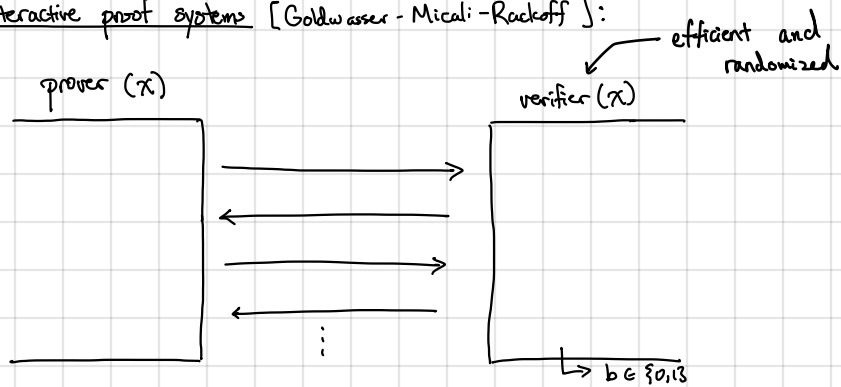Schnorr's protocol is actually an example of more general concept called zero-knowledge proofs:

Interactive proof systems [Goldwasser - Micali - Rackoff]:

efficient and randomized

prover $(x)$                              verifier $(x)$



$\hookrightarrow b \in \{0,1\}$

Interactive proof should satisfy completeness + soundness (as defined earlier)

Consider following example: Suppose prover wants to convince verifier that $N = pq$ where $p, q$ are prime (and secret).

prover $(N, p, q)$                          verifier $(N)$

$\pi = (p, q) \longrightarrow$

$\downarrow$

accept if $N = pq$ and reject otherwise

Proof is certainly complete and sound, but now verifier <u>also</u> learned the factorization of $N$... (may not be desirable if prover was trying to convince verifier that $N$ is a proper RSA modulus (for a cryptographic scheme) <u>without revealing</u> factorization in the process

$\hookrightarrow$ In some sense, this proof conveys <u>information</u> to the verifier [i.e., verifier learns something it did not know before seeing the proof]

<u>Zero-knowledge</u>: ensure that verifier does <u>not</u> learn anything (other than the fact that the statement is true)

<u>How do we define "zero-knowledge"</u>?  We will introduce a notion of a "simulator."

for a language $L$

<u>Definition</u>. An interactive proof system $\langle P, V \rangle$ is zero-knowledge if for all efficient (and possibly malicious) verifiers $V^*$, there exists an efficient simulator $S$ such that for all $x \in L$:

$$\underbrace{\text{View}_{V^*}(\langle P, V \rangle(x))}_{} \overset{c}{\approx} S(x)$$

random variable denoting the set of messages
sent and received by $V^*$ when interacting with the prover $P$ on input $x$

What does this definition mean?

$\text{View}_{V^*}(P \leftrightarrow V^*(x))$ : this is what $V^*$ sees in the interactive proof protocol with $P$

$S(x)$ : this is a function that only depends on the statement $x$, which $V^*$ already has

If these two distributions are indistinguishable, then anything that $V^*$ could have learned by talking to $P$, it could have learned just by invoking the simulator itself, and the simulator output only depends on $x$, which $V^*$ already knows

↳ In other words, anything $V^*$ could have learned (i.e., computed) after interacting with $P$, it could have learned _without_ ever talking to $P$!
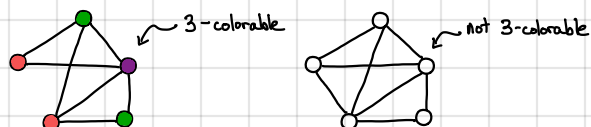
Very remarkable definition!

— can in fact be constructed from OWFs

**More remarkable** : Using cryptographic commitments, then every language $L \in IP$ has a zero-knowledge proof system.

↳ Namely, anything that can be proved can be proved in zero-knowledge!

We will show this theorem for NP languages. Here it suffices to construct a single zero-knowledge proof system for an NP-complete language. We will consider the language of graph 3-colorability.


— 3-colorable          — not 3-colorable

**3-coloring** : given a graph $G$, can you color the vertices so that no adjacent nodes have the same color?

— cryptographic analog of a sealed "envelope"

We will need a commitment scheme. A (non-interactive) commitment scheme consists of three algorithms (Setup, Commit, Open):

- Setup $\to \sigma$ : Outputs a common reference string (used to generate/validate commitments) $\sigma$
- Commit$(\sigma, m) \to (c, \pi)$ : Takes the CRS $\sigma$ and message $m$ and outputs a commitment $c$ and opening $\pi$
- Verify$(\sigma, m, c, \pi) \to 0/1$ : Checks if $c$ is a valid commitment to $m$ (given $\pi$)

**Typical setup** :

$\qquad$ Committer $\qquad\qquad\qquad\qquad\qquad\qquad$ Verifier

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $\sigma \leftarrow$ Setup

$\qquad\qquad\qquad\qquad\qquad \overset{\sigma}{\longleftarrow}$

$(c, \pi) \leftarrow \text{Commit}(\sigma, m) \quad \overset{c}{\longrightarrow}$

$\qquad\qquad\qquad\qquad\qquad$ (sometime later)

$\qquad\qquad\qquad\qquad\qquad \overset{m, \pi}{\longrightarrow}$

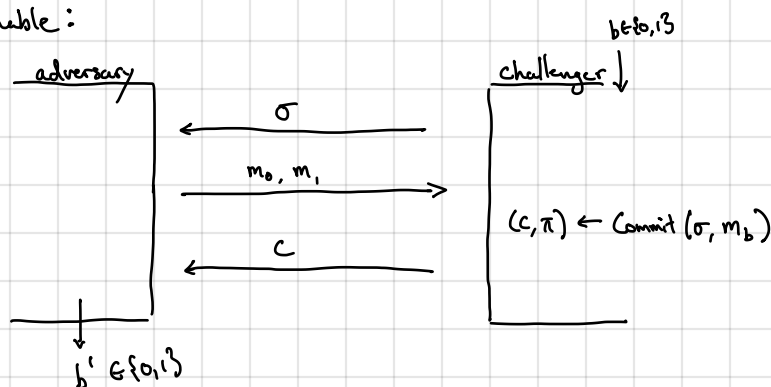$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ can check that Verify$(\sigma, m, c, \pi) = \underline{1}$

Requirements:

- **Correctness**: for all messages $m$:
$$\Pr\left[\sigma \leftarrow \text{Setup}, (c, \pi) \leftarrow \text{Commit}(\sigma, m); \text{Verify}(\sigma, c, m, \pi) = 1\right] = 1$$

- **Hiding**: for all common reference strings $\sigma \in \{0,1\}^n$ and all efficient $A$, following distributions are computationally indistinguishable:
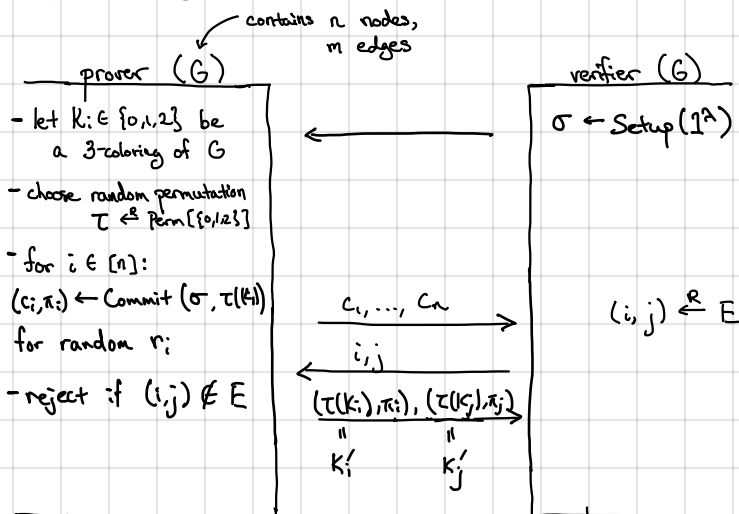


$b \in \{0,1\}$

adversary

challenger ↓

$\sigma$

$m_0, m_1$

$(c, \pi) \leftarrow \text{Commit}(\sigma, m_b)$

$c$

$b' \in \{0,1\}$

$$\left| \Pr[b' = 1 \mid b = 0] - \Pr[b' = 1 \mid b = 1] \right| = \text{negl}(\lambda)$$

- **Binding**: for all adversaries $A$, if $\sigma \leftarrow \text{Setup}$, then
$$\Pr\left[(m_0, m_1, c, \pi_0, \pi_1) \leftarrow A : m_0 \neq m_1 \text{ and } \text{Verify}(\sigma, c, m_0, \pi_0) = 1 = \text{Verify}(\sigma, c, m_1, \pi_1)\right] = \text{negl}$$

A ZK protocol for graph 3-coloring:



contains $n$ nodes, $m$ edges

prover $(G)$

verifier $(G)$

$\sigma \leftarrow \text{Setup}(1^\lambda)$

- let $k_i \in \{0,1,2\}$ be a 3-coloring of $G$

- choose random permutation $\tau \xleftarrow{\$} \text{Perm}[\{0,1,2\}]$

- for $i \in [n]$:
$(c_i, \pi_i) \leftarrow \text{Commit}(\sigma, \tau(k_i))$
for random $r_i$

- reject if $(i,j) \notin E$

$c_1, \ldots, c_n$

$(i,j) \xleftarrow{R} E$

$i, j$

$\underbrace{(\tau(k_i), \pi_i)}_{k_i'}, \underbrace{(\tau(k_j), \pi_j)}_{k_j'}$

↳ accept if $k_i' \neq k_j'$ and $k_i', k_j' \in \{0,1,2\}$
$\text{Verify}(\sigma, c_i, k_i', \pi_i) = 1 = \text{Verify}(\sigma, c_j, k_j', \pi_j)$

reject otherwise

**Intuitively**: Prover commits to a coloring of the graph

Verifier challenges prover to reveal coloring of a single edge

Prover reveals the coloring on the chosen edge and opens the entries in the commitment

**Completeness**: By inspection [if coloring is valid, prover can always answer the challenge correctly]

**Soundness**: Suppose $G$ is _not_ 3-colorable. Let $K_1, ..., K_n$ be the <span style="color:green">— except with prob. 1− negl.</span> coloring the prover committed to. If the commitment scheme is statistically binding, $c_1, ..., c_n$ _uniquely_ determine $K_1, ..., K_n$. Since $G$ is not 3-colorable, there is an edge $(i,j) \in E$ where $K_i = K_j$ or $i \notin \{0,1,2\}$ or $j \notin \{0,1,2\}$. <span style="color:green">[Otherwise, $G$ is 3-colorable with coloring $K_1, ..., K_n$.]</span> Since the verifier chooses an edge to check at random, the verifier will choose $(i,j)$ with probability $1/|E|$. Thus, if $G$ is not 3-colorable,

$$\Pr[\text{verifier rejects}] \geq \frac{1}{|E|}$$

Thus, this protocol provides soundness $1 - \frac{1}{|E|}$. We can repeat this protocol $O(|E|^2)$ times _sequentially_ to reduce soundness error to

$$\Pr[\text{verifier accepts proof of false statement}] \leq \left(1 - \frac{1}{|E|}\right)^{|E|^2} \leq e^{-|E|} = e^{-m} \quad \color{green}{\left[\text{since } 1+x \leq e^x\right]}$$

**Zero Knowledge**: We need to construct a simulator that outputs a valid transcript given only the graph $G$ as input.

Let $V^*$ be a (possibly malicious) verifier. Construct simulator $S$ as follows:

1. Run $V^*$ to get $\sigma^*$.

2. Choose $K_i \leftarrow \{0,1,2\}$ for all $i \in [n]$.
   Let $(c_i, \pi_i) \leftarrow \text{Commit}(\sigma^*, K_i)$     <span style="color:green">} Simulator does _not_ know coloring so it commits to a random one</span>
   Give $(c_1, ..., c_n)$ to $V^*$.

3. $V^*$ outputs an edge $(i,j) \in E$

4. If $K_i \neq K_j$, then $S$ outputs $(K_i, K_j, \pi_i, \pi_j)$.
   Otherwise, restart and try again (if fails $\lambda$ times, then abort)

Simulator succeeds with probability $2/3$ (over choice of $K_1, ..., K_n$). Thus, simulator produces a valid transcript with prob. $1 - \frac{1}{3^\lambda} = 1 - \text{negl}(\lambda)$ after $\lambda$ attempts. It suffices to show that simulated transcript is indistinguishable from a real transcript.

- _Real scheme_: prover opens $K_i, K_j$ where $k_i, k_j \overset{R}{\leftarrow} \{0,1,2\}$ [since prover randomly permutes the colors]
- _Simulation_: $K_i$ and $K_j$ sampled uniformly from $\{0,1,2\}$ and conditioned on $K_i \neq K_j$, distributions are identical

In addition, $(i,j)$ output by $V^*$ in the simulation is distributed correctly since commitment scheme is computationally-hiding (e.g. $V^*$ behaves essentially the same given commitments to a random coloring as it does given commitment to a valid coloring

If we repeat this protocol (for soundness amplification), simulator simulate one transcript at a time

**Summary**: Every language in NP has a zero-knowledge proof (assuming existence of PRGs)
                                                                                    ↰
                                                     PRGs imply commitments