

# Masters Project - Visual Hashes (temporary title)

Adrien Laydu

Sometime in 2022

## 1 Context [TO BE EXPANDED/CLARIFIED]

TrustID is an Identity Provider (IDP) that can have various uses, such as 2-Factor Authentication (2FA) or digital signatures. As such, TrustID has to offer ways to securely identify, verify, and authenticate users. To achieve this, TrustID has a strong onboarding procedure, either using a video onboarding process or by having users physically go to an onboarding office, as well as providing an authenticator application for the user's mobile phone. Once onboarded, a user that wants to access a service secured by TrustID enter their credentials ([USERNAME] and password) on the TrustID website. The TrustID server then sends a challenge to the authenticator application, which then uses the unlocking mechanism of the phone to authenticate the user: we assume only the targeted user can unlock their phone. The widespread use of biometrics as a way of unlocking phones (rather than traditional passwords or PINs) would ideally improve the security of this method drastically. However, unlocking processes typically offer one of the weaker methods as a fallback. Therefore, TrustID relies only on the ownership of the phone for security.

That brings us to our security issue: at time of the internship, When the user validates a [transaction on their mobile phone (i.e press the "accept" button of the authenticator app)], they have no way of being sure that the prompt on their phone corresponds to the actual transaction they wish to effectuate. This could lead to some vulnerabilities: for example, an impersonation attack is possible, as shown here (we take the example of digital signatures for illustration purposes):

Say an attacker A wants to sign a document we will call  $d_A$  while impersonating user U. U uses a weak password, which means A can log as U on the TrustID website.

- A logs as U on TrustID website.
- A waits for U to ask for a challenge for any document.
- When U asks for a challenge for some document  $d_U$ , A asks the server for a challenge for  $d_A$ .
- The server will then send two challenges to U's phone, one for  $d_A$  and one for  $d_U$ .
- Since U is unable to distinguish between the two challenges, they validate whichever challenge comes first.
- If the challenge for  $d_A$  came first, then the attacker successfully signed a document while impersonating U.

While not that impactful when used on digital signatures, the same attack could be used for other (much more sensitive) TrustID services, such as 2FA: the attack described above would make the second factor useless.

This project is aimed at solving the security issue presented above. The way we will achieve it is to present the user with a unique, recognizable fingerprint of the document they requested a challenge for. While this could be achieved with traditional hash functions, the fact that human users will have to compare fingerprints calls for a more user-friendly solution. In particular, we want the fingerprints to have the form of pictures, as we hypothesize they are less cumbersome to compare.

**Note:** in practice, A's and U's requests for challenges are augmented with metadata. Different metadata calls for different capabilities of A: a UserID would not call for anything more than assumed earlier. A timestamp (currently in use at time of writing) requires A to send their request simultaneously with U. A nonce would require A to guess the nonce. [AT IMPLEMENTATION TIME] We will see [later] how the choice of metadata alter the probability of a successful attack.

## 2 Visual Hashing

### 2.1 Definitions and notations

#### 2.1.1 Pixels

Because we work on image representation, we will work with pixels. We write  $\mathbb{P}$  the set of all pixels. We will often need separate categories of pixels: in particular, we write  $\mathbb{P}_2 \subset \mathbb{P}$  the set of pixel that may only take two values (namely "white" and "black"). We write  $\mathbb{P}_{\text{RGB}} \subset \mathbb{P}$  the set of pixels that are encoded by three color channels (Red, Green, Blue), each channel having an integer value between 0 and 255.

#### 2.1.2 Visual hash function

We define a visual hash function as a function  $\{0, 1\}^* \rightarrow \mathbb{P}^{m \times n}$  of the form:

$$\mathcal{H}(x; \lambda) = V(H(x; \lambda); \lambda)$$

Where  $\lambda$  is a security parameter,  $H : \{0, 1\}^* \rightarrow \{0, 1\}^\ell$  is a standard hash function (such as SHA-256) and  $V : \{0, 1\}^\ell \rightarrow \mathbb{P}^{m \times n}$  is a mapping from bits strings to arrays of pixels. The values of  $\ell, m, n$  depend on the value of  $\lambda$ .

Whenever it is not explicitly needed, we omit the security parameter  $\lambda$  in the notation. We call  $h := H(x)$  a *hash*, and  $F := V(h)$  a *fingerprint*. We note a visual hash function as  $\mathcal{H} = (H, V)$  or when contexts allows it, just  $\mathcal{H}$ .

#### 2.1.3 Visual indistinguishability

[TO BE CLARIFIED / EXPANDED IN DEDICATED SECTION]

For two fingerprints  $F_1, F_2 \in \mathbb{P}^{m \times n}$ , we write  $F_1 \approx_v F_2$  if  $F_1$  is *visually indistinguishable* from  $F_2$ , that is, if a user thinks  $F_1 = F_2$  by looking at them. Note that  $F_1 = F_2 \implies F_1 \approx_v F_2$ .

#### 2.1.4 Properties of visual hash functions

We adapt traditional hash functions properties to visual hash functions as follows [NOT SURE IF USEFUL YET]. When we say that an attack is "infeasible", we mean that an efficient [DEFINE EFFICIENT - POLY-TIME] adversary cannot realistically perform the attack because of its complexity. The properties are as follows:

- *First pre-image resistance*: we say a visual hash function  $\mathcal{H}$  is first pre-image resistant if given  $F \in \mathbb{P}^{m \times n}$ , it is infeasible to find  $x \in \{0, 1\}^*$  such that  $\mathcal{H}(x) \approx_v F$ .
- *Second pre-image resistance*: we say a visual hash function  $\mathcal{H}$  is second pre-image resistant if given  $y \in \{0, 1\}^*$ , it is infeasible to find  $x \neq y \in \{0, 1\}^*$  such that  $\mathcal{H}(x) \approx_v \mathcal{H}(y)$ .
- *Collision resistance*: we say a visual hash function  $\mathcal{H}$  is collision-resistant if it is infeasible to find  $x, y \in \{0, 1\}^*$  such that  $x \neq y$  and  $\mathcal{H}(x) \approx_v \mathcal{H}(y)$ .

In order to better categorize and analyze visual hash function, we define two additional properties:

- *Intermediate pre-image resistance*: we say a visual hash function  $\mathcal{H} = (H, V)$  is intermediate pre-image resistant if given  $F \in \mathbb{P}^{m \times n}$ , it is infeasible to find  $h \in \{0, 1\}^\ell$  such that  $V(h) \approx_v F$ .
- *Visual injectiveness*: we say a visual hash function  $\mathcal{H} = (H, V)$  is visually injective if for any  $h_1, h_2 \in \{0, 1\}^\ell$ ,  $h_1 \neq h_2 \implies V(h_1) \not\approx_v V(h_2)$ . Note that visual injectiveness implies that  $V$  is injective, because  $V(h_1) \not\approx_v V(h_2) \implies V(h_1) \neq V(h_2)$ .

## 3 Context of application

### 3.1 Protocol

We consider the protocol shown in fig. 1, in which a visual hash function  $\mathcal{H} = (H, V)$  is used. On the server side, **md** denotes some metadata such as a sessionID, timestamp, nonce, etc. [TO BE EXPANDED IN IMPLEMENTATION SECTION].

# VERIFY-TRANSACTION

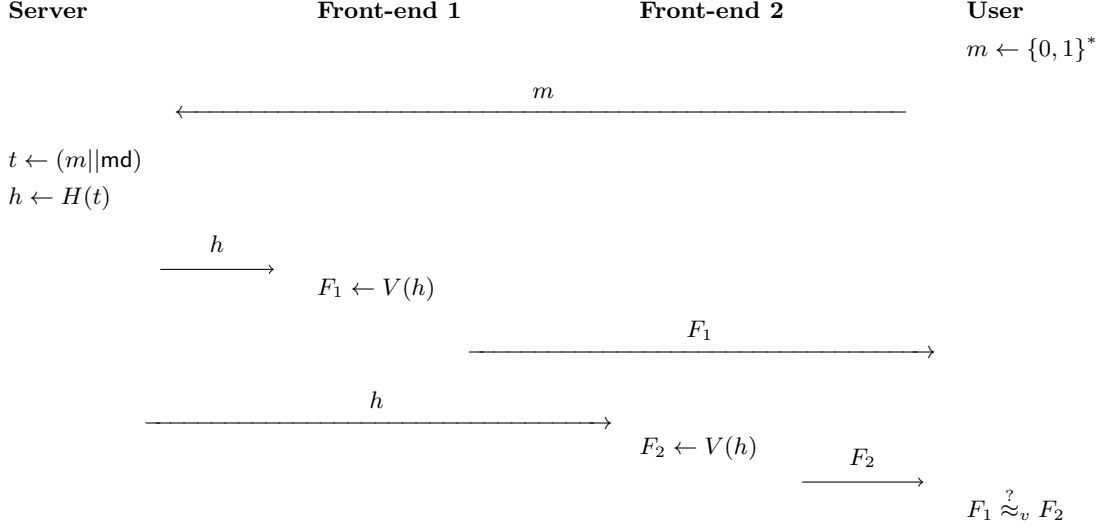


Figure 1: Protocol without adversary.

## 3.2 Threat models

### 3.2.1 A realistic TrustID threat

[GIVE IDEA OF TRANSACTION VALIDATION PROTOCOL (perhaps simplified version of Antoine’s UML sequence diagram)]. Following [that], we assume an adversary cannot impersonate the TrustID server. However, an attack by an external adversary is still possible. We define, for an adversary  $\mathcal{A}$ , the game  $\Gamma_{\text{TrustID}}$ :

**Game  $\Gamma_{\text{TrustID}}$**

---

$m \leftarrow \{0, 1\}^*$   
 $t \leftarrow (m || \text{md})$   
 $h \leftarrow H(t)$   
 $F \leftarrow V(h)$   
 $\mathcal{A}(t^*) \rightarrow m'$   
 $t' \leftarrow (m' || \text{md}')$   
 $h' \leftarrow H(t')$   
 $F' \leftarrow V(h')$   
**return**  $1_{F \approx_v F' \wedge m' \neq m}$

Where  $t^*$  denotes some part of  $t$ : depending on the metadata we choose,  $\mathcal{A}$  may know some bits of  $t$  (e.g, the adversary knows at which time the transaction is made, or the user’s ID). We assume however that the adversary doesn’t know  $m$ . We consider our scheme achieves the minimal required security if for any efficient adversary  $\mathcal{A}$ ,  $\text{Adv}_{\mathcal{A}}^{\Gamma_{\text{TrustID}}}(\lambda)$  is negligible in  $\lambda$ , where  $\text{Adv}_{\mathcal{A}}^{\Gamma_{\text{TrustID}}}(\lambda) = \Pr[\Gamma_{\text{TrustID}} \text{ returns } 1]$ . Figure 2 illustrates how the game translates to a concrete implementation.

While we could try and design some visual hash function that achieves this minimal security, we choose to go further and define a stronger threat model: this could alleviate our assumptions on the security of the rest of the TrustID architecture, as well as provide further insight. [Also it is fun]

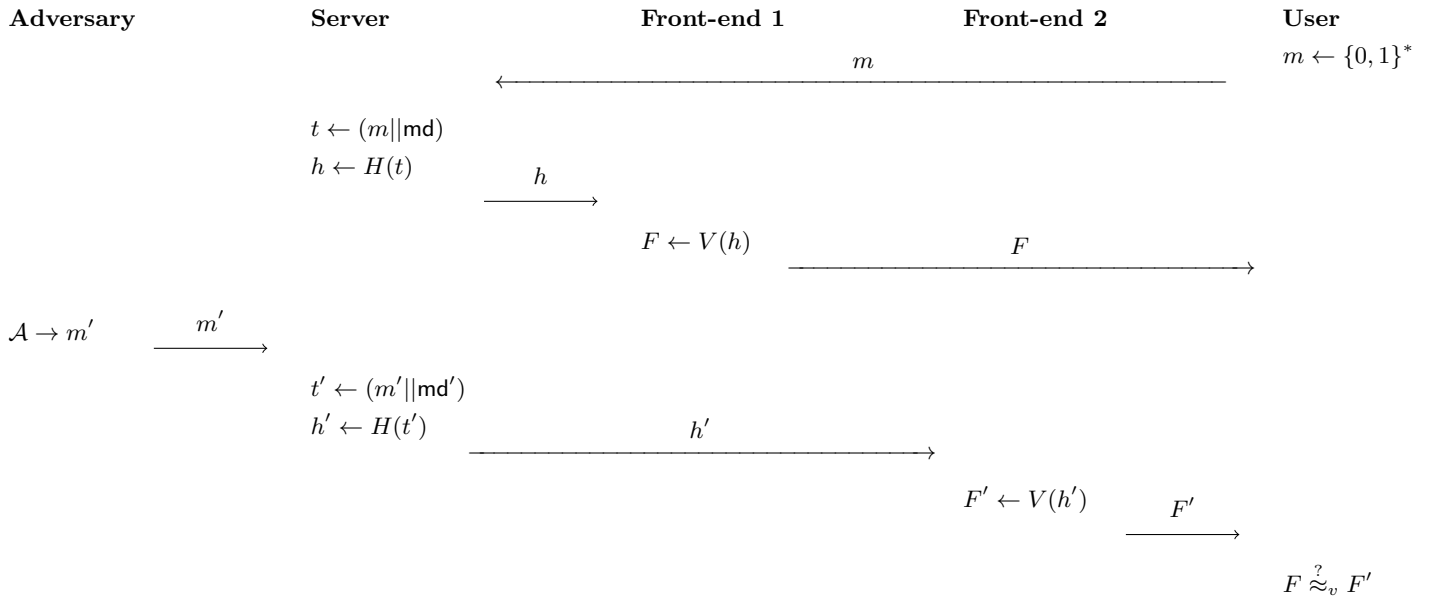


Figure 2: Protocol with a realistic TrustID adversary.

### 3.2.2 A stronger threat

We now define another game, with an adversary  $\mathcal{A}$ :

Game $\Gamma$
$m \leftarrow \{0, 1\}^*$
$t \leftarrow (m    \text{md})$
$h \leftarrow H(t)$
$F \leftarrow V(h)$
$\mathcal{A}(t, h, F) \rightarrow h'$
$F' \leftarrow V(h')$
<b>return</b> $1_{F \approx_v F'} \wedge h' \neq h}$

In words,  $\mathcal{A}$  wins if they can produce a hash  $h'$ , different from the true hash  $h$ , that yields a fingerprint  $V(h')$  that is visually indistinguishable from the true fingerprint  $F$ . We will consider our scheme is secure if for any efficient adversary  $\mathcal{A}$ ,  $\text{Adv}_{\mathcal{A}}^{\Gamma}(\lambda)$  is negligible in  $\lambda$ , where  $\text{Adv}_{\mathcal{A}}^{\Gamma}(\lambda) = \Pr[\Gamma \text{ returns } 1]$

We show that this security is indeed stronger than the minimal security we defined in the previous section. Concretely, we prove that,  $\text{Adv}_{\mathcal{A}}^{\Gamma}(\lambda)$  is negligible for all efficient adversaries  $\implies \text{Adv}_{\mathcal{A}}^{\Gamma_{\text{TrustID}}}(\lambda)$  is negligible for all efficient adversaries. We proceed by reduction: suppose  $\mathcal{A}_{\text{TrustID}}$  is an efficient adversary playing  $\Gamma_{\text{TrustID}}$  that wins with non-negligible probability. We build an adversary  $\mathcal{A}$  playing  $\Gamma$  that uses  $\mathcal{A}_{\text{TrustID}}$  to win as follows:

$\mathcal{A}(t, h, F)$
Extract $t^*$ from $t$ depending on what $\mathcal{A}_{\text{TrustID}}$ needs
$\mathcal{A}_{\text{TrustID}}(t^*) \rightarrow m'$
$t' \leftarrow (m'    \text{md}')$
$h' \leftarrow H(t')$
<b>return</b> $h'$

Now, if  $\mathcal{A}_{\text{TrustID}}$  wins at  $\Gamma_{\text{TrustID}}$ , then  $\mathcal{A}$  wins at  $\Gamma$ . Therefore,  $\text{Adv}_{\mathcal{A}}^{\Gamma}(\lambda) \geq \text{Adv}_{\mathcal{A}_{\text{TrustID}}}^{\Gamma_{\text{TrustID}}}(\lambda)$ , which was assumed to be non-negligible, which means  $\text{Adv}_{\mathcal{A}}^{\Gamma}(\lambda)$  is non-negligible. That concludes the proof: if there exists an efficient adversary

that breaks  $\Gamma_{\text{TrustID}}$ , then there exists an efficient adversary that breaks  $\Gamma$ . This is equivalent to the statement we wanted to prove.

Figure 3 illustrate how the threat is stronger: the adversary can bypass the server to submit a hash directly to the frontend.

#### VERIFY-TRANSACTION-STRONG

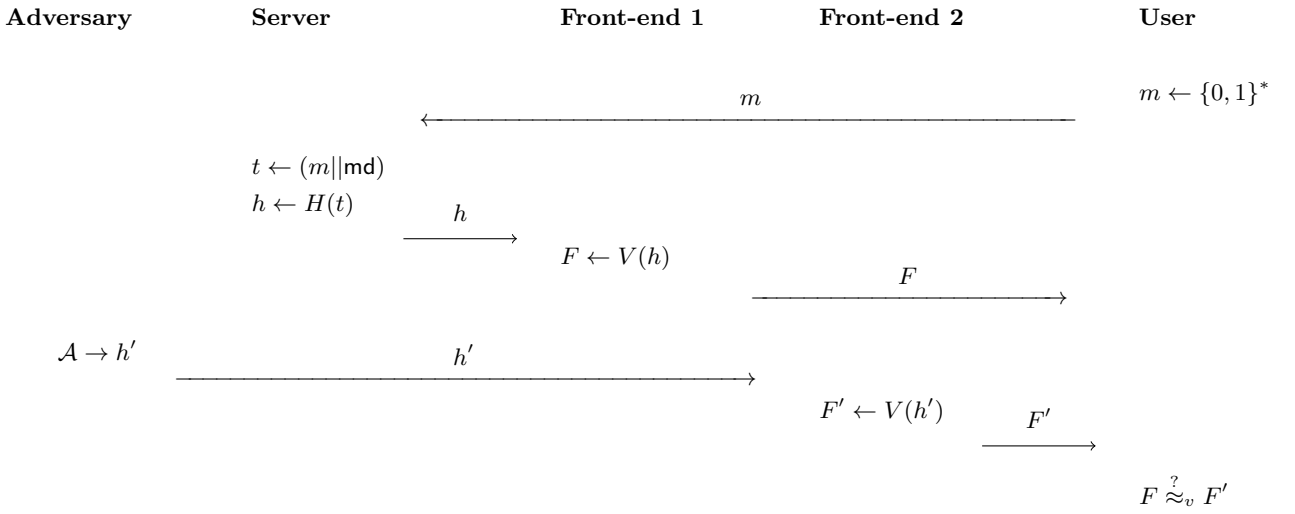


Figure 3: Protocol with a stronger adversary.

## 4 Security implications of visual hash function properties

We separate the different combinations of properties a visual hash function may have, to analyse their strengths and weaknesses in terms of security.

### 4.1 No visual injectiveness, no intermediate pre-image resistance

If a visual hash function  $\mathcal{H}$  is neither intermediate pre-image resistant nor visually injective, the adversary can create a forgery:

1. From  $F$ , find  $F' \neq F \in \text{Im}(V)$  such that  $F \approx_v F'$ .
2. From  $F'$ , get  $h'$  such that  $V(h') = F'$ .
3. Output  $h'$ .

This attack is efficient by assumption as long as the adversary can find  $F'$  efficiently. Therefore we need at least one of the two proprieties to keep the adversary from making forgeries, or somehow make it infeasible to find  $F'$  from  $F$ .

### 4.2 Visual injectiveness

If  $\mathcal{H} = (H, V)$  is visually injective, then the adversary never wins: indeed, by definition of visual injectiveness,  $h \neq h' \implies V(h) \not\approx_v V(h')$ .

Moreover, visual injectiveness implies that the only case where  $V(x_1) \approx_v V(x_2)$  with  $x_1 \neq x_2$  is when  $H(x_1) = H(x_2)$ . That is, in order to find a collision on  $\mathcal{H}$ , the adversary must find a collision on  $H$ . Theferore the following implication holds:

$$\mathcal{H} \text{ is visually injective} \wedge H \text{ is collision-resistant} \implies \mathcal{H} \text{ is collision-resistant} \quad (1)$$

Since there are well-documented hash functions that are considered collision-resistant, visual injectiveness would be the ultimate goal for our hash function.

### 4.3 Intermediate pre-image resistance

If  $\mathcal{H} = (H, V)$  is intermediate pre-image resistant but isn't visually injective, the attack described in section 4.1 is not feasible, because step 2 is infeasible by intermediate pre-image resistance. [BUT IS IT ENOUGH FOR SECURITY ?]

## 5 On visual indistinguishability

### 5.1 Evaluating the difference between fingerprints

We would like a metric to measure how well our scheme performs, so that we can analyse it formally. In particular, we need to define a distance between  $d$  between two fingerprints  $F_1, F_2 \in \mathbb{P}$  such that for some  $\delta > 0$  (ideally as small as possible)

$$d(F_1, F_2) \leq \delta \iff F_1 \approx_v F_2$$

#### 5.1.1 normalised Hamming distance

A natural choice for evaluating the differences between binary images (i.e, that are in  $\mathbb{P}_2^{m \times n}$ ) would be to compute the normalised Hamming distance of the two fingerprints

$$d_H(F_1, F_2) = \frac{1}{mn} \sum_{i=0}^{m-1} \sum_{j=0}^{n-1} \mathbb{I}(F_1[i][j] \neq F_2[i][j])$$

where  $F_1, F_2 \in \mathbb{P}^{m \times n}$ ,  $(\cdot)[i][j]$  denotes the pixel at position  $(i, j)$ , and  $\mathbb{I}$  is an indicator function.  $d_H$  ranges between 0 and 1 included, the extrema being obtained by having two pictures that differ in zero or all of their pixels, respectively. This metric would be useful because it relates directly to the Hamming distance between of the hashes functions.

Unfortunately, it seems this distance does not achieve the properties we stated above. We can refute the right direction ( $\implies$ ) very easily. Indeed, fig. 4 shows two fingerprint with the smallest non-zero Hamming distance ( $d_H = \frac{1}{mn}$ ,  $mn = 256$ ) that are clearly visually distinguishable by humans. Concerning the other direction ( $\impliedby$ ), we should proceed with caution: shifting bits can produce images that look similar despite having a large Hamming distance. Figure 5 shows fingerprints having a normalised Hamming distance of  $\frac{71}{256} = 27.73\%$  yet the results are not easy to distinguish.

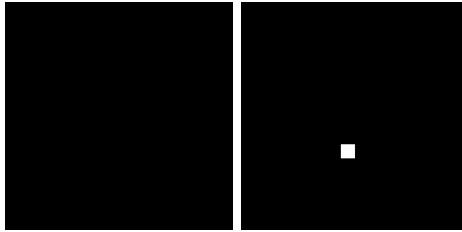


Figure 4: Two visually distinguishable fingerprints with low Hamming distance.

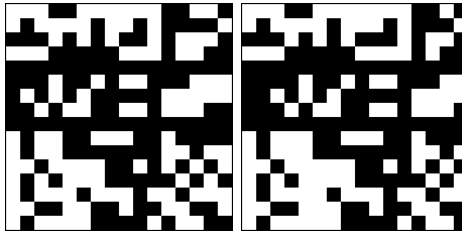


Figure 5: Two fingerprints that have a large Hamming distance, but are not trivial to distinguish.

### 5.1.2 HaarPSI distance

The Haar perceptual similarity index (HaarPSI) [ref] seems to be a relatively good measure to test whether a pair of fingerprints is distinguishable. The advantages it offers are a relatively low computation time compared to other PSI, which allows for automated parameter search, as well as a good correlation with human opinion scores [ref]. In contrast, the main downside is that although the results are re-linearized, the final output (a number between 0 and 1) is somewhat tricky to interpret: two clearly distinguishable images will have a HaarPSI around 0.1–0.2, and a HaarPSI over 0.6 yields a visual difference that is so small that we consider it imperceptible.

For our future scheme, we would of course like to minimize the HaarPSI distance between any 2 fingerprints. However, given the complexity of the HaarPSI formula, we cannot work our way from that to create a scheme. Instead, we pick a threshold value  $T_{Haar} = [0.5]$  and we consider that

$$\text{HaarPSI}(F_1, F_2) \geq T_{Haar} \implies F_1 \approx_v F_2$$

The value of  $T_{Haar}$  was picked arbitrarily: indeed, we would ideally determine the value experimentally would require user testing, a process for which this project has neither the ambition nor the resources. Note that the condition on  $T_{Haar}$  is necessary but not sufficient: indeed, we empirically found some pairs of fingerprints with a low PSI distance that are not easy to distinguish.

## 6 Existing schemes

### 6.1 Antoine’s

Antoine’s visual hash is a very simple scheme:  $H$  is the standard SHA-256 function, and  $V$  is a trivial bit-to-binary pixel mapping (a bit of value 1 gets mapped to a pixel of value 1 (white), a bit of value 0 gets mapped to a pixel of value 0 (black)). It is clear that this scheme is not intermediate pre-image resistant, as the trivial mapping can easily be reversed. Moreover, it is not sparse, as any combination of 256 binary pixels is a valid output: from a fingerprint  $F$  one can produce a fingerprint  $F'$  that differs from  $F$  in only one pixel, which means  $F \approx_v F'$  in general. Therefore, we can conduct the attack described in section 4.1 to create a forgery.

### 6.2 LifeHash

The LifeHash library [ref] employs the standard SHA-256 hash function, but it introduces some complexity in the definition of  $b$  using John Conway’s Game of Life [ref]. Let  $\mathcal{G} : \mathbb{P}^{m \times n} \rightarrow \mathbb{P}^{m \times n}$  be the result of a single generation of the game of life rules on the input. We suppose fingerprints support element-wise addition, as well as multiplication by a constant.

Algorithm 1 describes the working of the LifeHash visual hashing function.

```

Data:  $x \in \{0, 1\}^*$ 
Result:  $F \in \mathbb{P}_{RGB}^{m_n}$ 
begin
   $B \leftarrow \text{Antoine}(x)$ 
  Initialize  $G$  as the empty array
  while  $|G| < 150$  and Game of Life has not converged do
    Append  $B$  to  $G$ 
     $B \leftarrow \mathcal{G}(B)$  ▷ One iteration of Game of Life rules, 1 is alive, 0 is dead
  Initialize  $F$  as an all-white (255,255,255) fingerprint the same size as  $B$ 
  for  $i = 0, \dots, |G| - 1$  do
     $F = F - G[i] * \frac{255i}{|G|}$  ▷ Aggregate all generations
  Apply RGB mapping to  $F$ , using some of the first bits to pick a color scheme.
  Apply symmetry to  $F$ , using some of the first bits to pick a transformation scheme.
  return  $F$ 

```

**Algorithm 1:** LifeHash

Due to the unpredictable nature of the Game of Life, one would be tempted to think that this scheme is intermediate pre-image resistant: since one can construct a universal Turing machine in the Game of Life [ref], being able to predict for an arbitrary  $h$  the final state of the Game of Life (or even whether the game of life converges in 150 generations)

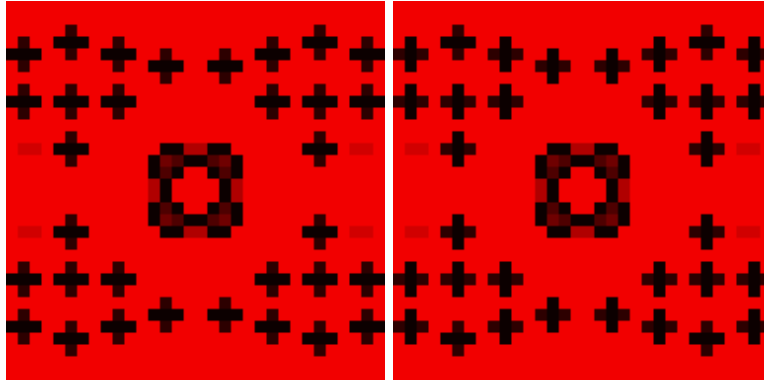


Figure 6: LifeHash outputs when  $H = \hat{h}_1$  (left) and  $H = \hat{h}_2$  (right)

can be reduced to the Halting problem [ref]. But while it may be intermediate resistant on complex inputs, the Game of Life has known, simple patterns that oscillate, which might cause problems if the input hash is simple (i.e. has a low Hamming weight). We describe an example of a forgery in order to illustrate that this scheme is neither intermediate pre-image resistant nor visually injective:

Let the initial state of the Game of Life  $B_1$  contain a known oscillator. Let the following state sequence when applying the Game of Life rules successively:

$$B_1 \xrightarrow{\mathcal{G}} B_2 \xrightarrow{\mathcal{G}} B_3 \xrightarrow{\mathcal{G}} \dots \xrightarrow{\mathcal{G}} B_1$$

Now, say we have two different input  $x_1$  and  $x_2$  such that  $\text{Antoine}(H(x_1)) = B_1$  and  $\text{Antoine}(H(x_2)) = B_2$ . Then the execution of LIFEHASH will go through the same sequence of Game of Life (although rotated one state left for the second), yielding visually similar results. Figure 6 shows the output picture on input  $x_1$  and  $x_2$  when

$$h(x_1) := \hat{h}_1 = 00000000000000e00e0e0000e00000000eee00000000000006e030002000e0008$$

$$h(x_2) := \hat{h}_2 = 00000000000004e404044404400044440444044400000040064030402000e0008$$

Both hashes contain the same period-3 oscillator (with the same initial state) in the center, and a collection of period-2 oscillators around it: 3 horizontally aligned pixels will flip to 3 vertically aligned pixels back and forth, thus giving a cross pattern on the fingerprint. Both need  $\text{lcm}(3, 2) = 6$  generations of the Game of Life to converge. The fingerprints differ only in the initial direction of the period-2 oscillators: on the left they start horizontally, whereas on the right they start vertically. Because each period-2 oscillators go through 3 entire cycles, the difference between the two setups becomes unnoticeable.

This pattern is really simple, but it is hard to determine whether other more complicated examples of oscillators combinations yield similar results.

### 6.3 PRG-based libraries

Some visual hash function use Pseudo-Random Generators (PRG) to create fingerprints. The idea is to seed a PRG with the input value, and let the PRG's randomness act to create fingerprints. Most implementations of such schemes use them directly on the messages (as opposed to hashes), which is not achievable in our case. However, they may have some desirable properties that we analyze in the following sections.

#### 6.3.1 Blockies

Blockies [ref] is a light-weight javascript library used to create icons. It is at the time of writing quite popular (more than 1500 weekly downloads according to the official website), following the gain in popularity of cryptocurrency exchange and the need to identify 42-hexadecimal-characters Ethereum addresses.

Blockies is based on a custom-made PRG (!) based on the method used by Java to calculate hashCodes



**Data:**  $x \in \{0, 1\}^*$

**Result:**  $F \in \mathbb{P}_{\text{RGB}}^{m \times n}$

**begin**

$\text{PRG}_{\text{blockies}}.\text{Seed}(x)$  //  $\text{PRG}_{\text{blockies}}$  returns values between 0 and 1

$\text{colorBack} \leftarrow (\text{PRG}_{\text{blockies}}.\text{next}() * 255, \text{PRG}_{\text{blockies}}.\text{next}() * 255, \text{PRG}_{\text{blockies}}.\text{next}() * 255)$

$\text{colorFront} \leftarrow (\text{PRG}_{\text{blockies}}.\text{next}() * 255, \text{PRG}_{\text{blockies}}.\text{next}() * 255, \text{PRG}_{\text{blockies}}.\text{next}() * 255)$

$\text{colorSpots} \leftarrow (\text{PRG}_{\text{blockies}}.\text{next}() * 255, \text{PRG}_{\text{blockies}}.\text{next}() * 255, \text{PRG}_{\text{blockies}}.\text{next}() * 255)$

    Initialize  $F$  as an all-zero fingerprint of size  $m \times n$

**foreach** pixel  $p$  in  $F$  such that  $p$  is in the left half of  $F$  **do**

$\text{val} \leftarrow \text{PRG}_{\text{blockies}}.\text{next}() * 2.3$

**if**  $\text{val} < 1$  **then**

$p \leftarrow \text{colorBack}$

// happens with prob  $\frac{1}{2.3} \approx 43.5\%$

**else if**  $1 \leq \text{val} < 2$  **then**

$p \leftarrow \text{colorFront}$

// happens with prob  $\frac{1}{2.3} \approx 43.5\%$

**else**

$p \leftarrow \text{colorSpots}$

// happens with prob  $\frac{0.3}{2.3} \approx 13\%$

    Set the right half of  $F$  to the mirrored left half of  $F$  **return**  $F$

**Algorithm 2:** PRG-based visual hash function

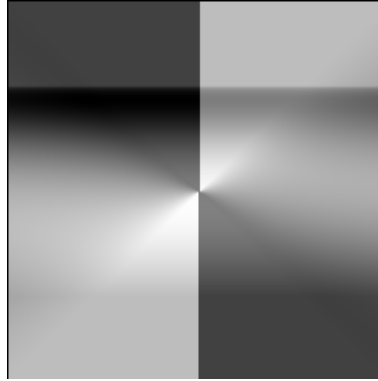


Figure 7: RandomArt visual hash.

### 6.3.2 RandomArt

[Ref] introduces another way of generating visual hashes using randomness: the idea is again to seed a PRG with the hash, then to use recursion to generate a random function to apply to each pixel. To visualise this, we construct an example. We choose a depth of 4, meaning we have 4 level of nesting. The following function was generated :

$f(x, y) = \text{SUB}(\text{SIN}(\text{SQRT}(\text{DIV}(y, 0.56))), \text{SIN}(\text{ADD}(\text{DIV}(y, x), \text{SQRT}(-0.35))))$

In mathematical notation :

$$f(x, y) = \sin\left(\sqrt{\frac{y}{0.56}}\right) - \sin\left(\frac{y}{x} + \sqrt{-0.35}\right)$$

Note that each function maps the interval  $[-1, 1]$  to itself, using the necessary pre-processing if needed, hence the square root of the negative value will not cause any error.

This function is then evaluated at each pixel. Figure 7 shows the resulting fingerprint. In the true scheme, we repeat the process for each color canal.

The issue with this scheme is that the time needed to generate the random function grows exponentially with depth. Also, similarly to the Blockies scheme, the security of RandomHash relies entirely on the security of the PRG. Also, because of the choice of primitives and their properties, different functions can produce the exact same image; for instance

$$\text{ADD}(\text{MULT}(x, y), \text{MULT}(0.5, x)) == \text{MULT}(x, \text{ADD}(y, 0.5))$$

because  $xy + 0.5x = x(y + 0.5)$  for all  $x, y$  by distributivity of addition. However, [ref] gives us another tool: the way they evaluated if their hashes were recognisable by humans. Indeed, they give without proof two criteria to assess the quality of their image:

- The fingerprints should have a high compression factor, as low compression is a sign of a chaotic image.
- The fingerprints' frequency spectra should be concentrated around low frequencies, as high frequencies typically indicate noise.

The second criterion is what inspired us in the design of our scheme.

## 7 Fourier scheme

The complexity of the RandomHash function generation algorithm grows exponentially with depth, which can incur high computing time. So we design a new scheme with the following idea: instead of generating a random function and hoping that its frequency spectrum is centered around low frequencies, we proceed "in reverse" and generate a frequency spectrum centered around low frequencies then apply a 2-dimensional inverse Fourier transform to create the image.

### 7.1 2D Discrete Fourier Transform

The 2D Discrete Fourier Transform (2DDFT) is defined as follows:

$$\mathcal{F}(u, v) = \frac{1}{MN} \sum_{x=0}^{M-1} \sum_{y=0}^{N-1} f(x, y) e^{-i2\pi(ux/M + vy/N)}$$

and the 2D Inverse Fourier Transform (2DIDFT):

$$f(x, y) = \sum_{u=0}^{M-1} \sum_{v=0}^{N-1} \mathcal{F}(u, v) e^{i2\pi(ux/M + vy/N)}$$

where  $M$  and  $N$  are the dimensions of the input image, and  $f(x, y)$  denotes the value of the pixel at position  $(x, y)$ . Although  $\mathcal{F}(u, v)$  is defined for any  $u, v \in \mathbb{Z}$ , since  $\mathcal{F}(u + M, v) = \mathcal{F}(u, v + N) = \mathcal{F}(u, v)$ , we restrict the input domain to  $[0, M - 1]$  for  $u$  and  $[0, N - 1]$  for  $v$  (similarly to how we would sometimes restrict the input domain to  $[0, 2\pi)$  for trigonometric functions).

### 7.2 Spectrum generation

What we would like to achieve is to map a hash to an array of coefficients in the frequency domain, and then apply the inverse 2D Transform on this frequency spectrum to get a fingerprint [UNCLEAR]. In this section we give an insight as to how we would generate this spectrum.

First, we have that  $\mathcal{F}(u, v)$  is a complex number, meaning each element in the output space of  $\mathcal{F}$  can be thought of as a pair of real numbers (either complex and imaginary parts in a Cartesian system or a modulus and phase in a polar system).

There is a symmetry in the output space: indeed, we have

$$\begin{aligned} \mathcal{F}(-u, -v) &= \frac{1}{MN} \sum_{x=0}^{M-1} \sum_{y=0}^{N-1} f(x, y) e^{-i2\pi(-ux/M - vy/N)} \\ &= \frac{1}{MN} \sum_{x=0}^{M-1} \sum_{y=0}^{N-1} f(x, y) e^{i2\pi(ux/M + vy/N)} \\ &= \mathcal{F}(u, v)^* \end{aligned}$$

where  $(\cdot)^*$  denotes the complex conjugate. That means we actually have to sample  $MN/2$  complex numbers in order to generate a  $M \times N$  spectrum.

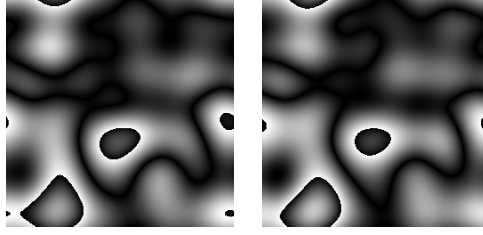


Figure 8: Two fingerprints with one bit flipped



Figure 9: The two fingerprints, filtered

Figure X shows a graphical representation of how we generate our spectrum. The darker squares represent a number with a greater module.

[SOMETHING ABOUT DISTANCE HERE PROBABLY]

Once the spectrum is generated, the 2DIDFT is applied. The resulting array is then scaled by the inverse of its mean in order to have a mean of 1 [bc very few have values > 1], and then multiplied by 255 to get an array of gray pixels. We apply the modulo 255 operator to values greater than 255 to generate the final pixel array.

### 7.2.1 Mapping bits to frequencies

We would like every bit to have roughly the same "importance" with respect to the final image. Therefore, we map words of 4 bits to frequencies according to the encoding depicted on figure [FIGURE ENCODING], in which flipping one bit changes both the phase and the module of the resulting frequency. One word is mapped to 2 frequencies, which are places according to [FIGURE SPECTRUM].

Figure 8- shows the resulting fingerprints with  $h = \dots$  and  $h' = \dots$ . We can see that the results are very hard to distinguish. Furthermore, the PSI is 0.6367, which is more than the threshold value, therefore, we consider those hashes to be visually indistinguishable.

### 7.2.2 Filtering

During testing, we noticed that the high-contrast zones due to overflows tend to stand out to the eye and thus help with the distinction of fingerprints [REF WOULD BE NICE]. Therefore, we decided to apply a simple filter to enforce very high contrast zones: namely, from an array  $a$  of real values issued from the 2DIDFT, we apply

$$f(a, x, y) = \begin{cases} 255 & \text{if } a(x, y) \geq 0.5 \\ 0 & \text{otherwise} \end{cases}$$

Figure 9- shows the two same  $h$  and  $h'$ . As expected, the two images are more easily distinguished. The PSI is 0.2135.

### 7.2.3 Number of functions

While filtering offers a clear improvement, there is still an issue: we have to map  $|h|$  bits to frequencies, which results in having to sample spectrum too far away from the axes. To address this problem, we instead perform the sampling procedure [3] times. That way we have 3 spectra array  $a_{red}, a_{green}, a_{blue}$ , each corresponding to one color channel. The color of a pixel at position is  $C(x, y) := (f(a_{red}, x, y), f(a_{green}, x, y), f(a_{blue}, x, y))$  where  $f$  is the same filter as

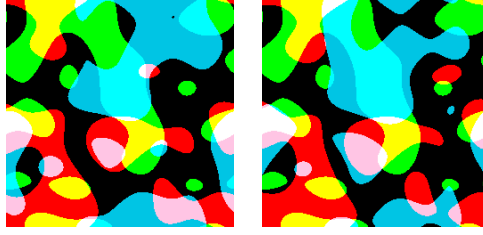


Figure 10: The two fingerprints, with 3 functions

Filtered values of $a_0(x, y), a_1(x, y), a_2(x, y)$	$C(x, y)$
0, 0, 0	<i>Colors</i> [0]
0, 0, 255	<i>Colors</i> [1]
0, 255, 0	<i>Colors</i> [2]
0, 255, 255	<i>Colors</i> [3]
255, 0, 0	<i>Colors</i> [4]
255, 0, 255	<i>Colors</i> [5]
255, 255, 0	<i>Colors</i> [6]
255, 255, 255	<i>Colors</i> [7]

Table 1: Color mapping

in the previous section. Figure 10 shows the two same hashes. We can see that flipping one bit changes one spectrum (in this case, the spectrum corresponding to the blue channel). The resulting images are therefore similarly close to the previous step (Haar-PSI = 0.2158).

#### 7.2.4 Colors

As we have just seen, a difference in colors helps with distinguishability. We would thus have a lot to gain picking a different set of colors for different hashes. As we would like to use colors more colors than the combinations of the pure colors (colors that have only one non-zero), we must slightly change the way we assign a color to a pixel: we must choose  $2^3$  colors to account for every possibility of the filtered values of  $a_0, a_1$  and  $a_2$ . Suppose we have an array named *Colors* of length 8, containing rgb encodings.

Note that the system described in section 7.2.3 corresponds to picking the color to be the filtered values directly. To choose the set of 8 colors that will make the output set of  $C(x, y)$ , we use parity bits:

Let the parity bits  $p_0$  to  $p_{15}$  be defined as

$$p_i = \begin{cases} b_i \oplus b_{12+i} \oplus b_{24+i} \oplus \dots \oplus b_{108+i} = \bigoplus_{k=0}^9 b_{12k+i} & 0 \leq i < 12 \\ b_{108+i} \oplus b_{112+i} & 12 \leq i < 16 \end{cases}$$

We compute the colors indices:

$$c_j = p_{2j} || p_{2j+1} \text{ for } j \in \{0, \dots, 7\}$$

where  $||$  denotes concatenation, therefore  $c_j \in \{0, 1, 2, 3\}$ . Let *Palette* be a fixed 32 colors palette. We fill *Colors* as :

$$Colors[j] = Palette[4j + c_j] \text{ for } j \in \{0, \dots, 7\}$$

Figure [FIGURE] depicts this process, and fig. 11 shows the two hashes. We can see that flipping one bit results in one color change. While in our case the results are quite obviously different (Haar PSI = 0.25 [COMMENT ON THIS]), there are cases where only a very small portion of the image is of a certain color, which would yield almost no difference should this color be changed.

An even stronger attack would be to flip 2 bits that correspond to the same parity bit: flipping both  $b_j$  and  $b_{j+12k}$  yields a fingerprint that differs only in the shapes of the patterns, which may not be significant enough.

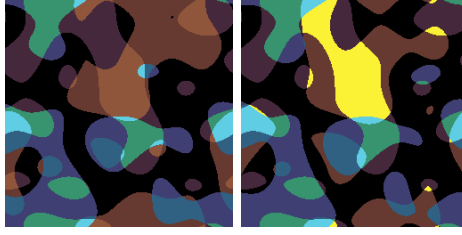


Figure 11: The two fingerprints, with 3 functions

### 7.2.5 Palette shift and permutation

As having only one color changed may result in very little zones switching colors, we would like to change more than one color if one bit is flipped. To achieve that, we apply a transformation to the palette before extracting the colors: namely, we shift it by  $h \bmod p$ , where  $h$  is interpreted as a 128-bit integer, and  $p$  is a prime number less than 32 [MAYBE JUST ODD NUMBER] to be determined later.

Once the colors are extracted, we re-order the array according to a permutation extracted from the hash. The set of possible permutations is a code that achieve a satisfactory maximal distance.

### 7.2.6 Symmetries

Humans tend to perform better at identifying symmetrical objects [ref]. Therefore, we would like our fingerprints to be symmetric. The way we do this is to first generate the fingerprint as we did previously, and then apply some symmetry according to a symmetry mode [DETAILS].

In order to determine the symmetry mode, we simply compute  $h \bmod q$ , where  $q \neq p$  is a small prime (in the order of 20).

To have the same the same symmetry, one would have to achieve  $k - \ell = 0 \bmod (|2| \bmod q)$ . Having the same palette shift AND same symmetry yields

$$\begin{cases} k - \ell = 0 \bmod (|2| \bmod p) \\ k - \ell = 0 \bmod (|2| \bmod q) \end{cases} \iff k - \ell = 0 \bmod \text{lcm}(|2| \bmod p, |2| \bmod q)$$

Picking  $p, q \in \{23, 29\}$  yields  $\text{lcm}(|2| \bmod p, |2| \bmod q) = \text{lcm}(11, 28) = 308$ . Therefore, if 2 bits are flipped, then either the palette is shifted differently, or the symmetry is not the same, or both. Either of these cases yields visually distinguishable results.

## 7.3 Worst cases

$$h + \sum_{\ell} 2^{\ell} - \sum_k 2^k = h \bmod p$$

If only two bits are changed we have

$$\begin{aligned} h + 2^k - 2^{k+12m} &= h \bmod p \\ h + 2^k(2^{12m} - 1) &= h \bmod p \\ 2^k(2^{12m} - 1) &= 0 \bmod p \\ 2^{12m} - 1 &= 0 \bmod p \\ 2^{12m} &= 1 \bmod p \\ 12m &= 0 \bmod (|2| \bmod p) \end{aligned}$$

We now analyse the worst case scenarios against our scheme. Let us assume  $n$  bits were flipped. Let us pick  $p = 29, q = 23$ . We have  $pq = 667$ ,  $|2| \bmod p = 28$ ,  $|q| \bmod q = 11$

- $n = 1$  : The palette is flipped, the shift is different, the symmetry is different.
- $n = 2$  : Either the shift is different, the symmetry is different, or both.
- $n > 2, n$  odd : The palette is flipped.

For an even  $n > 2$ , there are collisions in the sense that some pairs of hashes use the same palette, same colors and same symmetry.

- For  $n = 4$ , we have  $2^0 + 2^{84} + 2^{96} + 2^{108} \equiv 0 \pmod{667}$  which means flipping bits with indices  $0 + d, 84 + d, 96 + d, 108 + d$  with  $d \in \{0 \dots 11\}$  yields the same configuration if all the bits at these indices have the same value. We estimate the probability that  $h$  contains a 4 bits of same value at at least one of the possibilities for  $d$  to be  $1 - \left(1 - \frac{1}{8}\right)^{12} = 1 - \left(\frac{7}{8}\right)^{12} \approx \frac{4}{5}$ . (Note that this estimation assumes independent and uniform bits, which is not accurate but we just want to show that the probability has a high order of magnitude). Because that is a very high probability, we need patterns to change enough when 4 frequencies are changed. Empirically, this seems to be the case