

Sapling Security Proof

Ariel Gabizon

Daira Hopwood

Zcash

1 Signature schemes

When we say an algorithm \mathcal{A} is *efficient*, we mean it runs in time $\text{poly}(\lambda)$ for the “security parameter” λ .

Definition 1.1. Let \mathbb{G} be a group of prime order r . A signature scheme \mathcal{S} over \mathbb{G} in the random oracle model consists of algorithms $\mathcal{S} = (\text{sign}, \text{verifySig}, \mathcal{S} = (\mathcal{S}_{\text{sign}}, \mathcal{S}_{\mathcal{R}}))$ where $\text{sign}, \text{verifySig}$ are oracle machines with access to an oracle \mathcal{R} taking as input arbitrary strings and returning uniform elements of \mathbb{F}_r . Such that the following holds.

- The set of public/verification keys $\{\text{pk}\}$ is \mathbb{G} , and the set of private keys $\{\text{sk}\}$ is \mathbb{F}_r .
- For $\text{sk} \in \mathbb{F}_r$, the verification key of sk is $\text{pk} = \text{sk} \cdot g$ for a fixed generator $g \in \mathbb{G}$.
- We have the following “zero-knowledge” property: Fix any efficient \mathcal{A} . Suppose that \mathcal{A} interacts with \mathcal{S} with two types of queries
 1. Queries x , for an arbitrary string x that are answered according to $\mathcal{S}_{\mathcal{R}}$.
 2. Queries (pk, \mathbf{m}) , answered according to $\mathcal{S}_{\text{sign}}$.

Let π_1 be the distribution of the sequence of queries and replies to \mathcal{A} . Let π_2 be the distribution of the sequence of queries and replies to \mathcal{A} when

1. \mathcal{R} takes the place of \mathcal{S}_1
2. $\text{sign}^{\mathcal{R}}(\text{sk}, \mathbf{m})$ is returned instead of $\mathcal{S}_2(\text{pk}, \mathbf{m})$ where sk is the secret key corresponding to pk .

Then the distance between π_1 and π_2 is $\text{negl}(\lambda)$.

We say \mathcal{S} is unforgeable w.r.t key randomization if the following holds. Fix any efficient \mathcal{A} . A party \mathcal{O} chooses uniform $\text{sk} \in \mathbb{F}_r$ and sends $\text{pk} = \text{sk} \cdot g$ to \mathcal{A} . \mathcal{O} also initializes an empty set T . \mathcal{A} adaptively makes $\text{poly}(\lambda)$ queries of the form (α, \mathbf{m}) . \mathcal{O} replies with $\sigma := \text{sign}(\text{pk} + \alpha \cdot g, \mathbf{m})$ and adds $(\alpha, \mathbf{m}, \sigma)$ to T .

Finally \mathcal{A} outputs $(\alpha^*, \mathbf{m}^*, \sigma^*)$. Let $\text{pk}^* := \text{pk} + \alpha^* \cdot g$. Then the probability that

1. $\text{verifySig}(\text{pk}^*, \mathbf{m}^*, \sigma^*)$, and
2. $(\alpha^*, \mathbf{m}^*, \sigma^*) \notin T$

is $\text{negl}(\lambda)$.

We assume our group \mathbb{G} has a hard DL problem; meaning that for any efficient \mathcal{A} , given uniform $g, \text{sk} \cdot g \in \mathbb{G}$ the probability of outputting sk is $\text{negl}(\lambda)$.

We define the non-malleable version of Schnorr's signature scheme:

Schnorr:

Parameters: Group \mathbb{G} of prime order s . Non-zero $g \in \mathbb{G}$.

Signing: Given message \mathbf{m} and sk ,

- Choose random $a \in \mathbb{F}_r$ and let $R := a \cdot g$
- Compute $c := \mathcal{R}(R, \text{pk}, \mathbf{m})$
- Let $u := a + c \cdot \text{sk}$.
- Define $\text{sign}^{\mathcal{R}}(\text{sk}, \mathbf{m}) := (R, u)$.

Verifying: Given $\text{pk}, \mathbf{m}, \sigma = (R, u)$, $\text{verifySig}^{\mathcal{R}}(\text{pk}, \mathbf{m}, \sigma)$ accepts iff:

- Computing $c := \mathcal{R}(R, \text{pk}, \mathbf{m})$; we have $u \cdot g = R + c \cdot \text{pk}$.

Simulating:

- $\mathcal{S}_{\mathcal{R}}(x)$ checks if x has been queried before; if so answers consistently, otherwise answers uniformly in \mathbb{F}_r and records the answer.
- $\mathcal{S}_{\text{sign}}(\text{pk}, \mathbf{m})$: Choose uniform $c, u \in \mathbb{F}_r$. Let $x := (\text{pk}, \mathbf{m}, u \cdot g - c \cdot \text{pk})$. Check if $\mathcal{S}_{\mathcal{R}}(x)$ has been defined. If so, abort. Otherwise define $\mathcal{S}_{\mathcal{R}}(x) = c$ and return (c, u) .

Remark 1.2. At times when we wish to change the parameter g we work with from default to an element h , we will use it in the subscript, e.g. $\text{sign}_h^{\mathcal{R}}(\text{sk}, \mathbf{m})$.

We refer by $\text{Schnorr}' = (\text{sign}', \text{verifySig}')$ to the Schnorr scheme where pk is omitted from the computation of c .

Theorem 1.3. Schnorr is non-forgable w.r.t randomization.

Proof. Similarly to [1], we reduce to the non-forgability of standard Schnorr (where the public key is not part of the signature & without randomization) that was proven in [2].

Suppose we are given \mathcal{A} interacting with \mathcal{O} as described above, and finally outputting $(\alpha^*, \mathbf{m}^*, \sigma^*)$. We construct \mathcal{A}' that interacts with \mathcal{O}' which is a “standard” Schnorr oracle.

That is:

1. \mathcal{O}' begins by choosing a uniform $\text{sk} \in \mathbb{F}_r$
2. \mathcal{O}' computes $\text{pk} = \text{sk} \cdot g$ and sends pk to \mathcal{A}' . \mathcal{O}' initializes an empty set T' .
3. \mathcal{A}' sends queries \mathbf{m} to \mathcal{O}' and receives replies $\sigma = \text{sign}'_{\text{sk}}(\mathbf{m})$. \mathcal{O}' adds (\mathbf{m}, σ) to T' .

4. After all queries \mathcal{A}' outputs (\mathbf{m}^*, σ^*) .

\mathcal{A}' wins if

- $\text{verifySig}'(\text{pk}, \mathbf{m}^*, \sigma^*)$, and
- $(\mathbf{m}^*, \sigma^*) \notin T'$

\mathcal{A}' will simulate (\mathcal{A}) 's interaction with \mathcal{O} using \mathcal{O}' : Given a query (α, \mathbf{m}) of \mathcal{A} , \mathcal{A}' queries \mathcal{O}' with $\mathbf{m}' := (\text{pk} + \alpha \cdot g, \mathbf{m})$, to receive reply $\sigma' = (R, u')$ - *this is a Schnorr'-signature of \mathbf{m}' with sk , and we now convert this to a Schnorr-signature of \mathbf{m} with $\text{sk} + \alpha$* . Let $c := \mathcal{R}(R, \mathbf{m}') = \mathcal{R}(R, \text{pk} + \alpha \cdot g, \mathbf{m})$. It sends $\sigma := (R, u := u' + c\alpha)$ to \mathcal{A} .

We have

$$u \cdot g = u' \cdot g + c\alpha \cdot g = R + c \cdot \text{pk} + c\alpha \cdot g = R + c \cdot (\text{pk} + \alpha \cdot g).$$

So we have $\text{verifySig}(\text{pk} + \alpha \cdot g, \mathbf{m}, \sigma)$. Also R is uniformly distributed, thus \mathcal{A}' is answering (\mathcal{A}) 's queries with the same distribution \mathcal{O} would have.

Note that the mapping $F(\alpha, \mathbf{m}, \sigma) := (\mathbf{m}', \sigma')$ where $\mathbf{m}' := (\text{pk} + \alpha \cdot g, \mathbf{m})$, $\sigma' := (R, u - c\alpha)$ is injective.

Let T be the set of tuples $(\alpha, \mathbf{m}, \sigma)$ such that \mathcal{A} queried (α, \mathbf{m}) and \mathcal{A}' answered σ . We have $T' = \{F(x)\}_{x \in T}$.

When \mathcal{A} finally outputs $x^* = (\alpha^*, \mathbf{m}^*, \sigma^*)$; \mathcal{A}' outputs $F(x^*)$. As F is injective $x^* \notin T$ implies $F(x^*) \notin T'$.

Denote $(m', \sigma') := F(x^*)$. From [2]'s results on unforgeability of Schnorr', the probability that

- $\text{verifySig}'(\text{pk}, \mathbf{m}', \sigma')$, and
- $(\mathbf{m}', \sigma') \notin T'$

is $\text{negl}(\lambda)$. Noting that $\text{verifySig}'(\text{pk}, \mathbf{m}', \sigma') \equiv \text{verifySig}(\text{pk} + \alpha \cdot g, \mathbf{m}^*, \sigma^*)$, this means that the probability that

- $\text{verifySig}(\text{pk} + \alpha \cdot g, \mathbf{m}^*, \sigma^*)$, and
- $x^* \notin T$

is $\text{negl}(\lambda)$. This is exactly what we had to prove. \square

We state the following theorem that is almost implicit in [?]

Theorem 1.4 (Extractability of Schnorr). *There is an algorithm ξ with the following property. Fix any efficient \mathcal{A} and group element $\mathbf{g} \in \mathbb{G}$. Suppose that \mathcal{A} produces w.p. γ $(\text{pk}, \mathbf{m}, \sigma)$ such that $\text{verifySig}_{\mathbf{g}}^{\mathcal{R}}(\text{pk}, \mathbf{m}, \sigma)$. Then, given the internal randomness used by \mathcal{A} in the run, ξ produces w.p. $\gamma/2$ over (\mathcal{A}) 's randomness and its own randomness $s \in \mathbb{F}_r$ such that $\text{pk} = s \cdot \mathbf{g}$. Furthermore, ξ 's running time will be $P(\lambda, 1/\gamma)$ where P is a polynomial depending on the running time of \mathcal{A} .*

2 Description of Sapling

2.1 Basic components

Functions, and their requirements:

We do not explicitly state function domains and ranges; see the spec for more details. Whenever discussing a function in the properties below, we always think of an infinite sequence of functions indexed by the security parameter λ .

1. For any fixed values g, pk, v , and for any $\epsilon \geq 0$, $\mathbf{NC}(g, pk, v, rcm)$ is ϵ -close to uniform when rcm is ϵ -close to uniform.
2. \mathbf{NC} is collision resistant - i.e. the probability of finding $note, note'$ such that $\mathbf{NC}(note) = \mathbf{NC}(note')$ is $\text{negl}(\lambda)$.¹
3. For any fixed v and any $\epsilon \geq 0$, $\mathbf{VC}(v, rcv)$ is ϵ -close to uniform whenever rcv is ϵ -close to uniform.
4. \mathbf{VC} is collision-resistant.
5. **sighash** is collision-resistant.
6. **IVK** is collision-resistant.
7. **NF** is collision resistant; and has the structural property that $\mathbf{NF}(nk, note, pos) = \mathcal{R}(nk, \mathbf{NC}(note), pos)$ for some function \mathcal{R} (see another requirement for the indistinguishability property in Section 3.3).

Generators of \mathbb{G} We assume we are given generators g_{sig}, g_n, g_r, g_v that were sampled in a way that except w.p $\text{negl}(\lambda)$ an efficient \mathcal{A} cannot discover the discrete log relation between any two of them.

Statements:

OUT(cv, cm, epk): I know $note = (g, pk, v, rcm), rcv, esk$ such that

1. $cm = \mathbf{NC}(note)$.
2. $cv = \mathbf{VC}(v, rcv)$.
3. $epk = esk \cdot g$.
4. g has order greater than eight.

SPEND(rt, cv, nf, rk): I know $path, pos, note = (g, pk, v, rcm), cm, rcv, \alpha, ak, nsk$ such that

1. $cm = \mathbf{NC}(note)$.

¹A caveat here is that this is true when the rcm parameter is thought of as a field element; in the actual circuit it is received as a string of bits where some elements of \mathbb{F}_r have multiple representations; inspection of the proof shows that it suffices that CR w.r.t rcm as a field element; same story with rcv in \mathbf{VC} .

2. Either $v = 0$ (“dummy note”); or path is a merkle path from cm at position pos to rt .
3. $\text{rk} = \text{ak} + \alpha \cdot \text{g}_{\text{sig}}$.
4. Setting $\text{nk} := \text{nsk} \cdot \text{g}_{\text{n}}$, $\text{ivk} := \text{IVK}(\text{ak}, \text{nk})$; we have $\text{pk} = \text{ivk} \cdot \text{g}$.
5. $\text{nf} = \text{NF}(\text{nk}, \text{cm}, \text{pos})$

Components

A *note* is a tuple $\text{note} = (\text{g}, \text{pk}, v, \text{rcm})$ where

1. $\text{g}, \text{pk} \in \mathbb{G}$.
2. $v, \text{rcm} \in \mathbb{F}_r$.
3. $v \leq \text{MAX}$.

An *output base* $\text{output} = (\text{g}, \text{pk}, v)$ is the same as a note excluding the rcm component.

Remark 2.1. *It is convenient for us to define a note with g rather than its GH-preimage d as in the spec, as this is what’s given as input to the circuits; there are minor non-exploitable issues with this, see e.g. <https://github.com/zcash/zcash/issues/3490>.*

For $\text{ivk} \in \mathbb{F}_r$ we say *note belongs to ivk* if $\text{pk} = \text{ivk} \cdot \text{g}$.

An *input base*, usually denoted input , will consist of the values required to make an input in a Sapling transaction, except the spending key; namely $\text{input} = (\text{note}, \text{path}, \text{pos}, \text{pak})$ where

- note is a note
- path is a path in a merkle tree beginning from a leaf of value $\text{cm} = \text{NC}(\text{note})$.
- pos is the position of cm amongst the leaves of the Merkle tree (pos is redundant here as it can be derived from path , but convenient).
- pak is a proving key to make SNARK spend proofs about the note.

We say input is *consistent with* rt if path ends at rt .

A *transaction input*, usually denoted inp , is the final form an input appears in a transaction; inp consists of

1. A value commitment cv .
2. A nullifier nf .
3. A Merkle root rt of the tree containing the used note.
4. A public key rk that is (allegedly) a randomized version of the spent note’s proving key ak .
5. A SNARK proof π for the statement $\text{SPEND}(\text{rt}, \text{cv}, \text{nf}, \text{rk})$.

2.2 Methods

We use the convention that ℓ denotes the number of inputs in a transaction, and s the number of outputs.

makeinput($rt, \text{input} = (\text{note}, \text{path}, \text{pos}, \text{pak}), \text{rcv}, \alpha$)
 where input is an input base consistent with rt .

1. $\text{cm} = \mathbf{NC}(\text{note})$
2. $\text{nf} = \mathbf{NF}(\text{nk}, \text{note}, \text{pos})$
3. Define $\text{rk} := \text{ak} + \alpha \cdot \mathbf{g}$, $\text{cv} := \mathbf{v} \cdot \mathbf{g}_v + \text{rcv} \cdot \mathbf{g}_r$.
4. Let $\pi = \pi_{\text{spend}}(\text{cv}, \text{rt}, \text{nf}, \text{rk}; \text{note}, \text{pak}, \alpha, \text{path}, \text{pos})$.
5. Output $(\text{cv}, \text{rt}, \text{nf}, \text{rk}, \pi)$.

makeoutput($\text{note} = (\mathbf{g}, \text{pk}, \mathbf{v}, \text{rcm}), \text{rcv}$),

1. Choose random $\text{esk} \in \mathbb{F}_r$.
2. Let $\text{cv} := \mathbf{VC}(\mathbf{v}, \text{rcv}) = \mathbf{v} \cdot \mathbf{g}_v + \text{rcv} \cdot \mathbf{g}_r$.
3. Let $\text{note} = (\mathbf{g}, \text{pk}, \mathbf{v}, \text{rcm})$ and $\text{cm} := \mathbf{NC}(\text{note})$.
4. Let $\text{epk} = \text{esk} \cdot \mathbf{g}$.
5. Let $\text{enc} = \mathbf{ENC}_{\mathbf{KDF}(\text{esk} \cdot \text{pk}, \text{epk})}(\text{note})$
6. Let $\pi = \pi_{\text{output}}(\text{epk}, \text{cm}, \text{cv}; \text{note}, \text{rcv}, \text{esk})$.
7. Output $(\text{cv}, \text{cm}, \text{epk}, \pi, \text{enc})$

makerandomizedoutput($\text{note} = (\mathbf{g}, \text{pk}, \mathbf{v}), \text{rcv}$),

1. Choose random $\text{esk}, \text{rcm} \in \mathbb{F}_r$.
2. Let $\text{cv} := \mathbf{VC}(\mathbf{v}, \text{rcv}) = \mathbf{v} \cdot \mathbf{g}_v + \text{rcv} \cdot \mathbf{g}_r$.
3. Let $\text{note} = (\mathbf{g}, \text{pk}, \mathbf{v}, \text{rcm})$ and $\text{cm} := \mathbf{NC}(\text{note})$.
4. Let $\text{epk} = \text{esk} \cdot \mathbf{g}$.
5. Let $\text{enc} = \mathbf{ENC}_{\mathbf{KDF}(\text{esk} \cdot \text{pk}, \text{epk})}(\text{note})$
6. Let $\pi = \pi_{\text{output}}(\text{epk}, \text{cm}, \text{cv}; \text{note}, \text{rcv}, \text{esk})$.
7. Output $(\text{cv}, \text{cm}, \text{epk}, \pi, \text{enc})$

bindval($\text{raw}_{\text{tx}} = (\vec{\text{inp}}, \vec{\text{out}}, \mathbf{v}^{\text{bal}}), \vec{\text{rcv}}$)

1. Let $r := \sum_{i=1}^{\ell} \text{rcv}_i - \sum_{i=\ell+1}^{\ell+s} \text{rcv}_i$
2. Let $S := \sum_{i=1}^{\ell} \text{cv}_i - \sum_{i=\ell+1}^{\ell+s} \text{cv}_i - \mathbf{v}^{\text{bal}} \cdot \mathbf{g}_v$

3. Let $\sigma_{\text{bind}} := \text{sign}_{g_r}(r, \text{sighash}(\text{raw}_{\text{tx}}))$.

4. Output $\text{pre-tx} := (\text{raw}_{\text{tx}}, \sigma_{\text{bind}})$.

$\text{signtx}(\text{pre-tx} = (\text{raw}_{\text{tx}}, \sigma_{\text{bind}}), \vec{\text{ask}}, \vec{\alpha})$

1. For each $i \in [\ell]$, let $\sigma_i := \text{sign}_{g_{\text{sig}}}(\text{ask}_i + \alpha_i, \text{sighash}(\text{raw}_{\text{tx}}))$

2. Let $\vec{\sigma} := (\sigma_1, \dots, \sigma_\ell)$.

3. Output $(\text{raw}_{\text{tx}}, \vec{\sigma})$.

Given $(\text{rt}, v^{\text{bal}})$ we say $(\vec{\text{input}}, \vec{\text{output}})$ is *consistent* with $\text{rt}, v^{\text{bal}}$, if

- for each $j \in [\ell]$ input_j is consistent with rt , i.e. pak_j is from $\text{NC}(\text{note}_j)$ to rt ,
- $\sum_{j=1}^{\ell} v_j - \sum_{j=\ell+1}^{\ell+s} v_j = v^{\text{bal}}$.
- the positions $\{\text{pos}_j\}_{j \in [\ell]}$ are all distinct.

and

$\text{makerandomizedtx}(\text{rt}, v^{\text{bal}}, \vec{\text{input}}, \vec{\text{output}})$

where $\text{input}_j = (\text{note}_j, \text{pak}_j, \text{path}_j, \text{pos}_j)$, $\text{output}_j = (g_j, \text{pk}_j, v_j)$

1. Choose random $\vec{\text{rcv}} \in \mathbb{F}_r^s$.

2. For $j \in [\ell]$, $\text{inp}_j = \text{makeinput}(\text{rt}, \text{input}_j, \text{rcv}_j)$

3. For $j \in [s]$, $\text{out}_j = \text{makeoutput}(\text{output}_j, \text{rcv}_j)$

4. $\text{pre-tx} = \text{bindval}(\vec{\text{inp}}, \vec{\text{out}}, v^{\text{bal}})$.

5. Choose random $\vec{\alpha} \in \mathbb{F}_r^\ell$.

6. Output $\text{tx} = \text{signtx}(\text{pre-tx}, \text{ask}, \vec{\alpha})$

$\text{maketx}(\vec{\text{input}}, \vec{\text{output}}, \vec{\text{rcv}}, \text{ask}, \text{pak})$ where $\text{input}_j = (v_j, \text{note}_j, \text{pak}_j, \text{path}_j, \text{pos}_j)$, $\text{output}_j = (g_j, \text{pk}_j, v_j, \text{rcm}_j)$

1. Choose random $\vec{\alpha} \in \mathbb{F}_r^\ell$.

2. For $j \in [\ell]$, $\text{inp}_j = \text{makeinput}(\text{input}_j, \text{rcv}_j, \alpha_j, \text{pak})$

3. For $j \in [s]$, $\text{out}_j = \text{makeoutput}(\text{output}_j, \text{rcv}_j)$

4. Let $v^{\text{bal}} := \sum_{i=1}^{\ell} v_i - \sum_{j=\ell+1}^{\ell+s} v_j$.

5. $\text{pre-tx} = \text{bindval}(\vec{\text{inp}}, \vec{\text{out}}, v^{\text{bal}}, \vec{\text{rcv}})$.

6. Let $\text{tx} = \text{signtx}(\text{pre-tx}, \vec{\alpha}, \text{ask})$

$\text{verify-tx}(\text{L}, \text{tx})$

1. Suppose that $\text{tx} = (\text{raw}_{\text{tx}}, \vec{\sigma})$.
2. For each $\text{inp}_i = (\text{rt}, \text{cv}, \text{nf}, \text{rk}, \pi) \in \vec{\text{inp}}(\text{tx})$,
 - Check that $\text{nf} \notin \text{nf}(\text{L}) \cup \{\text{nf}(\text{inp}_1), \dots, \text{nf}(\text{inp}_{i-1})\}$.
 - Check that $\text{spendverify}(\text{rt}, \text{cv}, \text{nf}, \text{rk}; \pi)$.
 - Check that $\text{verifySig}_{\text{g}_{\text{sig}}}^{\mathcal{R}}(\text{rk}, \text{sighash}(\text{raw}_{\text{tx}}), \sigma_i)$
3. For each $\text{out} = (\text{cv}, \text{cm}, \text{epk}, \pi, \text{enc}) \in \vec{\text{out}}(\text{tx})$, check that $\text{outverify}(\text{cv}, \text{cm}, \text{epk}; \pi)$
4. Let $S := \sum_{i=1}^{\ell} \text{cv}_i - \sum_{i=\ell+1}^{\ell+s} \text{cv}_i - v^{\text{bal}} \cdot \text{g}_v$.
5. Check that $\text{verifySig}_{\text{g}_r}^{\mathcal{R}}(S, \text{sighash}(\text{raw}_{\text{tx}}), \sigma_{\text{bind}})$.

3 Non-Malleability of Sapling w.r.t. delegated spenders

We make the simplifying assumption when modelling non-malleability in this writeup; that *there is only one spending key* (ask, nsk) *of the honest signer involved, and all addresses are diversified addresses derived from this spending key.*

Modelling the adversary:

We wish to show that the delegated spender cannot create any new transactions of her own. We model this claim with the following non-malleability game: We model the honest signer as an oracle \mathcal{O} that \mathcal{A} interacts with as follows.

\mathcal{O} begins by choosing a new spending key $(\text{ask}, \text{nsk}) \leftarrow \mathcal{K}$ and sending the corresponding proof authorizing key $\text{pak} = (\text{ak}, \text{nsk})$ to \mathcal{A} . Where $\text{ak} = \text{ask} \cdot \text{g}$.

Afterwards, \mathcal{A} can make **sign-all-inputs** queries to \mathcal{O} , which intuitively correspond to asking for signatures on transactions whose inputs have spending key (ask, nsk) (though see remark).

Sign-all-inputs queries

1. \mathcal{A} sends $(\text{pre-tx} = (\text{raw}_{\text{tx}}, \sigma_{\text{bind}}), \vec{\alpha})$ to \mathcal{O} . Where $\text{raw}_{\text{tx}} = (\vec{\text{inp}}, \vec{\text{out}}, v^{\text{bal}})$
2. \mathcal{O} checks if $\text{spendverify}(\text{pub}_i, \pi_i)$ holds for each $i \in [\ell]$ and otherwise aborts.
3. \mathcal{O} computes for $i \in [\ell]$, $\sigma_i = \text{sign}_{\text{g}}(\text{ask} + \alpha_i, \text{sighash}(\text{raw}_{\text{tx}}))$.
4. Let $\vec{\sigma} := (\sigma_1, \dots, \sigma_{\ell})$. \mathcal{O} return $\text{tx} := (\text{raw}_{\text{tx}}, \sigma_{\text{bind}}, \vec{\sigma})$.

Remark 3.1. *The second item is another way of saying we assume \mathcal{A} can only ask \mathcal{O} for signatures of transactions with legitimate spend proofs. Otherwise the proof currently fails as we need to be able to extract the witness from each input.*

Terminology: We refer below to a transaction tx as $\text{tx} = (\text{raw}_{\text{tx}}, \sigma_{\text{bind}}, \vec{\sigma})$, where $\vec{\sigma}$ contains the ℓ input signatures and σ_{bind} is as described above in `maketx` that are added during `sign-all-inputs` and the signature σ_{bind} added in the last step of `maketx`.

Non-malleability says, \mathcal{A} should not be able to create a new valid transaction with inputs belonging to \mathcal{O} , even after seeing transactions of its choice with inputs of \mathcal{O} . New will mean that the `rawtx` part will be new. (If we had changed the signature scheme to sign in order and have each signature sign the previous ones we could have required that tx including the signature part must be different from all previous transactions).

The way we formalize “transaction with inputs of \mathcal{O} ” is that the transaction created by \mathcal{A} contains overlapping nullifiers with the transactions signed previously by \mathcal{O} ; precisely transactions that are outputs of `sign-all-inputs` queries.

Remark 3.2. *A somewhat odd thing about the construction with the delegated spender, is that valid transactions signed by \mathcal{O} , do not exactly correspond to transactions whose inputs \mathcal{O} knows the spending key of. We can only say \mathcal{O} and \mathcal{A} together know the spending key. For example, given (ak, nsk) , \mathcal{A} can choose random $s \in \mathbb{F}_r$, set $\text{ak}' := \text{ak} + s \cdot \mathbf{g}$. Now when \mathcal{A} wants to sign an input in address ak' , i.e. with some randomized key $\text{rk} = \text{ak}' + \alpha \mathbf{g} = \text{ak} + (s + \alpha) \cdot \mathbf{g}$, it can give \mathcal{O} the randomization $\alpha' = s + \alpha$.*

A way to avoid these oddities is to have \mathcal{O} only sign transactions where he recognizes the nullifiers as belonging to a note of his. For our purposes here, we get a stronger result without this restriction by showing non-malleability holds when \mathcal{O} signs any transaction.

Some more terminology Given a validly formatted transaction $\text{tx} = ((\vec{\text{inp}}, \vec{\text{out}}, \text{v}^{\text{bal}}), \sigma_{\text{bind}}, \vec{\sigma})$, we define

- $\text{nf}(\text{tx})$ to be the set of nullifiers appearing in one of its inputs; so $\text{nf}(\text{tx}) := \{\text{nf}(\text{inp})\}_{\text{inp} \in \vec{\text{inp}}}$.
- $\text{rk}(\text{tx})$ the set of randomized public keys appearing in inputs of tx , so $\text{rk}(\text{tx}) := \{\text{rk}(\text{inp})\}_{\text{inp} \in \vec{\text{inp}}}$.
- $\text{raw}(\text{tx}) := (\vec{\text{inp}}, \vec{\text{out}}, \text{v}^{\text{bal}})$. For a set T of validly formed transactions we define $\text{raw}(T) := \{\text{raw}(\text{tx})\}_{\text{tx} \in T}$.

Claim 3.3 (Non-malleability w.r.t delegated spenders). *Fix any efficient \mathcal{A} interacting with \mathcal{O} as described above. Let $T = \{\text{tx}'\}$ be the set of transactions that are replies of \mathcal{O} to \mathcal{A} ’s `sign-all-inputs` queries. The probability that \mathcal{A} manages to output a ledger L and transaction tx such that*

1. $\text{verify-tx}(L, \text{tx}) = \text{acc}$,
2. $\text{raw}(\text{tx})$ is not a prefix of an element of T .
3. $\text{nf}(\text{tx}) \cap \text{nf}(\text{tx}') \neq \emptyset$ for some $\text{tx}' \in T$.

is $\text{negl}(\lambda)$.

Proof. Let \mathcal{A} be an algorithm that after interacting with \mathcal{O} as described above outputs L, tx . Let ϵ be the probability that L, tx satisfy the above.

We construct \mathcal{A}' that receives a randomized forgery challenge for Schnorr as described in Definition 1.1, and with probability $\epsilon - \text{negl}(\lambda)$ either

- outputs a collision of **sighash**
- outputs a collision of **NF**,
- outputs a collision of **NC**,
- outputs a collision of **IVK**,
- Constructs a signature forgery for Schnorr w.r.t randomization.

Then, from CR of **sighash**, **NF**, **NC**, **IVK** and Theorem 1.3 the claim follows.

\mathcal{A}' works as follows:

1. \mathcal{A}' will receive a challenge \mathbf{ak} for the signature scheme Schnorr from a party \mathcal{O} .
2. \mathcal{A}' chooses random $\mathbf{nsk} \in \mathbb{G}$ and sends to \mathcal{A} the proof authorizing key $\mathbf{pak} = (\mathbf{nsk}, \mathbf{ak})$ *note that here we need to make a spending key that is not from the same seed \mathbf{sk} — ariel gabizon*
3. When \mathcal{A} makes a sign-all-inputs query $(\mathbf{raw}_{\mathbf{tx}}, \vec{\alpha})$ \mathcal{A}' first checks that the proofs in $\mathbf{raw}_{\mathbf{tx}}$ are valid (as \mathcal{O} does in the description of sign-all-inputs queries) and then answers with $\vec{\sigma}$ where $\sigma_i := \mathcal{S}_{\text{sign}}(\mathbf{pk} + \alpha \cdot g, \mathbf{m})$. If during invocations to $\mathcal{S}_{\text{sign}}$, $\mathcal{S}_{\mathcal{R}}$ is queried on a point on which \mathcal{A} queried \mathcal{R} , \mathcal{A}' aborts. (Note that the point queried by $\mathcal{S}_{\mathcal{R}}$ is $(R, \mathbf{pk}, \mathbf{m})$ for a uniform R chosen only during the execution of $\mathcal{S}_{\text{sign}}$, so the probability such a point was already queried is $\text{negl}(\lambda)$.)
4. When \mathcal{A}' makes a query to \mathcal{R} , \mathcal{A} answers according to \mathcal{R} unless the query has been answered according to $\mathcal{S}_{\mathcal{R}}$ during invocations of $\mathcal{S}_{\text{sign}}$ in sign-all-inputs queries; in which case \mathcal{A}' answers according to $\mathcal{S}_{\text{sign}}$. (This doesn't change the distribution of \mathcal{R} from the perspective of \mathcal{A} .)
5. When \mathcal{A} outputs L, \mathbf{tx} : \mathcal{A}' checks that it indeed satisfies the challenge - that is $\text{verify-tx}(L, \mathbf{tx})$; \mathbf{tx} contains an input \mathbf{inp} with $\mathbf{nf} = \mathbf{nf}(\mathbf{inp})$ being equal to $\mathbf{nf}(\mathbf{inp}')$ for some $\mathbf{inp}' \in \mathbf{tx}'$ for some $\mathbf{tx}' \in T$; appearing in one of the sign-all-inputs queries of \mathcal{A} ; and $\mathbf{raw}_{\mathbf{tx}} \notin \mathbf{raw}(T)$. If not \mathcal{A}' aborts.
6. \mathcal{A}' checks if $\mathbf{sighash}(\mathbf{raw}_{\mathbf{tx}}) = \mathbf{sighash}(\mathbf{raw}_{\mathbf{tx}}')$ for some $\mathbf{tx}' \in T$ with $\mathbf{raw}_{\mathbf{tx}'} \neq \mathbf{raw}_{\mathbf{tx}}$. If so it outputs $(\mathbf{raw}_{\mathbf{tx}}, \mathbf{raw}_{\mathbf{tx}}')$ as a collision of **sighash**.
7. Let $R := \{\mathbf{rk}_1, \dots, \mathbf{rk}_\ell\}$ be the randomized public keys in the inputs of \mathbf{tx} , and $R' := \{\mathbf{rk}'_1, \dots, \mathbf{rk}'_{\ell'}\}$ be the randomized public keys in the inputs of \mathbf{tx}' . \mathcal{A}' checks if $R' \cap R \neq \emptyset$, i.e. $\mathbf{rk}_i = \mathbf{rk}'_j$ for some i, j . In this case it outputs the forgery $(\alpha'_j, \mathbf{sighash}(\mathbf{raw}_{\mathbf{tx}}), \sigma_i)$.
8. Otherwise let ξ be the extractor guaranteed to exist for the combined party $\mathcal{A}', \mathcal{O}$ up to the point in step 5 where \mathcal{A} outputted \mathbf{tx} . Except with probability $\text{negl}(\lambda)$, ξ outputs a witness $\mathbf{w} = (\text{note}, \mathbf{pak} = (\mathbf{ak}, \mathbf{nsk}), \alpha, \text{path}, \text{pos})$. Similarly there is an extractor ξ' for the input \mathbf{inp}' in \mathbf{tx}' giving us a witness $\mathbf{w}' = (\text{note}', \mathbf{pak}' = (\mathbf{ak}', \mathbf{nsk}'), \alpha', \text{path}', \text{pos}')$.
9. Let $\mathbf{nk} := \mathbf{nsk} \cdot g, \mathbf{nk}' := \mathbf{nsk}' \cdot g$. We have

$$\mathbf{NF}(\mathbf{nk}, \mathbf{cm}, \text{pos}) = \mathbf{NF}(\mathbf{nk}', \mathbf{cm}', \text{pos}') = \mathbf{nf}$$

If $\mathbf{nk} \neq \mathbf{nk}'$, $\mathbf{cm} \neq \mathbf{cm}'$ or $\text{pos} \neq \text{pos}'$ \mathcal{A}' outputs $(\mathbf{nk}, \mathbf{cm}, \text{pos}), (\mathbf{nk}', \mathbf{cm}', \text{pos}')$ as a collision of **NF**.

10. Otherwise, we have $\mathbf{NC}(\text{note}) = \mathbf{NC}(\text{note}') = \text{cm}$. If $\text{note} \neq \text{note}'$, \mathcal{A}' outputs $\text{note}, \text{note}'$ as a collision of \mathbf{NC} .
11. Otherwise we have $\text{note} = \text{note}' = (g, \text{pk}, v, \text{rcm})$. Defining $\text{ivk} := \mathbf{IVK}(\text{ak}, \text{nk}), \text{ivk}' := \mathbf{IVK}(\text{ak}', \text{nk})$, we have $\text{pk} = \text{ivk} \cdot g = \text{ivk}' \cdot g$. Thus, $\text{ivk} = \text{ivk}'$. (Important here that ivk representation is unique and it is cause dfn of \mathbf{IVK} has $\text{mod } 2^{\ell_{\text{ivk}}=251}$.) If $\text{ak} \neq \text{ak}'$, \mathcal{A}' outputs $(\text{ak}, \text{nk}), (\text{ak}', \text{nk})$ as a collision of \mathbf{IVK} .
12. Otherwise, we have $\text{ak} = \text{ak}'$. Now, \mathcal{A} knows α^* such that $\text{rk}' = \text{ak}^* + \alpha^* \cdot g$, where ak^* is from the forgery challenger (as he used $(\alpha^*, \text{sighash}(\text{raw}_{\text{tx}}'))$ in the **sign-all-inputs** query for tx' for input inp'). And also $\text{rk}' = \text{ak}' + \alpha' \cdot g$. So $\text{ak} = \text{ak}' = \text{ak}^* + (\alpha^* - \alpha') \cdot g$. And $\text{rk} = \text{ak}^* + (\alpha^* - \alpha' + \alpha) \cdot g$. Thus, in this case \mathcal{A}' outputs $(\alpha^* - \alpha' + \alpha, \sigma, \text{sighash}(\text{raw}_{\text{tx}}))$ as a signature forgery.

□

3.1 Modelling the outside adversary

$\text{maketransaction}(\overrightarrow{\text{input}}, \overrightarrow{\text{output}}, v^{\text{bal}}, \text{ask})$:

1. For $i \in [\ell]$, check where input_i appears in the compute $\text{input}_i := \text{makeinput}(\text{inp}_i)$.
2. Check that $\sum_{i=1}^{\ell} v_i - \sum_{j=1}^s \text{ov}_j = v^{\text{bal}}$. If this is the case then $(\overrightarrow{\text{input}}, \overrightarrow{\text{output}}, v^{\text{bal}})$ is a *valid input* to maketransaction . Otherwise output rej and abort.
3. For $j \in [s]$, compute $\text{output}_j := \text{makeoutput}(\text{out}_j)$.
4. Choose uniform $\overrightarrow{\text{rcv}} \in \mathbb{F}_r^s$ and output $\text{sign-all-inputs}(\text{maketx}(\overrightarrow{\text{input}}, \overrightarrow{\text{output}}, v^{\text{bal}}, \overrightarrow{\text{rcv}}), \text{ask})$

\mathcal{O} begins by choosing a uniform spending key $(\text{ask}, \text{nsk}) \in \mathcal{K}$. And generates the corresponding keys $\text{ak} := \text{ask} \cdot g, \text{nk} := \text{nsk} \cdot g, \text{ivk} := \mathbf{IVK}(\text{ak}, \text{nk})$. \mathcal{A} and \mathcal{O} initialize an empty set T of diversified addresses. \mathcal{A} and \mathcal{O} initialize an empty set of current notes N . \mathcal{A} can make two kinds of queries.

Get new diversified address queries

- \mathcal{O} chooses $g \in \mathbb{G}$ according to the distribution GH of the group hash output.
- \mathcal{O} then outputs the diversified address $(g, \text{pk} := \text{ivk} \cdot g)$.
- \mathcal{A} and \mathcal{O} add (g, pk) to the set of diversified addresses T .

Make transaction queries

- \mathcal{A} chooses extended input notes $\text{input}_1, \dots, \text{input}_{\ell} \in N$.
- \mathcal{A} chooses output notes $\text{output}_1, \dots, \text{output}_s$, with $\text{output}_j = (g_j, \text{pk}_j, v_j, \text{rcm}_j)$, for $(g_j, \text{pk}_j) \in T$.
- \mathcal{A} Chooses uniform $\overrightarrow{\text{rcv}} \in \mathbb{F}_r^{s+\ell}$.
- \mathcal{A} sends $(\overrightarrow{\text{input}}, \overrightarrow{\text{output}}, \overrightarrow{\text{rcv}})$ to \mathcal{O} .

- \mathcal{O} returns $\text{tx} := \text{maketx}(\overrightarrow{\text{input}}, \overrightarrow{\text{output}}, \overrightarrow{\text{rcv}}, \text{ask}, \text{pak})$ to \mathcal{A} .
- \mathcal{A} and \mathcal{O} add output_j to \mathbf{N} for each j s.t. $(\mathbf{g}_j, \mathbf{pk}_j) \in T$.
- \mathcal{A} and \mathcal{O} remove the elements of $\overrightarrow{\text{input}}$ from \mathbf{N} .

3.2 Non-malleability w.r.t outside adversary

Claim 3.4. *Suppose \mathcal{A} interacts with \mathcal{O} as described above. Then it outputs a transaction $(\mathbf{L}, \text{tx}, \text{inp} \in \text{tx})$. The probability that there exists $\text{tx}' \in T, \text{inp}' \in \overrightarrow{\text{input}}(\text{tx}')$ such that*

1. $\text{raw}(\text{tx}) \neq \text{raw}(\text{tx}'), \forall \text{tx}' \in T$.
2. $\text{nf}(\text{inp}) = \text{nf}(\text{inp}')$.

is $\text{negl}(\lambda)$.

Proof. Reduce to Claim 3.3. We □

3.3 Indistinguishability w.r.t outside adversaries

For a sequence of random variables X_1, \dots, X_n it will be convenient in this section to denote $X_{<i} := (X_1, \dots, X_{i-1})$. Let us say that random variables X, Y are γ -independent if for any events A, B

$$|\Pr(X \in A \wedge Y \in B) - \Pr(X \in A) \cdot \Pr(Y \in B)| \leq \gamma.$$

We recall that the *statistical distance* between X and Y is the maximum over all events T of

$$|\Pr(X \in T) - \Pr(Y \in T)|.$$

We say X, Y are γ -close if they have statistical distance at most γ .

A calculation proves

Claim 3.5. *Suppose $X = (X_1, X_2), Y = (Y_1, Y_2)$ are such that*

- X_1 and Y_1 are on the same range, are γ_1 -independent and γ_1 -close.
- e.w.p γ_2 over the value (x_1, y_1) of (X_1, Y_1) , $(X|X_1 = x_1), (Y|Y_1 = y_1)$ are γ_3 -independent.
- e.w.p γ_2 over the value x_1 of X_1 , we have that $(X|X_1 = x_1), (Y|Y_1 = x_1)$ are γ_3 -close.

Then X, Y are $\gamma_1 + \gamma_2 + \gamma_3$ -independent and $\gamma_1 + \gamma_2 + \gamma_3$ -close.

Induction then shows that

Claim 3.6. *Suppose $t = \text{poly}(\lambda)$. Suppose random variables $X = (X_1, \dots, X_t), Y = (Y_1, \dots, Y_t)$ are such that for any $i \in [n]$,*

- e.w.p $\text{negl}(\lambda)$ over the value (x, y) of $(X_{<i}, Y_{<i})$,
 $(X_i|X_{<i} = x)$ and $(Y_i|Y_{<i} = y)$ are $\text{negl}(\lambda)$ -independent; and
- e.w.p $\text{negl}(\lambda)$ over the value x of $X_{<i}$, $(X_i|X_{<i} = x)$ and $(Y_i|Y_{<i} = x)$ are $\text{negl}(\lambda)$ -close.

Then X, Y are $\text{negl}(\lambda)$ -independent and $\text{negl}(\lambda)$ -close.

Below we use \mathcal{R}_{sig} to denote the random oracle used by the signature algorithm.

Theorem 3.7. *Assume that*

1. $\mathbf{NF}(\text{nk}, \mathbf{NC}(\text{note}), \text{pos}) = \mathcal{R}(\text{nk}, \mathbf{MPH}(\text{note}, \text{pos}))$ where \mathcal{R} is a random oracle and \mathbf{MPH} is a collision-resistant function²
2. \mathbf{KDF} and \mathcal{R}_{sig} are also random oracles.
3. $\mathbf{ENC}_K(m)$ produces a uniform output when K is uniform and m is fixed.
4. The SNARK we are using is witness indistinguishable - i.e. the proof distribution depends only on the public input and not on the witness.

Then, the probability of an efficient \mathcal{A} finding $\text{rt}, v^{\text{bal}}, \overrightarrow{\text{input}}, \overrightarrow{\text{output}}, \overrightarrow{\text{input}'}, \overrightarrow{\text{output}'}$ such that

- $|\overrightarrow{\text{input}}| = |\overrightarrow{\text{input}}'| = \ell, |\overrightarrow{\text{output}}| = |\overrightarrow{\text{output}}'| = s$.
- The positioned notes in $\overrightarrow{\text{input}}$ and $\overrightarrow{\text{input}}'$ are all distinct.
- $(\overrightarrow{\text{input}}, \overrightarrow{\text{output}})$ and $(\overrightarrow{\text{input}}', \overrightarrow{\text{output}}')$ are both consistent with $\text{rt}, v^{\text{bal}}$.
- The distributions of the random variables $D := \text{makerandomizedtx}(\text{rt}, v^{\text{bal}}, \overrightarrow{\text{input}}, \overrightarrow{\text{output}})$ and $D' := \text{makerandomizedtx}(\text{rt}, v^{\text{bal}}, \overrightarrow{\text{input}}', \overrightarrow{\text{output}}')$, over the randomness of the oracles $\mathcal{R}, \mathbf{KDF}$ and \mathcal{R}_{sig} , and the inner randomness of the signer, SNARK prover and the `makerandomizedtx` method, are not $\text{negl}(\lambda)$ -close and $\text{negl}(\lambda)$ -independent

is $\text{negl}(\lambda)$.

Proof. Let us denote by $(\overrightarrow{\text{inp}}, \overrightarrow{\text{out}}, \sigma_{\text{bind}}, \vec{\sigma})$ the output of `makerandomizedtx`($\text{rt}, v^{\text{bal}}, \overrightarrow{\text{input}}, \overrightarrow{\text{output}}$) and by $(\overrightarrow{\text{inp}}', \overrightarrow{\text{out}}', \sigma'_{\text{bind}}, \vec{\sigma}')$ the output of `makerandomizedtx`($\text{rt}, v^{\text{bal}}, \overrightarrow{\text{input}}', \overrightarrow{\text{output}}'$) when using independent inner randomness, but joint randomness for the oracles $\mathcal{R}, \mathcal{R}_{\text{sig}}, \mathbf{KDF}$.

We will consider D and D' as sequences of random variables $D = (X_1, \dots, X_m), D' = (Y_1, \dots, Y_m)$, and show that for every $i \in [m]$ they satisfy the conditions of Claim 3.6.

We begin with the inputs. Letting, for $i \in [\ell]$, $X_i = \text{inp}_i, Y_i = \text{inp}'_i$, the following claim shows those conditions hold for the first $i \in [\ell]$.

Claim 3.8. *E.w.p $\text{negl}(\lambda)$ over the randomness of \mathcal{A} , for each $i \in [\ell]$ $\text{inp}_i, \text{inp}'_i$ are identically distributed and independent given any fixing of $\text{inp}_{<i}, \text{inp}'_{<i}$.*

Proof. We show first that e.w.p. $\text{negl}(\lambda)$ over the randomness of \mathcal{A} , $\text{inp}_i, \text{inp}'_i$ are independent conditioned on any fixing of $\text{inp}_{<i}, \text{inp}'_{<i}$. $\text{inp}_1, \dots, \text{inp}_\ell, \text{inp}'_1, \dots, \text{inp}'_\ell$ are results of invocations of `makeinput` with independent randomness rcv, α and independent randomness of the SNARK prover.

²The requirement here may seem a bit odd; it models the fact that $\mathbf{NC}(\text{note})$ is a pedersen hash which is combined in \mathbf{NF} with a `pos`-multiple of an independent group generator, followed by an application of BLAKE-2 on the result prefixed with `nk`. In particular, BLAKE-2 takes the place of \mathcal{R} in the implementation.

Inspection shows the only opportunity for dependence amongst any two of them, even after conditioning on the value of the others, is having the random oracle \mathcal{R} queried at the same point during the invocations. \mathcal{R} is queried for the computation of \mathbf{NF} ; so this only happens if

$$(\mathbf{nk}_i, \mathbf{MPH}(\mathbf{note}_i, \mathbf{pos}_i)) = (\mathbf{nk}'_i, \mathbf{MPH}(\mathbf{note}'_i, \mathbf{pos}'_i)).$$

This implies $\mathbf{MPH}(\mathbf{note}_i, \mathbf{pos}_i) = \mathbf{MPH}(\mathbf{note}'_i, \mathbf{pos}'_i)$, but \mathcal{A} will only find such a collision w.p $\text{negl}(\lambda)$. When this doesn't happen \mathbf{inp}_i and \mathbf{inp}'_i are independent also given any fixing of the previous inputs.

Now to show they are identically distributed given a fixing of $\mathbf{inp}_{<i}, \mathbf{inp}'_{<i}$.

Suppose $\mathbf{inp}_i = (\mathbf{nf}, \mathbf{rt}, \mathbf{rk}, \mathbf{cv}, \pi)$, and $\mathbf{inp}'_i = (\mathbf{nf}', \mathbf{rt}', \mathbf{rk}', \mathbf{cv}', \pi')$. We show each element is identically distributed conditioned on any fixing of the previous ones.

- $\mathbf{nf} = \mathcal{R}(q)$ and $\mathbf{nf}' = \mathcal{R}(q')$ where $q = (\mathbf{nk}, \mathbf{MPH}(\mathbf{note}, \mathbf{pos}))$, $q' = (\mathbf{nk}', \mathbf{MPH}(\mathbf{note}', \mathbf{pos}'))$. These are both uniform *unless* one of the queries q, q' was already made to \mathcal{R} in a previous invocation; which would mean $\{(\mathbf{note}^*, \mathbf{pos}^*)\}_{(\mathbf{note}^*, \mathbf{pos}^*) \in \mathbf{inp}_{<i+1} \cup \mathbf{inp}'_{<i+1}}$ contains a collision of \mathbf{MPH} which \mathcal{A} can find only w.p $\text{negl}(\lambda)$.
- $\mathbf{rt} = \mathbf{rt}'$.
- $\mathbf{rk} = \mathbf{ak} + \alpha \cdot \mathbf{g}$, $\mathbf{rk}' = \mathbf{ak}' + \alpha' \cdot \mathbf{g}$. Are both uniform in \mathbb{G} because of the uniform choice of α, α' in makerandomizedtx .
- $\mathbf{cv} = \mathbf{v} \cdot \mathbf{g}_\mathbf{v} + \mathbf{rcv} \cdot \mathbf{g}_\mathbf{r}$, $\mathbf{cv}' = \mathbf{v}' \cdot \mathbf{g}'_\mathbf{v} + \mathbf{rcv}' \cdot \mathbf{g}'_\mathbf{r}$. Are both uniform in \mathbb{G} because of the uniform choices of $\mathbf{rcv}, \mathbf{rcv}' \in \mathbb{F}_r$ in the executions of makerandomizedtx .
- π, π' - When $(\mathbf{nf}, \mathbf{rt}, \mathbf{rk}, \mathbf{cv}) = (\mathbf{nf}', \mathbf{rt}', \mathbf{rk}', \mathbf{cv}')$, it follows from the witness indistinguishability of the SNARK that π and π' are identically distributed. They are independent for any fixing of the previous values, as given this fixing the value of π, π' depends only on the inner randomness of the SNARK prover.

□

We proceed with the elements of the the outputs. It will be convenient now to view each element in $\mathbf{out}_j, \mathbf{out}'_j$ as separate random variables X_i, Y_i , and show that

1. E.w.p $\text{negl}(\lambda)$ over the fixing of $X_{<i}$, they are identically distributed given this fixing of both $X_{<i}$ and $Y_{<i}$.
2. E.w.p $\text{negl}(\lambda)$ over the fixing of $X_{<i}, Y_{<i}$ they are independent given the fixing.

We show this for the different types of elements in $\mathbf{out}_j, \mathbf{out}'_j$:

- $\mathbf{cv} = \mathbf{v} \cdot \mathbf{g}_\mathbf{v} + \mathbf{rcv} \cdot \mathbf{g}_\mathbf{r}$, $\mathbf{cv}' = \mathbf{v}' \cdot \mathbf{g}'_\mathbf{v} + \mathbf{rcv}' \cdot \mathbf{g}'_\mathbf{r}$: are independent and uniform in \mathbb{G} because of the independent uniform choices of $\mathbf{rcv}, \mathbf{rcv}' \in \mathbb{F}_r$ in makerandomizedtx .
- $\mathbf{cm} = \mathbf{NC}(\mathbf{g}, \mathbf{pk}, \mathbf{v}, \mathbf{rcm})$, $\mathbf{cm}' = \mathbf{NC}(\mathbf{g}', \mathbf{pk}', \mathbf{v}', \mathbf{rcm}')$: are uniform and independent in \mathbb{G} because of the independent uniform choices of $\mathbf{rcm}, \mathbf{rcm}' \in \mathbb{F}_r$ in $\text{makerandomizedoutput}$.
- $\mathbf{epk} = \mathbf{esk} \cdot \mathbf{g}$, $\mathbf{epk}' = \mathbf{esk}' \cdot \mathbf{g}$ are uniform and independent in \mathbb{G} because of the independent and uniform choices of $\mathbf{esk}, \mathbf{esk}' \in \mathbb{F}_r$ in $\text{makerandomizedoutput}$.

- π, π' - Assuming the public inputs $(\text{epk}, \text{cm}, \text{cv}) = (\text{epk}', \text{cm}', \text{cv}')$, it follows from the witness indistinguishability of the SNARK that π and π' are identically distributed. They are independent for any fixing of the previous values, as given this fixing the value of π, π' depends only on the inner randomness of the SNARK prover.
- $\text{enc} = \text{ENC}_{\mathbf{KDF}(k)}((g, \text{pk}, v)), \text{enc}' = \text{ENC}_{\mathbf{KDF}(k')}((g', \text{pk}', v'))$ where $k := (\text{esk} \cdot \text{pk}, \text{epk})$ and $k' := (\text{esk}' \cdot \text{pk}', \text{epk}')$: Assuming $k \neq k'$, and moreover $\{k, k'\}$ are different from all the “key seeds” $\{k_j, k'_j\}$ used in previous outputs; we have that the encryption keys $\mathbf{KDF}(k), \mathbf{KDF}(k')$ are uniform and independent of all previous variables. And thus by the theorem’s assumption that \mathbf{KDF} is a random oracle enc, enc' are uniform and independent in this case. Thus there are at most ℓ values of the preceding esk and at most ℓ values of the preceding esk' that can prevent enc and enc' from being uniform and independent; which is a $\text{negl}(\lambda)$ -fraction of the possible values of the preceding values.

It is now left to deal with the signature elements. $\sigma_{\text{bind}}, \sigma'_{\text{bind}}, \{\sigma_i, \sigma'_i\}$.

The distribution of these elements is determined by the public key $\text{pk} = \text{rk}_i$ (or $\text{pk} = S$ the sum of value commitments in the case of σ_{bind}), the message $\mathbf{m} = \text{sighash}(\text{raw}_{\text{tx}})$ they are signing, the internal randomness of the signing algorithm and the reply of the random oracle \mathcal{R}_{sig} on the query point $(R, \text{pk}, \mathbf{m})$. Thus, given a fixing of previous variables, the only case where a dependence between σ_i or σ'_i could be created is if there as a collision between the signatures in the choice of R which happens w.p. $\text{negl}(\lambda)$. \square

3.4 Balance

The following claim states an adversary should not be able to create “money out of thin air”; or more specifically, extract more money from the shielded pool than was put in it. In Sapling, the value $v^{\text{bal}} = v^{\text{bal}}(\text{tx})$ in a transaction tx corresponds to the alleged difference of spend and output values (see Section 4.12 in the spec) and tx is thought of as having ; thus over-extracting from the pool corresponds to a constructing a ledger where the sum of all v^{bal} values is strictly positive.

Claim 3.9. *The probability that an efficient \mathcal{A} generates ledger $L = (\text{tx}_1, \dots, \text{tx}_n)$ such that*

$$\sum_{\text{tx} \in L} v^{\text{bal}}(\text{tx}) > 0$$

is $\text{negl}(\lambda)$.

Proof. Given \mathcal{A} that produces a ledger as in the claim statement w.p. γ , we construct an efficient \mathcal{A}' that w.p $\gamma/2 - \text{negl}(\lambda)$ produces a collision of **IVK**, **NC**, **treehash** or **VC**. It follows that $\gamma = \text{negl}(\lambda)$.

1. \mathcal{A}' begins by running \mathcal{A} and aborting if \mathcal{A} hasn’t output a ledger as in the claim.
2. Otherwise, given such a ledger L , \mathcal{A}' can apply an extractor for each SNARK proof in all inputs and outputs in all transactions. For each transaction input $\text{inp} \in \text{tx} \in L$, $\text{inp} = (\text{cv}, \text{nf}, \text{rt}, \text{rk}, \pi)$, the extractor except w.p. $\text{negl}(\lambda)$ outputs an input witness $\text{inpwit} = (\text{input} = (\text{note}, \text{path}, \text{pos}), \text{pak}, \text{rcv}, \alpha)$. We denote by posnote the *positioned note* corresponding to inp , $\text{posnote} := (\text{note}, \text{pos})$. Similarly for every transaction output in some tx in L , $\text{out} = (\text{cv}, \text{cm}, \text{epk}, \pi, \text{enc})$, the extractor outputs $\text{outwit} = (\text{note}, \text{esk}, \text{rcv})$. The value pos for the

output note can be deduced from when it was added to L, i.e., the location of **cm** in the commitment tree. So again we can define for each **out**, the corresponding positioned note **posnote** = (**note**, **pos**). For $i \in [n]$ let us denote respectively by $\mathcal{I}_i, \mathcal{O}_i$ the positioned input and output notes in **tx**_{*i*} with non-zero value³.

We also use the extractor from theorem 1.4 to find s such that $S = s \cdot \mathbf{g}_r$ where

$$S := \sum_{i=1}^{\ell} \mathbf{cv}_i - \sum_{i=\ell+1}^{\ell+s} \mathbf{cv}_i - \mathbf{v}^{\text{bal}} \cdot \mathbf{g}_v$$

is the public key in the value binding signature σ_{bind} .

If one of the extractor runs fails \mathcal{A}' aborts. Note that w.p. at least $\gamma/2 - \text{negl}(\lambda)$ \mathcal{A}' doesn't abort.

3. \mathcal{A}' checks if for some $i \in [n]$ and **inp** $\in \text{tx}_i$, **posnote**(**inp**) $\notin \mathcal{O}_{< i}$.

If so, let **tx** = **tx**_{*i*}. Let **rt** be the root of the tree used in the public input of **inp**; this is the tree T_j formed from $\{\text{tx}_1, \dots, \text{tx}_j\}$ for some $j < i$. Let **posnote** = (**g**, **pk**, **v**, **rcm**, **pos**) and **cm** = **NC**(**g**, **pk**, **v**, **rcm**). **inpwit** contains a path **path** from **cm** to **rt**. If **pos** is an index of a leaf in T_j , there exists an extended note **posnote'** that was inserted to this position when constructing the ledger and from **posnote'** we can derive a path **path'** from **cm'** = **NC**(**g'**, **pk'**, **v'**, **rcm'**) in position **pos** to **rt**. If **path** \neq **path'**, then going down from **rt** to the first difference between **path** and **path'** (ask Sean/Daira : is T always a full tree with zeroes on other leaves? No you have filler values for the empty subtrees, need to check this are values that are hard to find route to - their impossible to find rout to - have no preimage) this difference gives a collision of **treehash** that \mathcal{A}' can output.

Otherwise, we have **cm** = **cm'**. **note** must be different from **note'** because **posnote'** = (**note'**, **pos**) $\in \mathcal{O}_{< i}$ but (**note**, **pos**) $\notin \mathcal{O}_{< i}$.

Thus **note**, **note'** is a collision of **NC**. In this case, \mathcal{A}' outputs this collision and terminates.

Now suppose **pos** is not a position of a leaf in T_j . This means there is only a partial path **path'** in T_j from **rt** to a filler value with no preimage (see spec for details). So, similarly we follow **path** and **path'** to their first difference - a difference that must exist because of the filler value; and this gives us a collision of **treehash** that \mathcal{A}' outputs.

4. Now \mathcal{A}' checks if as a multiset $\mathcal{I} := \mathcal{I}_1 \cup \dots \cup \mathcal{I}_n$ contains a repetition. That is, there exists **posnote** = (**g**, **pk**, **v**, **rcm**, **pos**) such that for two distinct transaction inputs **inp** = (**cv**, **nf**, **rt**, **rk**, π), **inp'** = (**cv'**, **nf'**, **rt'**, **rk'**, π') in L; if the corresponding extracted witnesses are **inpwit** = (**input** = (**note**, **path**, **pos**), **pak**, **rcv**, α), **inpwit'** = (**input'** = (**note'**, **path'**, **pos'**), **pak'**, **rcv'**, α'); then (**note**, **pos**) = (**note'**, **pos'**) = **posnote**.

We show in this case that \mathcal{A}' can output a collision of **IVK**:

Let **cm** = **NC**(**g**, **pk**, **v**, **rcm**). Since **nf** \neq **nf'**, and **nf** = **NF**(**nk**, **cm**, **pos**), **nf'** = **NF**(**nk'**, **cm**, **pos**); we have **nk** \neq **nk'**.

³Sapling enables the creation of dummy notes with zero value, for which the spend statement doesn't check Merkle path validity, cf. Section 4.7.2 in the spec).

Also $\text{ivk} = \mathbf{IVK}(\text{ak}, \text{nk})$, $\text{ivk}' = \mathbf{IVK}(\text{ak}', \text{nk}')$, and $\text{pk} = \text{ivk} \cdot \mathbf{g} = \text{ivk}' \cdot \mathbf{g}$. So $\text{ivk} = \text{ivk}'$ (Check with Sean is ivk canonical - checked) And thus, \mathcal{A}' can output $(\text{ak}, \text{nk}), (\text{ak}', \text{nk}')$ as a collision of \mathbf{IVK} .

5. Let us denote by $\text{bal}(\text{tx})$ the (integer) sum of values in inputs of tx minus the sum of values in output of tx (notes meaning those output by the extractors); and by $\text{rcv}(\text{tx})$ the sum of values rcv in input witnesses of tx minus the sum of values rcv in output witnesses of tx . When reaching this point with no output we know that:

For each $i \in [n]$, $\mathcal{I}_i \subset \mathcal{O}_1 \cup \dots \cup \mathcal{O}_{i-1} \setminus (\mathcal{I}_1 \cup \dots \cup \mathcal{I}_{i-1})$.

This implies

$$\sum_{\text{tx} \in L} \text{bal}(\text{tx}) \leq 0.$$

We claim that we must have for some $\text{tx} \in L$, $\text{bal}(\text{tx}) \neq v^{\text{bal}}(\text{tx})$: Otherwise, we would have

$$\sum_{\text{tx} \in L} v^{\text{bal}}(\text{tx}) = \sum_{\text{tx} \in L} \text{bal}(\text{tx}) \leq 0,$$

contradicting the fact that \mathcal{A} has managed to output L with a positive sum of v^{bal} values.

Thus, let $\text{tx} = \text{tx}_i$ be such that $\text{bal}(\text{tx}) \neq v^{\text{bal}}(\text{tx})$. We show in the next step how \mathcal{A}' uses this to output a collision of \mathbf{VC} .

6. At this point, we know that $\text{bal}(\text{tx}) \neq v^{\text{bal}}(\text{tx})$. As both these values are in the open interval ⁴ $(-r/2, r/2)$, we have also $\text{bal}(\text{tx}) \neq v^{\text{bal}}(\text{tx}) \pmod{r}$. Suppose we are in this case with probability γ . We show how to find a collision of \mathbf{VC} with probability $\gamma/\text{poly}(\lambda)$. Since tx verifies, we know that $\text{verifySig}_{\mathbf{gr}}^{\mathcal{R}}(S, \text{sighash}(\text{raw}_{\text{tx}}), \sigma_{\text{bind}})$ for

$$S = \sum_{i=1}^{\ell} \text{cv}_i - \sum_{i=\ell+1}^{\ell+s} \text{cv}_i - v^{\text{bal}} \cdot \mathbf{g}_v = \left(\sum_{i=1}^{\ell} v_i - \sum_{i=\ell+1}^s v_i \right) \cdot \mathbf{g}_v + \left(\sum_{i=1}^{\ell} \text{rcv}_i - \sum_{i=\ell+1}^s \text{rcv}_i \right) \cdot \mathbf{g}_r - v^{\text{bal}} \cdot \mathbf{g}_v.$$

Using Theorem 1.4, we can with probability $\gamma/2$ we can use the forking lemma to rewind \mathcal{A} while altering the response of \mathcal{R} on the signature challenge in σ_{bind} , and find s such that $s \cdot \mathbf{g}_r = S$. Thus, we have $\mathbf{VC}(0, s) = S$.

Let $R := \sum_{i=1}^{\ell} \text{rcv}_i - \sum_{i=\ell+1}^s \text{rcv}_i$ and $v := \text{bal}(\text{tx}) - v^{\text{bal}}(\text{tx})$. We also have $\mathbf{VC}(v, R) = S$. Hence \mathcal{A}' can output $(0, s), (v, R)$ as a collision of \mathbf{VC} .

□

3.5 Spendability

Valid transaction bases: A sequence $\mathbf{x} = (\overrightarrow{\text{input}}, \overrightarrow{\text{output}}, v^{\text{bal}})$ is a *valid transaction base* if $v^{\text{bal}} = \sum v(\text{input}_i) - \sum v(\text{output}_j)$.

We review note encryption and decryption from the spec in our notation.

⁴See the spec for details: v^{bal} and v in each transaction input/output are at most 2^{64} in absolute value, so assuming less than, e.g., 2^{r-66} transaction inputs and outputs in any transaction, this is true.

Decrypting notes:

$\mathbf{dec}(\text{ivk}, \text{out} = (\text{cv}, \text{cm}, \text{epk}, \pi, \text{enc}))$

1. Let $K := \mathbf{KDF}(\text{epk} \cdot \text{ivk})$
2. Let $\text{np} = \mathbf{DEC}_K(\text{enc})$. If $\mathbf{DEC}()$ fails output rej .
3. Suppose $\text{np} = (\text{d}, \text{v}, \text{rcm}, \text{memo})$. If $\text{rcm} \geq r$ output rej .
4. Let $\text{g} := \mathbf{GH}(\text{d})$.
5. Let $\text{pk} := \text{g} \cdot \text{ivk}$. Let $\text{note} := (\text{g}, \text{pk}, \text{v}, \text{rcm})$.
6. Check that $\text{cm} = \mathbf{NC}(\text{note})$. Output rej if not.
7. Output note .

We define

$$\mathbf{dec}(\text{ivk}, \text{tx}) := \cup_{\text{out} \in \text{tx}} \mathbf{dec}(\text{ivk}, \text{out}),$$

$$\mathbf{dec}(\text{ivk}, \text{L}) := \cup_{\text{tx} \in \text{L}} \mathbf{dec}(\text{ivk}, \text{tx})$$

And also

$$\mathbf{nf}(\text{tx}) := \cup_{\text{inp} \in \overrightarrow{\text{inp}}(\text{tx})} \mathbf{nf}(\text{inp}), \mathbf{nf}(\text{L}) := \cup_{\text{tx} \in \text{L}} \mathbf{nf}(\text{tx})$$

In the spendability game \mathcal{A} tries to create a ledger where a note successfully decrypted with ivk cannot be spent. Formally, the game proceeds as follows.

1. We choose uniform $\text{sk} = (\text{ask}, \text{nsk})$; and give $\text{pak} = (\text{ask} \cdot \text{g}_{\text{sig}}, \text{nsk})$ to \mathcal{A} .
2. \mathcal{A} outputs a ledger L , a positioned note $(\text{note}, \text{pos})$, a set of output notes $\overrightarrow{\text{output}}$, and a set of incoming viewing keys $\overrightarrow{\text{ivk}}$.
3. We choose random $\overrightarrow{\text{rcv}} \in \mathbb{F}_r^{\ell+s}$ and compute $\text{tx} = \text{maketx}(\overrightarrow{\text{input}}, \overrightarrow{\text{output}}, \text{v}^{\text{bal}}, \text{ask})$.
4. Let $\text{ivk} := \mathbf{IVK}(\text{ak}, \text{nk})$. \mathcal{A} wins iff
 - (a) $\text{note} \in \mathbf{dec}(\text{ivk}, \text{L})$.
 - (b) $((\text{note}), \overrightarrow{\text{output}}, \text{v}^{\text{bal}})$ is a valid transaction base.
 - (c) For each $i \in [s]$, output_i belongs to ivk_i .
 - (d) $\text{verify-tx}(\text{L}, \text{tx})$.
 - (e) For some $i \in [s]$, $\mathbf{dec}(\text{ivk}_i, \text{out}_i)$ does not return output_i .

We wish to show that the success of any efficient \mathcal{A} in this game is $\text{negl}(\lambda)$.

Let $\text{nk} = \text{nsk} \cdot \text{g}_{\text{n}}$. Inspection of the protocol shows this exactly corresponds to the nullifier of note with nullifier key nk already appearing in the ledger. Thus, it suffices to prove the following.

Claim 3.10. *Fix any efficient \mathcal{A} . Suppose that \mathcal{A} is given uniformly chosen pak , and let $\text{ivk} := \mathbf{IVK}(\text{pak})$. The probability that \mathcal{A} generates a ledger L and positioned note $(\text{note}, \text{pos})$ such that*

1. $(\text{note}, \text{pos}) \in \mathbf{dec}(\text{ivk}, \text{L})$

2. $\mathbf{NF}(\mathbf{nk}, \mathbf{NC}(\text{note}), \text{pos}) \in \text{nf}(\mathbf{L})$

is $\text{negl}(\lambda)$.

Proof. Let γ be the probability that \mathcal{A} outputs \mathbf{L}, note satisfying the two properties in the claim. We construct an efficient \mathcal{A}' that receives a forgery challenge \mathbf{ak} of Schnorr and w.p. $\gamma - \text{negl}(\lambda)$ does one of the following.

- Output a collision of either \mathbf{NF} , \mathbf{NC} or \mathbf{IVK} .
 - Output a forgery w.r.t to randomization of Schnorr for the challenge \mathbf{ak} .
- \mathcal{A}' works as follows.

1. \mathcal{A}' receives a challenge \mathbf{ak} ; chooses random $\mathbf{nsk} \in \mathbb{F}_r$ and sends $\mathbf{pak} = (\mathbf{ak}, \mathbf{nsk})$ to \mathcal{A} .
2. \mathcal{A}' receives the output $(\mathbf{L}, \text{note}, \text{pos})$ of \mathcal{A} .
3. \mathcal{A}' checks that $\mathbf{L}, (\text{note}, \text{pos})$ satisfy the two properties in the claim; if not it aborts.
4. Let $\text{nf} := \mathbf{NF}(\mathbf{nk}, \mathbf{NC}(\text{note}), \text{pos})$. Fix the out, tx with $\text{out} \in \text{tx} \in \mathbf{L}$ such that $\mathbf{dec}(\text{ivk}, \text{out}) = (\text{note}, \text{pos})$. out contains a valid SNARK proof for $\text{SPEND}(\text{rt}, \text{cv}, \text{nf}, \text{rk})$ for some cv, rt . Apply the relevant extractor ξ relating to the snark proof to obtain e.w.p $\text{negl}(\lambda)$ a witness $\text{path}, \text{pos}', \mathbf{g}', \mathbf{pk}', \mathbf{v}', \text{rcm}', \text{cm}', \text{rcv}', \alpha, \mathbf{ak}', \mathbf{nsk}'$ for the statement.
5. Let $\mathbf{nk}' := \mathbf{nsk}' \cdot \mathbf{g}_n$. If $(\mathbf{nk}, \text{cm}, \text{pos}) \neq (\mathbf{nk}', \text{cm}', \text{pos}')$, \mathcal{A}' outputs $(\mathbf{nk}, \text{cm}, \text{pos}), (\mathbf{nk}', \text{cm}', \text{pos}')$ as a collision of \mathbf{NF} .
6. Otherwise, let $\text{note}' = (\mathbf{g}', \mathbf{pk}', \mathbf{v}', \text{rcm}')$. We have $\text{cm} = \mathbf{NC}(\text{note}) = \mathbf{NC}(\text{note}')$. If $(\mathbf{g}', \mathbf{pk}', \mathbf{v}') \neq (\mathbf{g}, \mathbf{pk}, \mathbf{v})$, \mathcal{A}' outputs $(\text{note}, \text{note}')$ as a collision of \mathbf{NC} .
7. Otherwise, we must have $\text{ivk}' = \text{ivk}$ (cause $\mathbf{g} \cdot \text{ivk} = \mathbf{g} \cdot \text{ivk}' = \mathbf{pk}$). Then $\text{ivk} = \mathbf{IVK}(\mathbf{ak}', \mathbf{nk})$ (by this stage we know $\mathbf{nk} = \mathbf{nk}'$). If $\mathbf{ak} \neq \mathbf{ak}'$, \mathcal{A}' outputs $(\mathbf{ak}, \mathbf{nk}), (\mathbf{ak}', \mathbf{nk})$ as a collision of \mathbf{IVK} .
8. Otherwise $\mathbf{ak} = \mathbf{ak}'$, and $\text{rk} = \mathbf{ak} + \alpha \cdot \mathbf{g}$. Let σ be the signature of raw_{tx} with public key rk in inp . and \mathcal{A}' outputs $(\alpha, \text{raw}_{\text{tx}}, \sigma)$ as a forgery of Schnorr with challenge \mathbf{ak} .

□

Remark 3.11. *Note that in the spendability and non-malleability property \mathcal{A} can choose what value nf to work with. It seems likely that in a weaker model where the values nf*

3.6 Indistinguishability of diversified addresses

El-gammal encryption private key \mathbf{sk} public key $[\mathbf{sk}] \cdot \mathbf{g}$ encryption of x : $[\mathbf{sk}]$

References

- [1] N. Fleischhacker, J. Krupp, G. Malavolta, J. Schneider, D. Schröder, and M. Simkin. Efficient unlinkable sanitizable signatures from signatures with re-randomizable keys. *IET Information Security*, 12(3):166–183, 2018.
- [2] D. Pointcheval and J. Stern. Security arguments for digital signatures and blind signatures. *J. Cryptology*, 13(3):361–396, 2000.