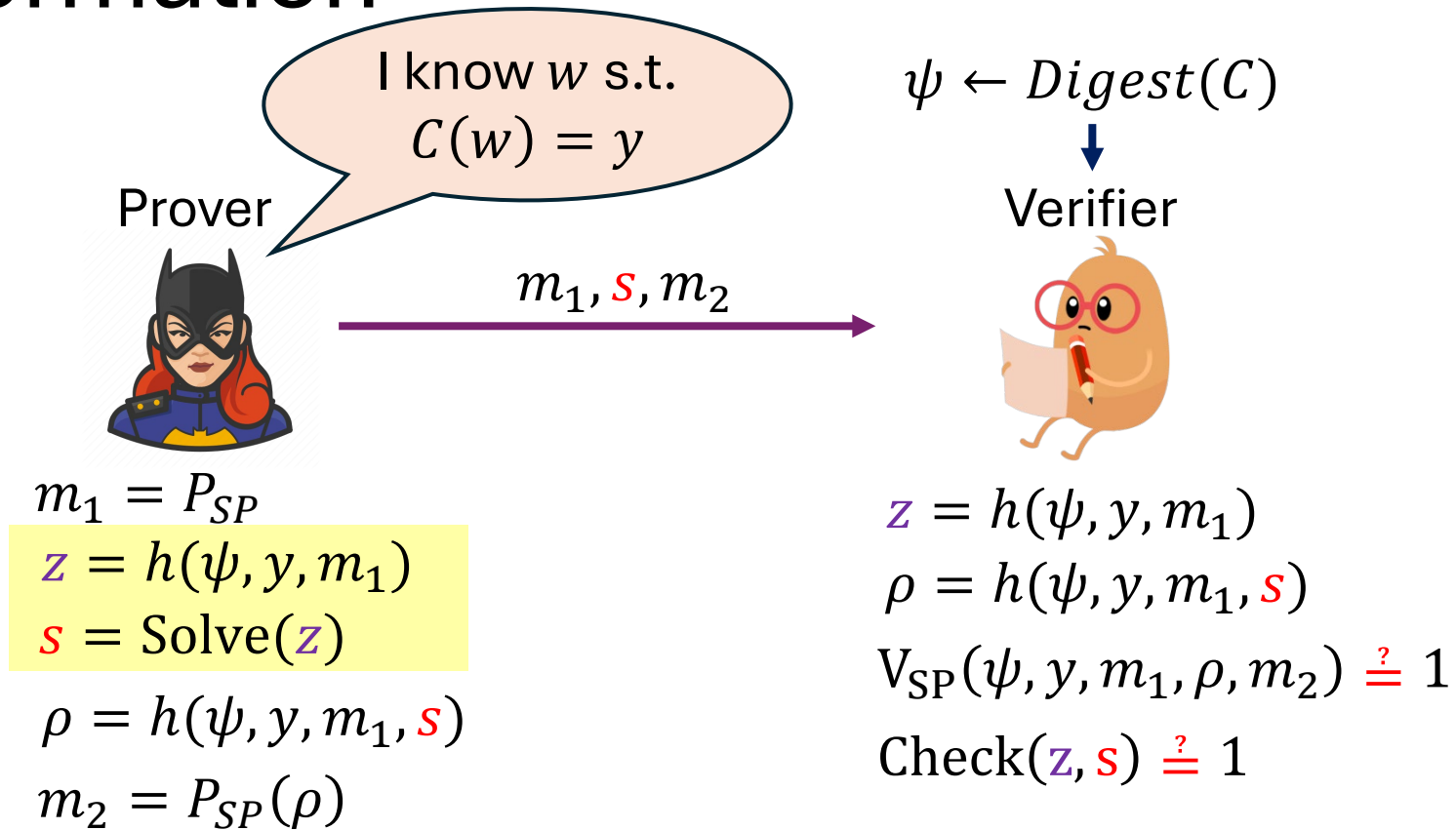
$\text{neg}(\lambda)$  **probability** with preprocessing



guessing a solution works  
with **probability**  $1/\ell$

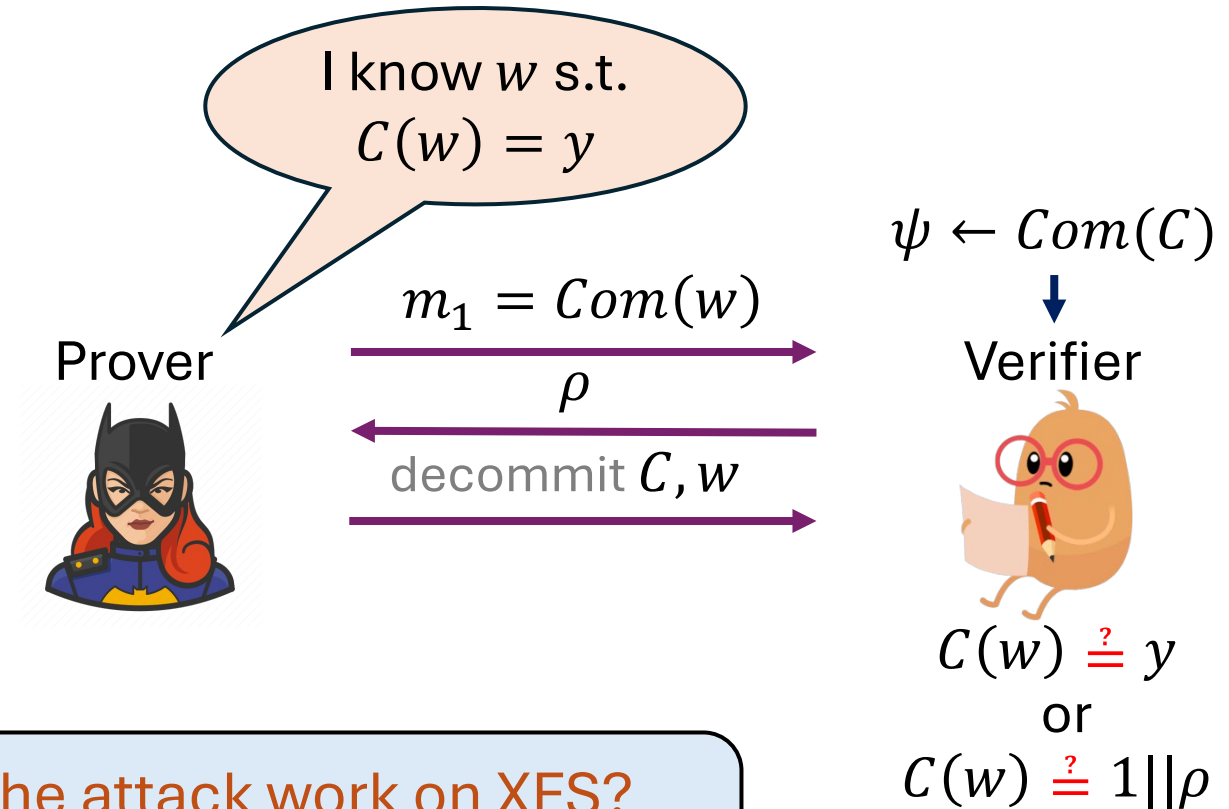
We give: two **constructs** of strong proof of work

# XFS Transformation



# XFS for the Toy Protocol

# XFS for Toy Protocol



Does the attack work on XFS?

Need to modify  $C$  to compute  $\rho$

# XFS for Toy Protocol

Fiat-Shamir hash function  $h$

**Attack:** circuit computes  $\rho$

$C(w)$ :

1. Parse  $w = \psi || y$
2.  $m_1 = Com(w)$
3.  $z = h(\psi, y, m_1)$
4.  $s = Solve(z)$
5. Output  $1 || h(\psi, y, m_1, s)$

Set  $y = 0^m$ ,  $\psi = Com(C)$ , and  $w = \psi || y$

**Observe:**  $z$  is computed after  $C, w$  are committed

$$z = h(\psi, y, m_1) = h(Com(C), y, Com(w))$$

Prover



I know  $w$  s.t.  
 $C(w) = y$

$m_1 = Com(w), s$   
decommit  $C, w$

Set PoW Hardness  $\ell > |C|$  so  
 $C$  can't compute  $Solve(z)$

Can  $w$  help solve the puzzle?

$\psi \leftarrow Com(C)$

Verifier



$z = h(\psi, y, m_1)$   
 $\rho = h(\psi, y, m_1, s)$   
Check( $z, s$ )  $\stackrel{?}{=} 1$

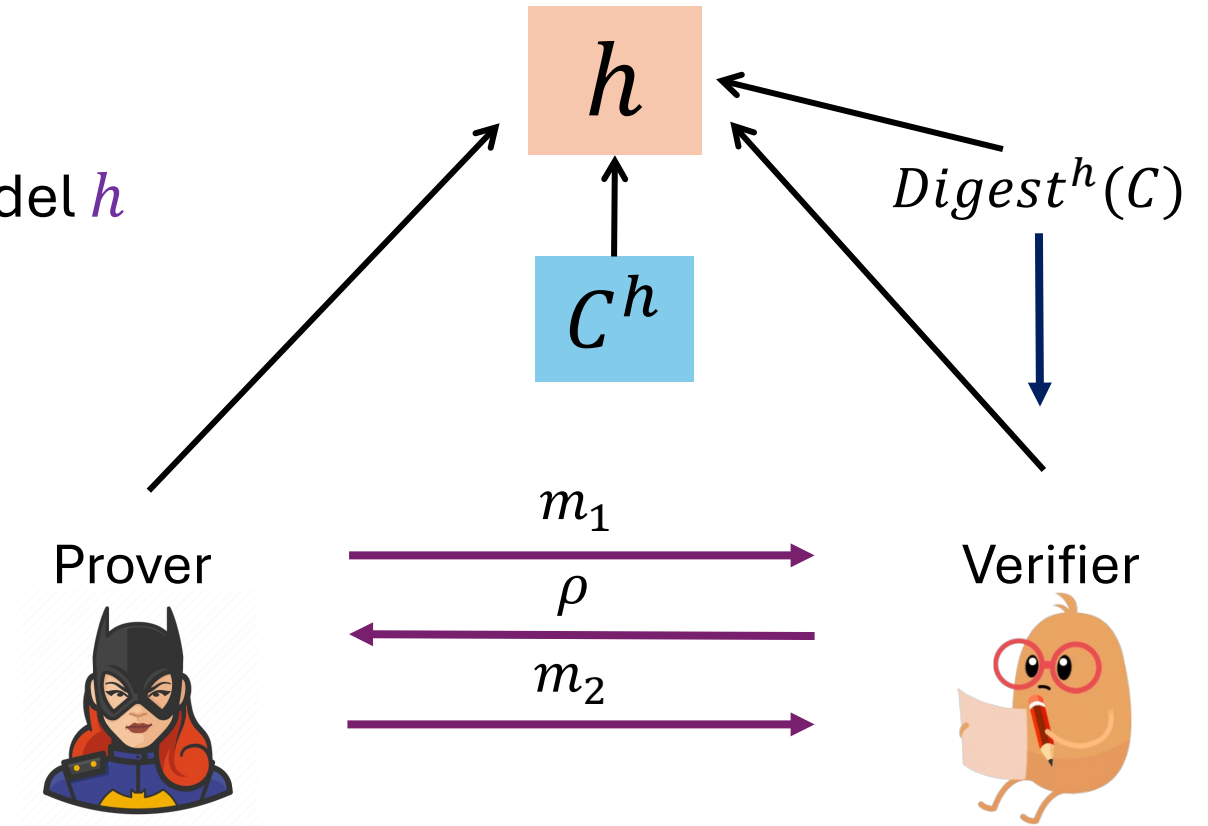
$C(w) \stackrel{?}{=} y$   
or  
 $C(w) \stackrel{?}{=} 1 || \rho$



# On the Security of XFS

# The Relativized World

- An **ideal** model with a random oracle model  $h$
- All parties have oracle access to  $h$
- This model **captures** the toy protocol



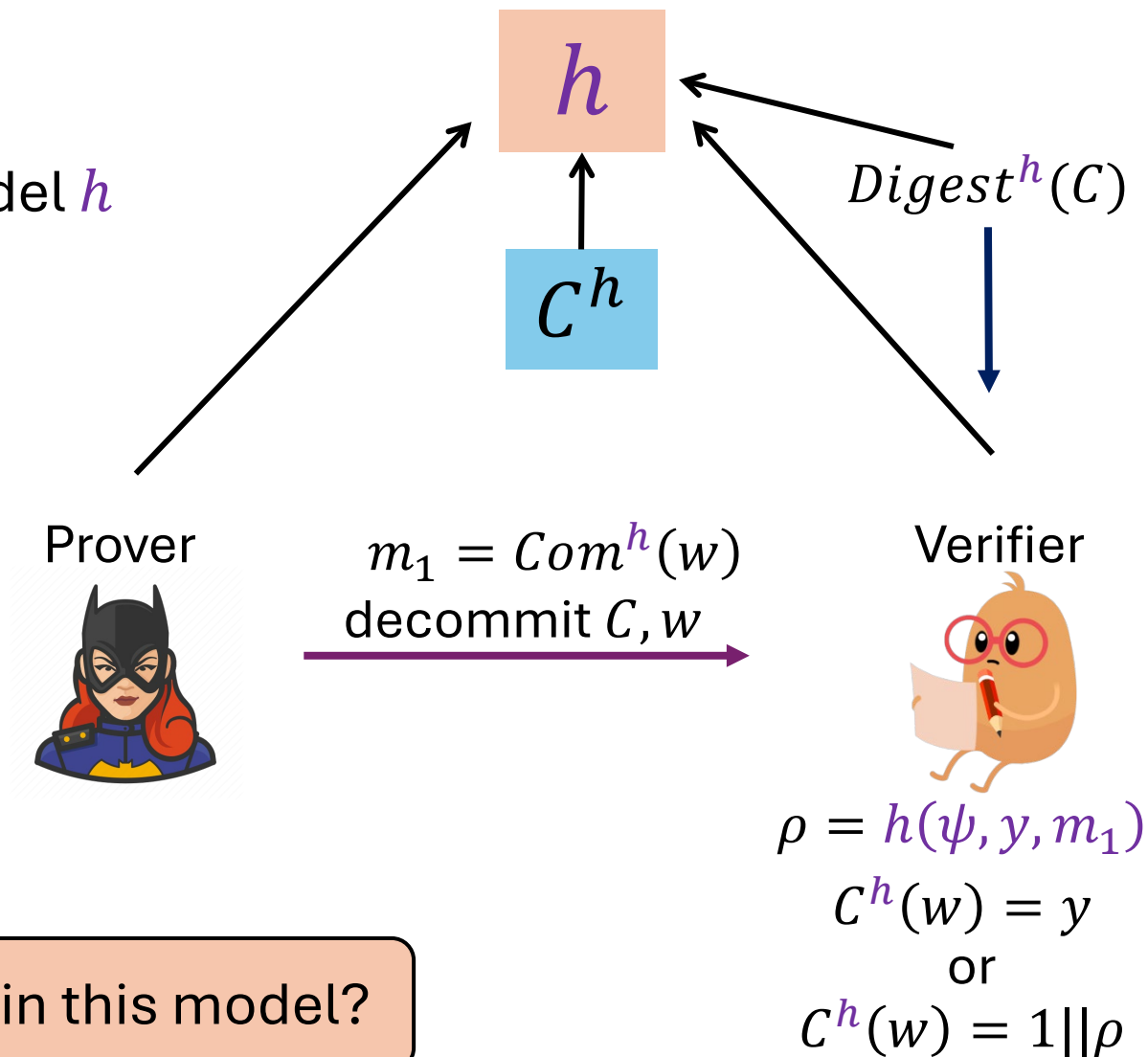
# The Relativized World

- An **ideal** model with a random oracle model  $h$
- All parties have oracle access to  $h$
- This model **captures** the toy protocol

$C^h(w)$ :

1. Parse  $w = \psi || y$
2.  $m_1 \leftarrow Com(w)$
3. Output  $1 || h(\psi, y, m_1)$

What about XFS in this model?



# Security in the Relativized World

## Theorem:

In the **relativized model**, the XFS\* transformation satisfies:

### Input:

1. sigma protocol with round-by-round knowledge error  $\kappa_{SP}$
2. strong PoW with error  $\epsilon_{PoW}$

**Output:** non-interactive protocol with knowledge error

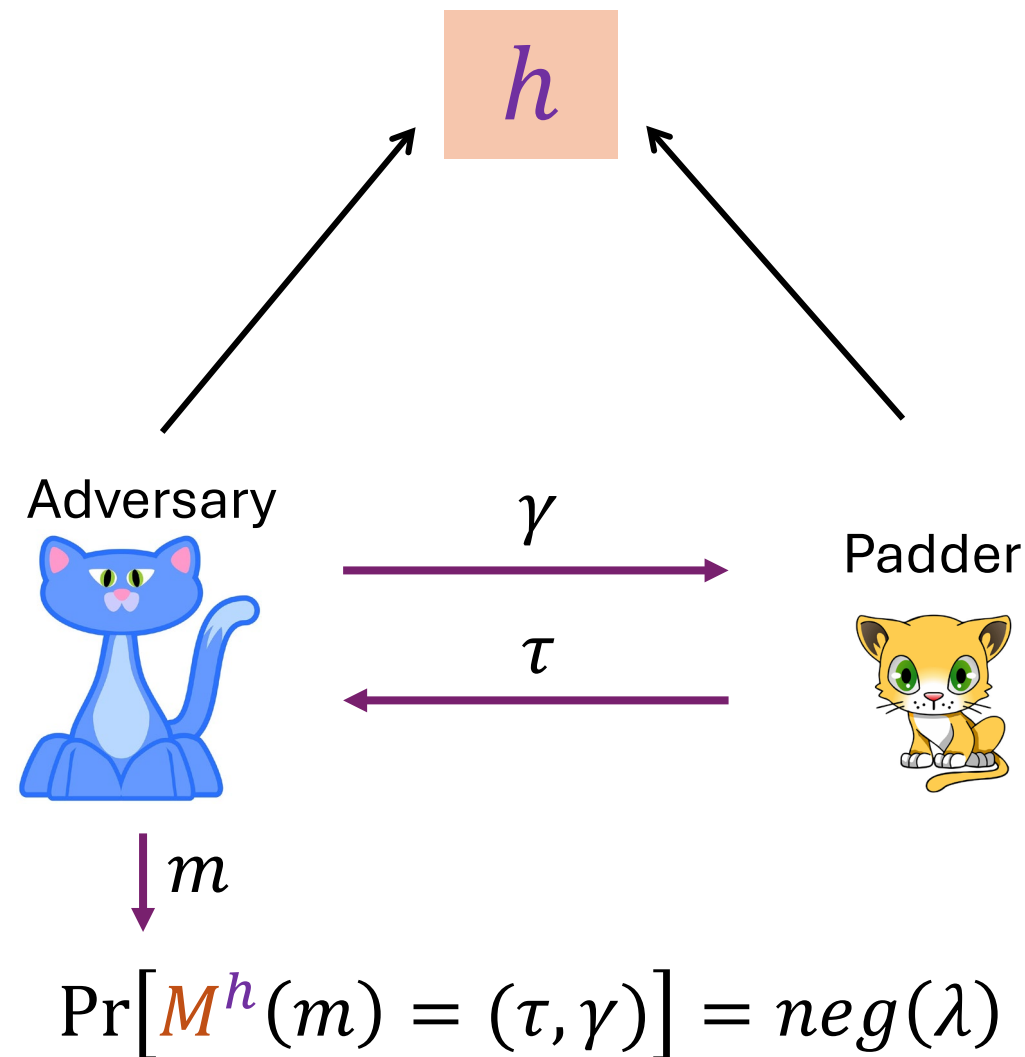
$$\kappa_{ARG} \leq O(t \cdot \kappa_{SP} + t \cdot \epsilon_{PoW})$$

where:

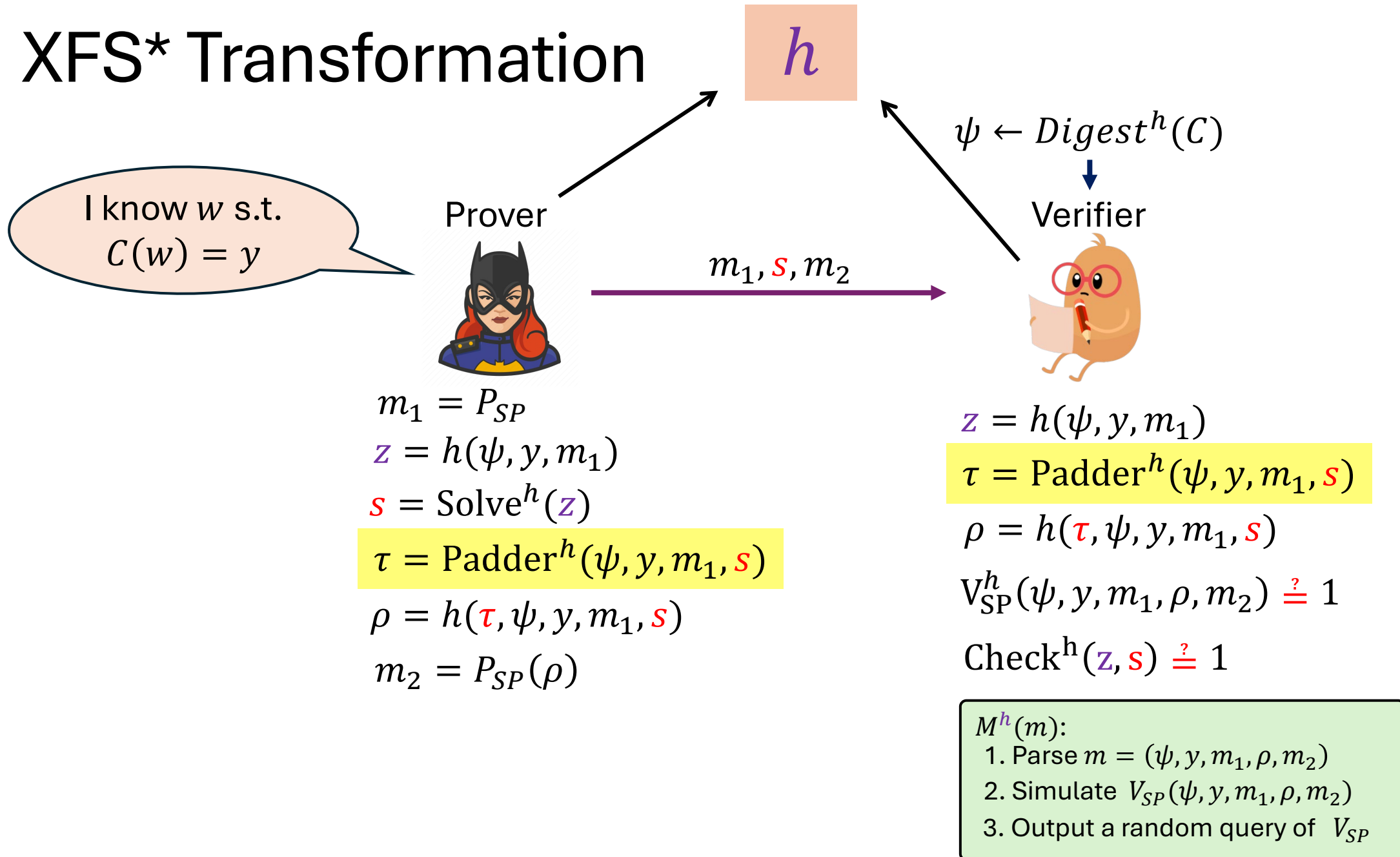
- $t$  is the query complexity of the malicious prover
- XFS\* is a **slight modification** of the transformation we saw

# Prefix Avoiding Padders

- Tool that facilitates our proof
  - Prevents the **verifier** from computing  $\rho$
- **Padder** outputs a prefix  $\tau$  (deterministically)
- For machine  $M^h$  the security game is:
- Example of padders:
  - $\text{Padder}(\gamma) = 0^{|M|}$
  - $\text{Padder}(\gamma) = h(M, \gamma)$  (extended to be long)
  - In practice could be trivial
- Used for  $M^h$  that simulates the verifier



# XFS\* Transformation



# Summary

We saw:

- FS **insecure** for arguments due to white-box attacks
- We propose **XFS** aimed to mitigate such attacks
- Uses **strong PoW** (ask me how to construct!)
- Heuristic proof of security in **relativized model**

Future work:

- Multi round version
- Security proofs in algebraic models
- Analyze with sponges

*Thank You!*

# Simple PoW construction

- Given a random puzzle  $z$ , and hardness  $\ell$ :
- Compute a Merkle tree of length  $\ell$ 
  - The  $i$ -th leaf is  $(z, i)$
- Hash the root to get small subset  $I \subseteq [\ell]$
- Open auth paths in  $I$
- Any algorithm that performs at most  $\ell/2$  hashes:
  - Can compute at most  $\ell/2$  leaves
  - Probability of opening all leaves in  $I$  is negligible