# Indistinguishability Obfuscation via Mathematical Proofs of Equivalence

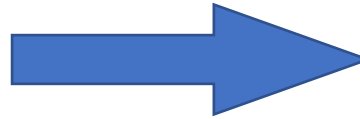**Abhishek Jain**

Johns Hopkins University

**Zhengzhong Jin**

Johns Hopkins University → MIT

# Indistinguishability Obfuscation (iO)

```
function main() {
  console.log('hello, world');
}
main()
```

Source code

iO

```
function _0x19e6(_0x4d301f,_0xcaab53){var _0x3a4e72=_0x3a4e();return
_0x19e6=function(_0x19e691,_0x5809f0){_0x19e691=_0x19e691-0x14e;var
_0x16ee0b=_0x3a4e72[_0x19e691];return
_0x16ee0b;},_0x19e6(_0x4d301f,_0xcaab53);}function _0x3a4e(){var _0x3f0a9d=
['log','199381NCGrSa','2328491tAiNSg','18mVqyqS','4cVQTsk','6PuGzwR','107410
32WsiTVO','104321yYIIVM','370911DTLqdw','10uRQffV','2024504eEkwnt','114dOcOh
j','hello,\x20world','2634710IatlOd'];_0x3a4e=function(){return
_0x3f0a9d;};return _0x3a4e();}(function(_0x3d9e47,_0x360e03){var
_0x3afd0b=_0x19e6,_0x2928d3=_0x3d9e47();while(!![]){try{var _0x33cc3a=-
parseInt(_0x3afd0b(0x15a))/0x1*(-parseInt(_0x3afd0b(0x158))/0x2)+-
parseInt(_0x3afd0b(0x15b))/0x3*(-parseInt(_0x3afd0b(0x157))/0x4)+-
parseInt(_0x3afd0b(0x152))/0x5+parseInt(_0x3afd0b(0x150))/0x6*
(parseInt(_0x3afd0b(0x154))/0x7)+-parseInt(_0x3afd0b(0x14f))/0x8*(-
parseInt(_0x3afd0b(0x156))/0x9)+parseInt(_0x3afd0b(0x14e))/0xa*
(parseInt(_0x3afd0b(0x155))/0xb)+-
parseInt(_0x3afd0b(0x159))/0xc;if(_0x33cc3a===_0x360e03)break;else
_0x2928d3['push'](_0x2928d3['shift']());}catch(_0x437e27){_0x2928d3['push']
(_0x2928d3['shift']());}}}(_0x3a4e,0x42c94));function main(){var
_0x29ace6=_0x19e6;console[_0x29ace6(0x153)](_0x29ace6(0x151));}main();
```
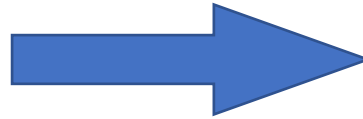
**"Unintelligible"**

# Indistinguishability Obfuscation (iO)

```javascript
1 function main() {
2   console.log('hello, world');
3 }
4 main()
```

Source code

iO →

```
function _0x19e6(_0x4d301f,_0xcaab53){var _0x3a4e72=_0x3a4e();return
_0x19e6=function(_0x19e691,_0x5809f0){_0x19e691=_0x19e691-0x14e;var
_0x16ee0b=_0x3a4e72[_0x19e691];return
_0x16ee0b;},_0x19e6(_0x4d301f,_0xcaab53);}function _0x3a4e(){var _0x3f0a9d=
['log','199381NCGrSa','2328491tAiNSg','18mVqyqS','4cVQTsk','6PuGzwR','107410
32WsiTVO','104321yYIIVM','370911DTLqdw','10uRQffV','2024504eEkwnt','114dOcOh
j','hello,\x20world','2634710IatlOd'];_0x3a4e=function(){return
_0x3f0a9d;};return _0x3a4e();}(function(_0x3d9e47,_0x360e03){var
_0x3afd0b=_0x19e6,_0x2928d3=_0x3d9e47();while(!![]){try{var _0x33cc3a=-
parseInt(_0x3afd0b(0x15a))/0x1*(-parseInt(_0x3afd0b(0x158))/0x2)+-
parseInt(_0x3afd0b(0x15b))/0x3*(-parseInt(_0x3afd0b(0x157))/0x4)+-
parseInt(_0x3afd0b(0x152))/0x5+parseInt(_0x3afd0b(0x150))/0x6*
(parseInt(_0x3afd0b(0x154))/0x7)+-parseInt(_0x3afd0b(0x14f))/0x8*(-
parseInt(_0x3afd0b(0x156))/0x9)+parseInt(_0x3afd0b(0x14e))/0xa*
(parseInt(_0x3afd0b(0x155))/0xb)+-
parseInt(_0x3afd0b(0x159))/0xc;if(_0x33cc3a===_0x360e03)break;else
_0x2928d3['push'](_0x2928d3['shift']());}catch(_0x437e27){_0x2928d3['push']
(_0x2928d3['shift']());}}}(_0x3a4e,0x42c94));function main(){var
_0x29ace6=_0x19e6;console[_0x29ace6(0x153)](_0x29ace6(0x151));}main();
```

"**Unintelligible**"

The obfuscated program *preserves the functionality* of the input program.
*(Produce the same output)*

# Indistinguishability Security (as a Game)

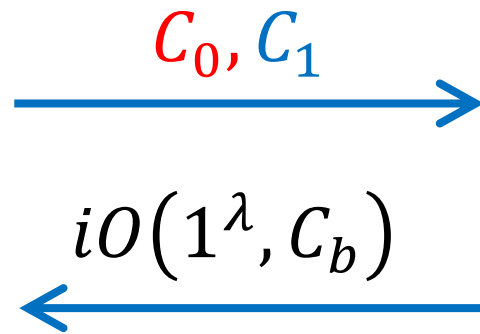# Indistinguishability Security (as a Game)



$C_0, C_1$

# Indistinguishability Security (as a Game)



$C_0, C_1$

$b \leftarrow \{0,1\}$

# Indistinguishability Security (as a Game)

$C_0, C_1$

$iO(1^\lambda, C_b)$

$b \leftarrow \{0,1\}$

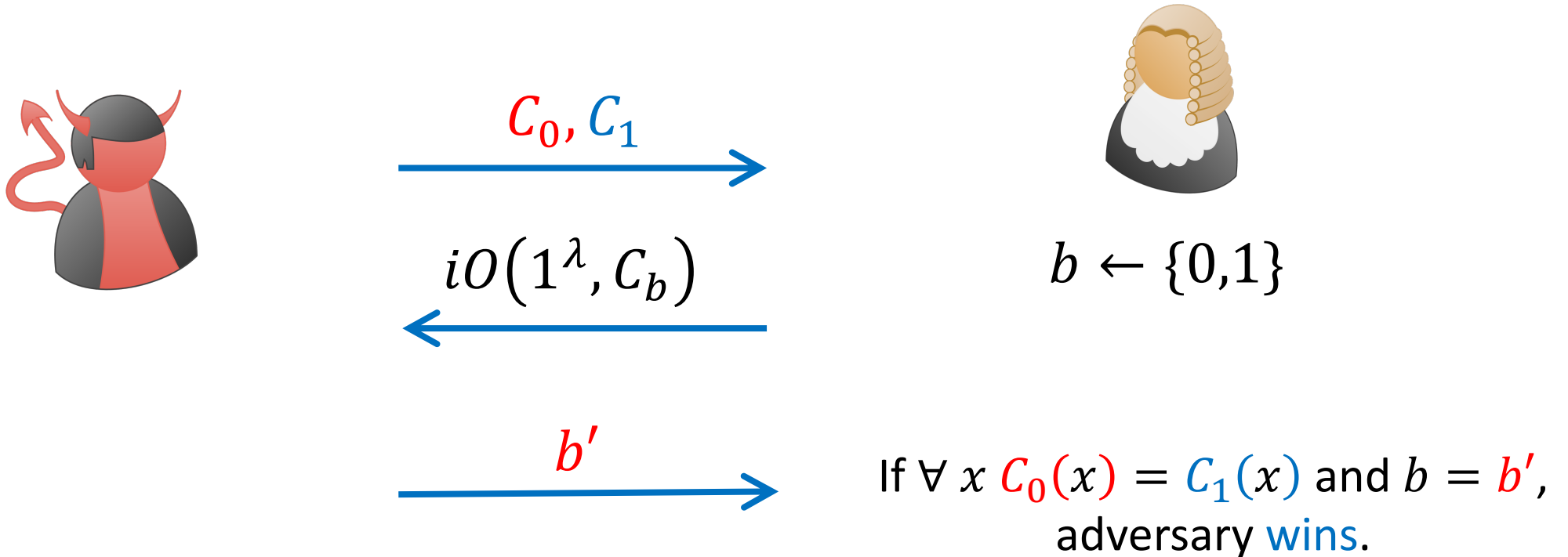# Indistinguishability Security (as a Game)

$$C_0, C_1$$
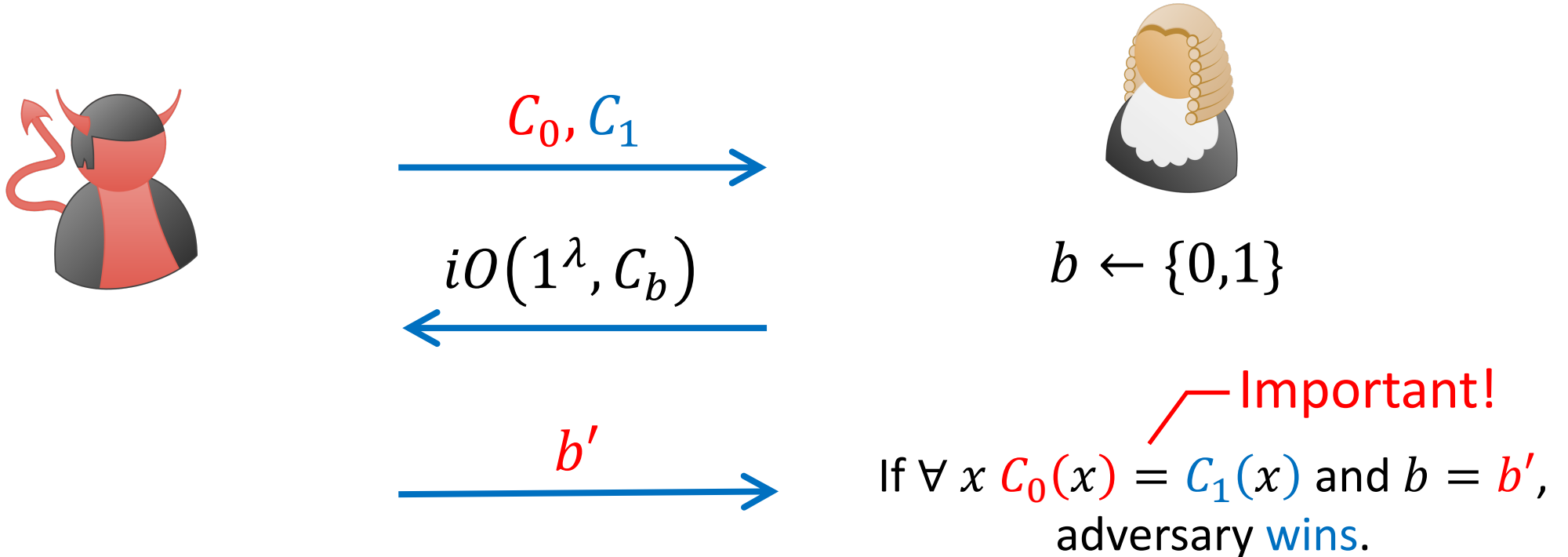
$$iO(1^\lambda, C_b)$$

$$b \leftarrow \{0,1\}$$

$$b'$$

# Indistinguishability Security (as a Game)



$C_0, C_1$

$iO(1^\lambda, C_b)$

$b \leftarrow \{0,1\}$

$b'$

If $\forall x\ C_0(x) = C_1(x)$ and $b = b'$, adversary wins.

# Indistinguishability Security (as a Game)



$C_0, C_1$

$iO(1^\lambda, C_b)$

$b \leftarrow \{0,1\}$

$b'$

If $\forall x \; C_0(x) = C_1(x)$ and $b = b'$, adversary wins.

$$\Pr[\text{\includegraphics{}} \; wins] \leq \frac{1}{2} + \text{negl}(\lambda)$$

# Indistinguishability Security (as a Game)



$C_0, C_1$

$iO(1^\lambda, C_b)$

$b \leftarrow \{0,1\}$

$b'$

Important!

If $\forall x$ $C_0(x) = C_1(x)$ and $b = b'$, adversary wins.

$\Pr[\text{wins}] \leq \dfrac{1}{2} + \mathrm{negl}(\lambda)$

# iO: Crypto "Complete" [Sahai-Waters'13,...]



Nash Equilibrium

Witness Encryption

Succinct Non-interactive
Zero-Knowledge Arg.

iO

...

Deniable
Encryption

Software
watermarking

# Can we build iO?

# Can we build iO?

A Long Line of Work:

[Garg-Gentry-Halevi-Raykova-Sahai-Waters'13][Pass-Seth-Telang'14]

[Gentry-Lewko-Sahai-Waters'15][Ananth-Jain'15][Bitansky-Vaikuntanathan'15]

[Lin'16][Lin-Vaikuntanathan'16][Lin-Pass-Karn Seth-Telang'16]

[Garg-Miles-Mukherjee-Sahai-Srinivasan-Zhandry'16][Ananth-Sahai'17][Lin'17]

[Lin-Tessaro'17][Agrawal'19][Jain-Lin-Matt-Sahai'19][Brakerski-Dottling-Malavolta'20]…

# Can we build iO?

A Long Line of Work:

[Garg-Gentry-Halevi-Raykova-Sahai-Waters'13][Pass-Seth-Telang'14]

[Gentry-Lewko-Sahai-Waters'15][Ananth-Jain'15][Bitansky-Vaikuntanathan'15]

[Lin'16][Lin-Vaikuntanathan'16][Lin-Pass-Karn Seth-Telang'16]

[Garg-Miles-Mukherjee-Sahai-Srinivasan-Zhandry'16][Ananth-Sahai'17][Lin'17]

[Lin-Tessaro'17][Agrawal'19][Jain-Lin-Matt-Sahai'19][Brakerski-Dottling-Malavolta'20]...

iO for *circuits* from well-founded assumptions

[Jain-Lin-Sahai'20]

**Question:** Can we build iO for **Turing machines**?

**Question:** Can we build iO for **Turing machines**?

**Why Turing machines?**

**Question:** Can we build iO for **Turing machines**?

**Why Turing machines?**
- Natural representation of programs

**Question:** Can we build iO for **Turing machines**?

**Why Turing machines?**
- Natural representation of programs

```
1  function main() {
2    console.log('hello, world');
3  }
4  main()
```

(Turing Machine)

**Question:** Can we build iO for **Turing machines**?

**Why Turing machines?**
- Natural representation of programs
- Support *any* input length

Circuit Model: input length is fixed

**Question:** Can we build iO for **Turing machines**?

**Why Turing machines?**

- Natural representation of programs
- Support *any* input length
- *Small* obfuscated program size



Obfuscated Turing Machine Size:
*Poly(input Turing Machine)*

# Prior Work: Only for *Bounded*-Input Length

[BGLPT'15][CHJV'15][KLW'15][GS'18]…

# Prior Work: Only for *Bounded*-Input Length

[BGLPT'15][CHJV'15][KLW'15][GS'18]…

Adversary for iO

# Prior Work: Only for *Bounded*-Input Length

[BGLPT'15][CHJV'15][KLW'15][GS'18]...



Adversary for iO

# Prior Work: Only for *Bounded*-Input Length

[BGLPT'15][CHJV'15][KLW'15][GS'18]…

# Prior Work: Only for *Bounded*-Input Length

[BGLPT'15][CHJV'15][KLW'15][GS'18]...



Adversary for iO

$2^{|input|}$-time
Reduction

Break

Assumptions

# Prior Work: Only for *Bounded*-Input Length

[BGLPT'15][CHJV'15][KLW'15][GS'18]...



Adversary for iO

$2^{|input|}$-time
Reduction

Break

Assumptions

**Assume $2^{\lambda^c}$-hardness of assumptions**
**& set $\lambda$ s.t. $2^{\lambda^c} > 2^{|input|}$**

# Prior Work: Only for *Bounded*-Input Length

[BGLPT'15][CHJV'15][KLW'15][GS'18]...



Adversary for iO

$2^{|input|}$-time
Reduction

Break

Assumptions

**Assume $2^{\lambda^c}$-hardness of assumptions
& set $\lambda$ s.t. $2^{\lambda^c} > 2^{|input|}$**

$\Rightarrow |input| < \lambda^c$

# Why $2^{|input|}$ Loss?

# Why $2^{|input|}$ Loss?

# Why $2^{|input|}$ Loss?



$C_0 \equiv C_1$

Real adversary

Reduction

Break

Assumptions

$C_0(x^*) \neq C_1(x^*)$

Fake adversary

Reduction

**Not** Break

Assumptions

# Why $2^{|input|}$ Loss?



Reduction

$C_0 \equiv C_1$

Real adversary

Break

Assumptions

Reduction needs to 'decide' the functionality equivalence.

Reduction

$C_0(x^*) \neq C_1(x^*)$

Fake adversary

**Not** Break

Assumptions

# Non-Falsifiability

# Non-Falsifiability



## iO Security

$C_0, C_1$

$iO(C_b)$

$b \leftarrow \{0,1\}$

$b'$

Check $\forall x \ C_0(x) = C_1(x)$

# Non-Falsifiability



iO Security

$C_0, C_1$

$iO(C_b)$

$b \leftarrow \{0,1\}$

$b'$

Check $\forall x \ C_0(x) = C_1(x)$
(inefficient checking)

# Non-Falsifiability



iO Security

$C_0, C_1$

$iO(C_b)$

$b'$

$b \leftarrow \{0,1\}$

**Non-Falsifiable**

Check $\forall x \ C_0(x) = C_1(x)$
(inefficient checking)

# Non-Falsifiability



iO Security

$C_0, C_1$

$iO(C_b)$

$b \leftarrow \{0,1\}$

$b'$

**Non-Falsifiable**

Check $\forall x \ C_0(x) = C_1(x)$
(inefficient checking)

Broader Perspective

**Non-Falsifiable** definitions appear in many other places,
e.g. proof systems. [Gentry-Wichs'10]

*This Talk:* How to overcome the non-falsifiability barrier?

*This Talk:* How to overcome the non-falsifiability barrier?

Prior Work

[Garg-Pandey-Srinivasan'16, Garg-Srinivasan'16, Garg-Pandey-Srinivasan-Zhandry'17][Liu-Zhandry'17]:
Require that "$\forall x \; C_0(x) = C_1(x)$" can be decided in **P**

# Observation: We *Prove* Equivalence in Math



Applications

↑

**iO**

# Observation: We *Prove* Equivalence in Math



Applications

iO

Security Proof of the Application

# Observation: We *Prove* Equivalence in Math



Applications

iO

Security Proof of the Application

# Observation: We *Prove* Equivalence in Math

# Observation: We *Prove* Equivalence in Math

Applications

**iO**

Security Proof of the Application

Build $C_0$ , $C_1$

Math proof of $\forall x \; C_0(x) = C_1(x)$

iO Security $\Rightarrow iO(C_0) \approx iO(C_1)$

# Observation: We *Prove* Equivalence in Math

# Observation: We *Prove* Equivalence in Math

# Our Approach

Short mathematical proof of "$\forall x \; C_0(x) = C_1(x)$"

**iO**

# Our Approach

**Short** mathematical proof of "$\forall x \; C_0(x) = C_1(x)$"

**iO**

# Our Approach

Short mathematical proof of "$\forall x \; C_0(x) = C_1(x)$"

**iO**

Assumptions

**Efficient** Reduction Algorithm
(this work)

## Our Result

$iO$ for any Turing machines $M_1, M_2$ with "$\forall x \; M_1(x) = M_2(x)$"

*provable in Cook's Theory PV,* based on well-founded assumptions.

# Cook's Theory $PV$ <span>[Cook'75]</span>

- Polynomial time reasoning

Polynomial-time Induction rule:
"If $\Phi(0)$ is true, and $\forall n, \Phi(n) \to \Phi(2n) \wedge \Phi(2n+1)$, then $\forall n\ \Phi(n)$."

Can define *any* **polynomial-time functions**, e.g.:
- Arithmetic: $+, -, \times, \div, \leq, <, \lfloor \cdot \rfloor, mod, \ldots$
- Logic Symbols: $\to, \neg, \wedge, \ldots$

# What Theorems Can $PV$ Prove?

# What Theorems Can $PV$ Prove?

[Prior work]

# What Theorems Can $PV$ Prove?

Prior work

- Correctness of "natural" algorithms in P

# What Theorems Can $PV$ Prove?

## Prior work

- Correctness of "natural" algorithms in P

- Basic Linear Algebra

# What Theorems Can $PV$ Prove?

**Prior work**

- Correctness of "natural" algorithms in P

- Basic Linear Algebra

- Combinatorial Theorems

# What Theorems Can $PV$ Prove?

Prior work

- Correctness of "natural" algorithms in P

- Basic Linear Algebra

- Combinatorial Theorems

  …

# What Theorems Can $PV$ Prove?

**Prior work**

- Correctness of "natural" algorithms in P

- Basic Linear Algebra

- Combinatorial Theorems

  ...

**This work**

Many crypto algorithms are "natural":

ElGamal Encryption

Regev's Encryption

Puncturable PRFs

...

# What Theorems Can *PV* Prove?

**Prior work**

- Correctness of "natural" algorithms in P

- Basic Linear Algebra

- Combinatorial Theorems

  ...

**This work**

Many crypto algorithms are "natural":

  ElGamal Encryption

  Regev's Encryption

  Puncturable PRFs

  ...

**Unprovable** Theorems (assume Factoring is hard)

- Fermat's Little Theorem

- Correctness for "Primes is in P"

# How to leverage mathematical proofs?

# How to leverage mathematical proofs?

*Overview of Techniques*

# What Information does a Proof Provide?

# What Information does a Proof Provide?

Mathematical Proofs Have Structures

| | | |
|---|---|---|
| 1. | $P \wedge Q$ | Premise |
| 2. | $P$ | Decomposing a conjunction (1) |
| 3. | $Q$ | Decomposing a conjunction (1) |
| 4. | $P \rightarrow \neg(Q \wedge R)$ | Premise |
| 5. | $\neg(Q \wedge R)$ | Modus ponens (3,4) |
| 6. | $\neg Q \vee \neg R$ | DeMorgan (5) |
| 7. | $\neg R$ | Disjunctive syllogism (3,6) |
| 8. | $S \rightarrow R$ | Premise |
| 9. | $\neg S$ | Modus tollens (7,8)    □ |

# What Information does a Proof Provide?

Mathematical Proofs Have Structures

| | | |
|---|---|---|
| 1. | $P \wedge Q$ | Premise |
| 2. | $P$ | Decomposing a conjunction (1) |
| 3. | $Q$ | Decomposing a conjunction (1) |
| 4. | $P \rightarrow \neg(Q \wedge R)$ | Premise |
| 5. | $\neg(Q \wedge R)$ | Modus ponens (3,4) |
| 6. | $\neg Q \vee \neg R$ | DeMorgan (5) |
| 7. | $\neg R$ | Disjunctive syllogism (3,6) |
| 8. | $S \rightarrow R$ | Premise |
| 9. | $\neg S$ | Modus tollens (7,8) |

- **Localness:** Each line is derived from $O(1)$ previous lines

# What Information does a Proof Provide?

Mathematical Proofs Have Structures

1. $P \wedge Q$ — Premise
2. $P$ — Decomposing a conjunction (1)
3. $Q$ — Decomposing a conjunction (1)
4. $P \rightarrow \neg(Q \wedge R)$ — Premise
5. $\neg(Q \wedge R)$ — Modus ponens (3,4)
6. $\neg Q \vee \neg R$ — DeMorgan (5)
7. $\neg R$ — Disjunctive syllogism (3,6)
8. $S \rightarrow R$ — Premise
9. $\neg S$ — Modus tollens (7,8)

- **Localness:** Each line is derived from $O(1)$ previous lines

- In Propositional Logic (Extended Frege): each line is also a *circuit*

# What Information does a Proof Provide?

Mathematical Proofs Have Structures

1. $P \wedge Q$      Premise
2. $P$      Decomposing a conjunction (1)
3. $Q$      Decomposing a conjunction (1)
4. $P \rightarrow \neg(Q \wedge R)$      Premise
5. $\neg(Q \wedge R)$      Modus ponens (3,4)
6. $\neg Q \vee \neg R$      DeMorgan (5)
7. $\neg R$      Disjunctive syllogism (3,6)
8. $S \rightarrow R$      Premise
9. $\neg S$      Modus tollens (7,8)

- **Localness:** Each line is derived from O(1) previous lines

- In Propositional Logic (Extended Frege): each line is also a *circuit*

Rest of the Talk: mainly focus on extended Frege ($\mathcal{EF}$), since PV-proof can be translated to $\mathcal{EF}$-Proof.

# Bypass $2^{|\text{input}|}$-Loss via $\mathcal{EF}$-Proofs

# Bypass $2^{|\text{input}|}$-Loss via $\mathcal{EF}$-Proofs

Hybrid Argument

$iO(C_0)$

$iO(C_1)$

# Bypass $2^{|\text{input}|}$-Loss via $\mathcal{EF}$-Proofs

Hybrid Argument

$$iO(C_0) \approx iO(\qquad) \approx iO(\qquad) \,\ldots\, iO(C_1)$$

# Bypass $2^{|\text{input}|}$-Loss via $\mathcal{EF}$-Proofs



Hybrid Argument

$$iO(C_0) \approx iO(\quad) \approx iO(\quad) \;\ldots\; iO(C_1)$$

# Bypass $2^{|\text{input}|}$-Loss via $\mathcal{EF}$-Proofs



**Hybrid Argument**

$$iO(C_0) \approx iO(\qquad) \approx iO(\qquad) \quad \ldots \quad iO(C_1)$$

*"**Locally** equivalent", checkable in polynomial time*

# Bypass $2^{|\text{input}|}$-Loss via $\mathcal{EF}$-Proofs

We build iO for locally equivalent circuits with loss independent of |input|.

### Hybrid Argument

$$iO(C_0) \approx iO(\quad) \approx iO(\quad) \ \dots \ iO(C_1)$$

*"**Locally** equivalent", checkable in polynomial time*

# Bypass $2^{|\text{input}|}$-Loss via $\mathcal{EF}$-Proofs

We build iO for locally equivalent circuits
with loss independent of |input|.

### Hybrid Argument

$$iO(C_0) \approx iO(\quad\quad) \approx iO(\quad\quad) \ \dots \ iO(C_1)$$

*"**Locally** equivalent"*, checkable *in polynomial time*

Poly. size $\mathcal{EF}$-proof for $C_0(x) \equiv C_1(x)$

# Technical Details

- $\mathcal{EF}$-Proofs $\Rightarrow$ local equivalence
- iO for locally equivalent ckts
- iO for Turing machines

# Technical Details

- $\mathcal{EF}$-**Proofs** $\Rightarrow$ **local equivalence**
- iO for locally equivalent ckts
- iO for Turing machines

# Define Local Equivalence

# Define Local Equivalence

$C:$  $:C'$

# Define Local Equivalence

$C$:     $: C'$

$C$ and $C'$ are almost the same (with same topology), except for a **functionality equivalent _sub-circuit_** of size $O(\log n)$

 $\equiv$ 

# Define Local Equivalence

$C$:

$: C'$

$C$ and $C'$ are almost the same (with same topology), except for

a **functionality equivalent _sub-circuit_** of size $O(\log n)$

$\equiv$

(Sub-circuit: induced subgraph from a subset of gates)

# Alternative View: A Series of Local Changes

$\mathcal{EF}$ proof for $C_0(x) \equiv C_1(x)$



$C_0$ ... $C_1$

Locally Equivalent

# Alternative View: A Series of Local Changes

$\mathcal{EF}$ proof for $C_0(x) \equiv C_1(x)$



$C_0$ ... $C_1$

Local Change

# Alternative View: A Series of Local Changes

$\mathcal{EF}$ proof for $C_0(x) \equiv C_1(x)$



Local Change

<u>Simplification in This Talk:</u> Ignore topology & allow multi-arity gates

# Stage I: Grow $C_1$

Add Gates in $C_1$ *one-by-one*

$C_0$

$x$

$C_0$

$x$

$C_1$

$x$

# Stage I: Grow $C_1$



Add Gates in $C_1$ *one-by-one*

Set output as $C_0(x)$

$C_0$

$x$

$C_0$

$x$

$C_1$

$x$

# Stage I: Grow $C_1$



**Local Equivalence**
When a gate is added, its output is not used anywhere

# Stage II: Grow the Proof

$\mathcal{EF}$-Proof of $C_0(x) \leftrightarrow C_1(x)$: $\theta_1, \theta_2, \dots, \theta_\ell$

# Stage II: Grow the Proof

$\mathcal{EF}$-Proof of $C_0(x) \leftrightarrow C_1(x)$: $\theta_1, \theta_2, \ldots, \theta_\ell$

# Stage II: Grow the Proof

$\mathcal{EF}$-Proof of $C_0(x) \leftrightarrow C_1(x)$: $\theta_1, \theta_2, \dots, \theta_\ell$

# Stage II: Grow the Proof

$\mathcal{EF}$-Proof of $C_0(x) \leftrightarrow C_1(x)$: $\theta_1, \theta_2, \ldots, \theta_\ell$

# Stage II: Grow the Proof

$\mathcal{EF}$-Proof of $C_0(x) \leftrightarrow C_1(x)$: $\theta_1, \theta_2, \dots, \theta_\ell$



**Intuition**: $\theta_i$'s (i.e. lines of the proof) are "true", so the functionality is preserved.

# Stage II: Local Equivalence

$i$-th Step: Add $\theta_i$

**Before:** $C_0(x) \wedge \theta_1 \wedge \cdots \wedge \theta_{i-1}$

**After:** $C_0(x) \wedge \theta_1 \wedge \cdots \wedge \theta_{i-1} \wedge \theta_i$

# Stage II: Local Equivalence

$i$-th Step: Add $\theta_i$

**Before:** $\quad C_0(x) \land \theta_1 \land \cdots \land \theta_{i-1}$

**After:** $\quad C_0(x) \land \theta_1 \land \cdots \land \theta_{i-1} \land \theta_i$

$\theta_i$ is derived via Modus Ponens: $\quad p, p \to q \vdash q$

# Stage II: Local Equivalence

$i$-th Step: Add $\theta_i$

**Before:**

**After:**

$\underline{\theta_i \text{ is derived via Modus Ponens:}}$   $p, p \rightarrow q \vdash q$

# Stage II: Local Equivalence

$i$-th Step: Add $\theta_i$

**Before:** $C_0(x) \wedge p \wedge \cdots \wedge (p \rightarrow q) \wedge \cdots$

**After:** $C_0(x) \wedge p \wedge \cdots \wedge (p \rightarrow q) \wedge \cdots \wedge q$

$\underline{\theta_i}$ is derived via Modus Ponens: $\quad p, p \rightarrow q \vdash q$

# Stage II: Local Equivalence

$i$-th Step: Add $\theta_i$

**Before:** $C_0(x) \wedge p \wedge \cdots \wedge (p \to q) \wedge \cdots$

**After:** $C_0(x) \wedge p \wedge \cdots \wedge (p \to q) \wedge \cdots \wedge q$

$\theta_i$ is derived via Modus Ponens: $\quad p, p \to q \vdash q$

$$p \wedge (p \to q) \equiv p \wedge (p \to q) \wedge q$$

# Stage III: **Switch** $o_0$ to $o_1$

$$o_0 \wedge \theta_1 \wedge \cdots \wedge \theta_\ell$$

$o_0$

$o_1$

$C_0$

$C_1$

# Stage III: **Switch** $o_0$ to $o_1$



$o_0 \wedge \theta_1 \wedge \cdots \wedge \theta_\ell$

$o_0$

$C_0$

$o_1$

$C_1$

$o_1 \wedge \theta_1 \wedge \cdots \wedge \theta_\ell$

$o_0$

$C_0$

$o_1$

$C_1$

# Stage III: Switch $o_0$ to $o_1$



$o_0 \wedge \theta_1 \wedge \cdots \wedge \theta_\ell$

$o_0$

$C_0$

$o_1$

$C_1$

$o_0$

$C_0$

$o_1 \wedge \theta_1 \wedge \cdots \wedge \theta_\ell$

$o_1$

$C_1$

**Local Equivalence**

$\theta_\ell$ is "$o_0 \leftrightarrow o_1$" (A proof of $C_0(x) \leftrightarrow C_1(x)$ must end with $o_0 \leftrightarrow o_1$ )

# Stage III: Switch $o_0$ to $o_1$



**Local Equivalence**

$\theta_\ell$ is "$o_0 \leftrightarrow o_1$" (A proof of $C_0(x) \leftrightarrow C_1(x)$ must end with $o_0 \leftrightarrow o_1$)

$$o_0 \wedge (o_0 \leftrightarrow o_1) \equiv o_1 \wedge (o_0 \leftrightarrow o_1)$$

# Stage IV: Shrink the Proof

$$o_1 \wedge \theta_1 \wedge \cdots \wedge \theta_\ell$$

# Stage IV: Shrink the Proof



$o_1 \wedge \theta_1 \wedge \cdots \wedge \theta_\ell$

$o_0$

$o_1$

$C_0$

$C_1$

Delete $\theta_i$
*one-by-one*

# Stage IV: Shrink the Proof

$o_1 \wedge \theta_1 \wedge \cdots \wedge \theta_\ell$

$o_0$

$o_1$

$C_0$

$C_1$

Delete $\theta_i$
*one-by-one*

$o_0$

$o_1$

$o_1$

$C_0$

$C_1$

# Stage IV: Shrink the Proof



$o_1 \wedge \theta_1 \wedge \cdots \wedge \theta_\ell$

$o_0$

$o_1$

$C_0$

$C_1$

Delete $\theta_i$
*one-by-one*

$o_0$

$o_1$

$o_1$

$C_0$

$C_1$

**Local Equivalence:** Similar to "Grow the proof" Stage

# Stage V: Shrink $C_0$

# Stage V: Shrink $C_0$

# Stage V: Shrink $C_0$

# Stage V: Shrink $C_0$

$o_1$

$o_0$

$o_1$

$C_0$

$C_1$

Delete $C_0$
*gate-by-gate*

$o_1$

$o_1$

$C_1$

# Stage V: Shrink $C_0$



**Local Equivalence:** Similar to "Grow $C_1$" Stage

# Technical Details

- $\mathcal{EF}$-Proofs $\Rightarrow$ local equivalence
- **iO for locally equivalent ckts**
- iO for Turing machines

# Gate-by-Gate Obfuscation

# Gate-by-Gate Obfuscation

# Gate-by-Gate Obfuscation



$o$

$g$

$l$     $r$

Obfuscate

iO for **small** ckt

$iO(C_g)$

# Gate-by-Gate Obfuscation

$o$

$g$

$l$ $r$

$\square$ : encrypt and sign the wire values

Obfuscate

iO for **small** ckt

$iO(C_g)$

# Gate-by-Gate Obfuscation

$o$

$g$

$l$     $r$

⬇ Obfuscate

iO for ***small*** ckt

$iO(C_g)$

□ : encrypt and sign the wire values

$$C_g( \boxed{m_l} \; \boxed{m_r} )$$
___
Decrypt $m_l, m_r$
$m_o = g(w_l, w_r)$
**Output**: $\boxed{m_o}$

# Gate-by-Gate Obfuscation

$o$

$g$

$l$     $r$

Obfuscate

iO for **small** ckt

$iO(C_g)$

: encrypt and sign the wire values

$$C_g(\ \boxed{m_l}\ \ \boxed{m_r}\ )$$

Decrypt $m_l, m_r$

$m_o = g(w_l, w_r)$

**Output**: $\boxed{m_o}$

**Key Feature**: obfuscated circuit preserves the **topology** of the input circuit

# Prove Security w/o $2^{|input|}$ Loss

# Prove Security w/o $2^{|\text{input}|}$ Loss

$C$:

$C'$:

$C, C'$: Locally Equivalent
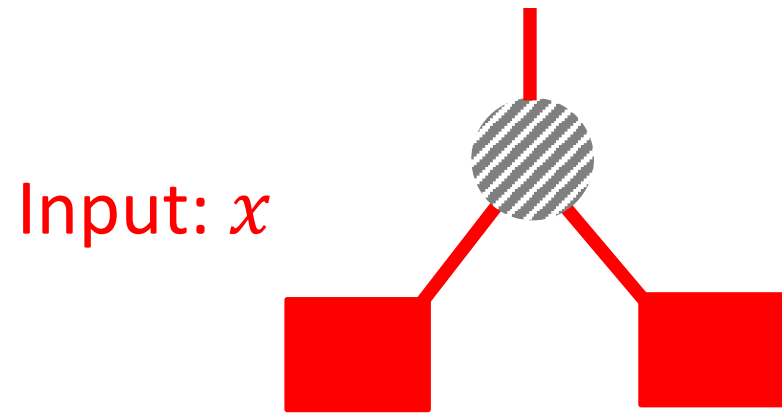
# Prove Security w/o $2^{|\text{input}|}$ Loss

$C$:

Gate-by-Gate
Obfuscation

$C'$:

$C, C'$: Locally Equivalent

# Prove Security w/o $2^{|\text{input}|}$ Loss

$C$:

Gate-by-Gate
Obfuscation

$C'$:

$C, C'$: Locally Equivalent

# Prove Security w/o $2^{|\text{input}|}$ Loss



$C$:

$C'$:

Gate-by-Gate Obfuscation

$C, C'$: Locally Equivalent

# Prove Security w/o $2^{|\text{input}|}$ Loss



$C$:

$C'$:

Gate-by-Gate Obfuscation

$C, C'$: Locally Equivalent

# Prove Security w/o $2^{|\text{input}|}$ Loss



$C$:

Gate-by-Gate Obfuscation

**Security Loss:** $2^{|subckt\ input|}$ (poly)

$C'$:

$C, C'$: Locally Equivalent

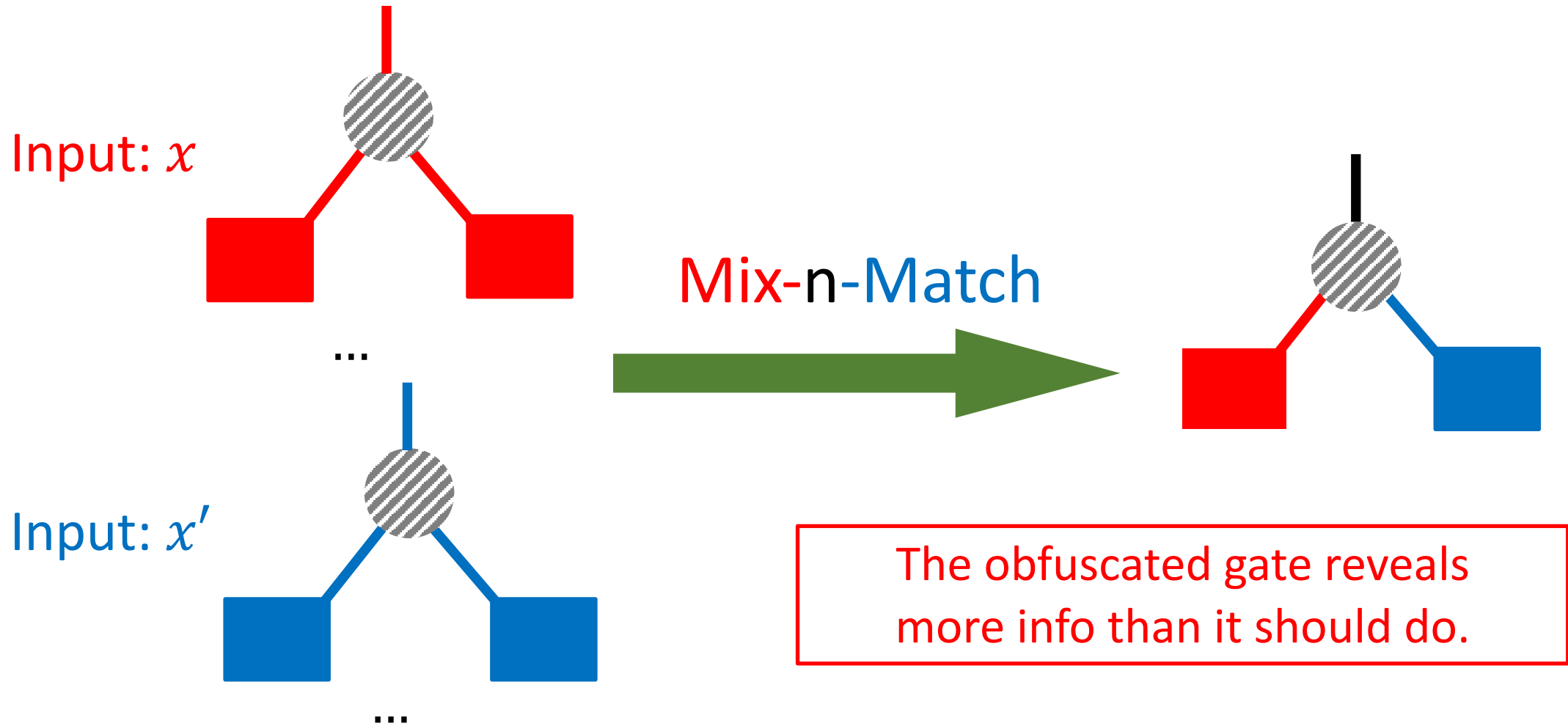# Challenge: Mix-and-Match Attack

# Challenge: Mix-and-Match Attack

Input: $x$

# Challenge: Mix-and-Match Attack

Input: $x$

...

Input: $x'$

...

# Challenge: Mix-and-Match Attack

Input: $x$

Input: $x'$

...

...

Mix-n-Match

# Challenge: Mix-and-Match Attack

Input: $x$

Input: $x'$

Mix-n-Match

# Challenge: Mix-and-Match Attack

Input: $x$

Mix-n-Match

Input: $x'$

The obfuscated gate reveals more info than it should do.

$$C_g(ct_l, ct_r, input)$$

Check consistency w.r.t input

....

# Idea 1: Replace Input with Dependent Wires

New Challenge
Too Long

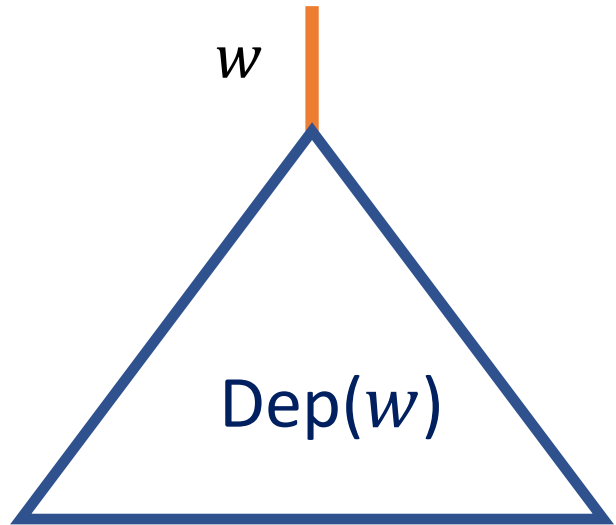$$C_g(ct_l, ct_r, input)$$
_____

Check consistency w.r.t input

....

# Idea 1: Replace Input with Dependent Wires

**New Challenge**
Too Long

$$\frac{C_g(ct_l, ct_r, input)}{}$$

Check consistency w.r.t input

....

Gate $g$ may not depend on the entire input
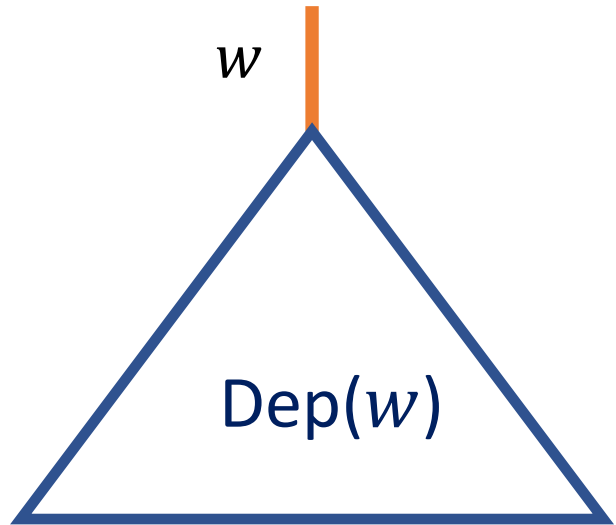(e.g. $NC^0$ circuits)

# Define Dependence

$w$

Dep($w$)

# Define Dependence



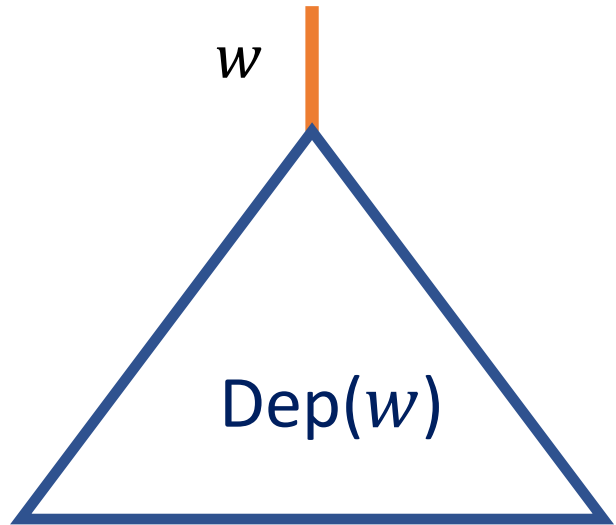$\mathrm{Dep}(w) := \{$ all wires that $w$ depends on $\}$

# Define Dependence

$w$



Dep($w$)

$\text{Dep}(w) := \{ \text{ all wires that } w \text{ depends on } \}$
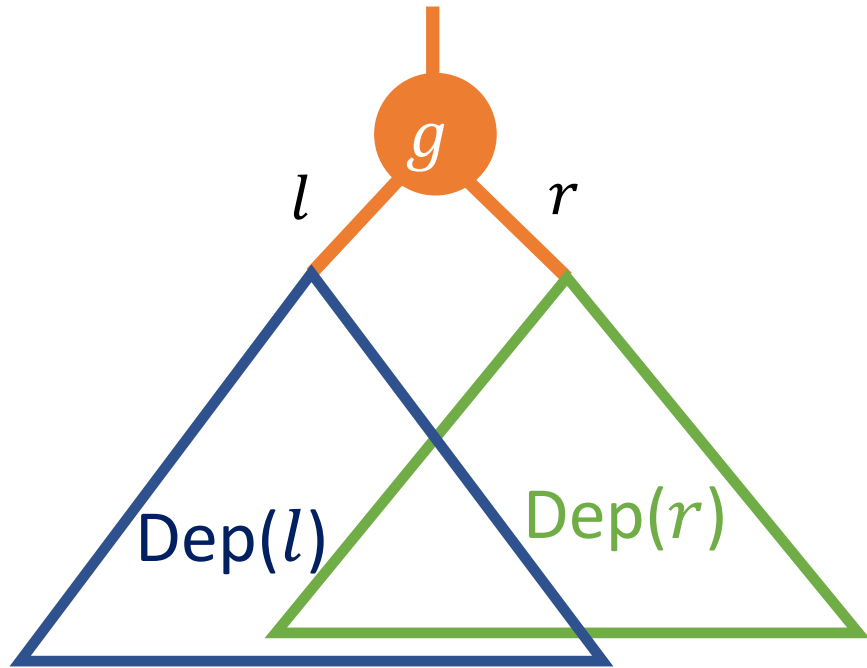
$CT_w := \{ciphertext\ of\ k\}_{k \in \text{Dep}(w)}$

# Define Dependence



$\text{Dep}(w) := \{ \text{ all wires that } w \text{ depends on } \}$

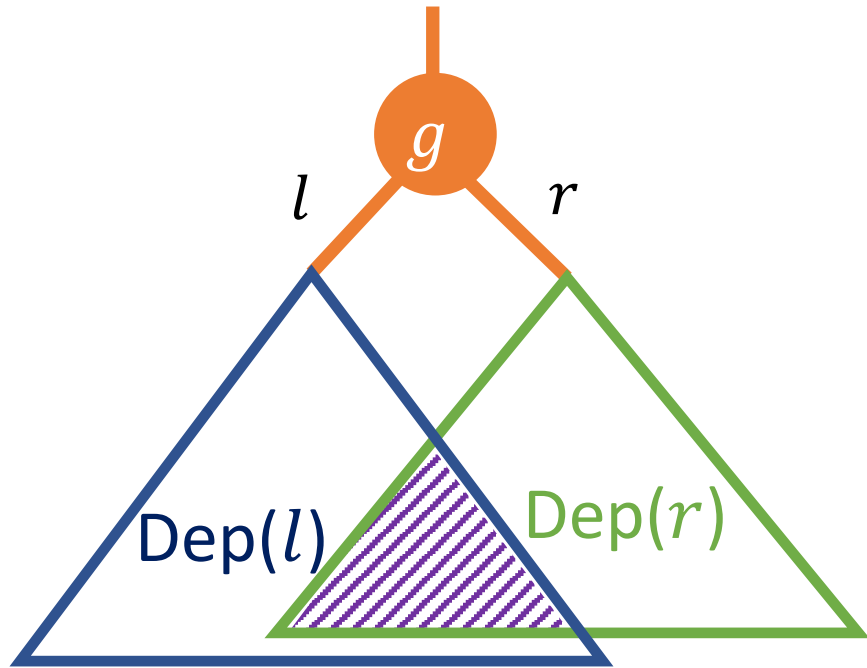$CT_w := \{ciphertext\ of\ k\}_{k \in \text{Dep}(w)}$

(An Index Set)

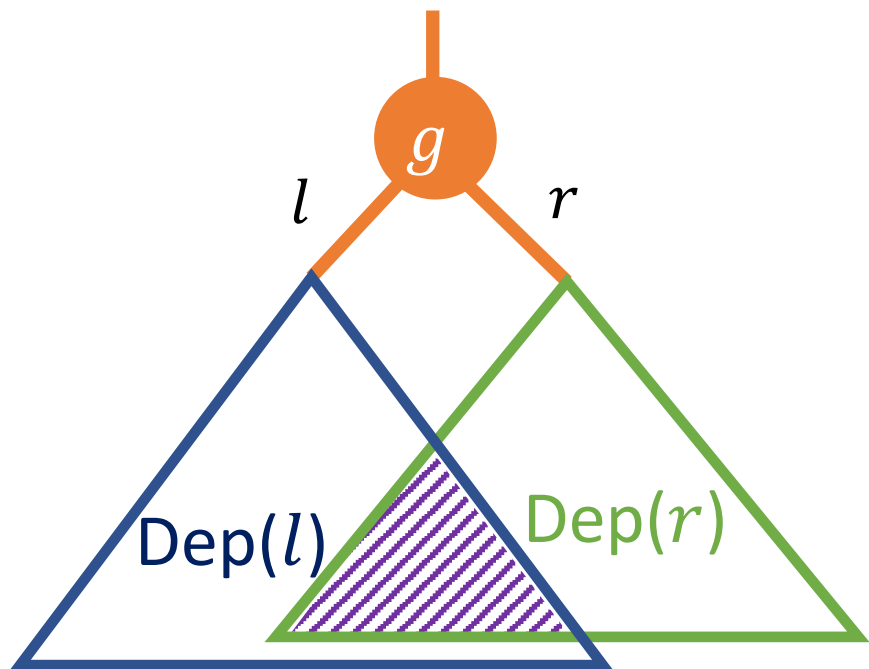Use $CT_l, CT_r$ in $C_g$



$$C_g(ct_l, ct_r, CT_l, CT_r)$$

...

Use $CT_l, CT_r$ in $C_g$



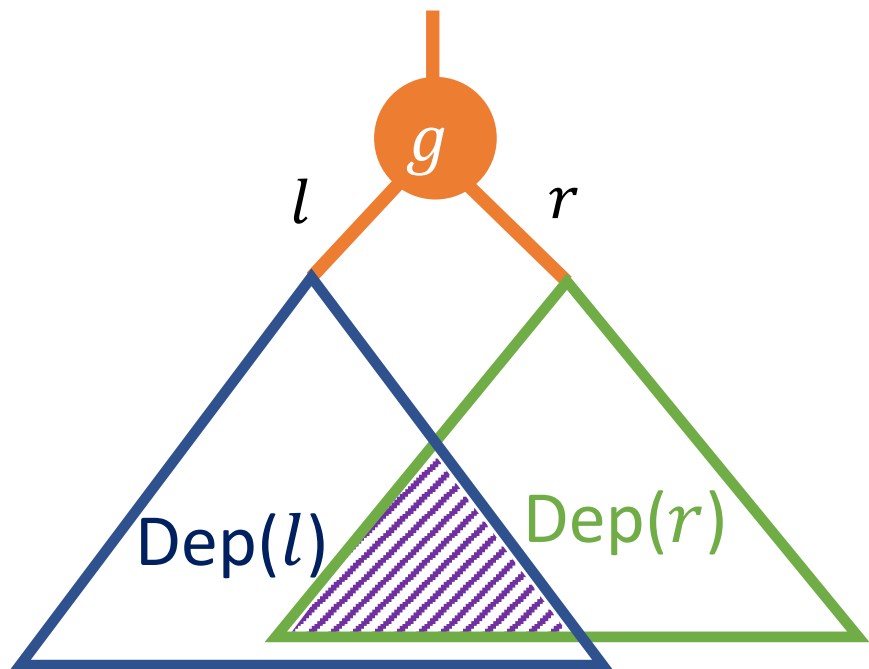$$C_g(ct_l, ct_r, CT_l, CT_r)$$

...

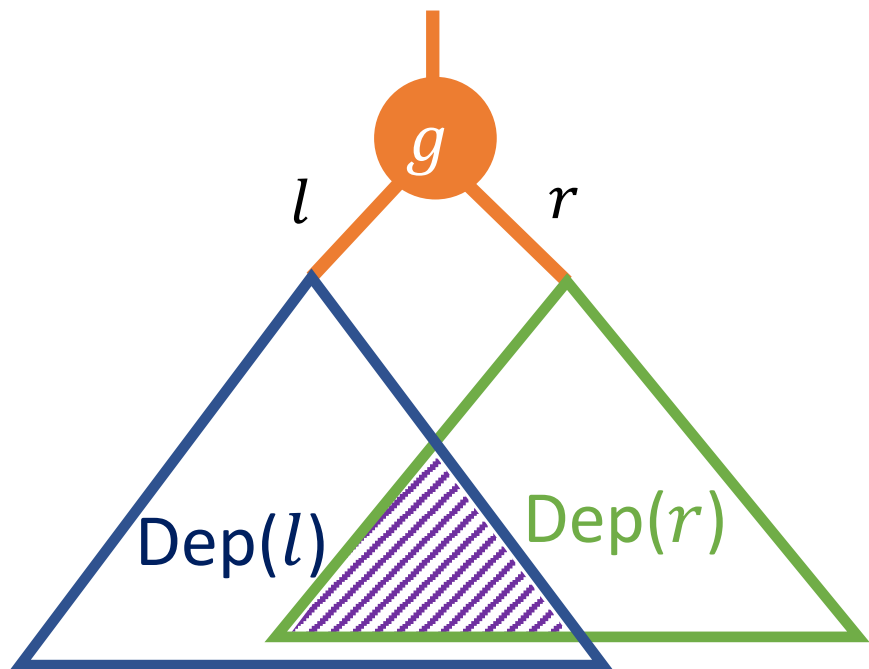Use $CT_l, CT_r$ in $C_g$



$$C_g(ct_l, ct_r, CT_l, CT_r)$$

**Consistency Check:**
$CT_l, CT_r$ contains same ciphertexts
in $Dep(l) \cap Dep(r)$

...

Use $CT_l, CT_r$ in $C_g$



$$C_g(ct_l, ct_r, CT_l, CT_r)$$

**Consistency Check:**
$CT_l, CT_r$ contains same ciphertexts
in $\text{Dep}(l) \cap \text{Dep}(r)$

...

# Idea 2: Hash $CT_l, CT_r$

$$C_g(ct_l, ct_r, CT_l, CT_r)$$
---
...

...

# Idea 2: Hash $CT_l, CT_r$

$$\frac{C_g(ct_l, ct_r, CT_l, CT_r)}{\cdots}$$

$$\cdots$$

Outside of $C_g$:

$h_l$      $h_r$

Hash      Hash

$CT_l$      $CT_r$

# Idea 2: Hash $CT_l, CT_r$

$$\frac{C_g(ct_l, ct_r, \qquad )}{\cdots}$$

$$\cdots$$

**Outside of $C_g$:**

$h_l$          $h_r$

**Hash**      **Hash**

$CT_l$          $CT_r$

# Idea 2: Hash $CT_l, CT_r$

$$C_g(ct_l, ct_r, \quad h_l, h_r \quad )$$

$$\cdots$$

$$\cdots$$

**Outside of $C_g$:**

$h_l$ $\qquad\qquad$ $h_r$

**Hash** $\qquad$ **Hash**

$CT_l$ $\qquad\qquad$ $CT_r$

# Idea 2: Hash $CT_l, CT_r$

# Idea 2: Hash $CT_l, CT_r$

$$C_g(ct_l, ct_r, \quad h_l, h_r \quad)$$

...

Check consistency of $CT_l$ and $CT_r$ **???**

...

**Outside of $C_g$:**

$h_l$               $h_r$

**Hash**        **Hash**

$CT_l$            $CT_r$

# Idea 3: Apply SNARGs

**Outside** $C_g$:
$$h_l = Hash(CT_l)$$
$$h_r = Hash(CT_r)$$

# Idea 3: Apply SNARGs

**Outside** $C_g$:
$$h_l = Hash(CT_l)$$
$$h_r = Hash(CT_r)$$

SNARGs (Succinct Cryptographic Proofs)

$\pi$ : prove $\exists$ consistent pre-images of $h_l, h_r$

Secure against poly-time adversary

# Idea 3: Apply SNARGs

**Outside** $C_g$:
$$h_l = Hash(CT_l)$$
$$h_r = Hash(CT_r)$$

SNARGs (Succinct Cryptographic Proofs)

$\pi$ : prove $\exists$ consistent pre-images of $h_l, h_r$

Secure against poly-time adversary

$$\frac{C_g(ct_l, ct_r, h_l, h_r, \pi)}{\text{Verify the proof } \pi}$$

…Decrypt, Compute, Re-encrypt…

# Idea 3: Apply SNARGs

**Outside** $C_g$:
$$h_l = Hash(CT_l)$$
$$h_r = Hash(CT_r)$$

SNARGs (Succinct Cryptographic Proofs)

$\pi$ : prove $\exists$ consistent pre-images of $h_l, h_r$
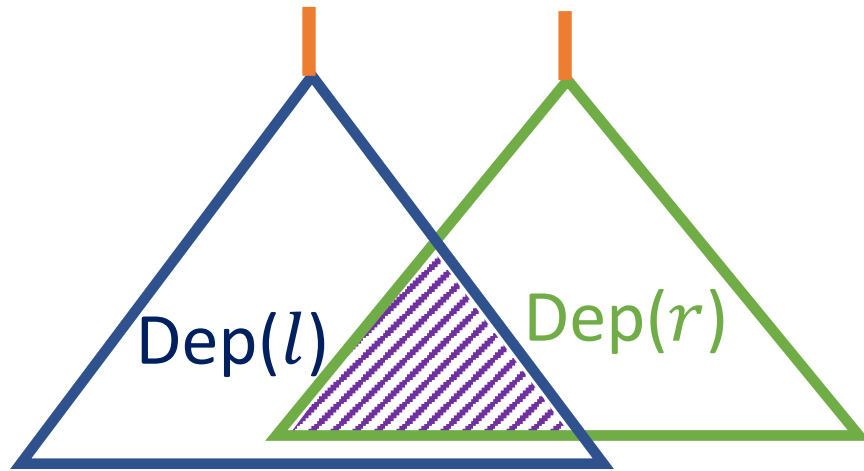
Secure against poly-time adversary

$$\frac{C_g(ct_l, ct_r, h_l, h_r, \pi)}{\text{Verify the proof } \pi}$$
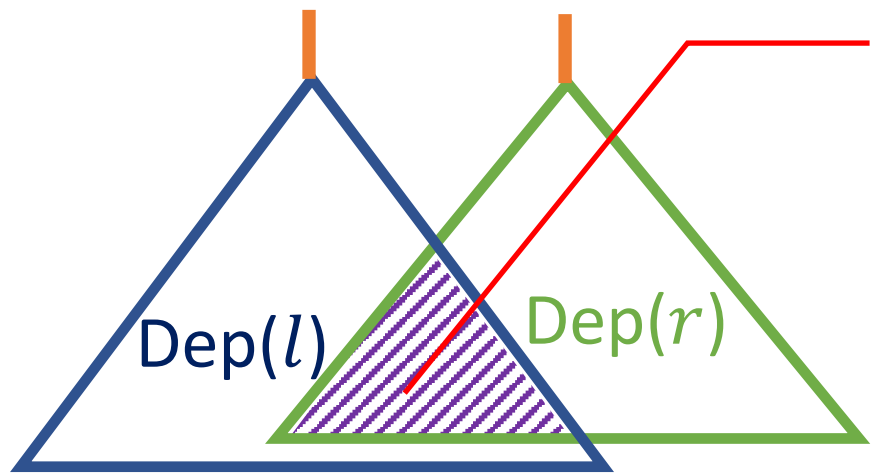
...Decrypt, Compute, Re-encrypt...

**New Challenge**: We need *statistical security* of SNARGs for iO.
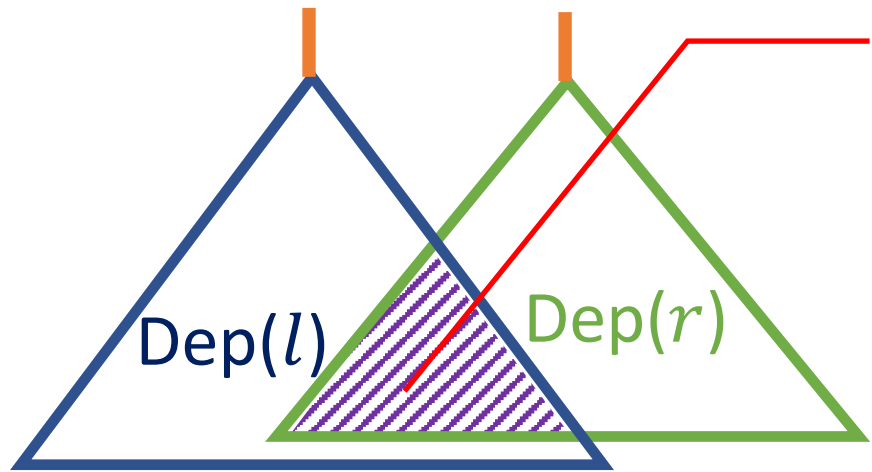
# We Use: iO-Friendly SNARGs

# We Use: iO-Friendly SNARGs
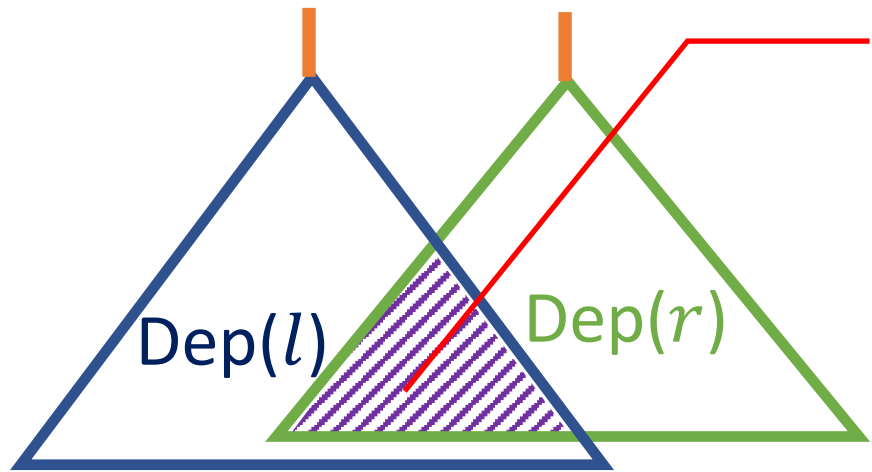
**Observation**: We only care about sub-circuit

$\text{Dep}(l)$   $\text{Dep}(r)$

# We Use: iO-Friendly SNARGs



**Observation**: We only care about sub-circuit

$\text{Dep}(l)$       $\text{Dep}(r)$

> Somewhere Statistical Soundness:
> If $CT_l$ and $CT_r$ are not *consistent in subcircuit*,
> then unbounded-time adversary can't cheat.

# We Use: iO-Friendly SNARGs



**Observation**: We only care about sub-circuit
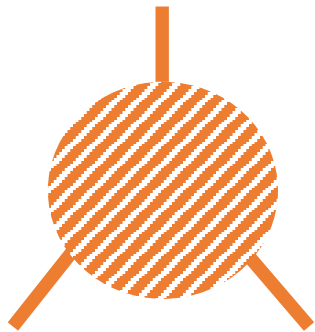
Dep($l$)  Dep($r$)

Somewhere Statistical Soundness:
If $CT_l$ and $CT_r$ are not *consistent in subcircuit*,
then unbounded-time adversary can't cheat.
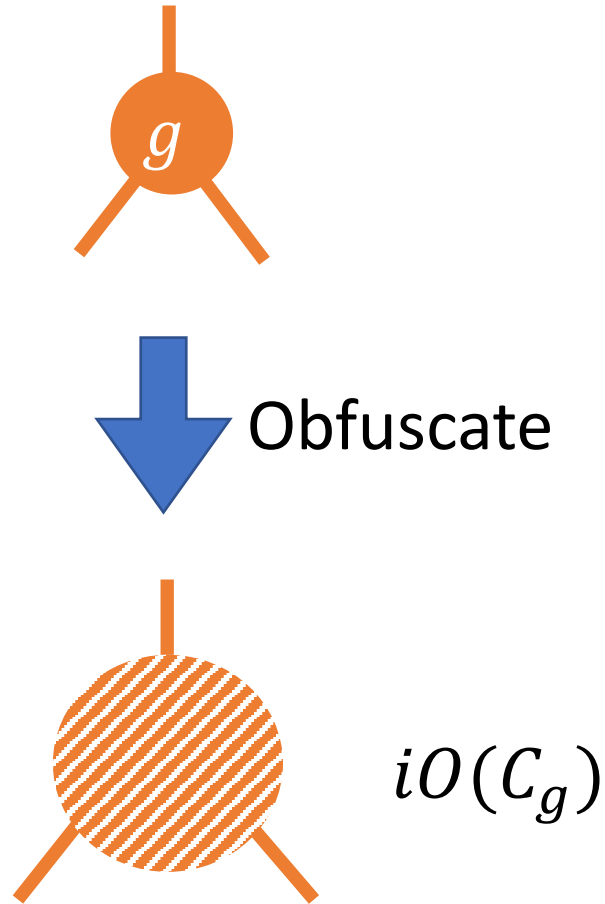
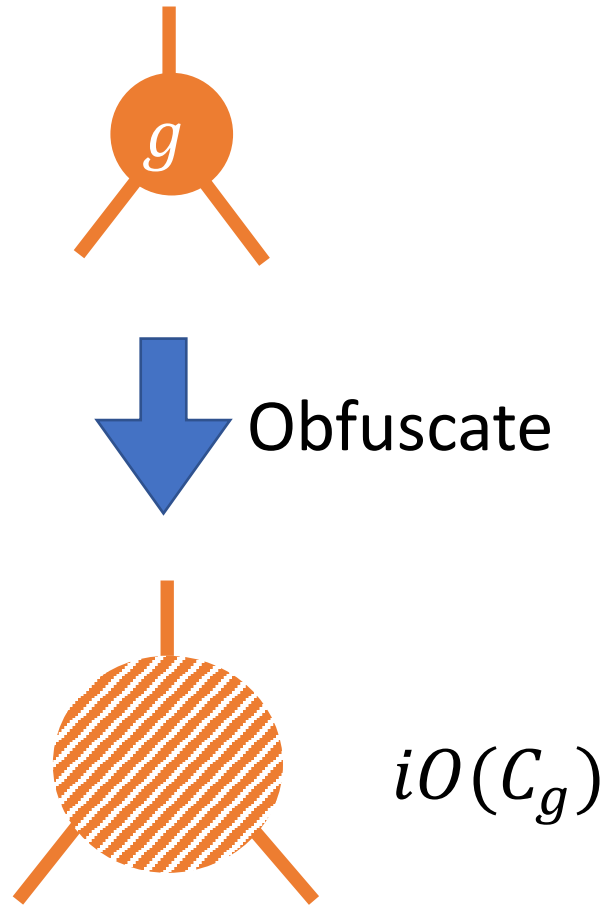Can be constructed from [CJJ'21]

# Summary

# Summary



**Outside** $C_g$:

$$h_l = \text{Hash}(CT_l)$$
$$h_r = \text{Hash}(CT_r)$$

$\pi$ : iO-friendly consistency proof for $h_l, h_r$

# Summary



**Outside** $C_g$:

$$h_l = \text{Hash}(CT_l)$$
$$h_r = \text{Hash}(CT_r)$$

$\pi$ : iO-friendly consistency proof for $h_l, h_r$

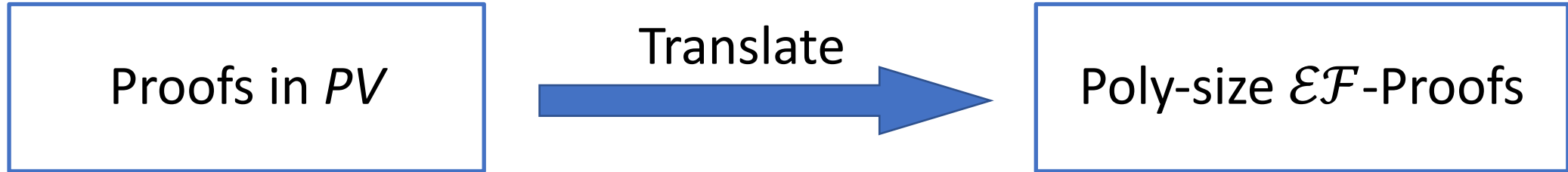$$C_g(ct_l, ct_r, h_l, h_r, \pi)$$
___

Verify the proof $\pi$
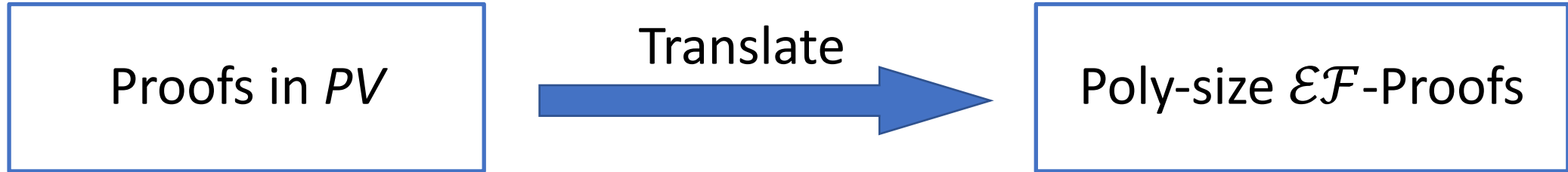
...Decrypt, Compute, Re-encrypt...

# Technical Details

- $\mathcal{EF}$-Proofs $\Rightarrow$ local equivalence
- iO for locally equivalent ckts
- **iO for Turing machines**

# Propositional Translation [Cook'75]
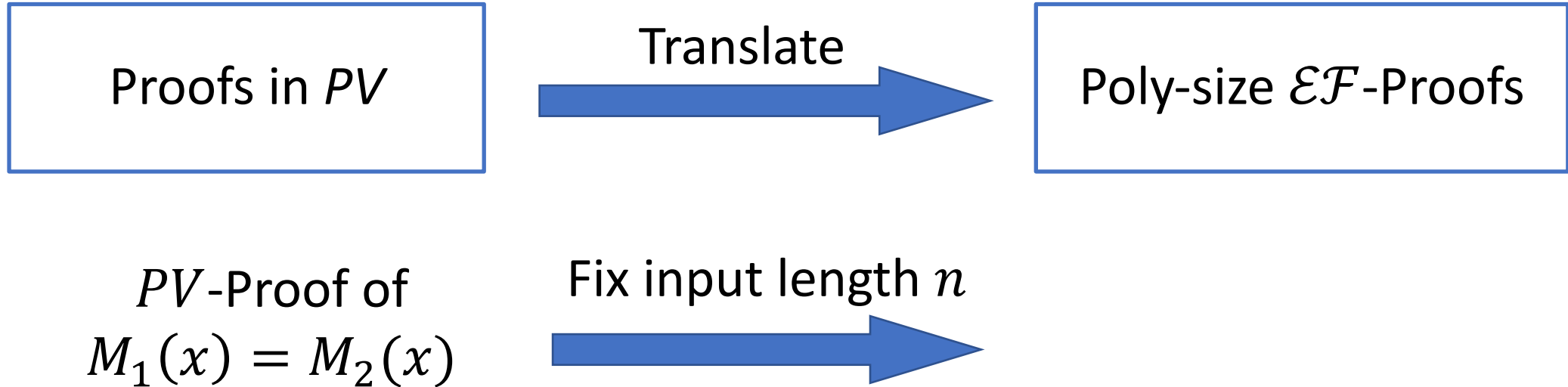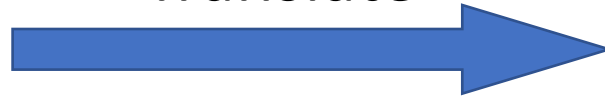
# Propositional Translation [Cook'75]

# Propositional Translation [Cook'75]

Proofs in $PV$ → Translate → Poly-size $\mathcal{EF}$-Proofs

$PV$-Proof of
$$M_1(x) = M_2(x)$$

# Propositional Translation [Cook'75]

| Proofs in $PV$ | Translate $\longrightarrow$ | Poly-size $\mathcal{EF}$-Proofs |

$PV$-Proof of
$M_1(x) = M_2(x)$

Fix input length $n$ $\longrightarrow$

# Propositional Translation [Cook'75]

| Proofs in $PV$ | Translate $\longrightarrow$ | Poly-size $\mathcal{EF}$-Proofs |

$PV$-Proof of
$M_1(x) = M_2(x)$

Fix input length $n$
$\longrightarrow$

$\mathcal{EF}$-Proof of
$C_{1,n}(x) \leftrightarrow C_{2,n}(x)$

# Propositional Translation [Cook'75]

Proofs in $PV$ → Translate → Poly-size $\mathcal{EF}$-Proofs

$PV$-Proof of
$M_1(x) = M_2(x)$

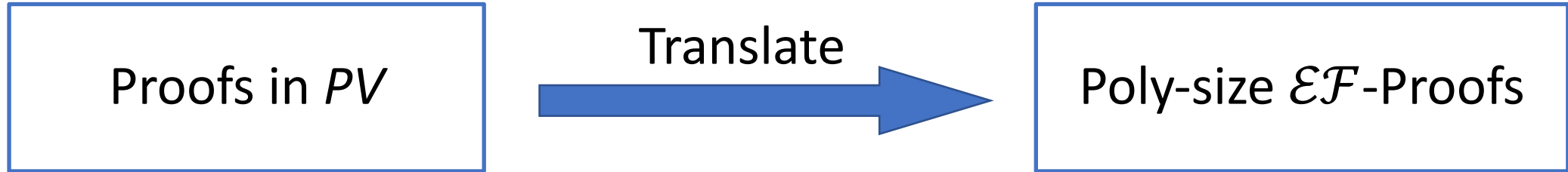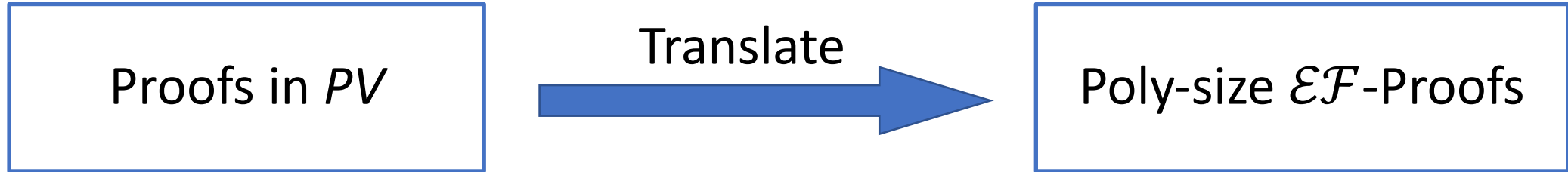Fix input length $n$ →

$\mathcal{EF}$-Proof of
$C_{1,n}(x) \leftrightarrow C_{2,n}(x)$

($C_{b,n}(x)$: Circuit that computes $M_b$ for input length $n$.)

# Propositional Translation [Cook'75]

| Proofs in $PV$ | Translate $\longrightarrow$ | Poly-size $\mathcal{EF}$-Proofs |
|---|---|---|

$PV$-Proof of
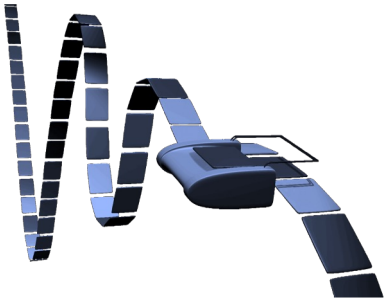$M_1(x) = M_2(x)$

Fix input length $n$ $\longrightarrow$

$\mathcal{EF}$-Proof of
$C_{1,n}(x) \leftrightarrow C_{2,n}(x)$

($C_{b,n}(x)$: Circuit that computes $M_b$ for input length $n$.)
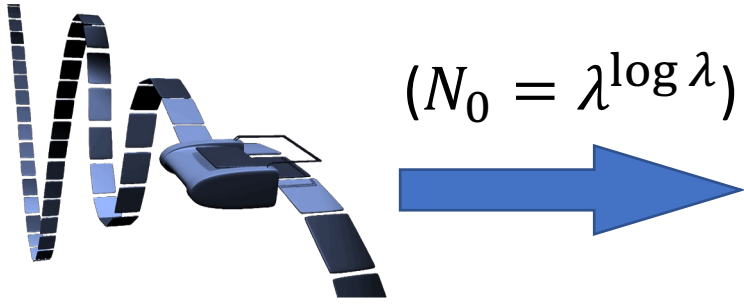
Apply $iO$ for locally equivalent circuits?

# iO for Turing Machines: Initial Attempt
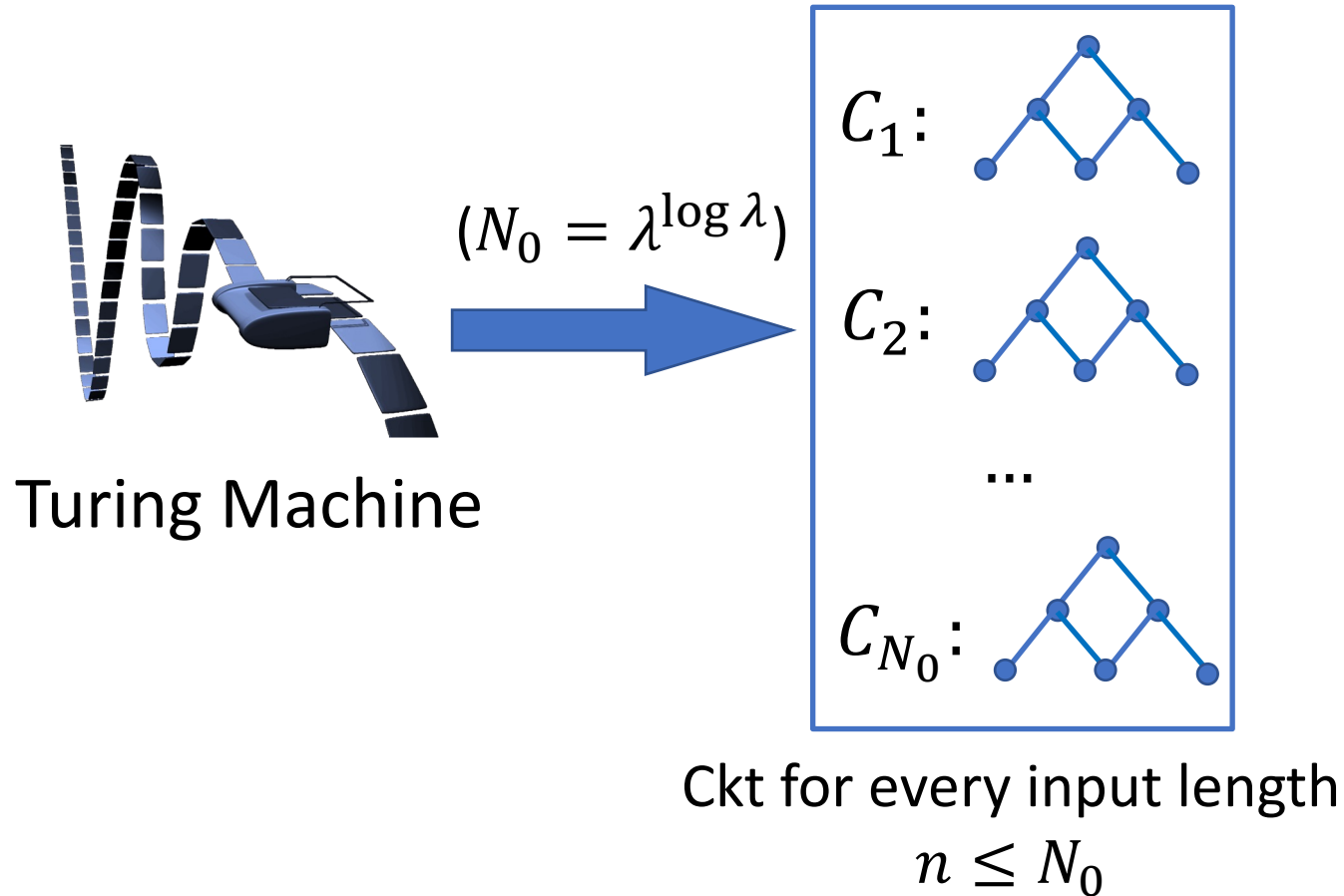
# iO for Turing Machines: Initial Attempt



Turing Machine

# iO for Turing Machines: Initial Attempt
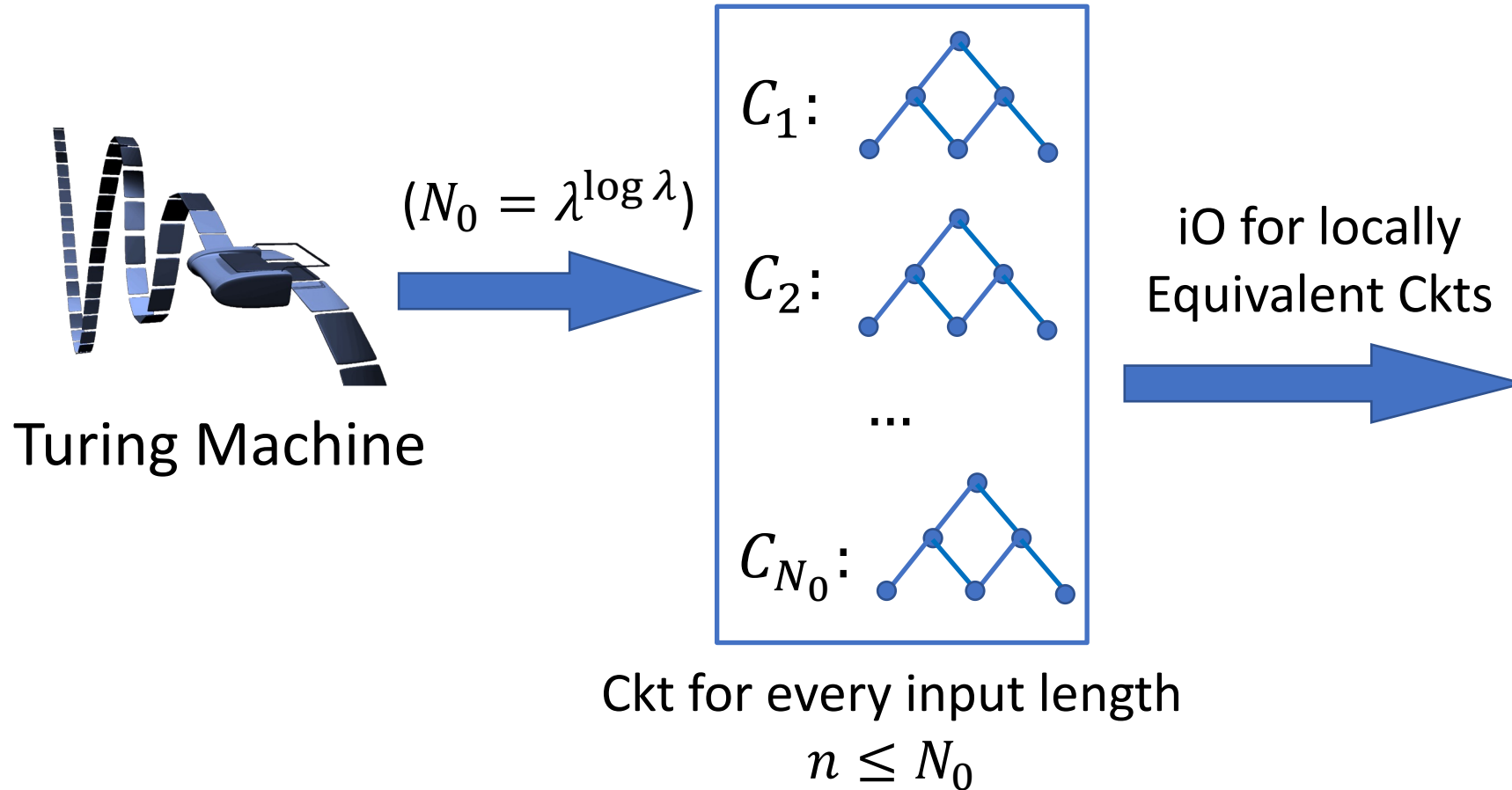


$$(N_0 = \lambda^{\log \lambda})$$

Turing Machine

# iO for Turing Machines: Initial Attempt



Turing Machine

$(N_0 = \lambda^{\log \lambda})$

$C_1:$

$C_2:$

$\ldots$

$C_{N_0}:$

Ckt for every input length
$n \leq N_0$

# iO for Turing Machines: Initial Attempt



$(N_0 = \lambda^{\log \lambda})$

Turing Machine

$C_1:$

$C_2:$

...

$C_{N_0}:$

iO for locally
Equivalent Ckts

Ckt for every input length
$n \leq N_0$

# iO for Turing Machines: Initial Attempt



$(N_0 = \lambda^{\log \lambda})$

$C_1$:

$C_2$:

...

$C_{N_0}$:

Turing Machine

iO for locally Equivalent Ckts

Ckt for every input length $n \leq N_0$

$\widetilde{C_1}$

$\widetilde{C_2}$

...

$\widetilde{C_{N_0}}$

# iO for Turing Machines: Initial Attempt



$(N_0 = \lambda^{\log \lambda})$

$C_1:$

$C_2:$

$\ldots$

$C_{N_0}:$

Turing Machine

Ckt for every input length
$n \leq N_0$

iO for locally
Equivalent Ckts

$\widetilde{C_1}$

$\widetilde{C_2}$

$\ldots$

$\widetilde{C_{N_0}}$

**Obfuscated Turing Machine**

# Leverage Uniform Description



$C_1$:

...

$C_n$:

...

$C_{N_0}$:

Small ckt

$[M]$ ← $n$

← $i$

Compute description of $i$-th gate of $C_n$

# Leverage Uniform Description

$C_1:$

...

$C_n:$ $i$

...

$C_{N_0}:$

Small ckt

$[M] \leftarrow n$

$\leftarrow i$

Compute description of $i$-th gate of $C_n$

**Recall: Obfuscation of $C_n$**

# Leverage Uniform Description

# Leverage Uniform Description

$C_1$:

...

$C_n$:

$i$

...

$C_{N_0}$:

Small ckt

$[M]$ ← $n$

← $i$

Compute description of $i$-th gate of $C_n$

**Recall: Obfuscation of $C_n$**

$iO(C_{g_1})$

$iO(C_{g_2})$

...

**Idea**: Generate $C_g$ on-the-fly using $[M]$
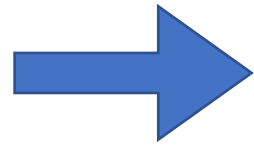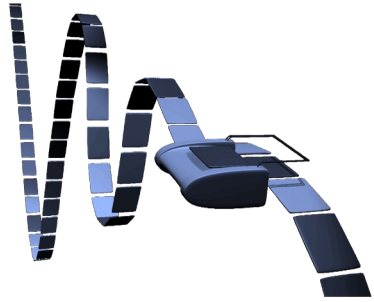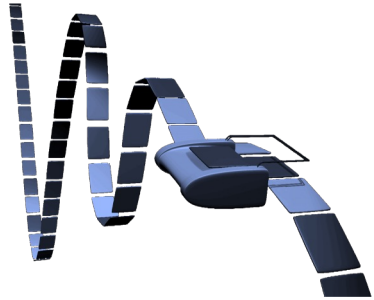
# Efficient Construction

# Efficient Construction
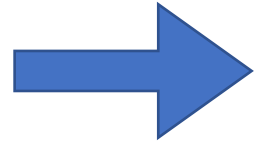


Turing Machine

# Efficient Construction



Turing Machine

# Efficient Construction
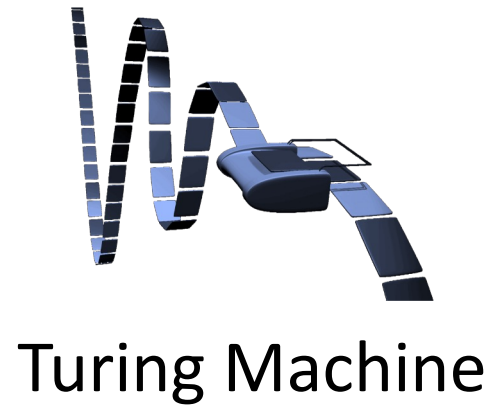


Turing Machine

$\longrightarrow$

$UC$

"Uniform Version" $C_g$

# Efficient Construction



Turing Machine

$UC$

"Uniform Version" $C_g$

$\underline{UC(n, i, input')}$

- Get description of $i$-th gate of $C_n$:

$$g \leftarrow [M] \leftarrow \begin{array}{c} n \\ i \end{array}$$

- Compute $C_g(input')$

# Efficient Construction



Turing Machine

$UC$

"Uniform Version" $C_g$

iO for Ckt

$\underline{UC(n, i, input')}$

- Get description of $i$-th gate of $C_n$:

  $g \leftarrow [M] \leftarrow n$

  $[M] \leftarrow i$

- Compute $C_g(input')$

# Efficient Construction



Turing Machine

$UC$

"Uniform Version" $C_g$

iO for Ckt

$iO(UC)$

Obfuscated Turing Machine

$\underline{UC(n, i, input')}$

- Get description of $i$-th gate of $C_n$:

$g \leftarrow [M] \leftarrow n$
$\quad\quad\quad \leftarrow i$

- Compute $C_g(input')$

# Summary & Future Directions

# Summary & Future Directions

Inference Rules in
**Logic systems** for
Proving Equivalence

# Summary & Future Directions

Inference Rules in
**Logic systems** for
Proving Equivalence

# Summary & Future Directions

Inference Rules in
**Logic systems** for
Proving Equivalence

⟷

Techniques to argue
Indistinguishability for iO

# Summary & Future Directions

Inference Rules in
**Logic systems** for
Proving Equivalence

$\mathcal{EF}$ / $PV$

Techniques to argue
Indistinguishability for iO

# Summary & Future Directions

Inference Rules in
**Logic systems** for
Proving Equivalence

$\mathcal{EF} \ / \ PV$

Techniques to argue
Indistinguishability for iO

# Summary & Future Directions

Inference Rules in
**Logic systems** for
Proving Equivalence

⟷

Techniques to argue
Indistinguishability for iO

$\mathcal{EF}$ / $\mathcal{PV}$

⟷

Local Equivalence

# Summary & Future Directions

Inference Rules in
**Logic systems** for
Proving Equivalence

⟷

Techniques to argue
Indistinguishability for iO

$\mathcal{EF}$ / $PV$

⟷

Local Equivalence

$ZFC$
(**Z**ermelo-**F**raenkel set theory
with axiom of **C**hoice)

⟷

**?**

# Thank you!

# Q & A