

Round Optimal Concurrent Non-Malleability from Polynomial Hardness

Dakshita Khurana
UCLA
dakshita@cs.ucla.edu

Abstract

Non-malleable commitments are a central cryptographic primitive that guarantee security against man-in-the-middle adversaries, and their exact round complexity has been a subject of great interest. Pass (TCC 2013, CC 2016) proved that non-malleable commitments with respect to commitment are impossible to construct in less than three rounds, via black-box reductions to polynomial hardness assumptions. Obtaining a matching positive result has remained an open problem so far.

While three-round constructions of non-malleable commitments have been achieved, beginning with the work of Goyal, Pandey and Richelson (STOC 2016), current constructions require super-polynomial assumptions.

In this work, we settle the question of whether three-round non-malleable commitments can be based on polynomial hardness assumptions. We give constructions based on polynomial hardness of Decisional Diffie-Hellman assumption or Quadratic Residuosity or N^{th} Residuosity, together with ZAPs. Our protocols also satisfy concurrent non-malleability.

1 Introduction

Non-malleable commitments are a fundamental primitive in cryptography, that help prevent man-in-the-middle attacks. A man-in-the-middle (MIM) adversary participates simultaneously in multiple protocol executions, using information obtained in one execution to breach security of the other execution. To counter such adversaries, the notion of non-malleable commitments was introduced in a seminal work of Dolev, Dwork and Naor [DDN91]. From their inception, non-malleable commitments have been instrumental to building various several important non-malleable protocols, including but not limited to non-malleable proof systems and round-efficient constructions of secure multi-party computation.

A commitment scheme is a protocol between two parties, a committer \mathcal{C} and receiver \mathcal{R} , where the committer has an input message m . Both parties engage in an interactive probabilistic commitment protocol, and the receiver’s view at the end of this phase is denoted by $\text{com}(m)$. Later in an opening phase, the committer and receiver interact again to generate a transcript, that allows the receiver to verify whether the message m was actually committed to, during the commit phase. A cryptographic commitment must be binding, that is, with high probability over the randomness of the experiment, no probabilistic polynomial time committer can claim to have used a different message $m' \neq m$ in the commit phase. In short, the commitment cannot be later opened to any message $m' \neq m$. A commitment must also be hiding, that is, for any pair of messages (m, m') , the distributions $\text{com}(m)$ and $\text{com}(m')$ should be computationally indistinguishable. Very roughly, a commitment scheme is *non-malleable* if for every message m , no MIM adversary, intercepting a commitment protocol $\text{com}(m)$ and modifying every message sent during this protocol arbitrarily, is able to efficiently *generate* a commitment to a message \tilde{m} related to the original message m .

Round Complexity. The study of the round complexity of non-malleable commitments has been the subject of a vast body of research over the past 25 years. The original construction of non-malleable commitments of [DDN91] was conceptually simple, but it required logarithmically many rounds. Subsequently, Barak [Bar02], Pass [Pas04], and Pass and Rosen [PR05] constructed constant-round protocols relying on non-black box techniques. Wee [Wee10], and then Goyal [Goy11], Lin and Pass [LP] and Goyal, Lee, Ostrovsky and Visconti [GLOV12] then gave various (increasingly round-optimized) constant-round black-box constructions of non-malleable commitments assuming sub-exponentially hard one-way functions, and one-way functions respectively.

In more recent years, there has been noteworthy progress in understanding the exact amount of interaction necessary for non-malleable commitments, in the plain model. Pass [Pas13] showed an impossibility for constructing non-malleable commitments using 2 rounds of communication or less, via a black-box reduction to any “standard” polynomial intractability assumption. Goyal, Richelson, Rosen and Vald [GRRV14] constructed four round non-malleable commitments in the standard model based on the existence of one-way functions. Even more recently, Goyal, Pandey and Richelson [GPR16] constructed three round non-malleable commitments (matching the lower bound of [Pas13]) using quasi-polynomially hard injective one-way functions, by exploiting properties of non-malleable codes. Ciampi, Ostrovsky, Siniscalchi and Visconti [COSV16] showed how to bootstrap the result of [GPR16] to obtain concurrent non-malleable commitments in three rounds assuming sub-exponential one-way functions. In fact, in the sub-exponential hardness regime, Khurana and Sahai [KS17] and concurrently Lin, Pass and Soni [LPS17] showed how to achieve two-round non-malleable commitments from DDH and from time-lock puzzles, respectively. All these works use complexity leveraging and therefore must inherently rely on super-polynomial hardness. This state of affairs begs the following fundamental question:

“Can we construct round optimal non-malleable commitments from polynomial assumptions?”

We answer this question in the affirmative, by giving an explicit construction of three-round non-malleable commitments, based on polynomial hardness of any one out of the Decisional Diffie-Hellman, Quadratic Residuosity or N^{th} residuosity assumptions. We additionally assume ZAPs, which can be built from trapdoor permutations [DN07], the decisional linear assumption on bilinear maps [GOS12] or indistinguishability obfuscation together with one-way functions [BP15]. Our construction additionally satisfies concurrent (many-many) non-malleability.

Informal Theorem 1. *Assuming polynomially hard DDH or QR or N^{th} -residuosity, together with ZAPs, there exist three-round concurrent non-malleable commitments.*

Related Work. Goyal, Khurana and Sahai [GKS16] recently constructed two-round non-malleable commitments with respect to opening, secure against synchronizing adversaries, from polynomial hardness of injective one-way functions. Their result is incomparable to ours because they achieve a weaker notion of security (non-malleability with respect to opening), in two rounds, but against only synchronizing adversaries.

2 Technical Overview

We now describe the key technical roadblocks that arise in constructing non-malleable commitments from polynomial hardness, and illustrate how we overcome these hurdles.

Proving non-malleability requires arguing that the value committed by a man-in-the-middle adversary remain independent of the value committed by the honest committer. This seems to inherently require extraction (as also implicit in [Pas13]): a reduction must successfully extract the value committed by the MIM and use this value to contradict an assumption. However, current protocols for non-malleable commitments, from polynomial assumptions, in three rounds [GPR16] suffer from a problem known as over-extraction. That is, they admit an extractor that sometimes extracts a valid value from the MIM even though the MIM committed to an invalid value. Non-malleable commitments built using such extractors suffer from “selective abort” attacks: a man-in-the-middle can choose to commit to invalid values depending upon the value in the honest commitment, and an over-extracting reduction may never be able to detect such cheating.

Non-synchronizing adversaries. Let us begin by considering a simple *non-synchronizing* man-in-the-middle (MIM) adversary that interacts with an honest committer \mathcal{C} in a left session, then tries to maul this message and commit to a related message when interacting with an honest receiver \mathcal{R} in a different (right) session. By non-synchronizing, we mean that this MIM completes the entire left execution before beginning the right session. Known protocols for achieving weaker notions of non-malleability from polynomial hardness (these include the three-round sub-protocol without the ZK argument from [GRRV14] which we will denote by Π , and the basic three-round protocol from [GPR16] which we will denote by Π') do not achieve non-malleability with respect to commitment, even in this restricted setting¹.

On the other hand, *any* extractable commitment is non-malleable in this restricted setting of non-synchronizing adversaries. The reason is simple: Suppose a non-synchronizing MIM managed to successfully maul the the honest commitment. For a fixed transcript of the honest commitment, a reduction can rewind the MIM and use the extractor of the commitment scheme to extract the

¹The basic protocol from [GPR16] however, does achieve non-malleability against synchronous adversaries.

value committed by the MIM. If this value is related to the value within the honest commitment, this can directly be used to contradict hiding of the honestly generated commitment.

The main technical goal of this paper is to find a way to bootstrap the basic schemes Π, Π' to obtain non-malleability against general synchronizing and non-synchronizing adversaries, while only relying on polynomial hardness.

Barrier I: Over-Extraction. A natural starting point, then, is to add extractability to the schemes Π, Π' , by using some variant of an AoK of committed values, and within three rounds.

We cannot rely on witness indistinguishable (WI) arguments of knowledge, since arguing hiding of the scheme would require allowing a committer to commit to *two* witnesses to invoke WI security. Moreover, all existing constructions of WI arguments with black-box proofs, involve a parallel repetition of constant-soundness arguments. Now, a malicious committer could commit to two different witnesses: and use one witness in some parallel executions of the WI argument, and a different witness in some others. In this situation, even though the commitment may be invalid, one cannot guarantee that an extractor will detect the invalidity of the commitment, and over-extraction is possible. This is a known problem with 3 round protocols based on one-one one-way functions.

On the other hand, very recently, new protocols have been constructed in situations unrelated to non-malleability, that do not suffer from over-extraction [JKKR17]. Assuming polynomial hardness of DDH or Quadratic Residuosity or N^{th} residuosity, [JKKR17] demonstrated how to achieve arguments of knowledge in three rounds, that do not over-extract and have a “weak” ZK property².

However, the protocols of [JKKR17] guarantee privacy only when proving statements that are chosen from a distribution, by a prover, exclusively in the third round. On the other hand, both schemes Π, Π' , and in fact most general non-malleable commitment schemes follow a commit-challenge-response structure, where cryptography is necessarily used in the first round. Thus, the statement being proved is already fully/partially decided in the first round, which are incompatible with the kind of statements that [JKKR17] allows proofs for. Thus ideally, we would either like to inject non-malleability into the scheme of [JKKR17], or we would like to give an argument of knowledge of the message committed in the first round of Π, Π' , that doesn’t overextract. The protocols of [JKKR17] are unlikely to directly help us achieve these objectives, because of their restriction to proving messages generated in the third round. However, before describing how we solve this problem, we describe another technical barrier.

Barrier II: Composing Non-Malleability with Extraction. Many state-of-the-art protocols for non-malleable commitments admit black-box proofs of security. Naturally then, security reductions for these protocols must rely on rewinding the adversary in order to prove non-malleability. This makes these protocols notoriously hard to compose with other primitives that rely on rewinding. More specifically, it is necessary to ensure that the knowledge extractor for the extractable commitment does not interfere with the rewinding strategies used in the proof of non-malleability, and vice-versa.

A relatively straightforward technique to get around this difficulty, used in [Goy11, LP, GLOV12, GRRV14] is to arrange the protocol such that the non-malleable component and the argument of knowledge appear in completely different rounds and do not overlap. A more challenging method that does not add rounds, that is also used in prior work [GRRV14], is to use “bounded-rewinding-secure” WIAoK’s while making careful changes to the non-malleable commitment scheme.

²Very roughly, this means that for every (malicious) PPT verifier and distinguisher \mathcal{D} , there exists a distinguisher-dependent simulator $\text{Sim}_{\mathcal{D}}$, that can generate a simulated proof.

Our Solution: First Attempt. Our first technical idea is to turn the problem of incompatibility between non-malleability and arguments of knowledge on its head, and try to use the same commitments to both argue non-malleability and perform knowledge-extraction. In other words, the only extractable primitive that we rely on will be a non-malleable commitment. This is explained in more detail below.

We will use non-malleable commitments with a weak extraction property. Very roughly, we will require the existence of a probabilistic “over”-extractor \mathcal{E} , that given a PPT (synchronizing) man-in-the-middle adversary and a transcript of an execution between the MIM and honest committer, successfully “extracts” the value committed by the MIM in the transcript unless the value is invalid, without having to rewind the honest execution (except with some tunable inverse-polynomial error). The weak extraction property is satisfied, even in the one-many setting (where the MIM participates in multiple right executions) by the protocol Π .

Note that this is not an extractable commitment (and in particular, does not imply non-malleability with respect to commitment), because \mathcal{E} is allowed to output a valid value even when the MIM committed to an incorrect/invalid value in the transcript. Thus, a MIM may cheat for example, by generating a commitment to an invalid value when the honest commitment is to 0, and to a valid value when the honest commitment is to 1.

Now in order to gain confidence in the correctness of the value we extract, our scheme will have the committer generate two non-malleable commitments in parallel, and give a WI argument that one of the two was correctly constructed. This argument will satisfy a specific type of security under rewinding, and can be constructed based on ZAPs and DDH in three rounds via [JKKR17]. For the purposes of this overview, even though we don’t actually require a non-interactive proof, assume that we use a non-interactive witness indistinguishable proof, NIWI [BOV07, GOS12]. Let ϕ_1 denote the protocol that results from committing to the message twice using the non-malleable commitment scheme Π , and giving a NIWI proof that one of the two was correctly computed.

This partial solution still leaves scope for over-extraction: how can we be sure that the extractor does not output any valid value even when a malicious committer could be committing to two different values within the non-malleable commitments and using both witnesses for the WI?

Second Attempt. Since protocol ϕ_1 also suffers from over-extraction, it may seem like we made no progress at all. However, note that the same protocol can be easily modified to a WIAoK (witness indistinguishable argument of knowledge): by committing to a witness twice using Π and proving via NIWI that one of the two non-malleable commitments is a valid commitment to a witness. Let us call the resulting protocol ϕ_2 . At a high level, the protocol ϕ_2 has the following properties:

- **Knowledge Extraction.** ϕ_2 is an argument of knowledge (which suffers from over-extraction).
- **Non-Malleability.** Weak non-malleability of Π implies a limited form of non-malleability of the protocol ϕ_2 .

Third Attempt. In order to prevent over-extraction, we will need to force any prover that generates a proof according to ϕ_2 to use a *unique* witness in ϕ_2 . We will now try to rely on three round “weak” zero-knowledge arguments of [JKKR17], which are secure when used to prove cryptographic statements chosen by the prover in the last round. These arguments in fact, also retain a limited type of security under rewinding, which will ensure that the simulation doesn’t interfere with extraction from the non-malleable commitment.

Assume again, for the purposes of this overview, that these arguments satisfy the standard notion of simulation for zero-knowledge, except that the statement to be proved, must be chosen in the last round. Let us denote them by **wzk**.

We will now use wzk to set up a trapdoor for ϕ_2 . This trapdoor will include a statistically binding commitment c_1 using a non-interactive statistically binding commitment scheme com , and a wzk argument that c_1 was generated correctly as a commitment to 1. The trapdoor statement will be that c_1 is a commitment to 0. This trapdoor statement will serve as the ‘other’ witness for ϕ_2 .

Given these building blocks, our actual commitment scheme will have the following structure:

- **Trapdoor:** The committer will generate commitment c_1 to 1, via com in the third round. In parallel, the committer will prove via wzk , that c_1 was correctly generated as a commitment to 1.
- **Actual Commitment:** The committer will also generate commitment c to input message m , via com , only in the third round. In parallel the committer will also run scheme ϕ_2 , proving that either c was correctly generated, or that c_1 was generated as a commitment to 0.

Note that the protocol ϕ_2 as described is not delayed-input: the non-malleable commitment Π requires an input (that is, the witness) in the first round, whereas the witness for the statement is only decided in the third round. However, we can just use one-time pads to get this delayed-input property from ϕ_2 .

A simple (informal) description that captures the essence of our final protocol, ϕ , is in [Figure 1](#). The scheme ϕ is opened up into its components: two non-malleable commitments and a WI argument. This scheme can be shown to be computationally hiding by the privacy properties of ϕ , wzk and com .

Extraction. We first argue that the scheme in [Figure 1](#) is an extractable commitment. We already discussed that there exists a knowledge extractor for ϕ_2 that extracts at least one out of γ_1, γ_2 : which can then be used to extract the randomness r via z_1, z_2 . All we need to argue is that this extractor does not over-extract. However, soundness of wzk already forces a computational committer to set c_1 as a commitment to 1, which means that there remains only one randomness (the randomness used for committing to m), that the committer can use in order to generate z_1 or z_2 in the WI. Extractability of this scheme is already enough to guarantee security against non-synchronizing adversaries, even if such adversaries simultaneously participate in many parallel executions.

Non-malleability. We must also argue that non-malleability is preserved, and in fact, the resulting scheme is concurrent non-malleable with respect to commitment, when instantiated with Π from [\[GRRV14\]](#).

While arguing non-malleability, some subtle technical issues arise that require careful analysis. For instance, the distinguisher-dependent simulation strategy of weak ZK if used naively, only guarantees that the view of the distinguisher remains indistinguishable under simulation. However, while arguing non-malleability, it is imperative to ensure that not just the view, but the joint distribution of the *view and the value committed* by the MIM remains indistinguishable under simulation. It is here that the over-extraction property of Π helps: in hybrids where we must argue non-malleability while also performing distinguisher-dependent simulation, we will use the extractor that is guaranteed by the weak non-malleability of Π , to extract the value committed by the MIM *without* having to rewind the left non-malleable commitment. This helps us guarantee that the joint distribution of the view and values committed by the MIM remains indistinguishable under simulation.

Our actual protocol is formalized in [Section 5](#) and is identical to the protocol described above, except the following modification: For technical reasons, in our actual protocol, instead of masking the randomness r' with γ , we mask it with $\text{PRF}(\gamma, \alpha)$ for randomly chosen α . The committer must

Inputs: Committer \mathcal{C} has input a message $m \in \{0, 1\}^n$, receiver \mathcal{R} has no input.

1.
 - \mathcal{C} samples $\gamma_1, \gamma_2 \xleftarrow{\$} \{0, 1\}^n$.
 - Next, \mathcal{C} sends the first message of wzk to \mathcal{R} .
 - Finally, \mathcal{C} sends the first message of $\Pi(\gamma_1), \Pi(\gamma_2)$.
2.
 - \mathcal{R} sends the second message of wzk to \mathcal{C} .
 - \mathcal{R} sends the second message for both executions of Π .
3.
 - \mathcal{C} computes and sends $c_1 = \text{com}(1; r)$ for $r \xleftarrow{\$} \{0, 1\}^n$.
 - \mathcal{C} sends the third message of wzk to \mathcal{R} , proving that c_1 is a commitment to 1.
 - \mathcal{C} computes and sends $c = \text{com}(m; r')$ for $r' \xleftarrow{\$} \{0, 1\}^n$.
 - \mathcal{C} sends the third message of $\Pi(\gamma_1), \Pi(\gamma_2)$.
 - \mathcal{C} sends $z_1 = (\gamma_1 \oplus r'), z_2 = (\gamma_2 \oplus r')$ to \mathcal{R} .
 - \mathcal{C} uses $(c, m, r', \gamma_1, z_1)$ as witness to prove using the WI that :
 - c is a valid commitment to some message m with randomness r' , and $\Pi(\gamma_1)$ is a valid non-malleable commitment to γ_1 and $z_1 = \gamma_1 \oplus r'$, OR
 - c is a valid commitment to some message m with randomness r' , and $\Pi(\gamma_2)$ is a valid non-malleable commitment to γ_2 and $z_2 = \gamma_2 \oplus r'$, OR
 - c_1 is a valid commitment to 0 with randomness r , and $\Pi(\gamma_1)$ is a valid non-malleable commitment to γ_1 and $z_1 = \gamma_1 \oplus r$, OR
 - c_1 is a valid commitment to 0 with randomness r , and $\Pi(\gamma_2)$ is a valid non-malleable commitment to γ_2 and $z_2 = \gamma_2 \oplus r$

Figure 1: A simplified description of the final non-malleable commitment scheme ϕ

also send α to the receiver. This is for similar reasons as [JKKR17]: the simulator for wzk sends *many* third protocol messages for the same fixed transcript of the first two messages, and we require security to hold even in this setting.

On Rewinding Techniques in the Proof. The weak ZK protocol of [JKKR17] that we use in this work, relies on the simulator rewinding the distinguisher. Because of this, our actual proof of security sometimes has two sequential rewindings happening within a three round protocol: one which rewinds to the end of the first round, and helps extract values committed in the MIM executions, and the second that rewinds (the distinguisher) to the end of the second round, in order to simulate the proof while maintaining an indistinguishable joint distribution of view and values with respect to a given distinguisher. This requires careful indistinguishability arguments that take such sequential rewindings into account, and can also be found in Section 5.

2.1 Organization

The rest of this paper is organized as follows. In Section 3, we will recall the preliminaries that will be of use in our constructions. In Section 4, we give definitions of non-malleable commitments. In Section 5, we describe our construction and provide a proof of non-malleability.

3 Preliminaries

In this section, we recall some preliminaries from [JKKR17] that will be useful in our constructions.

Definition 1 (Non-adaptive Distributional ϵ -Weak Zero Knowledge). *A delayed-input interactive argument (P, V) for a language L is said to be distributional ϵ -weak zero knowledge against non-adaptive verifiers if for every efficiently samplable distribution $(\mathcal{X}_n, \mathcal{W}_n)$ on R_L , i.e., $\text{Supp}(\mathcal{X}_n, \mathcal{W}_n) = \{(x, w) : x \in L \cap \{0, 1\}^n, w \in R_L(x)\}$, every non-adaptive PPT verifier V^* , every $z \in \{0, 1\}^*$, every PPT distinguisher \mathcal{D} , and every $\epsilon = 1/\text{poly}(n)$, there exists a simulator \mathcal{S} that runs in time $\text{poly}(n, \epsilon)$ such that:*

$$\left| \Pr_{(x,w) \leftarrow (\mathcal{X}_n, \mathcal{W}_n)} [\mathcal{D}(x, z, \text{view}_{V^*}[\langle P, V^*(z) \rangle](x, w)) = 1] - \Pr_{(x,w) \leftarrow (\mathcal{X}_n, \mathcal{W}_n)} [\mathcal{D}(x, z, \mathcal{S}^{V^*, D}(x, z)) = 1] \right| \leq \epsilon(n),$$

where the probability is over the random choices of (x, w) as well as the random coins of the parties.

Definition 2 (Weak Resetable Non-adaptive Distributional ϵ -Weak Zero Knowledge). *A three round delayed-input interactive argument (P, V) for a language L is said to be weak resetable distributional weak zero-knowledge, if for every efficiently samplable distribution $(\mathcal{X}_n, \mathcal{W}_n)$ on R_L , i.e., $\text{Supp}(\mathcal{X}_n, \mathcal{W}_n) = \{(x, w) : x \in L \cap \{0, 1\}^n, w \in R_L(x)\}$, every non-adaptive PPT verifier V^* , every $z \in \{0, 1\}^*$, every PPT distinguisher \mathcal{D} , and every $\epsilon = 1/\text{poly}(n)$, there exists a simulator \mathcal{S} that runs in time $\text{poly}(n, \epsilon)$ and generates a simulated proof for instance $x \stackrel{\$}{\leftarrow} \mathcal{X}_n$, such that over the randomness of sampling $(x, w) \leftarrow (\mathcal{X}_n, \mathcal{W}_n)$, $\Pr[b' = b] \leq \frac{1}{2} + \epsilon(n) + \text{negl}(n)$ in the following experiment, where the challenger C plays the role of the prover:*

- At the beginning, (C, V^*) receive the size of the instance, V^* receives auxiliary input z , and they execute the first 2 rounds. Let us denote these messages by τ_1, τ_2 .
- Next, (C, V^*) run $\text{poly}(n)$ executions, with the same fixed first message τ_1 , but different second messages chosen potentially maliciously by V^* . In each execution, C picks a fresh sample $(x, w) \leftarrow (\mathcal{X}_n, \mathcal{W}_n)$, and generates a proof for it according to honest verifier strategy.
- Next, C samples bit $b \stackrel{\$}{\leftarrow} \{0, 1\}$ and if $b = 0$, for $(x, w) \leftarrow (\mathcal{X}_n, \mathcal{W}_n)$ it generates an honest proof with first two messages τ_1, τ_2 , else if $b = 1$, for $x \stackrel{\$}{\leftarrow} \mathcal{X}_n$ it generates a simulated proof with first two messages τ_1, τ_2 using simulator \mathcal{S} that has oracle access to V^*, \mathcal{D} .
- Finally, V^* sends its view to a distinguisher \mathcal{D} that outputs b .

Remark 1. *The three message protocols in [JKKR17], based on DDH/QR/ N^{th} residuosity, along with ZAPs, satisfy weak resetable distributional ϵ -weak zero knowledge/strong WI against non-adaptive verifiers (refer to Remark 2 in [JKKR17], and Appendix A). In our protocols, we will always use weak zero-knowledge/strong witness-indistinguishable arguments in the “non-adaptive/delayed-input” setting, that is, to prove statements that are chosen by the prover only in the third round of the execution.*

Definition 3 (Resetable Reusable WI Argument). *We say that a two-message delayed-input interactive argument (P, V) for a language L is resetable reusable witness indistinguishable, if for every PPT verifier V^* , every $z \in \{0, 1\}^*$, $\Pr[b = b'] \leq \frac{1}{2} + \text{negl}(n)$ in the following experiment, where we*

denote the first round message function by $m_1 = \text{wi}_1(r_1)$ and the second round message function by $\text{wi}_2(x, w, m_1, r_2)$.

The challenger samples $b \xleftarrow{\$} \{0, 1\}$. V^* (with auxiliary input z) specifies $(m_1^1, x^1, w_1^1, w_2^1)$ where w_1^1, w_2^1 are (not necessarily distinct) witnesses for x^1 . V^* then obtains second round message $\text{wi}_2(x^1, w_b^1, m_1^1, r)$ generated with uniform randomness r . Next, the adversary specifies arbitrary $(m_1^2, x^2, w_1^2, w_2^2)$, and obtains second round message $\text{wi}_2(x^2, w_b^2, m_1^2, r)$. This continues $m(n) = \text{poly}(n)$ times for a-priori unbounded m , and finally V^* outputs b .

Remark 2. Note that ZAPs (more generally, any two-message WI) can be modified to obtain resettable reusable WI, by having the prover apply a PRF on the verifier message and the instance to compute randomness for the proof. This allows to argue, via a hybrid argument, that fresh randomness can be used for each proof, and therefore perform a hybrid argument so that each proof remains WI. In our construction, we will use resettable reusable ZAPs.

4 Definitions

4.1 Non-Malleability w.r.t. Commitment

Throughout this paper, we will use n to denote the security parameter, and $\text{negl}(n)$ to denote any function that is asymptotically smaller than $\frac{1}{\text{poly}(n)}$ for any polynomial $\text{poly}(\cdot)$. We will use PPT to describe a probabilistic polynomial time machine. We will also use the words “rounds” and “messages” interchangeably.

We follow the definition of non-malleable commitments introduced by Pass and Rosen [PR05] and further refined by Lin et al [LPV] and Goyal [Goy11] (which in turn build on the original definition of [DDN91]). In the real interaction, there is a man-in-the-middle adversary MIM interacting with a committer \mathcal{C} (where \mathcal{C} commits to value v) in the left session, and interacting with receiver \mathcal{R} in the right session. Prior to the interaction, the value v is given to \mathcal{C} as local input. MIM receives an auxiliary input z , which might contain a-priori information about v . Let $\text{MIM}_{\langle \mathcal{C}, \mathcal{R} \rangle}(\text{value}, z)$ denote a random variable that describes the value $\widetilde{\text{val}}$ committed by the MIM in the right session, jointly with the view of the MIM in the full experiment. In the simulated experiment, a simulator \mathcal{S} directly interacts with \mathcal{R} . Let $\text{Sim}_{\langle \mathcal{C}, \mathcal{R} \rangle}(1^n, z)$ denote the random variable describing the value $\widetilde{\text{val}}$ committed to by \mathcal{S} and the output view of \mathcal{S} . If the tags in the left and right interaction are equal, the value $\widetilde{\text{val}}$ committed in the right interaction, is defined to be \perp in both experiments.

Definition 4 (Non-malleable Commitments w.r.t. Commitment). *A commitment scheme $\langle \mathcal{C}, \mathcal{R} \rangle$ is said to be non-malleable if for every PPT MIM, there exists an expected PPT simulator \mathcal{S} such that the following ensembles are computationally indistinguishable:*

$$\{\text{MIM}_{\langle \mathcal{C}, \mathcal{R} \rangle}(\text{value}, z)\}_{n \in \mathbb{N}, v \in \{0, 1\}^n, z \in \{0, 1\}^*} \text{ and } \{\text{Sim}_{\langle \mathcal{C}, \mathcal{R} \rangle}(1^n, z)\}_{n \in \mathbb{N}, v \in \{0, 1\}^n, z \in \{0, 1\}^*}$$

4.2 Concurrent Non-Malleable Commitments

This setting considers an adversary that participates in multiple sessions with an honest committer, acting as receiver. The adversary simultaneously participates in multiple sessions with an honest receiver, acting as committer. In the left sessions, the MIM interacts with honest committer(s) obtaining commitments to values $m_1, m_2, \dots, m_{\text{poly}(n)}$ (say, from distribution val using tags $t_1, t_2, t_{\text{poly}(n)}$) of its choice. In the right session, \mathcal{A} interacts with \mathcal{R} attempting to commit to a sequence of related values $\tilde{m}_1, \dots, \tilde{m}_{\text{poly}(n)}$ again using identities $\tilde{t}_1, \dots, \tilde{t}_{\text{poly}(n)}$. If any of the right commitments are invalid, or undefined, their value is set to \perp . For any i such that $\tilde{t}_i = t_j$ for some j , set \tilde{m}_i (the value

committed using that tag) to \perp . Let $\text{MIM}_{\langle C, R \rangle}(\text{value}, z)$ denote a random variable that describes the values $\widetilde{\text{val}}$ committed by the MIM in the right sessions, jointly with the view of the MIM in the full experiment, when the value is the joint distribution of values committed in the left sessions. In a simulated execution, there is an expected polynomial time simulator that interacts with the MIM and generates a distribution Sim consisting of the views and values committed by the MIM. Then, the definitions of concurrent non-malleable commitment scheme w.r.t. commitment, replacement and opening are defined as above.

Definition 5 (Concurrent Non-malleable Commitments w.r.t. Commitment). *A commitment scheme $\langle C, R \rangle$ is said to be concurrently non-malleable if for every PPT MIM, there exists an expected PPT simulator \mathcal{S} such that the ensembles real and sim defined above are indistinguishable.*

Remark 3. *We will also consider the notion of non-malleability only against synchronizing adversaries. [GPR16] give a construction of a three-round non-malleable commitment scheme secure against synchronizing adversaries, based on polynomially hard injective one-way functions.*

Definition 6 (One-Many Weak Non-Malleable Commitments with respect to Synchronizing Adversaries). *A statistically binding commitment scheme $\langle C, R \rangle$ is said to be one-many weak non-malleable with respect to synchronizing adversaries, if there exists a probabilistic “over”-extractor \mathcal{E} parameterized by ϵ , that given a PPT synchronizing MIM which participates in one left session and $p = \text{poly}(n)$ right sessions, and given the transcript of a main-thread interaction τ , outputs a set of values v_1, v_2, \dots, v_p in time $\text{poly}(n, \frac{1}{\epsilon})$. These values are such that:*

- *For all $j \in [p]$, if the j^{th} commitment in τ is a commitment to a valid message m_j , then $v_j = m_j$ over the randomness of the extractor and the transcript, except with probability $\frac{\epsilon}{p}$.*
- *For all $j \in [p]$, if the j^{th} commitment in τ is a commitment to some invalid message (which we will denote by \perp), then v_j need not necessarily be \perp .*

Remark 4. *By the union bound, it is easy to see that by appropriately scaling to $\epsilon' = \frac{\epsilon}{p(n)}$, the values output by the extractor are correct for all sessions in which the MIM committed to valid messages in the transcript τ , except with probability ϵ .*

This formalization helps us to abstract out the exact properties satisfied by existing three-round schemes based on polynomial assumptions, which we can rely on for our bootstrapping protocol. We note that this is an alternative way of formalizing the requirement of “security against non-aborting adversaries” from [COSV17]. When invoking the security of non-malleable commitments in our proof, the adversary will always be forced (via appropriate proofs) to behave in a non-aborting way.

Instantiating one-many weak non-malleable commitments. The three-round sub-protocol in the non-malleable commitment scheme from [GRRV14] (their basic construction without the zero-knowledge argument of knowledge), based on one-way functions, is a one-many weak non-malleable commitment according to Definition 6. In fact, their proof of non-malleability proceeds via constructing such an extractor.

5 Non-Malleable Commitments w.r.t. Commitment

In this section, we describe a round-preserving way to transform (one-many) non-malleable commitments with respect to replacement to (one-many) non-malleable commitments with respect to commitment additionally assuming polynomial hardness of DDH and ZAPs.

5.1 Construction

Our construction of three round non-malleable commitments is described in Figure 2.

Let $\Pi^i = (\text{nmc}_1^i, \text{nmc}_2^i, \text{nmc}_3^i)$ for $i \in \{1, 2\}$ denote the three messages of, two independent instances (indexed by i) of a weak non-malleable commitment.

Let $\text{wi} = (\text{wi}_1, \text{wi}_2)$ denote the two messages of a reusable resettable ZAP for delayed-input statements. This can be instantiated via any ZAP, where the prover uses a PRF on verifier message and instance, to compute randomness for the proof.

Let $\text{wzk} = (\text{wzk}_1, \text{wzk}_2, \text{wzk}_3)$ denote the three messages of a weak resettable weak distributional ZK for delayed-input statements, against non-adaptive verifiers.

Let $\text{PRF}(K, r)$ denote the output of a pseudorandom function on key K and input r .

Let $\text{com}(\cdot)$ denote a non-interactive, statistically binding commitment scheme.

Tag: Let the tag for the interaction be $\text{tag} \in [n]$. Let n denote the security parameter.

Committer Input: A message $m \in \{0, 1\}^*$, along with tag tag .

1. **Committer Message:** Sample independent randomness $r_1, r_2, \gamma_1, \gamma_2$, and send $\text{nmc}_1^1(\gamma_1, r_1, \text{tag}), \text{nmc}_1^2(\gamma_2, r_2, \text{tag})$ together with wzk_1 .

2. **Receiver Message:** Send the second message for both non-malleable commitments $(\text{nmc}_2^1, \text{nmc}_2^2)$ for tag , to the prover together with $\text{wi}_1, \text{wzk}_2$.

3. **Committer Message:** Sample $r \xleftarrow{\$} \{0, 1\}^*$ and send $c = \text{com}(m; r)$ to \mathcal{R} .

Additionally, sample $\hat{r} \xleftarrow{\$} \{0, 1\}$ and send $c_1 = \text{com}(1; \hat{r})$. Along with c_1 , send wzk_3 proving that $\exists \hat{r}$ such that $c_1 = \text{com}(1; \hat{r})$.

Send $\text{nmc}_3^1(\gamma_1, r_1, \text{tag})$ and $\text{nmc}_3^2(\gamma_2, r_2, \text{tag})$ to \mathcal{R} .

Finally, sample $\{\alpha_1, \alpha_2\} \xleftarrow{\$} \{0, 1\}^{2n}$ and send $\delta_1 = \text{PRF}(\gamma_1, \alpha_1) \oplus r$ and $\delta_2 = \text{PRF}(\gamma_2, \alpha_2) \oplus r$. Send wi_2 proving (using witness Π^1) that:

- Either Π^1 is a valid non-malleable commitment to some γ_1 with randomness r_1 AND $r = \text{PRF}(\gamma_1, \alpha_1) \oplus \delta_1$ such that $(c = \text{com}(m; r) \text{ OR } c_1 = \text{com}(0; r))$
- Or, Π^2 is a valid non-malleable commitment to some γ_2 with randomness r_2 AND $r = \text{PRF}(\gamma_2, \alpha_2) \oplus \delta_2$ such that $(c = \text{com}(m; r) \text{ OR } c_1 = \text{com}(0; r))$

4. **Decommitment Phase:** The committer reveals the message m and randomness r . The verifier accepts if and only if c is a commitment to m using randomness r .

Figure 2: Non-Malleable Commitment Scheme ϕ

5.2 Proof of Security

We begin by proving that the scheme is statistically binding and computationally hiding. We note that computational hiding is in fact, implied by non-malleability: therefore as a warm up, we sketch the proof of hiding via a sequence of hybrid experiments without giving formal reductions. In Theorem 1, we prove formally that not only is the view of a receiver indistinguishable between

these hybrids, in fact, the joint distribution of the view *and values committed* by a MIM interacting with an honest committer remains indistinguishable between these hybrids.

Lemma 1. *The protocol in Figure 2 is a statistically binding, computationally hiding, commitment scheme.*

Proof. (Sketch) The statistical binding property follows directly from statistical hiding property of the underlying commitment scheme $\text{com}(\cdot)$.

The computational hiding property follows from the hiding of com and nmc , the weak zero-knowledge property of wzk , and the witness indistinguishability of wi . Here, we sketch a proof of computational hiding. Note that computational hiding is implied by non-malleability, therefore the proof of Theorem 1 can also be treated as a formal proof of hiding of the commitment scheme ϕ . Let $\langle \mathcal{C}_\phi(m; r), \mathcal{R} \rangle$ denote the transcript of an execution where the committer uses input message m and randomness R . We prove that $\langle \mathcal{C}_\phi(m_0; r), \mathcal{R} \rangle \approx_c \langle \mathcal{C}_\phi(m_1; r), \mathcal{R} \rangle$ for all m_0, m_1 , via the following sequence of hybrid experiments:

Hybrid $_{m_0}$: This hybrid corresponds to an interaction of \mathcal{C} and \mathcal{R} where \mathcal{C} uses input message m_0 , that is, $\langle \mathcal{C}_\phi(m_0; r), \mathcal{R} \rangle$.

Hybrid $_1$: In this hybrid, the challenger behaves identically to **Hybrid $_{m_0}$** , except that it generates nmc^2 as a non-malleable commitment to a different randomness γ'_2 than the (uniform) randomness γ_2 used to compute δ_2 . This hybrid is indistinguishable from **Hybrid $_0$** directly by the hiding of Π .

Hybrid $_{2, \mathcal{D}}$: In this hybrid, the challenger behaves identically to **Hybrid $_1$** , except that it outputs the transcript of an execution where the wzk argument is simulated³. Note that the challenger uses the simulation strategy of the weak zero-knowledge argument wzk , which executes the last message of the protocol multiple times, and learns the distinguisher's challenge to wzk . Each time, the simulation strategy samples fresh α_1, α_2 at random, and furthermore, learns by generating commitments to both m_0 and m_1 . However, the main transcript that is output still contains a commitment to m_0 , and is in fact identical to **Hybrid $_1$** except that it contains a simulated wzk argument. By the simulation security of wzk , for any distinguisher \mathcal{D} , there exists a distinguisher-dependent simulator/challenger such that **Hybrid $_{2, \mathcal{D}}$** is indistinguishable from **Hybrid $_1$** .

Hybrid $_{3, \mathcal{D}}$: In this hybrid, the challenger behaves identically to **Hybrid $_{2, \mathcal{D}}$** , except that it sets $c_1 = \text{com}(0; \hat{r})$ for some randomness \hat{r} , in the main output transcript. Note that this is possible because the challenger is generating a simulated proof in the output transcript. This hybrid is indistinguishable from **Hybrid $_{2, \mathcal{D}}$** by the computational hiding property of com .

Hybrid $_{4, \mathcal{D}}$: In this hybrid, the challenger behaves identically to **Hybrid $_{3, \mathcal{D}}$** except that in the output transcript, it sets $\delta_2 = \text{PRF}(\gamma_2, \alpha_2) \oplus \hat{r}$ where \hat{r} is the randomness used to generate $c_1 = \text{com}(0; \hat{r})$. Note that the committer is committing to a different value γ'_2 in the protocol Π^2 , thus the key γ_2 does not appear anywhere in the rest of the protocol. Therefore, this hybrid is indistinguishable from **Hybrid $_{3, \mathcal{D}}$** by the security of the PRF.

Hybrid $_{5, \mathcal{D}}$: In this hybrid, the challenger behaves identically to **Hybrid $_{4, \mathcal{D}}$** except that in all transcripts, it sets nmc^2 as a non-malleable commitment to the same randomness γ'_2 that is used to

³Note that in all hybrid experiments, we will actually use the extended simulation strategy of the weak ZK argument wzk – that is used for strong witness indistinguishability, and where the simulator takes into account both messages m_0 and m_1 during simulation.

compute δ_2 . This hybrid essentially “reverts” the cheating performed in **Hybrid**₁. Indistinguishability of this hybrid follows because of the hiding of Π^2 .

Note that the transcript output by the challenger in this experiment is such that Π^1 is a valid non-malleable commitment to γ_1 with randomness r_1 AND $r = \text{PRF}(\gamma_1, \alpha_1) \oplus \delta_1$ such that $c = \text{com}(m; r)$. Additionally, Π^2 is a valid non-malleable commitment to γ_2 with randomness r_2 AND $\hat{r} = \text{PRF}(\gamma_2, \alpha_2) \oplus \delta_2$ such that $c_1 = \text{com}(0; \hat{r})$.

Hybrid_{6, \mathcal{D}} : In this hybrid, the challenger behaves the same way as **Hybrid**_{5, \mathcal{D}} , except that it uses the second witness, (r_2, γ_2) , to generate the witness-indistinguishable proof wi in the output transcript. This hybrid is indistinguishable from **Hybrid**_{5, \mathcal{D}} by the reusable witness-indistinguishability of wi , that is, witness indistinguishability in the setting where multiple proofs are provided for different statements, using the same second message transcript.

Hybrid_{7, \mathcal{D}} : In this hybrid, the challenger behaves the same way as **Hybrid**_{6, \mathcal{D}} , except that it uses the second witness, r_2, γ_2 , to generate the witness-indistinguishable arguments wi all the “learning executions” of the simulation strategy, as well as in the output transcript. That is, in every message that the challenger sends, it uses the second witness instead of the first. This hybrid is indistinguishable from **Hybrid**_{6, \mathcal{D}} by the reusable witness-indistinguishability of wi .

Hybrid_{8, \mathcal{D}} : In this hybrid, the challenger behaves the same way as **Hybrid**_{7, \mathcal{D}} , except that in all transcripts, it sets nmc^1 as a non-malleable commitment to a *different* randomness γ'_1 than the one used to compute δ_1 . The view of a malicious receiver in this hybrid is indistinguishable from **Hybrid**_{7, \mathcal{D}} by the hiding of the non-malleable commitment Π^1 .

Hybrid_{9, \mathcal{D}} : In this hybrid, the challenger behaves the same way as **Hybrid**_{8, \mathcal{D}} , except that in the output transcript, it sets $\delta_1 \stackrel{\$}{\leftarrow} \{0, 1\}^*$, instead of setting $\delta_1 = \text{PRF}(\gamma_1, \alpha_1) \oplus r$. Note that the committer is committing to a different value γ'_1 in the protocol Π^1 , thus the key γ_1 does not appear in the rest of the protocol. Therefore, this hybrid is indistinguishable from **Hybrid**_{8, \mathcal{D}} by the security of the PRF.

Hybrid_{10, \mathcal{D}} : In this hybrid, the challenger behaves the same way as **Hybrid**_{10, \mathcal{D}} except that it replaces $\text{com}(m_0; r)$ with $\text{com}(m_1; r)$ in the output transcript. Note that at this point, r is not used anywhere else in the protocol, and hence the commitment can be obtained externally. This hybrid is indistinguishable from **Hybrid**_{9, \mathcal{D}} by computational hiding of the non-interactive commitment scheme com .

At this point, we have successfully indistinguishably switched to an experiment where the commitment is generated to message m_1 instead of m_0 in the main transcript output by the challenger. Computational hiding follows by repeating the above hybrids in reverse order, until in **Hybrid** _{m_1} , the challenger generates an honest commitment to message m_1 . \square

Lemma 2. *There exists a polynomial-time extractor that extracts the value committed by any adversary \mathcal{A} in any accepting transcript, with probability $1 - \text{negl}(n)$.*

Proof. For any accepting commitment transcript generated by a committer, because of adaptive soundness of wi , the i^{th} extractable commitment is generated as a valid extractable commitment to randomness r_i , such that $\text{PRF}(r_i, a_i) \oplus x_i$ yields a valid witness for wi , for some $i \in \{1, 2\}$. Furthermore, by soundness of wzk , c_1 is a commitment to 1, and by statistical binding of com , c_1 cannot be a commitment to 0. Thus, the only possible valid witness in wi , with overwhelming probability, must necessarily be a witness for c , which is the actual commitment to the message.

We now argue that this witness can be extracted by a polynomial time extractor. This follows roughly because of the (over)-extraction property of Π and the soundness of w_i , similar to [JKKR17]. Specifically, we consider a committer that generates an accepting transcript with probability $p > \frac{1}{\text{poly}(n)}$ for some polynomial $\text{poly}(\cdot)$. Then, within $\frac{n}{p}$ rewindings, such a committer generates an expected n accepting transcripts and, with overwhelming probability at least \sqrt{n} of the accepting transcripts (in the rewinding thread) produce a valid commitment using scheme Π for the same index i as the main thread, allowing for extraction from the over-extracting commitment. The extracted value can be used to compute r , that checked against c to ensure that r is the correct randomness that was used to compute c . We note that this commitment does not suffer from over-extraction, since by the soundness of wz and w_i , a malicious committer is always forced to use the unique witness corresponding to the commitment c . Furthermore, such an extractor extracts with error at most ϵ by running in time $\text{poly}(1/\epsilon)$. \square

Next, we directly prove concurrent non-malleability of the resulting scheme when instantiated with the basic protocol Π from [GRRV14]. We note that the scheme can also be instantiated with the protocol from [GPR16], yielding one-one non-malleability.

Theorem 1. *The protocol ϕ in Figure 2, when instantiated with the one-many weak non-malleable commitment Π from [GRRV14], is a concurrent non-malleable commitment with respect to commitment according to Definition 5.*

Proof. The proof of non-malleability against non-synchronizing adversaries, that complete the left execution before beginning a right (malicious) execution, follows directly because ϕ is an extractable commitment. In other words, given a non-synchronizing MIM adversary, there exists a reduction that runs an extractor to extract the value committed by the MIM from the right execution(s) by rewinding the adversary, and uses the view jointly with the values extracted from such a malleating adversary to directly break hiding of the commitment in the left execution. This is possible because of the non-synchronizing scheduling, it is possible to rewind the MIM's commitment and run the extractor of Lemma 2 without rewinding the honest commitment at all. This leads to a contradiction, ruling out the existence of any PPT MIM adversary that successfully mauls the honest commitment in the non-synchronizing setting.

Therefore, it only remains to argue non-malleability in the fully synchronizing setting (these arguments directly combine to argue security against adversaries that are synchronizing in some executions and non-synchronizing in others). We first note that it suffices to argue non-malleability against one-many adversaries, that participate in one left session and polynomially-many right sessions. By [LPV], security against such adversaries already implies concurrent non-malleability. Suppose the MIM opens $p = \text{poly}(n)$ sessions on the right, for some polynomial $p(\cdot)$.

We argue non-malleability against synchronizing adversaries via a sequence of hybrid experiments, relying on the non-malleability of Π , along with various properties of other primitives used in the protocol. These hybrids are all parameterized by an inverse polynomial error parameter ϵ , and sometimes require the challenger to run in time $\text{poly}(n, \frac{1}{\epsilon})$. Later, we will set ϵ to be significantly smaller than the advantage of any distinguisher between $\text{MIM}_{\langle C, R \rangle}(V_1, z)$ and $\text{MIM}_{\langle C, R \rangle}(V_2, z)$ (but ϵ will still be less than $\frac{1}{\text{poly}(\cdot)}$), thereby proving the lemma. We will use $\tilde{\phi}$ to denote message ϕ sent in the right execution, and a message ϕ sent during the left execution will just be denoted by ϕ .

5.2.1 Overview of Hybrid Experiments

Before describing the hybrid arguments in detail, we provide an overview. The sequence of experiments follows the same pattern as the proof of hiding, except that we now argue about the joint distribution of the view and values committed by the MIM. Moreover, when we say that the challenger rewinds and generates lookahead threads to learn γ or for the simulation of weak ZK, the challenger always generates multiple lookahead threads where some commit to value V_1 and some to V_2 (this is possible because the message is decided only in the last round), and takes the union of the values extracted using both V_1, V_2 .

In all these hybrids, the challenger will never generate simulated **wzk** proofs in any rewinding execution. The **wzk** proof will be carefully simulated only in the main transcript (in some of the hybrids). Thus, by soundness of the **wi**, the MIM will always commit to the witness for the commitment, by correctly generating a non-malleable commitment to at least one of the γ values, in any rewinding execution. Therefore, a rewinding extractor will correctly extract at least one γ value committed by the MIM, with high probability. Furthermore, when relying on the extractor of the non-malleable commitment scheme, we will again generate a transcript for the extractor that does not contain any simulated proofs – therefore, this extractor is guaranteed to correctly extract at least one of the γ values committed by the MIM.

The output of the first experiment, **Hybrid** $_{V_1}$ corresponds to the joint distribution of the view and values committed by the MIM on input an honest commitment to value V_1 .

Hybrid $_1$: In the first hybrid, the challenger changes the left execution by first sampling (γ_2, γ'_2) independently and uniformly at random. The value committed using the second non-malleable commitment Π^2 is γ'_2 , while the third message $\delta_2 = \text{PRF}(\gamma_2, \alpha_2) \oplus r$ is computed honestly using a different γ_2 . At this point, we invoke soundness of the **wi** and **wzk** to argue that the MIM must commit to at least one valid $\tilde{\gamma}_i^1$ or $\tilde{\gamma}_i^2$ in the main execution, for every $i \in [p(n)]$. Therefore, we can invoke the extractor for Π^2 , to extract the joint distribution of the values committed by the man-in-the-middle (MIM) in all right executions.

By the property of the non-malleable commitment, when the MIM commits to a valid value in the main execution, such an extractor will successfully extract at least one of the committed values $\tilde{\gamma}_i^1$ or $\tilde{\gamma}_i^2$ from the i^{th} right interaction, for all $i \in [p(n)]$. Because of soundness of **wi** and **wzk**, this extracted value will directly help recover the message committed by the MIM in this interaction. Since this extractor operates *without* rewinding the left execution, if the joint distribution of the view and values changes from **Hybrid** $_0$ to **Hybrid** $_1$, we obtain a contradiction to the hiding of Π .

Hybrid $_2$: In the next hybrid, the challenger modifies the left execution by generating an output view where the left execution contains a simulated weak ZK proof. When applied naively, the simulation guarantee is that the view of the MIM remains indistinguishable when provided a transcript with a simulated proof. However, there are no guarantees about the values committed by the MIM.

In order to ensure that the joint distribution of committed values remains indistinguishable, we modify the input to the distinguisher-dependent simulator. That is, we modify the experiment so that, the challenger first rewinds the MIM and extracts the joint distribution of values $\tilde{\gamma}$ committed by the MIM. Here, we rely on the fact that Π is stand-alone extractable (with over-extraction). Note that once extracted, these $\tilde{\gamma}$'s can be used to extract the messages committed by the MIM in any transcript with the same fixed first message, with overwhelming probability. The only situation in which the $\tilde{\gamma}_i^b$ extracted for some execution i does not help recover the message committed by the MIM from transcript τ with the same fixed first message, is if the MIM uses a different witness $\tilde{\gamma}_i^{1-b}$ in τ and uses $\tilde{\gamma}_i^b$ in all the rewinding executions. It is easy to observe this event occurs with

probability at most $\text{negl}(n)$.

Upon extracting these values, with the same fixed first message, the challenger begins running the simulation strategy of weak ZK to output a main transcript with a simulated proof. That is, the challenger uses the $\tilde{\gamma}$'s to extract the joint distribution of the values committed by the MIM from any right execution, and runs the distinguisher-dependent simulator on a distinguisher that obtains the joint distribution of the view output by the MIM, together with these extracted values. Now, by the guarantee of distinguisher-dependent simulation, we have that the joint distribution remains indistinguishable between **Hybrid₁** and **Hybrid₂**.

In our actual reduction, since we are first rewinding and then generating a simulated proof, we require a special type of weak resettable security of the weak ZK. Thus, for the proof to go through, it is crucial that we use a specific ordering to generate the lookahead threads for extracting the MIM's values, and the lookahead threads for simulation. Additional details can be found in the next section.

Hybrid₃ : In the next hybrid, the output transcript generated in the left execution, consists of a commitment $c_1 = \text{com}(0; \hat{r})$ for some randomness \hat{r} , instead of c_1 being a commitment to 1. This is allowed because the weak ZK proof is being simulated by this point. The joint distribution of the view and values committed do not change in this hybrid, because c_1 is non-interactive, and thus can be replaced in the main transcript, while rewinding the MIM and extracting the joint distribution of the values committed by the MIM in all right executions.

Hybrid₄ : In this next hybrid, the challenger sets $\delta_2 = \text{PRF}(\gamma_2, \alpha_2) \oplus \hat{r}$ (instead of $\oplus r$), where \hat{r} is the randomness used to generate c_1 . Since the PRF key γ_2 does not appear elsewhere in the protocol, the joint distribution of the view and values committed do not change in this hybrid. This is δ_2 can be replaced in the main transcript, while rewinding the MIM and extracting the joint distribution of the values committed by the MIM in all right executions.

Hybrid₅ : In this next hybrid, the challenger changes the non-malleable commitment Π^2 to commit to the same randomness γ_2 , that is used to compute δ_2 in all threads (instead of committing to a different γ'_2). In order to argue indistinguishability of the view and committed values, we now rely on the non-malleability of Π^2 . The challenger runs the extractor for Π^2 on a transcript that contains honestly generated **wzk** proofs: again by soundness, at least one of the $\tilde{\gamma}$ values committed by the MIM in every execution is a valid commitment in the main thread. Thus, the extractor outputs this value. Next, the challenger uses this extracted value to recover the joint distribution of messages from transcripts generated by the MIM. This helps the challenger generate an output transcript with a simulated **wzk** proof, such that the joint distribution of the view of the MIM and values committed remains indistinguishable.

Note that in this experiment, even though the left execution is rewound for distinguisher-dependent simulation, this rewinding happens after the first two rounds have been fixed: thus, the *non-malleable commitment* used in the left execution is never rewound, and can be obtained externally. If the joint distribution of view and values output by the extractor for Π changes in this hybrid, then this contradicts hiding of Π . The argument of indistinguishability for this hybrid again requires a delicate ordering to generate the lookahead threads for extracting the MIM's committed values, and the lookahead threads for simulation. Additional details can be found in the next section.

Hybrid₆, Hybrid₇ : By the end of these hybrids, the challenger will behave the same way as **Hybrid₅**, except that it will use the second witness γ_2 in all executions (in the main as well as lookahead threads). For the main thread, for which the witness is switched in **Hybrid₆**, the challenger will use witness $\gamma_2, \hat{r}, \delta_2, c_1$ to compute the **wi**. In the rewinding threads, for which the witness is

switched in **Hybrid₇**, the challenger will use witness γ_2, r, δ_2, c . The joint distribution of the view and value extracted remains indistinguishable because of the reusable resettable security of wi allows for switching the witness even when multiple proofs are given in the main as well as rewinding executions.

Hybrid₈ : In this hybrid, the challenger sets Π^1 as a non-malleable commitment to a different independently uniform randomness γ'_1 , than the randomness γ that is used to compute δ_1 in all executions. The joint distribution of view and values committed by the **MIM** remains indistinguishable by the non-malleability of Π . The proof follows in a similar manner as that of the indistinguishability of **Hybrid₅**.

Hybrid₉ : In this hybrid, the challenger behaves similar to the previous hybrid except setting δ_1 to uniformly at random, only in the output transcript. Since the key γ_1 no longer appears elsewhere in the protocol, indistinguishability of the view and committed values follows by security of the **PRF**.

Hybrid₁₀: In this hybrid, the challenger behaves similar to the previous hybrid, except in the output transcript, it sets c as a commitment to value V_2 instead of to value V_1 . This is allowed because the randomness used to compute c in the output transcript is not used elsewhere in the protocol. Indistinguishability of the view and values committed by the **MIM** in this execution, follows by hiding of the non-interactive commitment c .

At this point, the main transcript consists of a commitment to V_2 instead of to V_1 , while the lookahead transcripts are generated using both V_1 and V_2 . Now, following the same sequence of hybrids in reverse order, we get to a hybrid experiment where the challenger generates an honest commitment to V_2 in the left execution. Thus, the joint distribution of the view and values committed by the **MIM** remains indistinguishable between when the left commitment is to V_1 , versus to V_2 , which is the guarantee required by the definition of non-malleability.

5.2.2 Hybrid Experiments

We now formally describe the hybrid arguments required to prove non-malleability.

Hybrid_{V₁} : This hybrid corresponds to an interaction of the challenger and the **MIM** where the challenger uses input message V_1 in the honest interaction. Let $\text{MIM}_{(C,R)}(V_1, z)$ denote the joint distribution of the view and values committed by the **MIM** in this interaction.

Hybrid₁ : In this hybrid, the challenger behaves identically to **Hybrid_{V₁}**, except that it generates Π^2 as a non-malleable commitment to a different randomness γ'_2 chosen uniformly and independently at at random, from the randomness γ_2 that was used to compute δ_2 . Let $\text{MIM}_{(C,R)}(\text{value}, z)_{\text{Hybrid}_1}$ denote the joint distribution of the view and values committed by the **MIM** in this interaction, in all the right sessions.

Lemma 3. *For any PPT distinguisher \mathcal{D} with auxiliary information z , $|\Pr[\mathcal{D}(z, \text{MIM}_{(C,R)}(V_1, z)) = 1] - \Pr[\mathcal{D}(z, \text{MIM}_{(C,R)}(\text{value}, z)_{\text{Hybrid}_1}) = 1]| \leq \epsilon + \text{negl}(n)$.*

Proof. The proof of this lemma follows via a reduction to the weak non-malleability of the scheme Π . More specifically, given a distinguisher \mathcal{D} that distinguishes $\text{MIM}_{(C,R)}(\text{value}, z)_{\text{Hybrid}_1}$ and $\text{MIM}_{(C,R)}(V_1, z)$, we construct an adversary $\mathcal{A}^{\mathcal{D}}$ against the one-many non-malleability of Π .

The adversary \mathcal{A} participates in the experiment exactly as Hybrid_{V_1} , except that it samples $\gamma_2, \gamma'_2 \xleftarrow{\$} \{0, 1\}^*$ and submits these to an external challenger. It obtains externally, the messages of Π^2 , which are either a non-malleable commitment to γ_2 or to γ'_2 . It complete the third message of the protocol using γ_2 to compute δ_2 .

By the weak non-malleability of Π , there exists an extractor that runs in time $\text{poly}(\frac{1}{\epsilon})$ and extracts the values committed by the MIM in all the non-malleable commitments for all $j \in [p]$, *without* rewinding the honest execution. Further, this extractor has the property that it only extracts an incorrect value if the MIM is committing to \perp in the main thread in the honest execution, except with error ϵ .

However, in both Hybrid_{V_1} and Hybrid_1 , by the soundness of wi , the adversary is guaranteed to generate at least one out of the two non-malleable commitments (to $\tilde{\gamma}_1$ or $\tilde{\gamma}_2$) from each session, correctly in any execution, except with probability $\text{negl}(n)$. Moreover, by soundness of wzk , the extracted value from at least one of the non-malleable commitments generated by the MIM in each session, will correspond to a witness for the commitment c , and therefore directly help recover the value committed by the MIM in each right session.

\mathcal{A} then samples a random main thread execution, and then just runs this extractor to extract the values $\{\tilde{\gamma}_i^1, \tilde{\gamma}_i^2\}_{i \in [n]}$ committed by the MIM, and by soundness of wi and wzk , at least one is correctly extracted. Depending upon whether the challenge non-malleable commitment is to γ_2 or γ'_2 , the joint distribution of the view and value extracted by \mathcal{A} corresponds to either $\text{MIM}_{\langle C, R \rangle}(V_1, z)$ or $\text{MIM}_{\langle C, R \rangle}(\text{value}, z)_{\text{Hybrid}_1}$.

Therefore, if the joint distribution of the view and the values committed by the MIM changes by more than ϵ between these executions, it can be used to contradict the one-many weak non-malleability of Π . Thus, if

$$|\Pr[\mathcal{D}(z, \text{MIM}_{\langle C, R \rangle}(V_1, z)) = 1] - \Pr[\mathcal{D}(z, \text{MIM}_{\langle C, R \rangle}(\text{value}, z)_{\text{Hybrid}_1}) = 1]| \geq \epsilon + \frac{1}{\text{poly}}(n),$$

$$\text{then, } |\Pr[\mathcal{A}^{\mathcal{D}} = 1 | \gamma'] - \Pr[\mathcal{A}^{\mathcal{D}} = 1 | \gamma]| \geq \frac{1}{\text{poly}}(n).$$

This gives a contradiction, thus the distributions are indistinguishable upto ϵ error. \square

We note that in Hybrid_1 , soundness of the wi and wzk arguments in the left as well as right interactions is still maintained, thus a rewinding extractor always successfully extracts the value committed by the MIM.

Hybrid_{2, D} : In this hybrid, the challenger behaves similarly to Hybrid_1 , except that it outputs the transcript of an execution where the distinguisher-dependent weak zero-knowledge protocol wzk is simulated as follows:

- Run the execution until the MIM sends the first message for the right execution. With fixed first messages, ϕ_1 and $\tilde{\phi}_1^j$, run the rest of the protocol as follows.
- Send second messages $\tilde{\phi}_2^j$ for the right interactions, and wait for the MIM's response ϕ_2 . These will correspond to the first and second messages for the main transcript. Instead of completing the experiment by sending the third message, proceed to the next step.
- With the same fixed first messages, ϕ_1 and $\tilde{\phi}_1^j$, rewind the protocol $\text{poly}(1/\epsilon)$ times sending various second round challenge messages to the MIM on behalf of honest receiver. When the MIM sends a challenge for the left (honest) execution, complete the transcript as an honest commitment to V_1 , and wait for the MIM's response.

Use these rewinding executions to extract the value committed in at least one (or both) of the non-malleable commitments provided by the MIM adversary, for each session. Whenever the MIM completes a right execution (that is, it does not generate any invalid messages), by soundness of the ZK argument, we have that except with probability at most $\text{negl}(n)$, at least one of the non-malleable commitments were generated correctly in each execution. Thus, by the same argument as used in the Lemma 3, with overwhelming probability, the extractor runs in time $\text{poly}(\frac{1}{\sqrt{\epsilon}})$ and correctly extracts at least one of the values committed by the MIM using the non-malleable commitment in all executions, except with error $\sqrt{\epsilon}$.

For each right session $j \in [p]$, let us denote the values extracted by the challenger by $\tilde{\gamma}_1^j, \tilde{\gamma}_2^j$, where at least one value was correctly extracted except with failure probability $\sqrt{\epsilon}$. Moreover, if for any right execution the extractor successfully extracted only *one* value, then by a simple probabilistic argument, the MIM will continue to use the same value as witness for the w_i in other executions except with probability at most $\sqrt{\epsilon}$ (otherwise, if the MIM used a different value as witness for the w_i , then that value would also be extracted with significant probability). Therefore, $\{\tilde{\gamma}_1^j, \tilde{\gamma}_2^j\}_{j \in [p]}$ can be used to recover the value committed by the MIM from any transcript generated by the MIM with fixed first messages $\phi_1, \tilde{\phi}_1^j$, except with failure probability $\sqrt{\epsilon}$.

- After completing the previous step, with the first message transcript fixed, go back to the main transcript with messages $\phi_1, \tilde{\phi}_1^j, \phi_2^j, \phi_2$. These will remain fixed for the rest of the experiment. Since these were fixed before the rewindings, by the weak resettable weak ZK property of wzk , the simulation security of wzk holds with respect to the partial transcript $(\phi_1, \phi_2, \tilde{\phi}_1^j, \tilde{\phi}_2^j)$.

In particular, weak resettable security implies that indistinguishability between real and simulated view must hold even against a distinguisher that rewinds and obtains $\{\tilde{\gamma}_1^j, \tilde{\gamma}_2^j\}_{j \in [p]}$ – which can be used to extract the message committed in the string c by the MIM from any transcript generated by the MIM with fixed first messages $\phi_1, \tilde{\phi}_1^j$, except with error $\sqrt{\epsilon} + \text{negl}(n)$.

- Next, run the distinguisher-dependent simulation strategy \mathcal{S} of the weak zero-knowledge argument, with error $\sqrt{\epsilon}$, on the distinguisher \mathcal{D}' . \mathcal{D}' is given the view of the MIM, together with auxiliary information $\{\gamma_1^j, \gamma_2^j\}_{j \in [p]}$. On input the view of the MIM, it uses this information to extract the value committed by the MIM from all its executions. It then runs the distinguisher \mathcal{D} on the joint distribution of the view and the extracted values and mirrors the output of \mathcal{D} .

Recall, that the distinguisher-dependent simulation strategy \mathcal{S} of [JKKR17] generates several different third messages (corresponding to the same fixed messages $(\phi_1, \phi_2, \tilde{\phi}_1^j, \tilde{\phi}_2^j)$), while sampling fresh α_1, α_2 each time. Also note that the output transcript still contains a commitment to V_1 , and is in fact identical to Hybrid_1 except that it contains a simulated wzk proof.

Let $\text{MIM}_{\langle C, R \rangle}(\text{value}, z)_{\text{Hybrid}_{2, \mathcal{D}}}$ denote the joint distribution of the view and value committed by the MIM when interacting with an honest committer in this hybrid.

Lemma 4. *For any PPT distinguisher \mathcal{D} with auxiliary information z ,*

$$|\Pr[\mathcal{D}(z, \text{MIM}_{\langle C, R \rangle}(\text{value}, z)_{\text{Hybrid}_{2, \mathcal{D}}}) = 1] - \Pr[\mathcal{D}(z, \text{MIM}_{\langle C, R \rangle}(\text{value}, z)_{\text{Hybrid}_1}) = 1]| \leq \epsilon + \text{negl}(n).$$

Proof. This claim follows by the weak resettable security of distinguisher-dependent simulation: since $\text{MIM}_{\langle C, R \rangle}(\text{value}, z)_{\text{Hybrid}_{2, \mathcal{D}}}$ is the result of executing distinguisher-dependent simulation against distinguisher \mathcal{D}' , which itself runs the distinguisher \mathcal{D} on $\text{MIM}_{\langle C, R \rangle}(\text{value}, z)_{\text{Hybrid}_1}$. Note that the weak resettable security experiment for distinguisher-dependent simulation allows the adversary to obtain, in addition to a real/simulated main transcript, several “lookahead” transcripts, where all

lookahead transcripts contain honestly generated proofs, that all use the same first message of the argument.

In other words, we consider a reduction that first fixes the first two messages of the honest and MIM execution corresponding to the main thread. Next, it generates multiple lookahead threads, as allowed by the security experiment of weak resettable **wzk**, using these threads to extract the values committed by the MIM. It generates all messages on its own according to **Hybrid**₁, except that it obtains the honestly generated **wzk** proofs for these threads externally from a challenger for weak resettable weak ZK.

Finally, the challenger flips a bit b , and if $b = 0$, it outputs an honestly generated weak ZK argument for the main transcript. On the other hand, if $b = 1$, it outputs a simulated argument (with error at most $\sqrt{\epsilon} \cdot \sqrt{\epsilon} = \epsilon$), while simulating the output of distinguisher \mathcal{D} on input the view and values extracted from the MIM. The reduction obtains this proof from the challenger and uses it to complete the main transcript. Note that if $b = 0$, the experiment corresponds to running \mathcal{D} on $\text{MIM}_{\langle C, R \rangle}(\text{value}, z)_{\text{Hybrid}_{1, \mathcal{D}}}$ and if $b = 1$, the experiment corresponds to running the distinguisher \mathcal{D} on $\text{MIM}_{\langle C, R \rangle}(\text{value}, z)_{\text{Hybrid}_{2, \mathcal{D}}}$. Thus, if $|\Pr[\mathcal{D}(z, \text{MIM}_{\langle C, R \rangle}(\text{value}, z)_{\text{Hybrid}_{2, \mathcal{D}}}) = 1] - \Pr[\mathcal{D}(z, \text{MIM}_{\langle C, R \rangle}(\text{value}, z)_{\text{Hybrid}_{1, \mathcal{D}}}) = 1]| > \epsilon + \text{negl}(n)$, this gives a distinguisher against the weak resettable simulation security of the weak ZK argument according to [Definition 2](#), which is a contradiction. \square

Hybrid_{3, \mathcal{D}} : In this hybrid, the challenger behaves identically to **Hybrid**_{2, \mathcal{D}} , except that it sets $c_1 = \text{com}(0; \hat{r})$ for some randomness \hat{r} , in the main transcript (instead of generating it as a commitment to 1). Note that this is possible because the challenger is generating a simulated proof in the output transcript, for the fact that c_1 is a commitment to 1. Let $\text{MIM}_{\langle C, R \rangle}(\text{value}, z)_{\text{Hybrid}_{3, \mathcal{D}}}$ denote the joint distribution of the view and values committed by the MIM when interacting with the challenger in this hybrid.

Lemma 5. *For any PPT distinguisher \mathcal{D} with auxiliary information z ,*

$$|\Pr[\mathcal{D}(z, \text{MIM}_{\langle C, R \rangle}(\text{value}, z)_{\text{Hybrid}_{2, \mathcal{D}}}) = 1] - \Pr[\mathcal{D}(z, \text{MIM}_{\langle C, R \rangle}(\text{value}, z)_{\text{Hybrid}_{3, \mathcal{D}}}) = 1]| \leq \text{negl}(n).$$

Proof. This hybrid is indistinguishable from **Hybrid**₂ by the computational hiding property of the non-interactive commitment scheme **com**. More formally, consider a reduction \mathcal{R} that behaves identically to **Hybrid**_{2, \mathcal{D}} , first extracting $\{\tilde{\gamma}_1^j, \tilde{\gamma}_2^j\}_{j \in [p]}$. Next, it obtains the commitment c_1 (only for the main output transcript and not for any of the rewinding executions), externally, as either a commitment to 0 or a commitment to 1, and uses this to complete the main transcript. It then uses the extracted values $\{\tilde{\gamma}_1^j, \tilde{\gamma}_2^j\}_{j \in [p]}$ to recover the values committed by the MIM in the main transcript. It outputs the joint distribution of the transcript and the values committed by the MIM to distinguisher \mathcal{D} . Then given a distinguisher \mathcal{D} where:

$$|\Pr[\mathcal{D}(z, \text{MIM}_{\langle C, R \rangle}(\text{value}, z)_{\text{Hybrid}_{2, \mathcal{D}}}) = 1] - \Pr[\mathcal{D}(z, \text{MIM}_{\langle C, R \rangle}(\text{value}, z)_{\text{Hybrid}_{3, \mathcal{D}}}) = 1]| \geq \frac{1}{\text{poly}(n)}$$

The reduction mirrors the output of this distinguisher such that:

$$|\Pr[\mathcal{R} = 1 | c_1 = \text{com}(1; r)] - \Pr[\mathcal{R} = 1 | c_1 = \text{com}(0; r)]| \geq \frac{1}{\text{poly}(n)}$$

This is a contradiction to the hiding of **com**. \square

Hybrid_{4, \mathcal{D}} : In this hybrid, the challenger behaves identically to **Hybrid**_{3, \mathcal{D}} except that in the output transcript, it sets $\delta_2 = \text{PRF}(\gamma_2, \alpha_2) \oplus \hat{r}$ where \hat{r} is the randomness used to generate $c_1 = \text{com}(0; \hat{r})$.

Note that the committer is using PRF key γ'_2 in the protocol Π^2 , thus the key γ_2 does not appear anywhere else in the rest of the protocol.

Let $\text{MIM}_{\langle C,R \rangle}(\text{value}, z)_{\text{Hybrid}_{4,\mathcal{D}}}$ denote the joint distribution of the view and value committed by the MIM when interacting with an honest committer in this hybrid.

Lemma 6. *For any PPT distinguisher \mathcal{D} with auxiliary information z ,*

$$|\Pr[\mathcal{D}(z, \text{MIM}_{\langle C,R \rangle}(\text{value}, z)_{\text{Hybrid}_{4,\mathcal{D}}}) = 1] - \Pr[\mathcal{D}(z, \text{MIM}_{\langle C,R \rangle}(\text{value}, z)_{\text{Hybrid}_{3,\mathcal{D}}}) = 1]| \leq \text{negl}(n).$$

Proof. This hybrid is indistinguishable from $\text{Hybrid}_{3,\mathcal{D}}$ by the security of the PRF. More formally, consider a reduction \mathcal{R} that behaves identically to $\text{Hybrid}_{3,\mathcal{D}}$ except that for all rewinding (recall that the distinguisher is rewound several times) transcripts generated during distinguisher-dependent simulation, it samples fresh α_2 each time and obtains $\text{PRF}(\gamma_2, \alpha_2) \oplus \hat{r}$ externally from a PRF challenger.

Then, for the main output transcript it obtains the value δ_2 externally as either $\text{PRF}(\gamma_2, \alpha_2) \oplus \hat{r}$, or $\text{PRF}(\gamma_2, \alpha_2) \oplus r$, where r is the randomness used generate commitment c in the left execution, and \hat{r} is the randomness used to generate commitment c_1 . It uses the externally obtained δ_2 to complete the main transcript. It then uses the extracted values $\{\tilde{\gamma}_1^j, \tilde{\gamma}_2^j\}_{j \in [p]}$ to obtain the values committed by the MIM in the main transcript. It outputs the joint distribution of the transcript and the values committed by the MIM to distinguisher \mathcal{D} .

Given a distinguisher \mathcal{D} where:

$$|\Pr[\mathcal{D}(z, \text{MIM}_{\langle C,R \rangle}(\text{value}, z)_{\text{Hybrid}_{4,\mathcal{D}}}) = 1] - \Pr[\mathcal{D}(z, \text{MIM}_{\langle C,R \rangle}(\text{value}, z)_{\text{Hybrid}_{3,\mathcal{D}}}) = 1]| \geq \frac{1}{\text{poly}(n)}$$

In this case, the reduction can mirror the output of this distinguisher to directly contradict the security of the PRF. \square

Hybrid_{5, \mathcal{D}} : In this hybrid, the challenger behaves identically to $\text{Hybrid}_{4,\mathcal{D}}$ except that it sets Π^2 as a non-malleable commitment to the same randomness γ_2 that is used to compute δ_2 , for all executions.

This hybrid essentially “reverts” the changes performed in Hybrid_1 . Note that the challenger in this hybrid, first extracts the values committed via the non-malleable commitments provided by the MIM, and then rewinds the *distinguisher* multiple times – however, the first two messages of the protocol are fixed at the time of rewinding the distinguisher. In particular, for fixed nmc_1^2 and nmc_2^2 , the challenger gives the same response nmc_3^2 for all the third messages it generates while/before simulating wzk argument.

Since the main thread transcript output in this hybrid consists of a simulated proof, indistinguishability of this hybrid is the most interesting to argue. We prove that it follows by the weak non-malleability of Π^2 . It is important, for the proof of non-malleability to go through, that the witness used by the prover in the proof of WI in this hybrid, is always the randomness used to compute Π^1 and never the randomness used to compute Π^2 – because the messages of Π^2 will be obtained externally. Moreover, recall that the proof of non-malleability of the weak non-malleable commitment scheme Π requires a simulator-extractor to “cheat” in the scheme Π^2 in rewinding executions.

Note that the challenger in this hybrid, fixes the first two rounds for the output transcript. Then, with the same fixed first round, it attempts to extract the values $(\tilde{\gamma}_1^j, \tilde{\gamma}_2^j)$ committed by the MIM in the non-malleable commitments in all executions. It then rewinds the *distinguisher* multiple times – at this point the first two messages of the protocol are fixed. Note that the transcript output by the challenger in this experiment is such that Π^1 is a valid non-malleable commitment to γ_1 with randomness r_1 AND $r = \text{PRF}(\gamma_1, \alpha_1) \oplus \delta_1$ such that $c = \text{com}(m; r)$ (and this is the witness used

in wi). Additionally, Π^2 is also a valid non-malleable commitment to γ_2 with randomness r_2 AND $\hat{r} = \text{PRF}(\gamma_2, \alpha_2) \oplus \delta_2$ such that $c_1 = \text{com}(0; \hat{r})$. However, the witness used in wi is always Π^1 .

Let $\text{MIM}_{(C,R)}(\text{value}, z)_{\text{Hybrid}_{5,\mathcal{D}}}$ denote the joint distribution of the view and value committed by the MIM when interacting with an honest committer in this hybrid.

Lemma 7. *For any PPT distinguisher \mathcal{D} with auxiliary information z ,*

$$|\Pr[\mathcal{D}(z, \text{MIM}_{(C,R)}(\text{value}, z)_{\text{Hybrid}_{5,\mathcal{D}}}) = 1] - \Pr[\mathcal{D}(z, \text{MIM}_{(C,R)}(\text{value}, z)_{\text{Hybrid}_{4,\mathcal{D}}}) = 1]| \leq \epsilon + \text{negl}(n).$$

Proof. Recall that the challenger strategy in both $\text{Hybrid}_{5,\mathcal{D}}$ and $\text{Hybrid}_{4,\mathcal{D}}$ is as follows: The challenger first generates and fixes the first two messages of the main transcript $\phi_1, \tilde{\phi}_1^j, \tilde{\phi}_2^j, \phi_2$. It then rewinds the MIM multiple times with the same fixed first message but different second round messages, to extract $\tilde{\gamma}_1^j, \tilde{\gamma}_2^j$ for all $j \in [n]$. Finally, it runs the distinguisher-dependent simulation strategy with partial transcript $\phi_1, \tilde{\phi}_1^j, \tilde{\phi}_2^j, \phi_2$ to output a main transcript with a simulated proof.

The main difference between $\text{Hybrid}_{4,\mathcal{D}}$ and $\text{Hybrid}_{5,\mathcal{D}}$ is that the committer commits to γ_2' using Π^2 in $\text{Hybrid}_{4,\mathcal{D}}$, and uses a different γ_2 for the rest of the protocol, whereas in $\text{Hybrid}_{5,\mathcal{D}}$, $\gamma_2' = \gamma_2$. However, both hybrids involve the challenger rewinding the MIM (and consequently rewinding the left session) several times in order to extract $\tilde{\gamma}_1^j, \tilde{\gamma}_2^j$ for $j \in [n]$. In this rewinding situation, invoking weak one-malleability of Π^2 requires care.

Our first observation is that by the weak non-malleability of Π , there exists an extractor that runs in time $\text{poly}(\frac{1}{\epsilon})$ and extracts the values committed by the MIM in all the non-malleable commitments for all $j \in [p]$, *without rewinding the left execution*. The reduction to one-many weak non-malleability of Π uses this extractor and proceeds as follows:

- The reduction begins by fixing the first two messages in the left and right executions in the main thread. For these messages, it obtains an externally generated non-malleable commitment to either $\gamma_2' = \gamma_2$ or γ_2' chosen uniformly at random independent of γ_2 . The former corresponds to $\text{Hybrid}_{5,\mathcal{D}}$ and the latter to $\text{Hybrid}_{4,\mathcal{D}}$.

Instead of rewinding the MIM providing honestly generated transcripts in the left interaction as is done in $\text{Hybrid}_{5,\mathcal{D}}$ and $\text{Hybrid}_{4,\mathcal{D}}$, we will now consider two sub-hybrids, $\text{Hybrid}_{4,a,\mathcal{D}}$ and $\text{Hybrid}_{5,a,\mathcal{D}}$ where the reduction uses the extractor \mathcal{E} for the non-malleable commitment to extract the values committed by the MIM without rewinding the left interaction. We will show that the view and values extracted from these sub-hybrids will remain identical to the view and value extracted via rewinding in $\text{Hybrid}_{4,\mathcal{D}}$ and $\text{Hybrid}_{5,\mathcal{D}}$, respectively. This will essentially follow because of correctness of extractor \mathcal{E} , and because of soundness of wi and wzk in the interactions from which extraction occurs. We will also directly give a reduction proving that the joint distribution of the views and values extracted must be indistinguishable between these sub-hybrids.

- Recall that \mathcal{E} extracts the values committed by the MIM in a main transcript, without rewinding the messages sent in the non-malleable commitment in the left interaction (the extractor \mathcal{E} may still rewind the MIM, only in all such rewinds it will not need to rewind the left non-malleable commitment, indeed it will suffice to generate “fake” third round messages for the non-malleable commitment to γ_2 – please refer to [GRRV14] for details on the extraction procedure). It is important to note that the wzk simulation strategy requires that the MIM’s committed values be extracted first, therefore we cannot generate a simulated wzk argument without first extracting all values $\tilde{\gamma}_1^j, \tilde{\gamma}_2^j$ committed by the MIM.

- Thus, in sub-hybrids $\text{Hybrid}_{i,a,\mathcal{D}}$ for $i \in \{4, 5\}$, the challenger just runs extractor \mathcal{E} to extract the values $\{\tilde{\gamma}_1^j, \tilde{\gamma}_2^j\}_{j \in [n]}$, instead of rewinding the left execution. \mathcal{E} extracts the value committed in a main transcript without rewinding the left execution. Thus, first the challenger generates a special main transcript for the extractor \mathcal{E} as follows. It generates $\phi_1, \tilde{\phi}_1^j, \tilde{\phi}_2^j, \phi_2$ the same way as $\text{Hybrid}_{4,\mathcal{D}}$, and then completes the third message by generating an honest commitment to V_1 , that is, giving an honestly generated \mathbf{wzk} argument and using γ_1 as witness for the \mathbf{wi} ⁴. It waits for the MIM to generate the third messages for the right executions, and now feeds the transcript of the interaction to \mathcal{E} (if the MIM aborts, the challenger just repeats again with the same fixed first two messages, $\text{poly}(1/\epsilon)$ times). Whenever \mathcal{E} requests to rewind the MIM, the challenger rewinds the MIM, except that it obtains the messages for the left commitment Π^2 in all rewinding executions from \mathcal{E} . Further, recall that \mathcal{E} has the property that it only extracts an incorrect value when the MIM is committing to \perp in the honest execution, except with error ϵ , however, this is not true except with probability $1 - \text{negl}(n)$, by soundness of \mathbf{wi} and \mathbf{wzk} . The MIM waits for \mathcal{E} to output the extracted values $\{\tilde{\gamma}_1^j, \tilde{\gamma}_2^j\}$. Next, the MIM repeats this again with same fixed first two messages, waiting for the extractor to output (potentially different) extracted values. Finally the challenger uses the union of these extracted values to complete the rest of the experiment according to $\text{Hybrid}_{4,\mathcal{D}}$.

Claim 1. *The joint distribution of the views and values committed by the MIM remain indistinguishable (with error at most $\epsilon + \text{negl}(n)$) between $\text{Hybrid}_{i,\mathcal{D}}$ and $\text{Hybrid}_{i,a,\mathcal{D}}$ for $i \in \{4, 5\}$.*

Proof. Note that the special main transcript provided to \mathcal{E} to facilitate extraction in the sub-hybrids, is distributed identically to the transcripts provided in the lookahead executions for extraction in $\text{Hybrid}_{4,\mathcal{D}}$ and $\text{Hybrid}_{5,\mathcal{D}}$. Additionally, in all these executions, the challenger always provides honestly generated proofs, thus the soundness of \mathbf{wi} and \mathbf{wzk} provided by the MIM is guaranteed in all these executions. Therefore, the adversary is guaranteed to generate at least one out of the two non-malleable commitments from each session correctly in any non-aborting execution, except with probability $\text{negl}(n)$.

Moreover, by soundness of \mathbf{wzk} , the extracted value from at least one of the non-malleable commitments generated by the MIM in the j^{th} session, will correspond to a witness for the commitment to $\tilde{\gamma}_j^1$ or $\tilde{\gamma}_j^2$, directly allowing to recover the message committed by the MIM in each non-aborting right session (if only one $\tilde{\gamma}_j$ was extracted, w.h.p. the MIM continues to use the same witness). By correctness of extraction from \mathcal{E} and because of soundness of \mathbf{wi} and \mathbf{wzk} in all rewinding executions as well as the special main execution, the joint distribution of views and value extracted via rewinding in $\text{Hybrid}_{i,\mathcal{D}}$ is ϵ -indistinguishable from the distribution when \mathcal{A} extracts using \mathcal{E} in $\text{Hybrid}_{i,a,\mathcal{D}}$ for $i \in \{4, 5\}$. \square

- Next, keeping the first two messages of the transcript τ fixed, the challenger outputs a main transcript with a simulated weak ZK argument, where the simulation strategy runs on the distinguisher that obtains input the view of the MIM as well as the value extracted in the previous step, in a similar manner to $\text{Hybrid}_{4,\mathcal{D}}$.

If the joint distribution of the view and values committed by the MIM between $\text{Hybrid}_{4,a,\mathcal{D}}$ and $\text{Hybrid}_{5,a,\mathcal{D}}$ are more than ϵ -distinguishable, there exists a reduction to the hiding of the non-malleable commitment Π^2 , which obtains the messages of Π^2 externally to generate the first two round messages. In response to the MIM's challenge for the left execution, it obtains the

⁴Note that the actual transcript that is output by the experiment must contain a simulated \mathbf{wzk} argument: the transcript with the honest \mathbf{wzk} argument is only generated to facilitate extraction.

third message of Π^2 externally, and uses it to generate the special main transcript for \mathcal{E} . Next, it runs the extractor \mathcal{E} , which does not need to rewind Π^2 in the left execution. Once it obtains $\{\tilde{\gamma}_j^1, \tilde{\gamma}_j^2\}_{j \in [p]}$ from \mathcal{E} , it proceeds to run the distinguisher-dependent simulation strategy. In this step, since the first two messages for the main transcript have already been fixed, the challenger can use the same third message Π_2^3 that it obtained externally, to complete the second non-malleable commitment in the left execution, in all third messages it generates in order to simulate the wzk argument by rewinding the distinguisher.

Therefore, if the joint distribution of the view and the values committed by the MIM changes by more than ϵ between $\text{Hybrid}_{4,a,\mathcal{D}}$ and $\text{Hybrid}_{5,a,\mathcal{D}}$, it can be used directly to contradict the hiding of Π^2 . That is, if

$$|\Pr[\mathcal{D}(z, \text{MIM}_{\langle C,R \rangle}(\text{value}, z)_{\text{Hybrid}_{5,a,\mathcal{D}}}) = 1] - \Pr[\mathcal{D}(z, \text{MIM}_{\langle C,R \rangle}(\text{value}, z)_{\text{Hybrid}_{4,a,\mathcal{D}}}) = 1]| \geq \epsilon + \frac{1}{\text{poly}}(n),$$

$$\text{then, } |\Pr[\mathcal{A}^{\mathcal{D}} = 1 | \gamma_2 = \gamma'_2] - \Pr[\mathcal{A}^{\mathcal{D}} = 1 | \gamma_2 \neq \gamma'_2]| \geq \frac{1}{\text{poly}}(n).$$

This gives a contradiction, thus the distributions are indistinguishable upto at most ϵ -error. \square

Hybrid_{6,D} : In this hybrid, the challenger behaves the same way as **Hybrid_{5,D}**, except that it uses the second witness, r_2, γ_2 , to generate the witness-indistinguishable argument wi in the output transcript.

Lemma 8. *For any PPT distinguisher \mathcal{D} with auxiliary information z ,*

$$|\Pr[\mathcal{D}(z, \text{MIM}_{\langle C,R \rangle}(\text{value}, z)_{\text{Hybrid}_{6,\mathcal{D}}}) = 1] - \Pr[\mathcal{D}(z, \text{MIM}_{\langle C,R \rangle}(\text{value}, z)_{\text{Hybrid}_{5,\mathcal{D}}}) = 1]| \leq \epsilon + \text{negl}(n).$$

Proof. The proof of this lemma relies on the reusable resettable witness indistinguishability of wi .

The reduction R samples all messages for the experiment according to **Hybrid_{5,D}**, except that it obtains WI proofs for all lookahead (rewinding) executions externally from the challenger, by providing the first witness to the challenger. In this experiment, note that some executions rewind the MIM to the end of the first round, thus proofs for these executions are provided with respect to new verifier messages generated by the MIM. Some other executions (corresponding to weak ZK simulation strategy) rewind the MIM to the end of the second round: thus different statements are proved in these executions, corresponding to the same verifier message from the MIM, that is fixed before the end of the second round. Thus, this experiment exactly corresponds to the security game of resettable reusable WI.

For the main/output transcript generated during distinguisher-dependent simulation, R samples all messages except the WI proof according to **Hybrid_{5,D}**. Note that the statement being proved in this transcript has two valid witnesses, $w_1 = (r_1, \gamma_1$ randomness r and commitment c) and $w_2 = (r_2, \gamma_2$, randomness \hat{r} and commitment c_1), which are sampled by the reduction R together with the adversary. R forwards the verifier message wi_1 to the challenger, together with both witnesses, and obtains wi_2 that is generated using either witness w_1 or w_2 . The reduction uses this externally generated proof to complete the experiment. If w_1 was used, the experiment is identical to **Hybrid_{5,D}**, otherwise it is identical to **Hybrid_{6,D}**.

Note that in the experiment, R behaves according to **Hybrid_{5,D}** or **Hybrid_{6,D}**: that is, it first extracts $\{\tilde{\gamma}_1^j, \tilde{\gamma}_2^j\}_{j \in [p]}$. It then uses the extracted values $\{\tilde{\gamma}_1^j, \tilde{\gamma}_2^j\}_{j \in [p]}$ to obtain the values committed by the MIM in the main transcript. It outputs the joint distribution of the transcript and the values committed by the MIM to distinguisher \mathcal{D} . Given a distinguisher \mathcal{D} where:

$$|\Pr[\mathcal{D}(z, \text{MIM}_{\langle C,R \rangle}(\text{value}, z)_{\text{Hybrid}_{6,\mathcal{D}}}) = 1] - \Pr[\mathcal{D}(z, \text{MIM}_{\langle C,R \rangle}(\text{value}, z)_{\text{Hybrid}_{5,\mathcal{D}}}) = 1]| \geq \frac{1}{\text{poly}(n)}$$

In this case, the reduction mirrors the output of this distinguisher to directly contradict the security of wi . Thus, the joint distribution in this hybrid is indistinguishable from $\text{Hybrid}_{5,\mathcal{D}}$ by the resettable reusable witness-indistinguishability of wi . \square

Hybrid $_{7,\mathcal{D}}$: In this hybrid, the challenger behaves the same way as **Hybrid $_{6,\mathcal{D}}$** , except that it uses the second witness, r_2, γ_2 , to generate the witness-indistinguishable arguments wi in all the lookahead executions. That is, in every message sent by the challenger, it uses the second witness instead of the first. This hybrid is indistinguishable from **Hybrid $_{6,\mathcal{D}}$** by the resettable reusable witness-indistinguishability of wi .

Lemma 9. *For any PPT distinguisher \mathcal{D} with auxiliary information z ,*

$$|\Pr[\mathcal{D}(z, \text{MIM}_{\langle C,R \rangle}(\text{value}, z)_{\text{Hybrid}_{7,\mathcal{D}}}) = 1] - \Pr[\mathcal{D}(z, \text{MIM}_{\langle C,R \rangle}(\text{value}, z)_{\text{Hybrid}_{6,\mathcal{D}}}) = 1]| \leq \text{negl}(n).$$

Proof. The proof of this lemma follows similarly to that of [Lemma 8](#), by relying on the resettable reusable witness-indistinguishability of wi . In this experiment, note that some executions rewind the MIM to the end of the first round, thus proofs for these executions are provided with respect to new verifier messages generated by the MIM. Some other executions (corresponding to weak ZK simulation strategy) rewind the MIM to the end of the second round: thus different statements are proved in these executions, corresponding to the same verifier message from the MIM, that is fixed before the end of the second round. Thus, this experiment exactly corresponds to the security game of resettable reusable WI.

That is, the reduction obtains WI proofs externally from the challenger by providing both witnesses $w_1 = (r_1, \gamma_1, \text{randomness } r \text{ and commitment } c)$ and $w_2 = (r_2, \gamma_2, \text{randomness } r \text{ and commitment } c)$. The challenger sends proofs that are all generated either using witness w_1 or all using witness w_2 . The reduction completes the rest of the protocol according to **Hybrid $_{6,\mathcal{D}}$** , except using the externally generated proofs in the left execution. If the challenger used witness w_1 , the game corresponds to **Hybrid $_{6,\mathcal{D}}$** otherwise it corresponds to **Hybrid $_{7,\mathcal{D}}$** .

Note that in the experiment, R behaves according to **Hybrid $_{6,\mathcal{D}}$** or **Hybrid $_{7,\mathcal{D}}$** : that is, it first extracts $\{\tilde{\gamma}_1^j, \tilde{\gamma}_2^j\}_{j \in [p]}$. It then uses the extracted values $\{\tilde{\gamma}_1^j, \tilde{\gamma}_2^j\}_{j \in [p]}$ to obtain the values committed by the MIM in the main transcript. It outputs the joint distribution of the transcript and the values committed by the MIM to distinguisher \mathcal{D} . Given a distinguisher \mathcal{D} where:

$$|\Pr[\mathcal{D}(z, \text{MIM}_{\langle C,R \rangle}(\text{value}, z)_{\text{Hybrid}_{7,\mathcal{D}}}) = 1] - \Pr[\mathcal{D}(z, \text{MIM}_{\langle C,R \rangle}(\text{value}, z)_{\text{Hybrid}_{6,\mathcal{D}}}) = 1]| \geq \frac{1}{\text{poly}(n)}$$

In this case, the reduction mirrors the output of this distinguisher to directly contradict the resettable reusable security of wi . \square

We note that the changes made in **Hybrid $_{7,\mathcal{D}}$** and **Hybrid $_{6,\mathcal{D}}$** can be collapsed into a single hybrid experiment relying on resettable reusable security of WI, however we keep them separate for additional clarity – since the witness used in the main transcript refers to Π^2 and the randomness for $c_1 = \text{com}(0; \hat{r})$ while the witness used in the lookahead transcripts refer to Π^2 and the randomness for $c = \text{com}(V_1; r)$. Note that at this point, the value γ_1 committed using the first non-malleable commitment Π^1 is not used as a witness in any of the WI proofs.

Hybrid $_{8,\mathcal{D}}$: In this hybrid, the challenger behaves the same way as **Hybrid $_{7,\mathcal{D}}$** , except that in all transcripts, it sets Π^1 as a non-malleable commitment to a *different* randomness γ'_1 than the one used to compute δ_1 .

Lemma 10. For any PPT distinguisher \mathcal{D} with auxiliary information z ,

$$|\Pr[\mathcal{D}(z, \text{MIM}_{\langle C, R \rangle}(\text{value}, z)_{\text{Hybrid}_{8, \mathcal{D}}} = 1)] - \Pr[\mathcal{D}(z, \text{MIM}_{\langle C, R \rangle}(\text{value}, z)_{\text{Hybrid}_{7, \mathcal{D}}} = 1)]| \leq \epsilon + \text{negl}(n).$$

Proof. The proof of this lemma is exactly the same as that of Lemma 7. The joint distribution of the view and value committed by a malicious receiver in $\text{Hybrid}_{8, \mathcal{D}}$ is ϵ -indistinguishable from $\text{Hybrid}_{7, \mathcal{D}}$ by the non-malleability of the commitment Π^1 . \square

$\text{Hybrid}_{9, \mathcal{D}}$: In this hybrid, the challenger behaves the same way as $\text{Hybrid}_{8, \mathcal{D}}$, except that in the output transcript, it sets $\delta_1 \stackrel{\$}{\leftarrow} \{0, 1\}^*$, instead of setting $\delta_1 = \text{PRF}(\gamma_1, \alpha_1) \oplus r$. Note that the committer is using PRF key γ'_1 in the protocol Π^1 , thus the key γ_1 does not appear in the rest of the protocol.

Lemma 11. For any PPT distinguisher \mathcal{D} with auxiliary information z ,

$$|\Pr[\mathcal{D}(z, \text{MIM}_{\langle C, R \rangle}(\text{value}, z)_{\text{Hybrid}_{9, \mathcal{D}}} = 1)] - \Pr[\mathcal{D}(z, \text{MIM}_{\langle C, R \rangle}(\text{value}, z)_{\text{Hybrid}_{8, \mathcal{D}}} = 1)]| \leq \text{negl}(n).$$

Proof. The proof of this lemma is the same as that of Lemma 6, by relying on the security of the PRF. \square

$\text{Hybrid}_{10, \mathcal{D}}$: In this hybrid, the challenger behaves the same way as $\text{Hybrid}_{9, \mathcal{D}}$ except that it replaces $c = \text{com}(V_1; r)$ with $c = \text{com}(V_2; r)$ in the output transcript. Note that in this transcript, the randomness r is not used elsewhere in the protocol.

Lemma 12. For any PPT distinguisher \mathcal{D} with auxiliary information z ,

$$|\Pr[\mathcal{D}(z, \text{MIM}_{\langle C, R \rangle}(\text{value}, z)_{\text{Hybrid}_{10, \mathcal{D}}} = 1)] - \Pr[\mathcal{D}(z, \text{MIM}_{\langle C, R \rangle}(\text{value}, z)_{\text{Hybrid}_{9, \mathcal{D}}} = 1)]| \leq \text{negl}(n).$$

Proof. This hybrid is indistinguishable from $\text{Hybrid}_{9, \mathcal{D}}$ because of computational hiding of the non-interactive commitment scheme com . More formally, consider a reduction \mathcal{R} that behaves identical to $\text{Hybrid}_{9, \mathcal{D}}$ except that it obtains the commitment c (only for the main output transcript and not for any of the rewinding executions), externally, as either a commitment to V_1 or a commitment to V_2 . This is allowed because by the end of $\text{Hybrid}_{9, \mathcal{D}}$, the randomness used to generate this commitment is not used anywhere else in the protocol.

Note that in the experiment, the reduction it first extracts $\{\tilde{\gamma}_1^j, \tilde{\gamma}_2^j\}_{j \in [p]}$. It then uses the extracted values $\{\tilde{\gamma}_1^j, \tilde{\gamma}_2^j\}_{j \in [p]}$ to obtain the values committed by the MIM in the main transcript. It outputs the joint distribution of the transcript and the values committed by the MIM to distinguisher \mathcal{D} . Then given a distinguisher \mathcal{D} where:

$$|\Pr[\mathcal{D}(z, \text{MIM}_{\langle C, R \rangle}(\text{value}, z)_{\text{Hybrid}_{9, \mathcal{D}}} = 1)] - \Pr[\mathcal{D}(z, \text{MIM}_{\langle C, R \rangle}(\text{value}, z)_{\text{Hybrid}_{10, \mathcal{D}}} = 1)]| \geq \frac{1}{\text{poly}(n)}$$

The reduction mirrors the output of this distinguisher such that:

$$|\Pr[\mathcal{R} = 1 | c = \text{com}(V_1; r)] - \Pr[\mathcal{R} = 1 | c = \text{com}(V_2; r)]| \geq \frac{1}{\text{poly}(n)}$$

This is a contradiction to the hiding of com . \square

At this point, we have successfully switched (with distinguishing advantage at most $5\epsilon + \text{negl}(n)$) to an experiment where the commitment is generated to message V_2 instead of V_1 in the transcript output by the challenger. However, note that the wzk argument is still being simulated in this hybrid. Also note that throughout these hybrids, lookahead threads for extraction are generated

according to both values V_1 and V_2 . Non-malleability follows by repeating the above hybrids in reverse order, until in Hybrid_{V_2} , the challenger generates an honest commitment to message V_2 . The hybrids are at most $10\epsilon + \text{negl}(n)$ -distinguishable.

The proof of one-many non-malleability can then be completed by setting 10ϵ to be less than the distinguishing advantage of the given distinguisher \mathcal{D} , and arriving at a contradiction. By invoking [LPV], this completes the proof of concurrent non-malleability. \square

Acknowledgements

We are extremely grateful to Vipul Goyal for helpful discussions about the properties of the protocol in [GRRV14], Yael Kalai for several useful discussions regarding [JKKR17] and Amit Sahai for very valuable feedback about this writeup.

References

- [Bar02] Boaz Barak. Constant-Round Coin-Tossing with a Man in the Middle or Realizing the Shared Random String Model. In *FOCS 2002*, pages 345–355, 2002. 2
- [BOV07] Boaz Barak, Shien Jin Ong, and Salil P. Vadhan. Derandomization in cryptography. *SIAM J. Comput.*, 37(2):380–400, 2007. 5
- [BP15] Nir Bitansky and Omer Paneth. Zaps and non-interactive witness indistinguishability from indistinguishability obfuscation. In Yevgeniy Dodis and Jesper Buus Nielsen, editors, *Theory of Cryptography - 12th Theory of Cryptography Conference, TCC 2015, Warsaw, Poland, March 23-25, 2015, Proceedings, Part II*, volume 9015 of *Lecture Notes in Computer Science*, pages 401–427. Springer, 2015. 3
- [COSV16] Michele Ciampi, Rafail Ostrovsky, Luisa Siniscalchi, and Ivan Visconti. Concurrent non-malleable commitments (and more) in 3 rounds. In Matthew Robshaw and Jonathan Katz, editors, *Advances in Cryptology - CRYPTO 2016 - 36th Annual International Cryptology Conference, Santa Barbara, CA, USA, August 14-18, 2016, Proceedings, Part III*, volume 9816 of *Lecture Notes in Computer Science*, pages 270–299. Springer, 2016. 2
- [COSV17] Michele Ciampi, Rafail Ostrovsky, Luisa Siniscalchi, and Ivan Visconti. 4-round concurrent non-malleable commitments from one-way functions. In *CRYPTO 2017*, 2017. 10
- [DDN91] Danny Dolev, Cynthia Dwork, and Moni Naor. Non-Malleable Cryptography (Extended Abstract). In *STOC 1991*, 1991. 2, 9
- [DN07] Cynthia Dwork and Moni Naor. Zaps and their applications. *SIAM J. Comput.*, 36(6):1513–1543, 2007. 3
- [GKS16] Vipul Goyal, Dakshita Khurana, and Amit Sahai. Breaking the three round barrier for non-malleable commitments. In *FOCS*, 2016. 3
- [GLOV12] Vipul Goyal, Chen-Kuei Lee, Rafail Ostrovsky, and Ivan Visconti. Constructing non-malleable commitments: A black-box approach. In *FOCS*, 2012. 2, 4

- [GOS12] Jens Groth, Rafail Ostrovsky, and Amit Sahai. New techniques for noninteractive zero-knowledge. *J. ACM*, 59(3):11:1–11:35, 2012. 3, 5
- [Goy11] Vipul Goyal. Constant Round Non-malleable Protocols Using One-way Functions. In *STOC 2011*, pages 695–704. ACM, 2011. 2, 4, 9
- [GPR16] Vipul Goyal, Omkant Pandey, and Silas Richelson. Textbook non-malleable commitments. In Daniel Wichs and Yishay Mansour, editors, *Proceedings of the 48th Annual ACM SIGACT Symposium on Theory of Computing, STOC 2016, Cambridge, MA, USA, June 18-21, 2016*, pages 1128–1141. ACM, 2016. 2, 3, 10, 14
- [GRRV14] Vipul Goyal, Silas Richelson, Alon Rosen, and Margarita Vald. An algebraic approach to non-malleability. In *FOCS 2014*, pages 41–50, 2014. 2, 3, 4, 6, 10, 14, 22, 27
- [JKKR17] Abhishek Jain, Yael Kalai, Dakshita Khurana, and Ron Rothblum. Distinguisher-dependent simulation in two rounds and its applications. 2017. <http://eprint.iacr.org/2017/330>. 4, 5, 7, 8, 14, 19, 27, 28, 36
- [KS17] Dakshita Khurana and Amit Sahai. How to achieve non-malleability in one or two rounds. *Electronic Colloquium on Computational Complexity (ECCC)*, 24:100, 2017. 2
- [LP] Huijia Lin and Rafael Pass. Constant-round Non-malleable Commitments from Any One-way Function. In *STOC 2011*, pages 705–714. 2, 4
- [LPS17] Huijia Lin, Rafael Pass, and Pratik Soni. Two-round and non-interactive concurrent non-malleable commitments from time-lock puzzles. Cryptology ePrint Archive, Report 2017/273, 2017. <http://eprint.iacr.org/2017/273>. 2
- [LPV] Huijia Lin, Rafael Pass, and Muthuramakrishnan Venkitasubramaniam. Concurrent Non-malleable Commitments from Any One-Way Function. In *TCC 2008*, pages 571–588. 9, 14, 27
- [Pas04] Rafael Pass. Bounded-Concurrent Secure Multi-Party Computation with a Dishonest Majority. In *Proceedings of the 36th Annual ACM Symposium on Theory of Computing, STOC '04*, pages 232–241, 2004. 2
- [Pas13] Rafael Pass. Unprovable security of perfect NIZK and non-interactive non-malleable commitments. In *TCC*, pages 334–354, 2013. 2, 3
- [PR05] Rafael Pass and Alon Rosen. New and improved constructions of non-malleable cryptographic protocols. In *STOC 2005*, pages 533–542, 2005. 2, 9
- [Wee10] Hoeteck Wee. Black-Box, Round-Efficient Secure Computation via Non-malleability Amplification. In *Proceedings of the 51th Annual IEEE Symposium on Foundations of Computer Science*, pages 531–540, 2010. 2

A Proof of Weak Resettable Security of [JKKR17] Protocol

In this section, we expand the proof sketch in [JKKR17], proving that the protocol in Figure 6, [JKKR17] satisfies weak resettable weak ZK in the distributional setting according to Definition 2, against non-adaptive verifiers. We describe the protocol again in Figure 3. This protocol is modified as described in [JKKR17], where the randomness for ZAP and OT is computed via a PRF on the verifier/receiver message.

Distributional Weak Zero-Knowledge Argument**Prover Input:** Distribution $(\mathcal{X}, \mathcal{W})$.**Verifier Input:** Distribution $(\mathcal{X}, \mathcal{W})$, language L .

- **Prover Message:** Pick $r_1, r_2, r'_1, r'_2 \xleftarrow{\$} \{0, 1\}^*$, send $c_1 = \text{com}(r_1; r'_1), c_2 = \text{com}(r_2; r'_2)$ using non-interactive statistically binding commitment com .
- **Verifier Message:** Pick challenge $e \xleftarrow{\$} \{0, 1\}^n$ for the Σ -protocol, and for $i \in [n]$, send $o_1 = \text{OT}_{1,i}(e_i)$ in parallel. Here, each e_i is encrypted with a fresh OT instance. Additionally send $\tilde{r}_1, \tilde{r}_2 \xleftarrow{\$} \{0, 1\}^*$, and send wi_1 as the first message of ZAP.
- **Prover Message:** Send r_1, r_2 . Sample $(x, w) \xleftarrow{\$} (\mathcal{X}, \mathcal{W})$ and send x .

Sample $K \xleftarrow{\$} \{0, 1\}^*$ and compute $r_{\text{wi}} \parallel r_{\text{OT}} = \text{PRF}(K, \text{wi}_1 \parallel o_1 \parallel x)$.

Use r_{wi} to compute and send wi_2 as the second message of ZAP proving that $\exists r'_1$ such that $c_1 = \text{com}(r_1; r'_1)$ OR $\exists r'_2$ such that $c_2 = \text{com}(r_2; r'_2)$.

Set $\text{pk}_1 = r_1 \oplus \tilde{r}_1, \text{pk}_2 = r_2 \oplus \tilde{r}_2$ as public keys for a dense cryptosystem. Define $\text{commit}(M; R) = \text{enc}_{\text{pk}_1}(M; s_1), \text{enc}_{\text{pk}_2}(M; s_2)$ where $R = s_1 \parallel s_2$. This is decommitted by revealing R .

For $i \in [n]$, define $a_i = \text{commit}(h_i)$. Send $a_i, \text{OT}_{2,i}(z_i^0, z_i^1)$ in parallel, where $\text{OT}_{2,i}$ are computed using randomness r_{OT} . Note that the decommitment information in z_i^0, z_i^1 corresponding to any commitment in the Σ -protocol, only consists of the randomness R used to generate the commitment using commit .

- **Verifier Output:** The verifier V recovers z_i as the output of OT_i for $i \in [n]$, and outputs **accept** if and only if wi is an accepting transcript and $(a_i, e_i, z_i)_{i \in [n]}$ is an accepting transcript of the underlying Σ -protocol, according to the commitment scheme commit .

Figure 3: Three Round Argument System for NP

Theorem 2 (Distributional Weak Zero-Knowledge). *The protocol in Figure 3 is distributional weak zero-knowledge against malicious PPT verifiers.*

Proof. Fix any weak resetting PPT V^* , any distinguisher \mathcal{D} , any distribution $(\mathcal{X}, \mathcal{W}, \mathcal{Z})$, and any $\epsilon > 0$. We construct a simulator Sim_ϵ that obtains non-uniform advice $z, p_\epsilon = \text{poly}(1/\epsilon)$ random instance-witness samples $(x_1^*, w_1^*), (x_2^*, w_2^*), \dots (x_{p_\epsilon}^*, w_{p_\epsilon}^*)$ from the distribution $(\mathcal{X}, \mathcal{W})$. Or, if the distribution $(\mathcal{X}, \mathcal{W})$ is efficiently samplable, Sim_ϵ samples $(x_1^*, w_1^*), (x_2^*, w_2^*), \dots (x_{p_\epsilon}^*, w_{p_\epsilon}^*)$ on its own using the sampler for $(\mathcal{X}, \mathcal{W})$.

At a high level, the simulator uses these instances to approximately-learn the verifier's challenge string e (call this approximation e_{ch}), and then generates a transcript corresponding to a random $x \sim \mathcal{X}$, by using the honest-verifier ZK simulation strategy of the underlying Σ -protocol, corresponding to verifier challenge e_{ch} .

We now describe this sequence of hybrid experiments, where hybrid $\text{Hybrid}_{\text{Sim}_\epsilon}$ corresponds to our simulator Sim_ϵ .

Proof via Hybrid Experiments.

Hybrid₀ : This hybrid corresponds to an honest prover in the real world, behaving according to [Definition 2](#). That is, the challenger C and verifier V^* execute the first two rounds where C behaves according to honest prover strategy.

Next, (C, V^*) run $\text{poly}(n)$ executions, with the same fixed first message τ_1 , but different second messages chosen potentially maliciously by V^* . In each execution, C picks a fresh sample $(x, w) \leftarrow (\mathcal{X}_n, \mathcal{W}_n)$, and generates a proof according to honest verifier strategy. Next, C again samples $(x, w) \leftarrow (\mathcal{X}_n, \mathcal{W}_n)$ and for $i \in [n]$ and outputs an honestly generated proof with first two messages τ_1, τ_2 .

Hybrid_{0,ε} :⁵ This hybrid corresponds to a challenger behaving according to [Definition 2](#), except instead of generating randomness for ZAP and OT via PRF, the challenger samples fresh randomness for every thread where the verifier sends a different second round message. That is, the challenger C and verifier V^* execute the first two rounds where C behaves according to honest prover strategy.

Next, (C, V^*) run $\text{poly}(n)$ executions, with the same fixed first message τ_1 , but different second messages chosen potentially maliciously by V^* . In each execution, C picks a fresh sample $(x, w) \leftarrow (\mathcal{X}_n, \mathcal{W}_n)$, and generates a proof according to honest verifier strategy, but using randomness sampled uniformly and independently at random. Next, C samples $(x, w) \xleftarrow{\$} (\mathcal{X}, \mathcal{W})$ and for $i \in [n]$ and outputs an honestly generated proof with uniform randomness and first two messages τ_1, τ_2 .

Lemma 13. *For all PPT distinguishers \mathcal{D}_V that obtain the view of the verifier,*

$$|\Pr[\mathcal{D}_V(\text{Hybrid}_0) = 1] - \Pr[\mathcal{D}_V(\text{Hybrid}_{0,\epsilon}) = 1]| \leq \text{negl}(n)$$

Proof. The proof of security of [Lemma 13](#) follows by security of the PRF. Given a distinguisher that distinguishes the view of the verifier between both experiments, we construct a reduction to the security of the PRF. For each execution, whenever the verifier sends a fresh message, the reduction obtains $r_{\text{wi}} \| r_{\text{OT}}$ externally from the challenger, as either outputs of the PRF or uniformly chosen randomness. Then, if there exists a distinguisher \mathcal{D}_V where $|\Pr[\mathcal{D}_V(\text{Hybrid}_0) = 1] - \Pr[\mathcal{D}_V(\text{Hybrid}_{0,\epsilon}) = 1]| \geq \frac{1}{\text{poly}(n)}$ for some polynomial $\text{poly}(\cdot)$, the reduction mirrors the output of this distinguisher and breaks security of the PRF. \square

Hybrid_{1,ε} :

This hybrid is indexed by a small error parameter $\epsilon = \frac{1}{\text{poly}(n)}$ for some polynomial $\text{poly}(\cdot)$, and proceeds as follows. The challenger C and verifier V^* execute the first two rounds where C behaves according to [Hybrid_{0,ε}](#).

Next, (C, V^*) run $\text{poly}(n)$ executions, with the same fixed first message τ_1 , but different second messages chosen potentially maliciously by V^* . In each execution, C picks a fresh sample $(x, w) \leftarrow (\mathcal{X}_n, \mathcal{W}_n)$, and generates a proof according to [Hybrid_{0,ε}](#).

Next, C samples $(x, w) \xleftarrow{\$} (\mathcal{X}, \mathcal{W})$, and for fixed first two rounds, it does the following.

1. Run the algorithm in [Figure 4](#) parameterized by $I = 1$ with oracle access to the distinguisher \mathcal{D} , and error parameter ϵ , to obtain guess $e_{\text{ch},1}$ for the first bit of the verifier challenge.
2. Next, compute $a_1 = f_1(x, w, r_1), z_1^0 = f_2(x, w, r_1, e_{\text{ch},1}), z_1^1 = f_2(x, w, r_1, e_{\text{ch},1})$.
3. For $i \in [2, n]$, compute (a_i, z_i^0, z_i^1) honestly.

⁵This hybrid does not actually depend on ϵ and is only denoted as [Hybrid_{0,ε}](#) for notational convenience.

4. Send prover message according to [Figure 3](#) using the a_i, z_i computed for $i \in [n]$.

Hybrid $_{I,\epsilon}$ for $I \in [2, n]$:

This hybrid is indexed by a small error parameter $\epsilon = \frac{1}{\text{poly}(n)}$ for some polynomial $\text{poly}(\cdot)$, and proceeds as follows. The challenger C and verifier V^* execute the first two rounds where C behaves according to **Hybrid** $_{0,\epsilon}$.

Next, (C, V^*) run $\text{poly}(n)$ executions, with the same fixed first message τ_1 , but different second messages chosen potentially maliciously by V^* . In each execution, C picks a fresh sample $(x, w) \leftarrow (\mathcal{X}_n, \mathcal{W}_n)$, and generates a proof according to **Hybrid** $_{0,\epsilon}$.

Next, C samples $(x, w) \leftarrow^{\$} (\mathcal{X}, \mathcal{W})$, and for fixed first two round transcript, it does the following.

1. Run the algorithm in [Figure 4](#) parameterized by I with oracle access to the verifier V , distinguisher \mathcal{D} , and error parameter ϵ , to obtain guess e_{ch} for the first I bits of the verifier challenge.
2. Next, for $i \in [I]$, compute $a_i = f_1(x, w, r_i)$, $z_i^0 = f_2(x, w, r_i, e_{\text{ch},i})$, $z_i^1 = f_2(x, w, r_i, e_{\text{ch},i})$.
3. For $i \in [I + 1, n]$, compute (a_i, z_i^0, z_i^1) honestly.
4. Send prover message according to [Figure 3](#) using the a_i, z_i computed for $i \in [n]$.

Lemma 14. For all $I \in [0, n - 1]$,

$$|\Pr[\mathcal{D}_{\mathcal{V}} = 1 | \text{Hybrid}_{I,\epsilon}] - \Pr[\mathcal{D}_{\mathcal{V}} = 1 | \text{Hybrid}_{I+1,\epsilon}]| \leq \frac{\epsilon}{n+1}$$

Proof. The only difference between **Hybrid** $_{I,\epsilon}$ and **Hybrid** $_{I+1,\epsilon}$ is that in **Hybrid** $_{I+1}$, $e_{\text{ch},I+1}$ is computed according to the algorithm in [Figure 4](#) and the challenger sets $a_{I+1} = f_1(x, w, r_{I+1})$, $z_{I+1}^0 = z_{I+1}^1 = f_2(x, w, r_{I+1}, e_{\text{guess},I+1})$, and then sends prover message according to [Figure 3](#).

For the fixed prover first message and fixed verifier message (which fixes OT_1), for $i \in [n]$ and a fixed prefix $e_{\text{pre}} = e_{\text{ch},[I]}$, denoting the first I bits of e_{ch} ,

- Let $\mathcal{D}_{e_{\text{pre}},0,x}$ denote the actual distribution output by the distinguisher when the challenger samples random $(x, w) \leftarrow^{\$} (\mathcal{X}, \mathcal{W})$,
 - For $j \leq I$, sets $a_j = f_1(x, w, r_j)$, $z_j^0 = z_j^1 = f_2(x, w, r_j, e_j = e_{\text{pre},j})$, and using these sends prover message according to [Figure 3](#). Here, $e_{\text{pre},j}$ denotes the j^{th} bit of e_{pre} .
 - For $j = I + 1$, sets $a_j = f_1(x, w, r_j)$, $z_j^0 = z_j^1 = f_2(x, w, r_j, e_j = 0)$, and using these sends prover message according to [Figure 3](#).
 - For $j \in [I + 2, n]$, sets $a_j = f_1(x, w, r_j)$, $z_j^0 = f_2(x, w, r_j, e_j = 0)$, $z_j^1 = f_2(x, w, r_j, e_j = 1)$, and using these sends prover message according to [Figure 3](#).

We will abuse notation and also use $\mathcal{D}_{e_{\text{pre}},0,x}$ to denote the probability that the distinguisher outputs 1 in this situation.
- Let $\mathcal{D}_{e_{\text{pre}},1,x}$ denote the actual distribution output by the distinguisher when the challenger samples random $(x, w) \leftarrow^{\$} (\mathcal{X}, \mathcal{W})$ and fresh randomness r ,
 - For $j \leq I$, sets $a_j = f_1(x, w, r_j)$, $z_j^0 = z_j^1 = f_2(x, w, r_j, e_j = e_{\text{pre},j})$, and using these sends prover message according to [Figure 3](#).

Algorithm $\mathcal{M}^{V, \mathcal{D}_V}$ to approximate the verifier's challenge upto the I^{th} bit.

- Set $p = n^2/\epsilon^3, i = 1, e_{\text{ch}} = \perp$. For fixed verifier message r ,
- While $i \leq I$, repeat:
 - Set $\mathcal{D}_0 = 0$ and for $j \in [p]$, repeat:
 1. For $k < i$, sample fresh randomness r_k and set $a_k = f_1(x_j^*, w_j^*, r_k), z_k^0 = z_k^1 = f_2(x_j^*, w_j^*, r_k, e = e_{\text{ch},k})$.
 2. Sample fresh r_i , set $a_i = f_1(x_j^*, w_j^*, r_i), z_i^0 = z_i^1 = f_2(x_j^*, w_j^*, a, \mathbf{e} = \mathbf{0}, r_i)$.
 3. For $k \in [i + 1, n]$, sample fresh randomness r_k and honestly set $a_k = f_1(x_j^*, w_j^*, r_k), z_k^0 = f_2(x_j^*, w_j^*, a, e = 0, r_k), z_k^1 = f_2(x_j^*, w_j^*, a, e = 1, r_k)$
 4. Using (a, z) computed above, send prover message according to [Figure 3](#), together with the instance x_j^* .
Set $\mathcal{D}_0 = \mathcal{D}_0 + \frac{1}{p}$ if the output of the distinguisher $\mathcal{D}_V = 1$ (w.l.o.g., we assume that the distinguisher \mathcal{D}_V outputs either 0 or 1).
 - Set $\mathcal{D}_1 = 0$ and for $j \in [p]$, repeat:
 1. For $k < i$, sample fresh randomness r_k and set $a_k = f_1(x_j^*, w_j^*, r_k), z_k^0 = z_k^1 = f_2(x_j^*, w_j^*, r_k, e = e_{\text{ch},k})$.
 2. Sample fresh r_i , set $a_i = f_1(x_j^*, w_j^*, r_i), z_i^0 = z_i^1 = f_2(x_j^*, w_j^*, a, \mathbf{e} = \mathbf{1}, r_i)$.
 3. For $k \in [i + 1, n]$, sample fresh randomness r_k and honestly set $a_k = f_1(x_j^*, w_j^*, r_k), z_k^0 = f_2(x_j^*, w_j^*, a, e = 0, r_k), z_k^1 = f_2(x_j^*, w_j^*, a, e = 1, r_k)$
 4. Using (a, z) computed above, send prover message according to [Figure 3](#), together with the instance x_j^* .
Set $\mathcal{D}_1 = \mathcal{D}_1 + \frac{1}{p}$ if the output of the distinguisher $\mathcal{D}_V = 1$.
 - Set $\mathcal{D}_w = 0$ and for $j \in [p]$, repeat:
 1. For $k < i$, sample fresh randomness r_k and set $a_k = f_1(x_j^*, w_j^*, r_k), z_k^0 = z_k^1 = f_2(x_j^*, w_j^*, r_k, e = e_{\text{ch},k})$.
 2. For $k \in [i, n]$, sample fresh randomness r_k and honestly set $a_k = f_1(x_j^*, w_j^*, r_k), z_k^0 = f_2(x_j^*, w_j^*, a, e = 0, r_k), z_k^1 = f_2(x_j^*, w_j^*, a, e = 1, r_k)$.
 3. Using (a, z) computed above, send prover message according to [Figure 3](#), together with the instance x_j^* .
Set $\mathcal{D}_w = \mathcal{D}_w + \frac{1}{p}$ if the output of the distinguisher $\mathcal{D}_V = 1$.
 - If $|\mathcal{D}_1 - \mathcal{D}_w| \leq |\mathcal{D}_0 - \mathcal{D}_w|$, set $e_{\text{ch},i} = 1$, else set $e_{\text{ch},i} = 0$.
 - Set $i = i + 1$ and go to beginning of the while loop.
- Output e_{ch} .

Figure 4: Approximately Learning the Verifier's Challenge

- For $j = I + 1$, sets $a_j = f_1(x, w, r_j), z_j^0 = z_j^1 = f_2(x, w, r_j, e_j = 1)$, and using these sends prover message according to [Figure 3](#).
- For $j \in [I + 2, n]$, sets $a_j = f_1(x, w, r_j), z_j^0 = f_2(x, w, r_j, e_j = 0), z_j^1 = f_2(x, w, r, e_j = 1, r_j)$, and using these sends prover message according to [Figure 3](#).

We will abuse notation and also use $\mathcal{D}_{e_{\text{pre}},1,x}$ to denote the probability that the distinguisher outputs 1 in this situation.

- Let $\mathcal{D}_{e_{\text{pre}},w,x}$ denote the actual distribution output by the distinguisher when the challenger samples random $(x, w) \xleftarrow{\$} (\mathcal{X}, \mathcal{W})$ and fresh randomness r ,
 - For $j \leq I$, sets $a = f_1(x, w, r_j)$, $z_j^0 = z_j^1 = f_2(x, w, r_j, e_j = e_{\text{pre},j})$, and using these sends prover message according to [Figure 3](#).
 - For $j \in [I + 1, n]$, sets $a = f_1(x, w, r_j)$, $z_j^0 = f_2(x, w, r_j, e_j = 0)$, $z_j^1 = f_2(x, w, r_j, e_j = 1)$, and using these sends prover message according to [Figure 3](#).

We will abuse notation and also use $\mathcal{D}_{e_{\text{pre}},w,x}$ to denote the probability that the distinguisher outputs 1 in this situation.

Claim 2. *Either of the following statements is true:*

- For any prefix $e_{\text{pre}} \in \{0, 1\}^I$, $e |\Pr[\mathcal{D}_{e_{\text{pre}},0,x} = 1] - \Pr[\mathcal{D}_{e_{\text{pre}},w,x} = 1]| \leq \text{negl}(n)$
- For any prefix $e_{\text{pre}} \in \{0, 1\}^I$, $e |\Pr[\mathcal{D}_{e_{\text{pre}},1,x} = 1] - \Pr[\mathcal{D}_{e_{\text{pre}},w,x} = 1]| \leq \text{negl}(n)$

Proof. This claim follows from security of the OT. Assume, for contradiction, that there exist \mathcal{V} and $\mathcal{D}_{\mathcal{V}}$ for which the claim is not true. We will use them to break receiver security of the OT. Consider a reduction \mathcal{R} that first generates all the rewinding transcripts in the same way as [Hybrid \$_{0,\epsilon}\$](#) . For the final transcript, obtains the OT receiver message from \mathcal{V} and forwards this message to the OT challenger.

The reduction also picks $(x, w) \xleftarrow{\$} (\mathcal{X}, \mathcal{W})$, $r \xleftarrow{\$} \{0, 1\}^*$ and sets $a_{I+1} = f_1(x, w, r)$, $z_{I+1}^0 = f_2(x, w, r, e = 0)$, $z_{I+1}^1 = f_2(x, w, r, e = 1)$, and sends (z_{I+1}^0, z_{I+1}^1) to the OT challenger.

The OT challenger generates either the real message $\text{OT}_2(z_{I+1}^0, z_{I+1}^1)$ corresponding to verifier input, or a simulated message $\text{OT}_2(z^*, z^*)$, for some $z^* \in \{z_0, z_1\}$. The reduction sets all other (a^i, z_0^i, z_1^i) for $i \neq (I + 1)$ according to [Hybrid \$_I\$](#) , and generates sender message accordingly.

Then, the output of distinguisher $\mathcal{D}_{\mathcal{V}}$ on input the simulated message is either distributed identically to $\mathcal{D}_{e_{\text{pre}},0,x}$ or $\mathcal{D}_{e_{\text{pre}},1,x}$ (depending upon whether z^* is 0 or 1). The reduction mirrors the output of $\mathcal{D}_{\mathcal{V}}$ and it holds that, $\Pr[\mathcal{D}_{\mathcal{V}} = 1 | \text{real OT message}] - \Pr[\mathcal{D}_{\mathcal{V}} = 1 | \text{simulated OT message}] \geq \frac{1}{\text{poly}(n)}$ for some polynomial $\text{poly}(\cdot)$, for both $z^* = z_{I+1}^0$ and $z^* = z_{I+1}^1$, which is a contradiction. \square

This claim establishes that for any prefix pre , *at least one* of the distributions $\mathcal{D}_{e_{\text{pre}},0,x}$ and $\mathcal{D}_{e_{\text{pre}},1,x}$ is negligibly close to $\mathcal{D}_{e_{\text{pre}},w,x}$.

If both $\mathcal{D}_{e_{\text{pre}},0,x}$ and $\mathcal{D}_{e_{\text{pre}},1,x}$ are $\epsilon/(n + 1)$ -close to $\mathcal{D}_{e_{\text{pre}},w,x}$, then for any value of $e_{\text{ch},I+1} \in \{0, 1\}$, $|\Pr[\mathcal{D}_{\mathcal{V}} = 1 | \text{Hybrid}_{I,\epsilon}] - \Pr[\mathcal{D}_{\mathcal{V}} = 1 | \text{Hybrid}_{I+1,\epsilon}]| \leq \epsilon/(n + 1)$ and we are done.

Therefore, for the rest of this lemma, we restrict ourselves to the case where one and only one out of $\mathcal{D}_{e_{\text{pre}},0,x}$ and $\mathcal{D}_{e_{\text{pre}},1,x}$ is $\frac{\epsilon}{n+1}$ -close to $\mathcal{D}_{e_{\text{pre}},w,x}$. In particular, this also implies that $|\mathcal{D}_{e_{\text{pre}},0,x} - \mathcal{D}_{e_{\text{pre}},1,x}| > \frac{\epsilon}{n+1}$.

If the challenger could “magically” set $e_{\text{ch},I+1}$ to 0 if $\mathcal{D}_{e_{\text{pre}},0,x}$ was close to $\mathcal{D}_{e_{\text{pre}},w,x}$, and to 1 if $\mathcal{D}_{e_{\text{pre}},0,x}$ was close to $\mathcal{D}_{e_{\text{pre}},w,x}$, then again we would have that

$$|\Pr[\mathcal{D}_{\mathcal{V}} = 1 | \text{Hybrid}_{I,\epsilon}] - \Pr[\mathcal{D}_{\mathcal{V}} = 1 | \text{Hybrid}_{I+1,\epsilon}]| \leq \epsilon/(n + 1)$$

Unfortunately, the challenger cannot magically know which distributions are close, and will therefore have to approximate these distributions to obtain an answer. We now bound the probability that the challenger’s approximation $e_{\text{ch},I}$ is incorrect conditioned on $|\mathcal{D}_{e_{\text{pre}},0,x} - \mathcal{D}_{e_{\text{pre}},1,x}| > \frac{\epsilon}{n+1}$, i.e., we show:

Claim 3.

$$\Pr\left[(e_{\text{ch},I} = b) \mid (\mathcal{D}_{e_{\text{pre},1,x}} - \mathcal{D}_{e_{\text{pre},0,x}} > \frac{\epsilon}{n+1}) \wedge (|\mathcal{D}_{\text{correct},w} - \mathcal{D}_{\text{correct},b,w}| > \frac{\epsilon}{n+1})\right] \leq \text{negl}(n)$$

Proof. We note that for the $(I+1)^{\text{th}}$ iteration of Figure 4, \mathcal{D}_0 just consists of p random samples of a distribution with mean $\mathcal{D}_{e_{\text{pre},0,x}}$, \mathcal{D}_1 just consists of p random samples of a distribution with mean $\mathcal{D}_{e_{\text{pre},1,x}}$, and \mathcal{D}_w just consists of p random samples of a distribution with mean $\mathcal{D}_{e_{\text{pre},w,x}}$.

Then, using a simple Chernoff bound, we have:

- $\Pr[(\mathcal{D}_0 > \mathcal{D}_{e_{\text{pre},0,x}}(1+\alpha)) \vee (\mathcal{D}_0 < \mathcal{D}_{e_{\text{pre},0,x}}(1-\alpha))] \leq 2 \exp^{-\frac{\alpha^2 p \mathcal{D}_0}{2}}$
- $\Pr[(\mathcal{D}_1 > \mathcal{D}_{e_{\text{pre},1,x}}(1+\alpha)) \vee (\mathcal{D}_1 < \mathcal{D}_{e_{\text{pre},1,x}}(1-\alpha))] \leq 2 \exp^{-\frac{\alpha^2 p \mathcal{D}_1}{2}}$
- $\Pr[(\mathcal{D}_w > \mathcal{D}_{e_{\text{pre},w,x}}(1+\alpha)) \vee (\mathcal{D}_w < \mathcal{D}_{e_{\text{pre},w,x}}(1-\alpha))] \leq 2 \exp^{-\frac{\alpha^2 p \mathcal{D}_w}{2}}$

Setting $\alpha = \frac{\epsilon}{2n}$, and since $p = \frac{n^2}{\epsilon^3}$, by a simple union bound we have that

$$\Pr\left[\left(|\mathcal{D}_{e_{\text{pre},0,x}} - \mathcal{D}_0| > \frac{\epsilon}{2n}\right) \vee \left(|\mathcal{D}_{e_{\text{pre},1,x}} - \mathcal{D}_1| > \frac{\epsilon}{2n}\right) \vee \left(|\mathcal{D}_{e_{\text{pre},w,x}} - \mathcal{D}_w| > \frac{\epsilon}{2n}\right)\right]$$

$\leq 6 \exp^{-\frac{1}{8\epsilon}}$. Since ϵ will always be set to $\frac{1}{\text{poly}(n)}$ for some polynomial $\text{poly}(\cdot)$,

$$\Pr\left[\left(|\mathcal{D}_{e_{\text{pre},0,x}} - \mathcal{D}_0| > \frac{\epsilon}{2n}\right) \vee \left(|\mathcal{D}_{e_{\text{pre},1,x}} - \mathcal{D}_1| > \frac{\epsilon}{2n}\right) \vee \left(|\mathcal{D}_{e_{\text{pre},w,x}} - \mathcal{D}_w| > \frac{\epsilon}{2n}\right)\right]$$

$\leq \text{negl}(n)$.

Recall that one of $\mathcal{D}_{e_{\text{pre},0,x}}$ and $\mathcal{D}_{e_{\text{pre},w,x}}$ is at least $\epsilon/(n+1)$ -far from $\mathcal{D}_{e_{\text{pre},b,w}}$, and the other is at most $\text{negl}(n)$ -far. The bit $e_{\text{ch},I}$ is estimated via $\mathcal{D}_0, \mathcal{D}_1, \mathcal{D}_w$ which each have error at most $\frac{\epsilon}{2n}$, from the corresponding $\mathcal{D}_{e_{\text{pre},0,x}}, \mathcal{D}_{e_{\text{pre},1,x}}, \mathcal{D}_{e_{\text{pre},w,x}}$. Thus,

$$\Pr\left[e_{\text{ch},I} = b \mid (|\mathcal{D}_{e_{\text{pre},1,x}} - \mathcal{D}_{e_{\text{pre},0,x}}| > \epsilon/(n+1)) \wedge (|\mathcal{D}_{e_{\text{pre},w,x}} - \mathcal{D}_{e_{\text{pre},b,w}}| > \epsilon/(n+1))\right] \leq \text{negl}(n).$$

□

This completes the proof of the lemma. □

Hybrid_{Sim,ε} : This hybrid is similar to the interaction of the simulator with the verifier and distinguisher. It is indexed by a small error parameter $\epsilon = \frac{1}{\text{poly}(n)}$ for some polynomial $\text{poly}(\cdot)$, and proceeds as follows.

The simulator Sim_ϵ and verifier V^* execute the first two rounds where Sim_ϵ behaves according to **Hybrid_{0,ε}**. Next, $(\text{Sim}_\epsilon, V^*)$ run $\text{poly}(n)$ executions, with the same fixed first message τ_1 , but different second messages chosen potentially maliciously by V^* . In each execution, Sim_ϵ picks a fresh sample $(x, w) \leftarrow (\mathcal{X}_n, \mathcal{W}_n)$, and generates a proof according to **Hybrid_{0,ε}**.

Next, Sim_ϵ samples $(x, w) \leftarrow^{\$} (\mathcal{X}, \mathcal{W})$, and for the first two round transcript that was fixed in the beginning, it does the following.

1. Run the algorithm in Figure 4 parameterized by n with oracle access to the distinguisher \mathcal{D} , and error parameter ϵ , to obtain guess e_{ch} for the entire verifier challenge (all n bits).

2. Next, for $i \in [n]$, compute (without using the witness), $a_i = f_1(x, w, e_{\text{ch},i}, r_i), z_i^0 = z_i^1 = f_2(x, w, e_{\text{ch},i}, r_i)$ and send prover message according to [Figure 3](#).

Lemma 15. $\left| \Pr[\mathcal{D}_{\mathcal{V}}(\text{Hybrid}_{n,\epsilon}) = 1] - \Pr[\mathcal{D}_{\mathcal{V}}(\text{Hybrid}_{\text{Sim},\epsilon}) = 1] \right| \leq \text{negl}(n)$

Proof. Assume, for contradiction, that there exist \mathcal{V} and $\mathcal{D}_{\mathcal{V}}$ for which the claim is not true.

We will describe a sequence of sub-hybrids where we use $\mathcal{V}, \mathcal{D}_{\mathcal{V}}$ to break hiding of the commitment scheme `com` or the IND-CPA security of the dense public-key cryptosystem.

Hybrid_a : In this sub-hybrid, the challenger \mathcal{C} first conducts the experiment identically to **Hybrid_{n,ε}**, except it obtains a public key pk_2 externally and instead of opening the second commitment to the correct value r_2 , it sets the opening $r'_2 := \text{pk}_2 \oplus \tilde{r}_2$. Note that it uses this value r'_2 in all lookahead/rewinding executions as well as the main transcript. Also note that \tilde{r}_2 is fixed in the first two rounds, at the beginning of the execution. Also, in this hybrid, the challenger uses r_1 as witness for w_i .

The view of a verifier in this experiment remains computationally indistinguishable from the view in **Hybrid_{n,ε}** because of hiding of the commitment scheme `com`. If there exists a distinguisher $\mathcal{D}_{\mathcal{V}}$ that distinguishes the view in **Hybrid_a** from the view in **Hybrid_{n,ε}**, the reduction can just mirror the output of this distinguisher, to distinguish an experiment where the commitment to r_2 is correctly opened to r_2 , from one where it is opened to a different, uniformly random r'_2 .

Hybrid_b : This next sub-hybrid is identical to **Hybrid_a**, except *after* computing e_{ch} using [Figure 4](#), the challenger changes the second set of encryptions enc_{pk_2} in the final transcript, to be computed without using the witness, using the honest-verifier ZK strategy for the Σ -protocol. Since the randomness used to compute the changed commitments is never revealed, the view of a verifier in this experiment remains computationally indistinguishable from the view in **Hybrid_a** by the IND-CPA security of the dense public key encryption scheme. That is, the reduction obtains the public key pk_2 externally along with the parts of `commit` that consist of encryptions corresponding to pk_2 which are not opened when the challenge is e_{ch} . If there exists a distinguisher $\mathcal{D}_{\mathcal{V}}$ that distinguishes the view in **Hybrid_b** from the view in **Hybrid_a**, the reduction can just mirror the output of this distinguisher, to break IND-CPA security of the encryption scheme.

Hybrid_c : In this next sub-hybrid, the challenger behaves the same was as **Hybrid_b** except that it opens r_2 honestly (instead of setting it as $\text{pk}_2 \oplus \tilde{r}_2$ for externally obtained pk_2). The view of a verifier in this experiment remains computationally indistinguishable from the view in **Hybrid_b** because of hiding of the commitment scheme `com`. The indistinguishability argument is identical to that between **Hybrid_a** and **Hybrid_{n,ε}**.

Hybrid_d : In this next sub-hybrid, the challenger uses r_2 as witness for w_i instead of using r_1 . Since the ZAP is computed with fresh randomness, the view of a verifier in this experiment remains computationally indistinguishable from the view in **Hybrid_c**. Thus, if there exists a distinguisher $\mathcal{D}_{\mathcal{V}}$ that distinguishes the view in **Hybrid_c** from the view in **Hybrid_d**, the reduction can just mirror the output of this distinguisher, to break WI of the ZAP.

Next, the challenger changes the first set of encryptions enc_{pk_1} to be computed without using the witness, corresponding to verifier challenge e_{ch} for the Σ -protocol. Note that this is possible because r_2 is correctly opened and used as a witness for the WI in all these sub-hybrids. The view

remains indistinguishable by the same sequence of hybrid arguments as Hybrid_a to Hybrid_c . This corresponds to the simulation strategy, and therefore proves the lemma. \square

Suppose the distinguisher \mathcal{D}_V has a distinguishing advantage ϵ between Hybrid_0 and $\text{Hybrid}_{\text{Sim}_\epsilon}$, then it necessarily has advantage at least $\epsilon/(n+1)$ in distinguishing one consecutive pair of hybrids between Hybrid_0 and $\text{Hybrid}_{\text{Sim}_\epsilon}$, which is a contradiction. This completes our proof. \square

Strong WI and Witness Hiding. The proof of strong witness indistinguishability and witness hiding even against weak resetting verifiers, follows from distributional weak ZK with extended simulation, identically to [JKKR17].