

Automated Security Proofs for Almost-Universal Hash for MAC Verification*

Martin Gagné¹, Pascal Lafourcade², and Yassine Lakhnech²

¹ Department of Computer Science, Saarland University, Germany

² Université Grenoble 1, CNRS, VERIMAG, France

Abstract. Message authentication codes (MACs) are an essential primitive in cryptography. They are used to ensure the integrity and authenticity of a message, and can also be used as a building block for larger schemes, such as chosen-ciphertext secure encryption, or identity-based encryption. MACs are often built in two steps: first, the ‘front end’ of the MAC produces a short digest of the long message, then the ‘back end’ provides a mixing step to make the output of the MAC unpredictable for an attacker. Our verification method follows this structure. We develop a Hoare logic for proving that the front end of the MAC is an almost-universal hash function. The programming language used to specify these functions is fairly expressive and can be used to describe many block-cipher and compression function-based MACs. We implemented this method into a prototype that can automatically prove the security of almost-universal hash functions. This prototype can prove the security of the front-end of many CBC-based MACs (DMAC, ECBC, FCBC and XCBC to name only a few), PMAC and HMAC. We then provide a list of options for the back end of the MAC, each consisting of only two or three instructions, each of which can be composed with an almost-universal hash function to obtain a secure MAC.

1 Introduction

Message authentication codes (MACs) are among the most common primitives in symmetric key cryptography. They ensure the integrity and provenance of a message, and they can be used, in conjunction with chosen-plaintext (CPA) secure encryption, to obtain chosen-ciphertext (CCA) secure encryption. Given the importance of this primitive, it is important that their proofs of security be the object of close scrutiny. The study of the security of MACs is, of course, not a new field. Bellare et al. [5] were the first to prove the security of CBC-MAC for fixed-length inputs. Following this work, a myriad of new MACs secure for variable-length inputs were proposed ([4,7,8,9,17]). None of these protocols’ proofs have been verified by any means other than human scrutiny. Automated proofs can provide additional assurance of the correctness of these security proofs by providing an independent proof of complex schemes. This paper presents a method for automatically proving the security of MACs based on block ciphers and hash functions.

Contributions: To analyze the security of MACs, we first decompose the MAC algorithms into two parts: a ‘front-end’, whose work is to compress long input messages

* This work was partially supported by ANR project ProSe and Minalogic project SHIVA.

into small digests, and a ‘back-end’, usually a mixing step, which obfuscates the output of the front-end. We present a Hoare logic to prove that the front-ends of block-cipher based and hash based MACs are almost-universal hash functions in the ideal cipher model and random oracle model respectively. We then make a list of operations which, when composed with an almost-universal hash function, yield a secure MAC. We can then attest the security of MACs by first proving the security of the front end using our logic, and then by manually verifying that the back end of the MAC belongs to our list.

Our result differs significantly from previous works that used Hoare logic to generate proofs of cryptographic protocols (such as [12, 15]) because those results proved the security of encryption schemes. Proving the security of MACs proved to be singularly more challenging: the security of encryption schemes could be simply proven by showing that the ciphertext is indistinguishable from a random value, whereas the unforgeability property required of MACs cannot, to our knowledge, be captured by their predicates. As a result, we have to consider the simultaneous execution of the program, define a dedicated semantics to capture these executions, and introduce appropriated predicates that keep track of equality and inequality of values between the two executions.

In contrast to the previous results that only deal with schemes that had fixed-length inputs, we are able to analyze for-loops, which allows us to prove the security of protocols that can take arbitrary strings as an input. We describe two heuristics that can be used to discover stable loop invariants and apply them to one example. These heuristics successfully find stable invariants for all the hash functions analyzed in this paper.

Finally, we implemented our method into a prototype [14] that can be used to verify the security of the front-end of several well-known MACs, such as HMAC [4], DMAC [17], ECB, FCBC and XCBC [8] and PMAC [9], and could be used to verify the security of other hash functions based on the same primitives. We also give a predicate filter that enables us to discard unnecessary predicates, which speeds up our implementation and facilitates the discovery of loop invariants. Our prototype goes through the programs from beginning to end, instead of the more common backward approach, to avoid an exponential blowout in the number of possibilities to examine, due to the many choices of rules that can cause certain predicates caused by the presence of the logical or connector in our Hoare logic.

Related Work: The idea of using Hoare logic to automatically produce proofs of security for cryptographic protocols is not new. Courant et al. [12] presented a Hoare logic to prove the security of asymmetric encryption schemes in the random oracle model. A Hoare logic was also used by Gagné et al. [15] to verify proofs of security of block cipher modes of encryption. Also worth mentioning is the paper by Corin and Den Hartog [11], which presented a Hoare-style proof system for game-based cryptographic proofs.

Fournet et al. [13] developed a framework for modular code-based cryptographic verification. However, their approach considers interfaces for MACs. In a way, our work is complementary to theirs, as our result, coupled with theirs, could enable a more complete verification of systems.

In [1], the authors introduce a general logic for proving the security of cryptographic primitives. This framework can easily be extended using external results, such as [12], to add to its power. Our result could also be added to this framework to further extend it.

Other tools, such as Cryptoverif [10] and EasyCrypt [3,2], can be used to verify the security of cryptographic schemes, but they are not as convenient as our method for proving the security of MACs. Cryptoverif does not support loop constructs, which are an important part of our result, and is generally used for proving the security of higher level protocols, assuming the security of primitives such as MACs. As for EasyCrypt, it relies on a game-based approach and requires human assistance to enter the sequence of games. Our result is complementary to these approaches. Integrating our method to these tools would enable a more complete analysis of cryptographic protocols and remove the need for human assistance when analysing MACs.

Outline: In Section 2, we introduce cryptographic background. The following section introduces our grammar, semantics and assertion language. In Section 4, we present our Hoare logic and method for proving the security of almost-universal hash functions, and we discuss our implementation of this logic and treatment of loops in Section 5. We then obtain a secure MAC by combining these with one of the back-end options described in Section 6. Finally, we conclude in Section 7.

2 Cryptographic Background

In this section, we introduce a few notational conventions, and we recall a few cryptographic concepts.

Notation and Conventions

We assume that all variables range over domains whose cardinality is exponential in the security parameter η and that all programs have length polynomial in η . We say that a function $f : \mathbb{N} \rightarrow \mathbb{R}$ is *negligible* if, for any polynomial p , there exists a positive integer n_0 such that for all $n \geq n_0$, $f(n) \leq \frac{1}{|p(n)|}$.

For a probability distribution \mathcal{D} , we denote by $x \stackrel{\$}{\leftarrow} \mathcal{D}$ the operation of sampling a value x according to distribution \mathcal{D} . If S is a finite set, we denote by $x \stackrel{\$}{\leftarrow} S$ the operation of sampling x uniformly at random among the values in S .

MAC Security

A message authentication code ensures the authenticity of a message m by computing a small tag τ , which is sent together with the message to the intended receiver. Upon receiving the message and the tag, the receiver recomputes the tag τ' using the message and his own copy of the key, and he accepts the message as authentic if $\tau = \tau'$. More formally:

Definition 1 (MAC). *A message authentication code is a triple of polynomial-time algorithms (K, MAC, V) , where $K(1^\eta)$ takes a security parameter 1^η and outputs a secret key sk , $MAC(sk, m)$ takes a secret key and a message m , and outputs a tag, and $V(sk, m, tag)$ takes a secret key sk , a message m and a tag, and outputs a bit: 1 for a correct tag, 0 otherwise.*

We say that a MAC is secure, or *unforgeable* if it is impossible to compute a new valid message-tag pair for anybody who does not know the secret key, even when given access to oracles that can compute and verify the MACs. This way, when one receives

a valid message-tag pair, he can be certain that the message was sent by someone who possesses a copy of his secret key.

Definition 2 (Unforgeability [5]). A MAC (K, Mac, V) is unforgeable under a chosen-message attack (UNF-CMA) if for every polynomial-time algorithm \mathcal{A} that has oracle access to the MAC and verification algorithm and whose output message m^* is different from any message it sent to the Mac oracle, the following probability is negligible

$$\Pr[sk \xleftarrow{\$} K(1^\eta); (m^*, tag^*) \xleftarrow{\$} \mathcal{A}^{Mac(sk, \cdot), V(sk, \cdot, \cdot)} : V(sk, m^*, tag^*) = 1]$$

A standard method for constructing MACs is to apply a pseudo-random function, or some other form of ‘mixing’ step, to the output of an almost-universal hash function [18, 19]. We assume that a MAC is constructed in this way.

Definition 3 (Almost-Universal Hash). A family of functions $\mathcal{H} = \{h_i\}$ indexed with key $i \in \{0, 1\}^\eta$ is a family of almost-universal hash functions if for any two distinct strings M and M' , $\Pr_{h_i \in \mathcal{H}}[h_i(M) = h_i(M')]$ is negligible, where the probability is taken over the choice of h_i in \mathcal{H} .

It is much easier to work with this definition than with the unforgeability definition because of the absence of an adaptive adversary, and the collision probability is taken over all possible choices of key.

Block Cipher Security

Many MAC constructions are based on block cipher, so we quickly recall the definition of block ciphers and their security definition.

A block cipher is a family of permutations $\mathcal{E} : \{0, 1\}^{K(\eta)} \times \{0, 1\}^\eta \rightarrow \{0, 1\}^\eta$ indexed with a key $k \in \{0, 1\}^{K(\eta)}$ where $K(\eta)$ is a polynomial. A block cipher is *secure* if, for a randomly sampled key, the block cipher is indistinguishable from a permutation sampled at random from the set of all permutations of $\{0, 1\}^\eta$. However, since random permutations of $\{0, 1\}^\eta$ and random functions from $\{0, 1\}^\eta$ to $\{0, 1\}^\eta$ are statistically close, and that random functions are often more convenient for proof purposes, it is common to assume that secure block ciphers are pseudo-random functions.

Definition 4 (Pseudo-Random Functions). Let $P : \{0, 1\}^{K(\eta)} \times \{0, 1\}^\eta \rightarrow \{0, 1\}^\eta$ be a family of functions and let \mathcal{A} be an algorithm that takes an oracle and returns a bit. The prf-advantage of \mathcal{A} is defined as follows.

$$\text{Adv}_{\mathcal{A}, P}^{\text{prf}} = \left| \Pr[k \xleftarrow{\$} \{0, 1\}^{K(\eta)}; \mathcal{A}^{P(k, \cdot)} = 1] - \Pr[R \xleftarrow{\$} \Phi_\eta; \mathcal{A}^{R(\cdot)} = 1] \right|$$

where Φ_η is the set of all functions from $\{0, 1\}^\eta$ to $\{0, 1\}^\eta$. We say that P is a family of pseudo-random functions if for every polynomial-time adversary \mathcal{A} , $\text{Adv}_{\mathcal{A}, P}^{\text{prf}}$ is a negligible function in η .

Since all the schemes in this paper require only one key for the block cipher, to simplify the notation, we write only $\mathcal{E}(m)$ instead of $\mathcal{E}(k, m)$, but it is understood that a key was selected at the initialization of the scheme, and remains the same throughout.

Random Oracle Model

For MACs that make use of a hash function, we assume that the hash function behaves like a random oracle. That is, we assume that the hash function is picked at random among all possible functions from the given domain and range, and that every algorithm participating in the scheme, including all adversaries, has oracle access to this random function. This is a fairly common assumption to provide a heuristic argument for the security of cryptographic protocols [6].

Indistinguishable Distributions

Given two distribution ensembles $X = \{X_\eta\}_{\eta \in \mathbb{N}}$ and $X' = \{X'_\eta\}_{\eta \in \mathbb{N}}$, an algorithm \mathcal{A} and $\eta \in \mathbb{N}$, we define the *advantage* of \mathcal{A} in distinguishing X_η from X'_η as the following quantity:

$$\text{Adv}(\mathcal{A}, \eta, X, X') = \left| \Pr[x \stackrel{\$}{\leftarrow} X_\eta : \mathcal{A}(x) = 1] - \Pr[x \stackrel{\$}{\leftarrow} X'_\eta : \mathcal{A}(x) = 1] \right|.$$

We say that X and X' are *indistinguishable*, denoted by $X \sim X'$, if $\text{Adv}(\mathcal{A}, \eta, X, X')$ is negligible as a function of η for every probabilistic polynomial-time algorithm \mathcal{A} .

3 Model

In this section, we introduce the grammar for the programs describing almost-universal hash function. We present the semantics of each commands, and introduce the assertion language that will be used in for our Hoare logic.

3.1 Grammar

We consider the language defined by the BNF grammar below, where p and q are positive integers.

$$\begin{aligned} \text{cmd} ::= & x := \mathcal{E}(y) \mid x := \mathcal{H}(y) \mid x := y \mid x := y \oplus z \mid x := y \parallel z \mid x := \rho(i, y) \\ & \mid \text{for } l = p \text{ to } q \text{ do: } [\text{cmd}_l] \mid \text{cmd}_1; \text{cmd}_2 \end{aligned}$$

We refer to individual instructions as *commands* and to lists of commands as *programs*. Each command has the following effect:

- $x := \mathcal{E}(y)$ denotes application of the block cipher \mathcal{E} to the value of y and assigning the result to x .
- $x := \mathcal{H}(y)$ denotes the application of the hash function \mathcal{H} to the value of y and assigning the result to x .
- $x := y$ denotes the assignment to x of the values of y .
- $x := y \oplus z$ denotes the assignment to x of the xor or the values of y and z .
- $x := y \parallel z$ denotes the assignment to x of the concatenation of the values of y and z .
- $x := \rho(i, y)$ denotes the computation of the function ρ on input i (an integer) and the value of y and assigning the result to x .
- $c_1; c_2$ is the sequential composition of c_1 and c_2 .

- for $l = p$ to q do: $[\text{cmd}_l]$ denotes the successive execution of $\text{cmd}_p; \text{cmd}_{p+1}; \dots; \text{cmd}_q$ when $p \leq q$. If $p > q$, the command has no effect.

The function ρ is used to process the *tweak* in a common construction for *tweakable block ciphers* [16]. A fixed-input-length almost-universal function is often sufficient, but exact implementations vary from one scheme to the next, and we want to allow for the possibility of functions that have additional properties. When a scheme uses a function ρ , the properties of the function ρ required for the proof will be added to the initial conditions of the verification procedure using the predicates of Section 3.3. We do not any other assumptions about ρ other than it is a function with fixed output length.

Definition 5 (Generic Hash Function). *A generic hash function $Hash$ on message blocks m_1, \dots, m_n with output c_n , is represented by a tuple $(\mathcal{F}_\mathcal{E}, \mathcal{F}_\mathcal{H}, Hash(m_1 \parallel \dots \parallel m_n, c_n) : \text{var } \mathbf{x}; \text{cmd})$, where $\mathcal{F}_\mathcal{E}$ is a family of pseudorandom permutations (usually a block cipher), $\mathcal{F}_\mathcal{H}$ is a family of cryptographic hash functions, and $Hash(m_1 \parallel \dots \parallel m_n, c_n) : \text{var } \mathbf{x}; \text{cmd}$ is the program of the hash function, where \mathbf{x} is the set of all the variables in the program that are neither input variables m_i , output variable c_n , or the special variable k (used to hold a secret key), and the program cmd is in the language described by our grammar.*

The secret key sk of the generic hash is a combination of the value of the special variable k and the choice of the block cipher \mathcal{E} in the family $\mathcal{F}_\mathcal{E}$.

We assume that, prior to executing the MAC, the message has been padded using some unambiguous padding scheme, so that all the message blocks m_1, \dots, m_n are of equal and appropriate length for the scheme, usually the input length of the block cipher. We also assume that each variable in the program cmd is assigned at most once, as it is clear that any program obtained from our language can be transformed into an equivalent program with this property, and that the input variables m_1, \dots, m_n never appear on the left side of any command since these variables already hold a value before the execution of the program. For simplicity of exposition, we henceforth assume that all the programs in this paper satisfy these assumptions.

We present to the right the program for $Hash_{CBC}$, the hash function that is used as a running example in this paper. We give the program for other hash functions that can be verified with our method in the full version of this paper [14].

$$\begin{aligned}
 & Hash_{CBC}(m_1 \parallel \dots \parallel m_n, c_n) : \\
 & \text{var } i, z_2, \dots, z_n, c_1, \dots, c_{n-1}; \\
 & c_1 := \mathcal{E}(m_1); \\
 & \text{for } i = 2 \text{ to } n \text{ do:} \\
 & \quad [z_i := c_{i-1} \oplus m_i; c_i := \mathcal{E}(z_i)]
 \end{aligned}$$

3.2 Semantics

In our analysis, we consider the execution of a program on two inputs simultaneously. These simultaneous executions will enable us to keep track of the probability of equality and inequality of strings between the two executions, thereby allowing us to prove that the function is almost-universal.

Each command is a function that takes a configuration and outputs a configurations. A *configuration* γ is a tuple $(S, S', \mathcal{E}, \mathcal{H}, \mathcal{L}_\mathcal{E}, \mathcal{L}_\mathcal{H})$ where S and S' are states, \mathcal{E} is a

block cipher, \mathcal{H} is a hash function (that will be modeled as a random oracle), and $\mathcal{L}_\mathcal{E}$ and $\mathcal{L}_\mathcal{H}$ are sets of strings.

A *state* is a function $S : \mathbf{Var} \rightarrow \{0, 1\}^* \cup \perp$, where \mathbf{Var} is the full set of variables in the program, that assigns bitstrings to variables (the symbol \perp is used to indicate that no value has been assigned to the variable yet). A configuration contains two states, one for each execution of the program.

The set $\mathcal{L}_\mathcal{E}$ records the values for which the functions \mathcal{E} was computed. The set is common for both executions of the program. Every time a command of the type $x := \mathcal{E}(y)$ is executed in the program, we add $S(y)$ and $S'(y)$ to $\mathcal{L}_\mathcal{E}$ if they are not already present. We define $\mathcal{L}_\mathcal{H}$ for the hash function \mathcal{H} similarly.

Let Γ denote the set of configurations and $\text{DIST}(\Gamma)$ the set of distributions on configurations. The semantics is given below, where $S\{x \mapsto v\}$ denotes the state which assigns the value v to the variable x , and behaves like S for all other variables and \circ denotes function composition. The semantic function $\text{cmd} : \Gamma \rightarrow \Gamma$ of commands can be lifted in the usual way to a function $\text{cmd}^* : \text{DIST}(\Gamma) \rightarrow \text{DIST}(\Gamma)$ by point-wise application of cmd . By abuse of notation we also denote the lifted semantics by $\llbracket \text{cmd} \rrbracket$.

$$\begin{aligned}
 \llbracket x := \mathcal{E}(y) \rrbracket(S, S', \mathcal{E}, \mathcal{H}, \mathcal{L}_\mathcal{E}, \mathcal{L}_\mathcal{H}) &= \\
 & (S\{x \mapsto \mathcal{E}(S(y))\}, S'\{x \mapsto \mathcal{E}(S'(y))\}, \mathcal{E}, \mathcal{H}, \mathcal{L}_\mathcal{E} \cup \{S(y), S'(y)\}, \mathcal{L}_\mathcal{H}) \\
 \llbracket x := \mathcal{H}(y) \rrbracket(S, S', \mathcal{E}, \mathcal{H}, \mathcal{L}_\mathcal{E}, \mathcal{L}_\mathcal{H}) &= \\
 & (S\{x \mapsto \mathcal{H}(S(y))\}, S'\{x \mapsto \mathcal{H}(S'(y))\}, \mathcal{E}, \mathcal{H}, \mathcal{L}_\mathcal{E}, \mathcal{L}_\mathcal{H} \cup \{S(y), S'(y)\}) \\
 \llbracket x := y \rrbracket(S, S', \mathcal{E}, \mathcal{H}, \mathcal{L}_\mathcal{E}, \mathcal{L}_\mathcal{H}) &= (S\{x \mapsto S(y)\}, S'\{x \mapsto S'(y)\}, \mathcal{E}, \mathcal{H}, \mathcal{L}_\mathcal{E}, \mathcal{L}_\mathcal{H}) \\
 \llbracket x := y \oplus z \rrbracket(S, S', \mathcal{E}, \mathcal{H}, \mathcal{L}_\mathcal{E}, \mathcal{L}_\mathcal{H}) &= \\
 & (S\{x \mapsto S(y) \oplus S(z)\}, S'\{x \mapsto S'(y) \oplus S'(z)\}, \mathcal{E}, \mathcal{H}, \mathcal{L}_\mathcal{E}, \mathcal{L}_\mathcal{H}) \\
 \llbracket x := y || z \rrbracket(S, S', \mathcal{E}, \mathcal{H}, \mathcal{L}_\mathcal{E}, \mathcal{L}_\mathcal{H}) &= \\
 & (S\{x \mapsto S(y) || S(z)\}, S'\{x \mapsto S'(y) || S'(z)\}, \mathcal{E}, \mathcal{H}, \mathcal{L}_\mathcal{E}, \mathcal{L}_\mathcal{H}) \\
 \llbracket x := \rho(i, y) \rrbracket(S, S', \mathcal{E}, \mathcal{H}, \mathcal{L}_\mathcal{E}, \mathcal{L}_\mathcal{H}) &= \\
 & (S\{x \mapsto \rho(i, S(y))\}, S'\{x \mapsto \rho(i, S'(y))\}, \mathcal{E}, \mathcal{H}, \mathcal{L}_\mathcal{E}, \mathcal{L}_\mathcal{H}) \\
 \llbracket \text{for } l = p \text{ to } q \text{ do: } [\text{cmd}_l] \rrbracket \gamma &= \begin{cases} \llbracket \text{cmd}_q \rrbracket \circ \llbracket \text{cmd}_{q-1} \rrbracket \circ \dots \circ \llbracket \text{cmd}_p \rrbracket \gamma & \text{if } p \leq q \\ \gamma & \text{otherwise} \end{cases} \\
 \llbracket c_1; c_2 \rrbracket &= \llbracket c_2 \rrbracket \circ \llbracket c_1 \rrbracket
 \end{aligned}$$

The set of initial distributions $\text{DIST}_0(\mathbb{H})$, where $\mathbb{H} = (\mathcal{F}_\mathcal{E}, \mathcal{F}_\mathcal{H}, \text{Hash}(m_1 || \dots || m_n, c_n) : \mathbf{var } x; \text{cmd})$ is a generic hash, contains all the following distributions:

$$\begin{aligned}
 \mathcal{D}_0^{(M, M')} &= [\mathcal{E} \stackrel{\$}{\leftarrow} \mathcal{F}_\mathcal{E}(1^n); \mathcal{H} \stackrel{\$}{\leftarrow} \mathcal{F}_\mathcal{H}(1^n); u \stackrel{\$}{\leftarrow} \{0, 1\}^n : \\
 & (S\{k \mapsto u, m_1 || \dots || m_n \mapsto M\}, S'\{k \mapsto u, m_1 || \dots || m_n \mapsto M'\}, \mathcal{E}, \mathcal{H}, \emptyset, \emptyset)]
 \end{aligned}$$

where M and M' are any two n block messages and k is a variable holding a secret string needed in some MACs (among our examples, $\text{Hash}_{\text{PMAC}}$ and $\text{Hash}_{\text{HMAC}}$ need it). Note that $\mathcal{F}_\mathcal{E}$, $\mathcal{F}_\mathcal{H}$, the domain \mathbf{Var} of the states and the length n of the input messages are defined in \mathbb{H} . These distributions capture the initial situation of Definition 3 where the variables m_i contain the blocks of M and M' in S and S' respectively.

The set $\text{DIST}(\mathbb{H})$ is obtained by executing a program on one of the initial distributions. It contains all the distributions of the form $\llbracket \text{cmd} \rrbracket X_0$, where $X_0 \in \text{DIST}_0(\mathbb{H})$ and cmd is a program.

A notational convention. It is easy to see that commands never modify \mathcal{E} or \mathcal{H} . Therefore, we can, without ambiguity, write $(\hat{S}, \hat{S}', \mathcal{L}'_{\mathcal{E}}, \mathcal{L}'_{\mathcal{H}}) \stackrel{\S}{\leftarrow} \llbracket c \rrbracket(S, S', \mathcal{L}_{\mathcal{E}}, \mathcal{L}_{\mathcal{H}})$ instead of $(\hat{S}, \hat{S}', \mathcal{E}, \mathcal{H}, \mathcal{L}'_{\mathcal{E}}, \mathcal{L}'_{\mathcal{H}}) \stackrel{\S}{\leftarrow} \llbracket c \rrbracket(S, S', \mathcal{E}, \mathcal{H}, \mathcal{L}_{\mathcal{E}}, \mathcal{L}_{\mathcal{H}})$.

3.3 Assertion Language

Like [15], our assertion languages deals with block ciphers, so it stands to reason that some of our predicates will be similar to theirs. However, the definition of all the predicates has to be adapted to our new semantics with two simultaneous executions. We also need additional predicates to describe equality or inequality of strings between the two executions, that will allow us to capture the definition of almost-universal hash functions. We first give an intuitive description of our predicates, then we define them all formally.

Empty: means that the probability that $\mathcal{L}_{\mathcal{E}}$ contains an element is negligible.

Eq(x, y): means that the probability that $S(x) \neq S'(y)$ is negligible.

Uneq(x, y): means that the probability that $S(x) = S'(y)$ is negligible.

E($\mathcal{E}; x; V$): means that the probability that the value of x is either in $\mathcal{L}_{\mathcal{E}}$ or equal to that of a variable in V is negligible.

H($\mathcal{H}; x; V$): means that the probability that the value of x is either in $\mathcal{L}_{\mathcal{H}}$ or equal to that of a variable in V is negligible.

Ind($x; V; V'$): means that no adversary has non-negligible probability to distinguish whether he is given results of computations performed using the value of x or a random value, when he is given the values of the variables in V and the values of the variables in V' from the parallel execution. In addition to variables in **Var**, the set V can contain special symbols $\ell_{\mathcal{E}}$ or $\ell_{\mathcal{H}}$. When the symbol $\ell_{\mathcal{E}}$ is present, it means that, in addition to the other variables in V , the distinguisher is also given the values in $\mathcal{L}_{\mathcal{E}}$, similarly for $\ell_{\mathcal{H}}$.

Our Hoare logic is based on statements from the following language.

$$\begin{aligned} \varphi &::= \varphi \wedge \varphi \mid \varphi \vee \varphi \mid \psi \\ \psi &::= \text{Ind}(x; W; V') \mid \text{Eq}(x, y) \mid \text{Uneq}(x, y) \mid \text{Empty} \mid \text{E}(\mathcal{E}; x; V) \mid \text{H}(\mathcal{H}; x; V) \end{aligned}$$

where $x, y \in \mathbf{Var}$ and $V, V' \subseteq \mathbf{Var}$, and $W \subseteq \mathbf{Var} \cup \{\ell_{\mathcal{E}}, \ell_{\mathcal{H}}\}$. We refer to the statements produced by this grammar as *formulas*.

We introduce a few notational shortcuts that will help in formally defining our predicates. For any set $V \subseteq \mathbf{Var}$, we denote by $S(V)$ the multiset resulting from the application of S on each variable in V . Also, for a set $W \subseteq \mathbf{Var} \cup \{\ell_{\mathcal{E}}\}$ with $\ell_{\mathcal{E}} \in W$, we use $S(W)$ as a shorthand for $S(W \setminus \{\ell_{\mathcal{E}}\}) \cup \mathcal{L}_{\mathcal{E}}$, and similarly for $\ell_{\mathcal{H}}$. For a set $V \subseteq \mathbf{Var} \cup \{\ell_{\mathcal{E}}, \ell_{\mathcal{H}}\}$ and an element $x \in \mathbf{Var} \cup \{\ell_{\mathcal{E}}, \ell_{\mathcal{H}}\}$, we write V, x as a shorthand for $V \cup \{x\}$ and $V - x$ as a shorthand for $V \setminus \{x\}$.

We define that a distribution X satisfies φ , denoted $X \models \varphi$ as follows:

- $X \models \varphi \wedge \varphi'$ iff $X \models \varphi$ and $X \models \varphi'$
- $X \models \varphi \vee \varphi'$ iff $X \models \varphi$ or $X \models \varphi'$
- $X \models \text{Empty}$ iff $\Pr[(S, S', \mathcal{L}_{\mathcal{E}}, \mathcal{L}_{\mathcal{H}}) \stackrel{\S}{\leftarrow} X : \mathcal{L}_{\mathcal{E}} \neq \emptyset]$ is negligible

- $X \models \text{Eq}(x, y)$ iff $\Pr[(S, S', \mathcal{L}_\mathcal{E}, \mathcal{L}_\mathcal{H}) \stackrel{\$}{\leftarrow} X : S(x) \neq S'(y)]$ is negligible
- $X \models \text{Uneq}(x, y)$ iff $\Pr[(S, S', \mathcal{L}_\mathcal{E}, \mathcal{L}_\mathcal{H}) \stackrel{\$}{\leftarrow} X : S(x) = S'(y)]$ is negligible
- $X \models \text{E}(\mathcal{E}; x; V)$ iff $\Pr[(S, S', \mathcal{L}_\mathcal{E}, \mathcal{L}_\mathcal{H}) \stackrel{\$}{\leftarrow} X : \{S(x), S'(x)\} \cap (\mathcal{L}_\mathcal{E} \cup S(V - x) \cup S'(V - x)) \neq \emptyset]$ is negligible¹
- $X \models \text{H}(\mathcal{H}; x; V)$ iff $\Pr[(S, S', \mathcal{L}_\mathcal{E}, \mathcal{L}_\mathcal{H}) \stackrel{\$}{\leftarrow} X : \{S(x), S'(x)\} \cap (\mathcal{L}_\mathcal{H} \cup S(V - x) \cup S'(V - x)) \neq \emptyset]$ is negligible
- $X \models \text{Ind}(x; V; V')$ iff the two following formulas hold:

$$\begin{aligned} & [(S, S', \mathcal{L}_\mathcal{E}, \mathcal{L}_\mathcal{H}) \stackrel{\$}{\leftarrow} X : (S(x), S(V - x) \cup S'(V'))] \sim \\ & \quad [(S, S', \mathcal{L}_\mathcal{E}, \mathcal{L}_\mathcal{H}) \stackrel{\$}{\leftarrow} X; u \stackrel{\$}{\leftarrow} \mathcal{U} : (u, S(V - x) \cup S'(V'))] \\ & [(S, S', \mathcal{L}_\mathcal{E}, \mathcal{L}_\mathcal{H}) \stackrel{\$}{\leftarrow} X : (S'(x), S'(V - x) \cup S(V'))] \sim \\ & \quad [(S, S', \mathcal{L}_\mathcal{E}, \mathcal{L}_\mathcal{H}) \stackrel{\$}{\leftarrow} X; u \stackrel{\$}{\leftarrow} \mathcal{U} : (u, S'(V - x) \cup S(V'))] \end{aligned}$$

We now present a few lemmas that show useful relations and properties of our predicates. In all these lemmas, it is assumed that \mathbb{H} is any generic hash. The proof of these lemmas is in the full version of this paper [14].

Lemma 1. *The following relations are true for any sets V_1, V_2, V_3, V_4 and variables x, y with $x \neq y$*

1. $\text{Ind}(x; V_1; V_2) \Rightarrow \text{Ind}(x; V_3; V_4)$ if $V_3 \subseteq V_1$ and $V_4 \subseteq V_2$
2. $\text{H}(\mathcal{H}; x; V_1) \Rightarrow \text{H}(\mathcal{H}; x; V_2)$ if $V_2 \subseteq V_1$
3. $\text{E}(\mathcal{E}; x; V_1) \Rightarrow \text{E}(\mathcal{E}; x; V_2)$ if $V_2 \subseteq V_1$
4. $\text{Ind}(x; V_1, \ell_\mathcal{H}; \emptyset) \Rightarrow \text{H}(\mathcal{H}; x; V_1)$
5. $\text{Ind}(x; V_1, \ell_\mathcal{E}; \emptyset) \Rightarrow \text{E}(\mathcal{E}; x; V_1)$
6. $\text{Ind}(x; \emptyset; \{y\}) \Rightarrow \text{Uneq}(x, y) \wedge \text{Uneq}(y, x)$

Note that lines 4, 5 and 6 are particularly helpful because the predicate Ind is much easier to propagate than the other predicates.

We also show that, as a consequence of our definition of $\text{DIST}(\mathbb{H})$, we can always infer the following predicates on the message blocks. This lemma is useful for proving the rules corresponding to commands that introduce a new message block.

Lemma 2. *Let $X \in \text{DIST}(\mathbb{H})$. Then for any integer i , $1 \leq i \leq n$, $X \models \text{Eq}(m_i, m_i) \vee \text{Uneq}(m_i, m_i)$.*

The following formalizes the intuition that if a value can be computed in polynomial time from other values available, then adding this value does not give the adversary any useful information. In general, we say that an expression e is *constructible* from values in a set V if e can be computed in polynomial time from V . But for the purpose of the following lemma, it is sufficient to define constructible expressions as only single variables x , as well as $x \oplus y$ and $x \parallel y$ for any variables x and y .

¹ Since the variable x is removed from the set V when taking the probability, we always have $X \models \text{E}(\mathcal{E}; x; V)$ iff $X \models \text{E}(\mathcal{E}; x; V, x)$. This is to remove the trivial case that $\{S(x), S'(x)\} \cap (\mathcal{L}_\mathcal{E} \cup \{S(x), S'(x)\}) = \emptyset$ never holds, and to simplify the notation. The same is also used for predicates $\text{H}(\mathcal{H}; x; V)$ and $\text{Ind}(x; V; V')$.

Lemma 3. *For any any $X \in \text{DIST}(\mathbb{H})$, any sets of variables V , any expression e constructible from V , and any variable x, z such that $z \notin \{x\} \cup \text{Var}(e)$ if $X \models \text{Ind}(z; V; V')$ then $\llbracket x := e \rrbracket(X) \models \text{Ind}(z; V, x; V')$. We emphasize that here we use the notation $\text{Var}(e)$ (in its usual sense), that is to say, the variable z does not appear at all in e . Similarly, if $X \models \text{Ind}(z; V'; V)$, then $\llbracket x := e \rrbracket(X) \models \text{Ind}(z; V'; V, x)$.*

The following, which is useful for proving some of the rules dealing with the concatenation commands, shows that the value of any given variable always have the same length in each execution.

Lemma 4. *For any distribution $X \in \text{DIST}(\mathbb{H})$, any program cmd produced by our grammar any $(S, S', \mathcal{E}, \mathcal{H}, \mathcal{L}_{\mathcal{E}}, \mathcal{L}_{\mathcal{H}}) \stackrel{\S}{\leftarrow} \llbracket \text{cmd} \rrbracket X$ and any variable $v \in \text{Var}$, $|S(v)| = |S'(v)|$.*

4 Proving Almost-Universal Hash

Our main contribution is a Hoare logic for proving that a program is an almost-universal hash function. We require that the program be written in a way so that, on input $m_1 \parallel \dots \parallel m_n$, the program must assign values to variables c_1, \dots, c_n in such a way that the variable c_1 contains the output of the function on input m_1 , the variable c_2 contains the output of the function on input $m_1 \parallel m_2$ and so on. We model the security of an almost-universal hash function using our predicates as follows.

Proposition 1. *Let $\mathbb{H} = (\mathcal{F}_{\mathcal{E}}, \mathcal{F}_{\mathcal{H}}, \text{Hash}(m_1 \parallel \dots \parallel m_n, c_n) : \text{var } x; \text{cmd})$ be a generic hash function on n -block messages. Then, \mathbb{H} is an almost-universal hash function if, for every positive integer n , $\text{UNIV}(n)$ holds in the distribution obtained by executing the program on any distribution in $\text{DIST}_0(\mathbb{H})$, where*

$$\text{UNIV}(n) = \left(\bigwedge_{i=1}^{n-1} \text{Uneq}(c_n, c_i) \wedge \bigwedge_{i=1}^n \text{Eq}(m_i, m_i) \right) \vee \bigwedge_{i=1}^n \text{Uneq}(c_n, c_i)$$

The proof of this proposition is in the full version of this paper [14].

Hoare Logic Rules

We present a set of rules of the form $\{\varphi\} \text{cmd} \{\varphi'\}$, meaning that execution of command cmd in any distribution in $\text{DIST}(\mathbb{H})$ that satisfies φ leads to a distribution that satisfies φ' . Using Hoare logic terminology, this means that the triple $\{\varphi\} \text{cmd} \{\varphi'\}$ is valid.

Since the predicates $\text{Eq}(m_i, m_i)$ are useful only if the whole prefix of the two messages up to the i^{th} block are equal, so that keeping track of the equality or inequality of the message blocks after the first point at which the messages are different is unnecessary. For this reason, when we design our rules, we never produce the predicates $\text{Uneq}(m_i, m_i)$ even when they would be correct.

We group rules together according to their corresponding commands. In all the rules, unless indicated otherwise, we assume that $t \notin \{x, y, z\}$ and $x \notin \{y, z\}$. In addition, for all rules involving the predicate Ind , we assume that $\ell_{\mathcal{E}}$ and $\ell_{\mathcal{H}}$ can be among the elements in the set V . Since some of the rules (for example, rule (G5)) are valid only under certain slightly complex conditions, we use square brackets in the statement of

some conditions to remove any ambiguity about their meaning. The proofs of soundness of our rules are given in the full version of this paper [14].

We first introduce a few general rules for consequence, sequential composition, conjunction and disjunction. Let $\phi_1, \phi_2, \phi_3, \phi_4$ be any four formulas in our logic, and let $\text{cmd}, \text{cmd}_1, \text{cmd}_2$ be any three commands. These rules are standard, and their proof are omitted.

- (Csq) if $\phi_1 \Rightarrow \phi_2, \phi_3 \Rightarrow \phi_4$ and $\{\phi_2\}\text{cmd}\{\phi_3\}$, then $\{\phi_1\}\text{cmd}\{\phi_4\}$
- (Seq) if $\{\phi_1\}\text{cmd}_1\{\phi_2\}$ and $\{\phi_2\}\text{cmd}_2\{\phi_3\}$, then $\{\phi_1\}\text{cmd}_1; \text{cmd}_2\{\phi_3\}$
- (Conj) if $\{\phi_1\}\text{cmd}\{\phi_2\}$ and $\{\phi_3\}\text{cmd}\{\phi_4\}$, then $\{\phi_1 \wedge \phi_3\}\text{cmd}\{\phi_2 \wedge \phi_4\}$
- (Disj) if $\{\phi_1\}\text{cmd}\{\phi_2\}$ and $\{\phi_3\}\text{cmd}\{\phi_4\}$, then $\{\phi_1 \vee \phi_3\}\text{cmd}\{\phi_2 \vee \phi_4\}$

Initialization:

We find that the following predicates holds in any distribution $X \in \text{DIST}_0(\mathbb{H})$.

- (Init) $\{\text{Ind}(k; \text{Var}, \ell_{\mathcal{E}}, \ell_{\mathcal{H}}; \text{Var} - k) \wedge \text{Eq}(k, k) \wedge \text{Empty}\}$

We recall that k is a special variable holding a secret key. It is sampled at random before executing the program and is the same in both executions, so it is indistinguishable from a random value given any other value.

Generic preservation rules:

Rules (G1) to (G6) show how predicates are preserved by most of the commands when the predicates concern a variable other than that being operated on. For all these rules, we assume that t and t' can be y or z and cmd is either $x := \rho(i, y), x := y, x := y \parallel z, x := y \oplus z, x := \mathcal{E}(y)$, or $x := \mathcal{H}(y)$.

- (G1) $\{\text{Eq}(t, t')\} \text{cmd} \{\text{Eq}(t, t')\}$ even if $t = y$ or $t = z$
- (G2) $\{\text{Uneq}(t, t')\} \text{cmd} \{\text{Uneq}(t, t')\}$ even if $t = y$ or $t = z$
- (G3) $\{\mathcal{E}(t; V)\} \text{cmd} \{\mathcal{E}(t; V)\}$ provided $x \notin V$ and cmd is not $x := \mathcal{E}(y)$
- (G4) $\{\mathcal{H}(t; V)\} \text{cmd} \{\mathcal{H}(t; V)\}$ provided $x \notin V$ and cmd is not $x := \mathcal{H}(y)$
- (G5) $\{\text{Ind}(t; V; V')\} \text{cmd} \{\text{Ind}(t; V; V')\}$ provided $[\text{cmd} \text{ is not } x := \mathcal{E}(y) \text{ or } x := \mathcal{H}(y)],$
 $[x \notin V \text{ unless } x \text{ is constructible from } V - t]$ and $[x \notin V' \text{ unless } x \text{ is constructible from } V' - t]$
- (G6) $\{\text{Empty}\} \text{cmd} \{\text{Empty}\}$ provided cmd is not $x := \mathcal{E}(y)$

We note that, for rules (G3) to (G6), the straightforward preservation rule does not apply when the command is either of the form $x := \mathcal{E}(y)$ or $x := \mathcal{H}(y)$, because some predicates may no longer hold if the block cipher or the random oracle is computed more than once on any given point. Therefore, the preservation of these predicates for the block cipher and hash commands will have to be handled separately in rules (B4) to (B6) and (H3) to (H5). For rule (G5), in general, we say that the value of a variable x is *constructible* from the values of variables in V if there exists a deterministic polynomial-time algorithm that can compute the value of x from the values in V . In this case, it means that the variables in the right-hand side of cmd are all in V .

Function ρ :

- (P1) $\{\text{Eq}(y, y)\} x := \rho(i, y) \{\text{Eq}(x, x)\}$ for integer i

Since the details of the function ρ are not known in advance, we can infer only one rule, that ρ preserves equality, because it is a deterministic function.

Assignment:

Rules (A1) to (A8), for the assignment, are all straightforward, and follow simply from the simple fact that after the command, the value of x is equal to the value of y .

- (A1) $\{\text{true}\} x := m_i \{(\text{Eq}(m_i, m_i) \wedge \text{Eq}(x, x)) \vee \text{Uneq}(x, x)\}$
- (A2) $\{\text{Eq}(y, y)\} x := y \{\text{Eq}(x, x)\}$
- (A3) $\{\text{Uneq}(y, y)\} x := y \{\text{Uneq}(x, x)\}$
- (A4) $\{\text{Ind}(y; V; V')\} x := y \{\text{Ind}(x; V; V')\}$ if $x \notin V'$ unless $y \in V'$ and $y \notin V$
- (A5) $\{\text{E}(\mathcal{E}; y; V)\} x := y \{\text{E}(\mathcal{E}; x; V) \wedge \text{E}(\mathcal{E}; y; V)\}$ if $y \notin V$
- (A6) $\{\text{H}(\mathcal{H}; y; V)\} x := y \{\text{H}(\mathcal{H}; x; V) \wedge \text{H}(\mathcal{H}; y; V)\}$ if $y \notin V$
- (A7) $\{\text{E}(\mathcal{E}; t; V, y)\} x := y \{\text{E}(\mathcal{E}; t; V, x, y)\}$
- (A8) $\{\text{H}(\mathcal{H}; t; V, y)\} x := y \{\text{H}(\mathcal{H}; t; V, x, y)\}$

Concatenation:

Rules (C1) to (C6) propagate the predicates for the concatenation command.

- (C1) $\{\text{Eq}(y, y)\} x := y \| m_i \{(\text{Eq}(m_i, m_i) \wedge \text{Eq}(x, x)) \vee \text{Uneq}(x, x)\}$
- (C2) $\{\text{Eq}(y, y) \wedge \text{Eq}(z, z)\} x := y \| z \{\text{Eq}(x, x)\}$
- (C3) $\{\text{Uneq}(y, y)\} x := y \| z \{\text{Uneq}(x, x)\}$
- (C4) $\{\text{Ind}(y; V, y, z; V') \wedge \text{Ind}(z; V, y, z; V')\} x := y \| z \{\text{Ind}(x; V, x; V')\}$ provided $[y \neq z]$, $[x, y, z \notin V]$ and $[x \notin V' \text{ unless } y, z \in V']$
- (C5) $\{\text{Ind}(y; V, \ell_{\mathcal{E}}; V)\} x := y \| z \{\text{E}(\mathcal{E}; x; V)\}$
- (C6) $\{\text{Ind}(y; V, \ell_{\mathcal{H}}; V)\} x := y \| z \{\text{H}(\mathcal{H}; x; V)\}$

The most important rule for the concatenation is (C4), which states that the concatenation of two random strings results in a random string. Note that it is important for this rule that $y \neq z$, otherwise the string x consists of a string twice repeated, which can be distinguished easily from a random value. The condition $x \notin V'$ unless $y, z \in V'$ is similar to rule (G5), and follows from the constructibility of x from y and z . Rules (C5) and (C6) state that if a string is indistinguishable from a random value given all the values in the set of queries to the block cipher (or the hash function), then clearly it cannot be a prefix of one of the strings $\mathcal{L}_{\mathcal{E}}$. For rules (C1), (C3), (C5) and (C6), the roles of y and z , or y and m_i in the case of (C1), can be exchanged.

Xor operator:

Rules (X1) to (X4) describe the effect of the Xor operation.

- (X1) $\{\text{Eq}(y, y)\} x := y \oplus m_i \{(\text{Eq}(m_i, m_i) \wedge \text{Eq}(x, x)) \vee \text{Uneq}(x, x)\}$
- (X2) $\{\text{Ind}(y; V, y, z; V')\} x := y \oplus z \{\text{Ind}(x; V, x, z; V')\}$ provided $[y \neq z]$, $[y \notin V]$ and $[x \notin V' \text{ unless } y, z \in V']$
- (X3) $\{\text{Eq}(y, y) \wedge \text{Eq}(z, z)\} x := y \oplus z \{\text{Eq}(x, x)\}$
- (X4) $\{\text{Eq}(y, y) \wedge \text{Uneq}(z, z)\} x := y \oplus z \{\text{Uneq}(x, x)\}$

Rules (X2) is reminiscent of a one-time-pad encryption: if a value z is xor-ed with a random-looking value y , then the result is similarly random-looking provided the value of y is not given. Again, the condition $x \notin V'$ unless $y, z \in V'$ is similar to rule (G5), and follows from the constructibility of x from y and z . The other rules are propagation of the Eq and Uneq predicates. Due to the commutativity of the xor, the role of y and z , or y and m_i in the case of (X1), can be exchanged in all the rules above.

Block cipher:

Since block ciphers are modeled as random functions, that is, functions picked at random among all functions from $\{0, 1\}^n$ to $\{0, 1\}^n$, the output of the function for a point

on which the block cipher has never been computed is indistinguishable from a random value.

- (B1) $\{\text{Empty}\} x := \mathcal{E}(m_i) \{(\text{Uneq}(x, x) \wedge \text{Ind}(x; \text{Var}, \ell_{\mathcal{E}}, \ell_{\mathcal{H}}; \text{Var})) \vee (\text{Eq}(m_i, m_i) \wedge \text{Eq}(x, x) \wedge \text{Ind}(x; \text{Var}, \ell_{\mathcal{E}}, \ell_{\mathcal{H}}; \text{Var} - x))\}$
- (B2) $\{\text{E}(\mathcal{E}; y; \emptyset) \wedge \text{Uneq}(y, y)\} x := \mathcal{E}(y) \{\text{Ind}(x; \text{Var}, \ell_{\mathcal{E}}, \ell_{\mathcal{H}}; \text{Var})\}$
- (B3) $\{\text{E}(\mathcal{E}; y; \emptyset) \wedge \text{Eq}(y, y)\} x := \mathcal{E}(y) \{\text{Ind}(x; \text{Var}, \ell_{\mathcal{E}}, \ell_{\mathcal{H}}; \text{Var} - x) \wedge \text{Eq}(x, x)\}$
- (B4) $\{\text{E}(\mathcal{E}; y; \emptyset) \wedge \text{Ind}(t; V; V')\} x := \mathcal{E}(y) \{\text{Ind}(t; V, x; V', x)\}$ even if $t = y$, provided $\ell_{\mathcal{E}} \notin V$
- (B5) $\{\text{E}(\mathcal{E}; y; \emptyset) \wedge \text{Ind}(t; V, \ell_{\mathcal{E}}, y; V', y)\} x := \mathcal{E}(y) \{\text{Ind}(t; V, \ell_{\mathcal{E}}, x, y; V', x, y)\}$
- (B6) $\{\text{E}(\mathcal{E}; y; \emptyset) \wedge \text{E}(\mathcal{E}; t; V, y)\} x := \mathcal{E}(y) \{\text{E}(\mathcal{E}; t; V, y)\}$

This is expressed in rules (B1) to (B3), and also used in the proof of many other rules. Note that, when executing $x := \mathcal{E}(y)$ on a new value, if the values of y from the two executions are equal, then of course the values of x will be equal afterwards. However, if the values of y are not the same in the two executions, then the values of x will be indistinguishable from two *independent* random values afterwards.

Since the querying of a block cipher twice at any point is undesirable, we always require the predicate E as a precondition. We also have rules similar to (B2) to (B6), with the predicate $\text{E}(\mathcal{E}; y; \emptyset)$ replaced by the predicate Empty , since both imply that the value of y is not in $\mathcal{L}_{\mathcal{E}}$.

Hash Function:

We note that the distinguishing adversary, described in Section 2, does not have access to the random oracle. This is sufficient for our purpose since our goal is only to prove inequality of strings, not their indistinguishability from random strings. As a result, the rules for the hash function are essentially the same as those for the block cipher.

- (H1) $\{\text{H}(\mathcal{H}; y; \emptyset) \wedge \text{Uneq}(y, y)\} x := \mathcal{H}(y) \{\text{Ind}(x; \text{Var}, \ell_{\mathcal{E}}, \ell_{\mathcal{H}}; \text{Var})\}$
- (H2) $\{\text{H}(\mathcal{H}; y; \emptyset) \wedge \text{Eq}(y, y)\} x := \mathcal{H}(y) \{\text{Ind}(x; \text{Var}, \ell_{\mathcal{H}}; \text{Var} - x) \wedge \text{Eq}(x, x)\}$
- (H3) $\{\text{H}(\mathcal{H}; y; \emptyset) \wedge \text{Ind}(t; V; V')\} x := \mathcal{H}(y) \{\text{Ind}(t; V, x; V', x)\}$ even if $t = y$, provided $\ell_{\mathcal{H}} \notin V$
- (H4) $\{\text{H}(\mathcal{H}; y; \emptyset) \wedge \text{Ind}(t; V, \ell_{\mathcal{H}}, y; V', y)\} x := \mathcal{H}(y) \{\text{Ind}(t; V, \ell_{\mathcal{H}}, x, y; V', x, y)\}$
- (H5) $\{\text{H}(\mathcal{H}; t; V, y)\} x := \mathcal{H}(y) \{\text{H}(\mathcal{H}; t; V, y)\}$

For loop:

- (F1) $\{\psi(p-1)\}$ for $l = p$ to q do: $\text{cmd}_l \{\psi(q)\}$ provided $\{\psi(l-1)\} \text{cmd}_l \{\psi(l)\}$ for $p \leq l \leq q$

The rule for the **For** loop simply states that if an indexed formula $\psi(i)$ is preserved through one iteration of the loop, then it is preserved through the entire loop. We discuss methods for finding such a formula in Section 5.

Combining our logic with Proposition 1, we obtain the following theorem.

Theorem 1. *Let $(\mathcal{F}_{\mathcal{E}}, \mathcal{F}_{\mathcal{H}}, \text{Hash}(m_1 \| \dots \| m_n, c_n) : \text{var } x; \text{cmd})$ describe the program to compute a hash function Hash on an n block message. Then, Hash is an almost-universal hash function if, for every positive integer n , $\{\text{init}\} \text{cmd} \{UNIV(n)\}$.*

The theorem is the consequence of Proposition 1 and of the soundness of our Hoare logic. We then say that a sequence of formulas $[\phi_0, \dots, \phi_n]$ is a proof that a program $[\text{cmd}_1, \dots, \text{cmd}_n]$ computes an almost-universal hash function if $\phi_0 = \text{true}$, $\phi_n \Rightarrow UNIV(n)$ and for all i , $1 \leq n$, $\{\phi_{i-1}\} \text{cmd}_i \{\phi_i\}$ holds.

5 Implementation

We chose to go forward through the program, instead of the more common approach of going backward from the end, after implementing both methods. Going backward through the program can require exploring multiple combinations of choices that all need to be explored when many rules can lead to the necessary predicate. The presence of the logical-or connector in our logic often resulted in an exponential number of possibilities at each step. As a result, our prototype for the forward method was able to find proofs much faster than an implementation of the backwards method.

We start at the beginning of the program and, at each command, apply every possible rule. Once done, we test if the predicate $UNIV(n)$ holds at the end of the program. One downside of this forward approach is that the application of every possible rule can be very time consuming because the formulas tend to grow after each command, which leads to more and more rules being applied at every step. For this reason, we need a way to filter out unneeded predicates, so that execution time remains reasonable.

5.1 Predicate Filter

We say that ϕ is a *predicate on x* if ϕ is either $\text{Eq}(x, y)$, $\text{Uneq}(x, y)$, $\text{E}(\mathcal{E}; x; V)$, $\text{H}(\mathcal{H}; x; V)$ or $\text{Ind}(x; V_1, V_2)$ (for some $y \in \text{Var}$ and $V_1, V_2 \subseteq \text{Var}$). We say that a predicate ϕ on variable x is *obsolete for program p* if x does not appear anywhere in p and if $\neg(\phi \Rightarrow \text{Uneq}(c_n, c_i))$ and $\neg(\phi \Rightarrow \text{Eq}(m_i, m_i))$ for any i , $1 \leq i \leq n$.² The following theorem shows that once a predicate is obsolete, it can be discarded.

Theorem 2. *If there exists a proof $[\phi_0, \dots, \phi_n]$ that a program $p = [\text{cmd}_1, \dots, \text{cmd}_n]$ computes an almost-universal hash function, then there also exists a proof $[\phi'_0, \dots, \phi'_n]$ that p computes an almost-universal hash function where for each i , $\phi_i \Rightarrow \phi'_i$ and each ϕ'_i does not contain any obsolete predicates for $[\text{cmd}_{i+1}, \dots, \text{cmd}_n]$.*

The theorem is a consequence of the fact that, in our logic, the rules for creating a predicate on x following the execution of command $x := e$ only have as preconditions predicates on the variables in e . As a result, we can always filter out obsolete predicates after processing each command.

Also, we note that the only commands that can make a predicate $\text{Eq}(m_i, m_i)$ appear are those of the form $x := e$ in which m_i appears in e . As a result, if we find that, for some integer l , the predicate $\text{Eq}(m_l, m_l)$ is not present in one of the conjunctions of the current formula (after transforming the formula in disjunctive normal form) and that the variable m_l is no longer present in the rest of the program, then there is no longer any chance that it will satisfy the conjunction with $\bigwedge_{j=1}^n \text{Eq}(m_j, m_j)$ from $UNIV(n)$. Therefore, we can also safely filter out all other predicates of the form $\text{Eq}(m_i, m_i)$ from that conjunction.

We also add a *heuristic filter* to speed up the execution of our method. We make the hypothesis that the predicate $\text{Ind}(c_n; V; \{c_1, \dots, c_{n-1}\})$ will be present at the end of the program, which is the case for all our examples, so that we can filter out $\text{Ind}(c_i; V; V')$

² Here, p will usually be the rest of the program after the program point at which the predicate ϕ holds.

if $i < n$ and c_i is no longer present in the remainder of the program. In addition to speeding up the program, filtering out these predicates greatly simplifies the construction of loop invariants discussed in the next section. If we fail to produce a proof while using the heuristic filter, we simply attempt again to find a proof without it.

5.2 Finding Loop Invariants

The programs describing the almost-universal hash function usually contains for loops. It is therefore necessary to have an automatic procedure to detect the formula $\psi(i)$ that allows us to apply rule (F1). We now show a heuristic that can be used to construct such an invariant, and illustrate how it works by applying them to $Hash_{CBC}$, described in Section 3.1. One could easily verify that it also works on $Hash_{CBC'}$, $Hash_{HMAC}$ and $Hash_{PMAC}$.

Once we hit a command "for $l = p$ to q do: $[cmd_l]$ ", we express the formula that holds before the loop is executed in the form $\varphi(p-1)$. The classical method for finding a stable invariant consists in processing the instructions cmd_l contained in the loop to find the formula $\psi(l)$ such that $\{\varphi(l-1)\} cmd_l \{\psi(l)\}$. If $\psi(l) \Rightarrow \varphi(l)$, then we have found a formula such that $\{\varphi(l-1)\} cmd_l \{\varphi(l)\}$ and we can apply rule (F1).

Unfortunately, for most loops, this simple process either does not yield a stable invariant, or gives a stable invariant too weak to produce a proof. We need a heuristic to construct stronger stable invariants. The heuristic we describe here is inspired from widening methods in abstract interpretation. We start with formula $\varphi(l-1)$, and process the program of the loop once to find formula $\psi_1(l)$ such that $\{\varphi(l-1)\} cmd_l \{\psi_1(l)\}$. Then, we repeat this starting with formula $\psi_1(l-1)$ to find formula $\psi_2(l)$ such that $\{\psi_1(l-1)\} cmd_l \{\psi_2(l)\}$. The idea is then to inspect formulas $\varphi(l)$, $\psi_1(l)$ and $\psi_2(l)$ for patterns that can be extrapolated. For example, we can try to identify a predicate $\gamma(l)$ such that: (i) $\gamma(l)$ appears in $\varphi(l)$, (ii) $\gamma(l-1) \wedge \gamma(l)$ appears in $\psi_1(l)$, (iii) $\gamma(l-2) \wedge \gamma(l-1) \wedge \gamma(l)$ appears in $\psi_2(l)$. We then use a new starting formula $\varphi'(l)$ which is just like $\varphi(l)$, except that the occurrence of $\gamma(l)$ in $\varphi(l)$ is replaced by $\bigwedge_{j=p-1}^{j=l} \gamma(j)$ in $\varphi'(l)$. Note that, by construction, $\varphi(p-1)$ is equal to $\varphi'(p-1)$, so we know that $\varphi'(p-1)$ is satisfied at the beginning of the loop.³

Example: We now apply this method to $Hash_{CBC}$. After processing command $c_1 := \mathcal{E}(m_1)$, we obtain the formula $\varphi(1) = (\text{Ind}(c_1; \text{Var}, \ell_{\mathcal{E}}; \text{Var} - c_1) \wedge \text{Eq}(m_1, m_1) \wedge \text{Eq}(c_1, c_1)) \vee \text{Ind}(c_1)$. Parameterizing this in terms of l , we obtain

$$\varphi(l) = (\text{Eq}(m_l, m_l) \wedge \text{Eq}(c_l, c_l) \wedge \text{Ind}(c_l; \text{Var}, \ell_{\mathcal{E}}; \text{Var} - c_l)) \vee \text{Ind}(c_l)$$

We recall that the two instructions in the loop of $Hash_{CBC}$ are the following: $z_i := c_{i-1} \oplus m_i$; $c_i := \mathcal{E}(z_i)$. After processing the program of the loop on $\varphi(l-1)$, we obtain the following.

$$\psi_1(l) = (\text{Eq}(m_{l-1}, m_{l-1}) \wedge \text{Eq}(m_l, m_l) \wedge \text{Eq}(c_l, c_l) \wedge \text{Ind}(c_l; \text{Var}, \ell_{\mathcal{E}}; \text{Var} - c_l)) \vee \text{Ind}(c_l)$$

³ We can similarly try to find patterns that appear only after the first iteration of the loop, that is, $\gamma(l)$ appears in $\psi_1(l)$ and $\gamma(l-1) \wedge \gamma(l)$ appears in $\psi_2(l)$, in which case $\bigwedge_{j=p}^{j=l} \gamma(j)$ is added in $\varphi'(l)$.

We get this by applying rules (G1), (X1) and (X2) for the first command and rules (G1), (B2) and (B3) for the second command. Note that $\psi_1(l) \Rightarrow \varphi(l)$, so we could use $\varphi(l)$ to apply rule (F1), but this would not yield a proof of $Hash_{CBC}$. We repeat the same process with $\psi_1(l-1)$ to obtain

$$\begin{aligned} \psi_2(l) = & (\text{Eq}(m_{l-2}, m_{l-2}) \wedge \text{Eq}(m_{l-1}, m_{l-1}) \wedge \text{Eq}(m_l, m_l) \wedge \\ & \text{Eq}(c_l, c_l) \wedge \text{Ind}(c_l; \text{Var}, \ell_{\mathcal{E}}; \text{Var} - c_l)) \vee \text{Ind}(c_l). \end{aligned}$$

This requires applying the same rules as before, but rule (G1) more often applied for each command. We find $\gamma(l) = \text{Eq}(m_l, m_l)$ and use

$$\varphi'(l) = \left(\left(\bigwedge_{i=1}^l \text{Eq}(m_i, m_i) \right) \wedge \text{Eq}(c_l, c_l) \wedge \text{Ind}(c_l; \text{Var}, \ell_{\mathcal{E}}; \text{Var} - c_l) \right) \vee \text{Ind}(c_l)$$

as our next attempt at finding a stable invariant. We find that $\varphi'(l)$ is a stable invariant for the loop. So we apply the rule (F1) to obtain that $\varphi'(n)$ holds at the end of the program, and we easily find that $\varphi'(n) \Rightarrow UNIV(n)$ for all positive integer n , thereby proving that $Hash_{CBC}$ computes an almost-universal hash function.

5.3 Prototype

We programmed an OCaml prototype of our method for proving that the front end of MACs are almost-universal hash functions. The program requires about 2000 lines of code, and can successfully produce proofs of security for all the examples discussed in this paper in less than one second on a personal workstation. Our prototype is available on [14].

6 Proving MAC Security

As mentioned in Section 2, we prove the security of MACs in two steps: first we show that the ‘compressing’ part of the MAC is an almost-universal hash function family, and then we show that the last section of the MAC, when applied to an almost-universal hash function, results in a secure MAC. The following shows how a secure MAC can be constructed from an almost-universal hash function. The proof can be found in [4,8,9], so we do not repeat them here.

Theorem 3. *Let $\mathcal{H} = \{h_i\}_{i \in \{0,1\}^\eta}$ and $\mathcal{H}' = \{h_i\}_{i \in \{0,1\}^\eta}$ be families of almost-universal hash function, $\mathcal{F}_{\mathcal{E}}$ be a family of block ciphers and \mathcal{G} be a random oracle. If $h \xleftarrow{\$} \mathcal{H}$, $h_{\mathcal{E}} \xleftarrow{\$} \mathcal{H}'$, $\mathcal{E}, \mathcal{E}_1, \mathcal{E}_2 \xleftarrow{\$} \mathcal{F}_{\mathcal{E}}$, \mathcal{G} is sampled at random from all functions with the appropriate domain and range and $k, k_1, k_2 \xleftarrow{\$} \{0,1\}^\eta$, then the following hold:*

- $MAC_1(m) = \mathcal{E}(h_i(m))$ is a secure MAC with key $sk = (i, k_{\mathcal{E}})$.⁴
- $MAC_2(m) = \mathcal{G}(k || h_i(m))$ is a secure MAC with key $sk = (i, k)$.
- $MAC_3(m) = \begin{cases} \mathcal{E}_1(h_i(m')) & \text{where } m' = \text{pad}(m) \text{ if } m \text{'s length is not a multiple of } \eta \\ \mathcal{E}_2(h_i(m)) & \text{if } m \text{'s length is a multiple of } \eta \end{cases}$ is a secure MAC with key $sk = (i, k_{\mathcal{E}_1}, k_{\mathcal{E}_2})$.

⁴ Here, $k_{\mathcal{E}}$ denotes the secret key associated with block cipher \mathcal{E} .

$$\begin{aligned}
 - \text{MAC}_4(m) = & \\
 & \begin{cases} \mathcal{E}(h_{\mathcal{E}}(m') \oplus k_1) \text{ where } m' = \text{pad}(m) \text{ if } m\text{'s length is not a multiple of } \eta \\ \mathcal{E}(h_{\mathcal{E}}(m) \oplus k_2) \text{ if } m\text{'s length is a multiple of } \eta \end{cases} \\
 & \text{is a secure MAC with key } sk = (k_{\mathcal{E}}, k_1, k_2)
 \end{aligned}$$

Combining Hash_{CBC} with MAC_1 and MAC_3 yields the message authentication codes DMAC and ECBC respectively, using $\text{Hash}_{CBC'}$ with MAC_3 and MAC_4 yields FCBC and XCBC, combining Hash_{PMAC} and MAC_4 yields a four key construction of PMAC and using Hash_{HMAC} with MAC_2 yields HMAC.

7 Conclusion

We presented a Hoare logic that can be used to automatically prove the security of constructions for almost-universal hash functions based on block ciphers and compression functions modeled as random oracles. We can then obtain a secure MAC by combining with a few operations, such as those presented in Section 6. Our method can be used to prove the security of DMAC, ECBC, FCBC, XCBC, a two-key variant of HMAC and a four-key variant of PMAC. Since the final step of the proof for the MACs is not integrated in the logic, we cannot prove the one key variants of HMAC or PMAC, nor can we prove CMAC or OMAC, which are one-key variants of XCBC. It is however relatively simple to derive the security of these one-key schemes by hand once the security of the multiple key variants has been proven. It remains an open problem to integrate this step into the logic.

It should be possible to extend our logic to prove exact reduction bounds for the security of the ϵ -universal hash function. This could be done by keeping track of exact security for each predicate to obtain a bound on the final invariant. We are also working on integrating our tool for verifying the security of MACs with the tool for verifying the security of encryption modes of operation of [15], to get a general tool for producing security proofs of symmetric modes of operation.

References

1. Barthe, G., Daubignard, M., Kapron, B., Lakhnech, Y.: Computational indistinguishability logic. In: Proceedings of the 17th ACM Conference on Computer and Communications Security, CCS 2010, pp. 375–386. ACM (2010)
2. Barthe, G., Grégoire, B., Heraud, S., Béguelin, S.Z.: Computer-aided security proofs for the working cryptographer. In: Rogaway, P. (ed.) CRYPTO 2011. LNCS, vol. 6841, pp. 71–90. Springer, Heidelberg (2011)
3. Barthe, G., Grégoire, B., Lakhnech, Y., Zanella Béguelin, S.: Beyond provable security verifiable IND-CCA security of OAEP. In: Kiayias, A. (ed.) CT-RSA 2011. LNCS, vol. 6558, pp. 180–196. Springer, Heidelberg (2011)
4. Bellare, M., Canetti, R., Krawczyk, H.: Keying hash functions for message authentication. In: Koblitz, N. (ed.) CRYPTO 1996. LNCS, vol. 1109, pp. 1–15. Springer, Heidelberg (1996)
5. Bellare, M., Kilian, J., Rogaway, P.: The security of cipher block chaining. In: Desmedt, Y.G. (ed.) CRYPTO 1994. LNCS, vol. 839, pp. 341–358. Springer, Heidelberg (1994)

6. Bellare, M., Rogaway, P.: Random oracles are practical: a paradigm for designing efficient protocols. In: CCS 1993: Proceedings of the 1st ACM Conference on Computer and Communications Security, pp. 62–73. ACM, New York (1993)
7. Black, J., Halevi, S., Krawczyk, H., Krovetz, T., Rogaway, P.: UMAC: Fast and secure message authentication. In: Wiener, M. (ed.) CRYPTO 1999. LNCS, vol. 1666, pp. 216–233. Springer, Heidelberg (1999)
8. Black, J., Rogaway, P.: CBC MACs for arbitrary-length messages: The three-key constructions. In: Bellare, M. (ed.) CRYPTO 2000. LNCS, vol. 1880, pp. 197–215. Springer, Heidelberg (2000)
9. Black, J., Rogaway, P.: A block-cipher mode of operation for parallelizable message authentication. In: Knudsen, L.R. (ed.) EUROCRYPT 2002. LNCS, vol. 2332, pp. 384–397. Springer, Heidelberg (2002)
10. Blanchet, B., Pointcheval, D.: Automated security proofs with sequences of games. In: Dwork, C. (ed.) CRYPTO 2006. LNCS, vol. 4117, pp. 537–554. Springer, Heidelberg (2006)
11. Corin, R., den Hartog, J.: A probabilistic hoare-style logic for game-based cryptographic proofs. In: Bugliesi, M., Preneel, B., Sassone, V., Wegener, I. (eds.) ICALP 2006, Part II. LNCS, vol. 4052, pp. 252–263. Springer, Heidelberg (2006)
12. Courant, J., Daubignard, M., Ene, C., Lafourcade, P., Lakhnech, Y.: Towards automated proofs for asymmetric encryption schemes in the random oracle model. In: Proceedings of the 15th ACM Conference on Computer and Communications Security, CCS 2008, Alexandria, USA (October 2008)
13. Fournet, C., Kohlweiss, M., Strub, P.: Modular code-based cryptographic verification. In: Chen, Y., Danezis, G., Shmatikov, V. (eds.) ACM-CCS 2011, pp. 341–350. ACM (2011)
14. Gagné, M., Lafourcade, P., Lakhnech, Y.: Full paper and OCaml implementation of our method. Computer Science Department, Saarland University, Germany (June 2013), <http://www.infsec.cs.uni-saarland.de/gagne/macChecker/macChecker.html>
15. Gagné, M., Lafourcade, P., Lakhnech, Y., Safavi-Naini, R.: Automated security proof for symmetric encryption modes. In: Datta, A. (ed.) ASIAN 2009. LNCS, vol. 5913, pp. 39–53. Springer, Heidelberg (2009)
16. Liskov, M., Rivest, R.L., Wagner, D.: Tweakable block ciphers. In: Yung, M. (ed.) CRYPTO 2002. LNCS, vol. 2442, pp. 31–46. Springer, Heidelberg (2002)
17. Petrank, E., Rackoff, C.: Cbc mac for real-time data sources. *Journal of Cryptology* 13, 315–338 (1997)
18. Wegman, M., Carter, J.L.: Universal classes of hash functions. *Journal of Computer and System Sciences* 18(2), 143–154 (1979)
19. Wegman, M., Carter, J.L.: New hash functions and their use in authentication and set equality. *Journal of Computer and System Sciences* 22(3), 265–279 (1981)