

Formalization of the Standard Uniform Random Variable

Osman Hasan and Sofiène Tahar

*Dept. of Electrical & Computer Engineering,
Concordia University
Montreal, Quebec, H3G 1M8, Canada*

Abstract

Continuous random variables are widely used to mathematically describe random phenomenon in engineering and physical sciences. In this paper, we present a higher-order logic formalization of the Standard Uniform random variable as the limit value of the sequence of its discrete approximations. We then show the correctness of this specification by proving the corresponding probability distribution properties within the HOL theorem prover and the proof steps have been summarized. The formalized Standard Uniform random variable can be transformed to formalize other continuous random variables, such as Uniform, Exponential, Normal, etc., by using various Non-uniform random number generation techniques. The formalization of these continuous random variables will enable us to perform error free probabilistic analysis of systems within the framework of a higher-order-logic theorem prover. For illustration purposes, we present the formalization of the Continuous Uniform random variable based on the formalized Standard Uniform random variable and then utilize it to perform a simple probabilistic analysis of roundoff error in HOL.

Key words: Continuous Probability Distributions, Probabilistic Analysis, Formal Verification, Theorem Proving, HOL Theorem Prover.

1 Introduction

The engineering and physical science approach to mastering the unavoidable elements of randomness and uncertainty is to use random variables to mathe-

Email addresses: o_hasan@ece.concordia.ca (Osman Hasan),
tahar@ece.concordia.ca (Sofiène Tahar).

matically describe random phenomenon. Random variables are basically functions that map random events to numbers. Every random variable gives rise to a probability distribution, which contains most of the important information about this random variable. The probability distribution of a random variable X can be uniquely described by its *cumulative distribution function* (CDF), which is defined by $F_X(x) = \mathbb{P}(X \leq x)$ for any real number x , where \mathbb{P} represents the probability. A distribution is called discrete if its CDF consists of a sequence of finite jumps, which means that it belongs to a random variable that can only attain values from a certain finite or countable set. Similarly, a distribution is called continuous if its CDF is continuous, which means that it belongs to a random variable that ranges over a continuous set of real numbers. A continuous set of real numbers, sometimes referred to as an interval, contains all real numbers between two limits. An interval can be open (a,b) corresponding to the set $\{x|a < x < b\}$, closed [a,b] which represents the set $\{x|a \leq x \leq b\}$, or half-open (a,b], [a,b). Continuous probability distributions are widely used to mathematically describe random phenomenon in engineering and scientific applications. For example, Continuous Uniform distribution is used to model quantization errors in computer arithmetic applications [25], Exponential distribution occurs in applications such as queuing theory to model *interarrival* and *service* times and Normal distribution is extensively used to model signals in data transmission and digital signal processing systems [24].

Today, simulation is the most commonly used computer based probabilistic analysis technique. Most simulation softwares provide a programming environment for defining functions that approximate random variables for probability distributions. The random elements in a given system are modeled by these functions and the system is analyzed using computer simulation techniques [5], such as the Monte Carlo Method [16], where the main idea is to approximately answer a query on a probability distribution by analyzing a large number of samples. Due to these approximations the results can be quite unreliable at times. It is a common occurrence that different software packages come up with different solutions to the same probabilistic problem. This unreliability mainly arises because of three major reasons. First, the primary source of randomness in these simulation based packages is pseudorandom numbers [14]. Secondly, functions for evaluating probability distributions are approximated using a variety of efficient but less accurate numerical methods [13]. Finally, like all computer based computations, roundoff and truncation errors also creep into the numerical computations conducted in these simulation based software packages. In [17], McCullough proposed a collection of intermediate-level tests for assessing the numerical reliability of a statistical package and uncovered flaws in most of the mainframe statistical packages, e.g., [18]. These unreliable results pose a serious problem in safety critical applications, such as space travel, military and medicine, where a mismatch between the predicted and the actual system performance may result in either inefficient usage of

the available resources or paying higher costs to meet some performance or reliability criteria unnecessarily. Another major limitation of simulation based probabilistic analysis is the enormous amount of CPU time requirement for attaining meaningful estimates. This approach generally requires hundreds of thousands of simulations to calculate the probabilistic quantities and becomes impractical when each simulation step involves extensive computations.

As an alternative to simulation techniques, we propose to use higher-order logic interactive theorem proving [7] for probabilistic analysis. Higher-order logic is a system of deduction with a precise semantics and can be used for the development of almost all classical mathematics theories. Interactive theorem proving is the field of computer science and mathematical logic concerned with computer based formal proof tools that require some sort of human assistance. We believe that probabilistic analysis can be performed by specifying the behavior of systems which exhibit randomness in higher-order logic and formally proving the intended probabilistic properties within the environment of an interactive theorem prover. The probabilistic analysis carried out in this way will be free from approximation and precision issues, and the accuracy of the results will be independent of the CPU time.

The foremost criteria for implementing a formalized probabilistic analysis framework is to be able to formalize random variables in higher-order logic. Hurd's PhD thesis [12] can be considered a pioneering work in this regard as it presents a methodology for the formalization of probabilistic algorithms in the higher-order-logic (HOL) theorem prover [8]. Random variables are basically probabilistic algorithms and Hurd formalized some discrete random variables in [12]. On the other hand, to the best of our knowledge no higher-order logic formalization of continuous random variables exists in the literature so far. In this paper, we show how to extend Hurd's formalization framework for the formalization of continuous random variables using Non-uniform random number generation methods [5]. The main advantage of this approach is that only one continuous random variable needs to be formalized from scratch, i.e., the Standard Uniform random variable. Other continuous random variables can then be formalized by using the definition of Standard Uniform random variable and formalizing the corresponding Non-uniform random number generation method.

The rest of the paper is organized as follows: In Section 2, we provide some preliminaries including a brief introduction to the HOL theorem prover and an overview of Hurd's methodology for the formalization of probabilistic algorithms in HOL. In Section 3, we formally specify the Standard Uniform random variable in the HOL theorem prover as the limit value of a discrete uniform random variable in the interval $[0, 1 - (\frac{1}{2})^n]$. In Section 4, we verify the correctness of the above specification by proving its corresponding probability distribution properties. In order to illustrate the fact that the formalized

Standard Uniform random variable can be utilized to formalize other continuous random variables using Non-uniform random number methods, we formally specify and verify the Continuous Uniform random variable in Section 5. Then in Section 6, we present the process of evaluating probabilistic quantities within the HOL theorem prover by considering a simplified probabilistic analysis example of roundoff error in a digital processor. We present a review of related work in the open literature in Section 7 and finally conclude the paper in Section 8.

2 Preliminaries

In this section we give a brief introduction to the HOL theorem prover and present an overview of Hurd’s methodology for the formalization of probabilistic algorithms proposed in [12]. The intent is to introduce the main ideas along with some notation that is going to be used in the coming sections.

2.1 HOL Theorem Prover

The HOL theorem prover is an interactive theorem prover which is capable of conducting proofs in higher-order logic. It utilizes the simple type theory of Church [3] along with Hindley-Milner polymorphism [19] to implement higher-order logic. HOL has been successfully used as a verification framework for both software and hardware as well as a platform for the formalization of pure mathematics. It supports the formalization of various mathematical theories including sets, natural numbers, real numbers, measure and probability. The HOL theorem prover includes many proof assistants and automatic proof procedures. The user interacts with a proof editor and provides it with the necessary tactics to prove goals while some of the proof steps are solved automatically by the automatic proof procedures.

In order to ensure secure theorem proving, the logic in the HOL system is represented in the strongly-typed functional programming language ML [21]. The ML abstract data types are then used to represent higher-order-logic theorems and the only way to interact with the theorem prover is by executing ML procedures that operate on values of these data types. Users can prove theorems using a natural deduction style by applying inference rules to axioms or previously generated theorems. The HOL core consists of only basic 5 axioms and 8 primitive inference rules, which are implemented as ML functions. Soundness is assured as every new theorem must be created from these basic axioms and primitive inference rules or any other pre-existing theorems/inference rules.

We selected the HOL theorem prover for the proposed formalization mainly because of its inherent soundness and ability to handle higher-order logic and in order to benefit from the built-in mathematical theories for measure and probability. The table below summarizes some of the HOL symbols used in this paper and their corresponding mathematical interpretation [8].

HOL Symbol	Standard Symbol	Meaning
\wedge	<i>and</i>	Logical <i>and</i>
\vee	<i>or</i>	Logical <i>or</i>
$\sim t$	$\neg t$	Not <i>t</i>
<i>num</i>	$\{0, 1, 2, \dots\}$	Natural data type
<i>real</i>	All Real numbers	Real data type
SUC <i>n</i>	$n + 1$	Successor of a <i>num</i>
$\lambda x.t$	$\lambda x.t$	Function that maps <i>x</i> to $t(x)$
$\{x P(x)\}$	$\{\lambda x.P(x)\}$	Set of all <i>x</i> that satisfy the condition <i>P</i>
(a, b)	a x b	A mathematical pair of two elements

Table 1
HOL Symbols

2.2 Verifying Probabilistic Algorithms in HOL

Hurd [12] proposed to formalize the probabilistic algorithms in higher-order logic by thinking of them as deterministic functions with access to an infinite Boolean sequence \mathbb{B}^∞ ; a source of infinite random bits. These deterministic functions make random choices based on the result of popping the top most bit in the infinite Boolean sequence and may pop as many random bits as they need for their computation. When the algorithms terminate, they return the result along with the remaining portion of the infinite Boolean sequence to be used by other programs. Thus, a probabilistic algorithm which takes a parameter of type α and ranges over values of type β can be represented in HOL by the function

$$\mathcal{F} : \alpha \rightarrow B^\infty \rightarrow \beta \times B^\infty$$

For example, a *Bernoulli*($\frac{1}{2}$) random variable that returns 1 or 0 with equal probability $\frac{1}{2}$ can be modeled as follows

```
⊢ bit = λs. (if shd s then 1 else 0, stl s)
```

where s is the infinite Boolean sequence and `shd` and `stl` are the sequence equivalents of the list operation *'head'* and *'tail'*. The function *bit* accepts the infinite Boolean sequence and returns a random number, which is either 0 or 1 together with a sequence of unused Boolean sequence, which in this case is the tail of the sequence. The above methodology can be used to model most probabilistic algorithms. All probabilistic algorithms that compute a finite number of values equal to 2^n , each having a probability of the form $\frac{m}{2^n}$: where m represents the hol data type *nat* and is always less than 2^n , can be modeled, using Hurd's framework, by well-founded recursive functions. The probabilistic algorithms that do not satisfy the above conditions but are sure to terminate can be modeled by the *probabilistic while loop* proposed in [12].

The probabilistic programs can also be expressed in the more general state-transforming monad where the states are the infinite Boolean sequences.

$$\begin{aligned} &\vdash \forall a, s. \text{unit } a \text{ } s = (a, s) \\ &\vdash \forall f, g, s. \text{bind } f \text{ } g \text{ } s = \text{let } (x, s') \leftarrow f(s) \text{ in } g \text{ } x \text{ } s' \end{aligned}$$

The `unit` operator is used to lift values to the monad, and the `bind` is the monadic analogue of function application. All the monad laws hold for this definition, and the notation allows us to write functions without explicitly mentioning the sequence that is passed around, e.g., function *bit* can be defined as

$$\vdash \text{bit_monad} = \text{bind } \text{sdest } (\lambda b. \text{if } b \text{ then unit } 1 \text{ else unit } 0)$$

where `sdest` gives the head and tail of a sequence as a pair (*shd* s , *stl* s).

Hurd [12] also formalized some mathematical measure theory in HOL in order to define a probability function \mathbb{P} from sets of infinite Boolean sequences to real numbers between 0 and 1. The domain of \mathbb{P} is the set \mathcal{E} of events of the probability. Both \mathbb{P} and \mathcal{E} are defined using the Carathéodory's Extension theorem, which ensures that \mathcal{E} is a σ -algebra: closed under complements and countable unions. The formalized \mathbb{P} and \mathcal{E} can be used to derive the basic laws of probability in the HOL prover, e.g., the additive law, which represents the probability of two disjoint events as the sum of their probabilities:

$$\begin{aligned} &\vdash \forall A B. A \in \mathcal{E} \wedge B \in \mathcal{E} \wedge A \cap B = \emptyset \Rightarrow \\ &\quad \mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B) \end{aligned}$$

The formalized \mathbb{P} and \mathcal{E} can also be used to prove probabilistic properties for probabilistic programs such as

$$\vdash \mathbb{P} \{s \mid \text{fst } (\text{bit } s) = 1\} = \frac{1}{2}$$

where the function `fst` selects the first component of a pair.

The measurability of a function is an important concept in probability theory and also a useful practical tool for proving that sets are measurable [2]. In Hurd’s formalization of probability theory, a set of infinite Boolean sequences, S , is said to be measurable if and only if it is in \mathcal{E} , i.e., $S \in \mathcal{E}$. Since the probability measure \mathbb{P} is only defined on sets in \mathcal{E} , it is very important to prove that sets that arise in verification are measurable. Hurd [12] showed that a function is guaranteed to be measurable if it accesses the infinite boolean sequence using only the `unit`, `bind` and `sdest` primitives and thus leads to only measurable sets.

Hurd formalized four discrete random variables and proved their correctness by proving the corresponding *Probability Mass Functions* (PMF) properties [12]. Because of their discrete nature, all these random variables either compute a finite number of values or are sure to terminate. Thus, they can be expressed using Hurd’s methodology by either well formed recursive functions or the probabilistic while loop [12]. On the other hand, continuous random variables always compute an infinite number of values and therefore would require all the random bits in the infinite Boolean sequence if they are to be represented using Hurd’s methodology. The corresponding deterministic functions cannot be expressed by either recursive functions or the probabilistic while loop and it is mainly for this reason that the specification of continuous random variables needs to be handled differently than their discrete counterparts.

3 Specification of Standard Uniform Random Variable in HOL

In this section, we present a formal specification of the Standard Uniform random variable in HOL. Standard Uniform random variable can be formalized in a number of different ways using the various formal semantics of probabilistic programs available in the literature of theoretical computer science. In order to minimize the effort and speed up the formalization process, our intent is to find a solution that enables us to build upon the existing work of Hurd [12].

Standard Uniform random variable is a continuous random variable for which the probability that it will belong to a subinterval of $[0,1]$ is proportional to the length of that subinterval. It can be characterized by the CDF as follows:

$$\mathbb{P}(X \leq x) = \begin{cases} 0 & \text{if } x < 0; \\ x & \text{if } 0 \leq x < 1; \\ 1 & \text{if } 1 \leq x. \end{cases} \quad (1)$$

It is a well known mathematical fact, see [6] for example, that a Standard

Uniform random variate can be modeled by an infinite sequence of random bits (informally coin flips) as follows

$$\sum_{k=0}^{\infty} \left(\frac{1}{2}\right)^{k+1} X_k \quad (2)$$

where, X_k denotes the outcome of the k^{th} random bit; *true* or *false* represented as 1 or 0 respectively. The mathematical relation of Equation (2) presents a sampling algorithm for the Standard Uniform random variable which is quite consistent with Hurd’s formalization methodology, i.e, it allows us to model the Standard Uniform random variable by a deterministic function with access to the infinite Boolean sequence. The specification of this sampling algorithm in higher-order logic is not very straight forward though. Due to the infinite sampling, it cannot be modeled by either of the approaches proposed in [12], i.e., a recursive function or the probabilistic while loop. We approach this problem by splitting the corresponding mathematical relation into two steps. The first step is to mathematically represent a discrete version of the Standard Uniform random variable

$$(\lambda n. \sum_{k=0}^{n-1} \left(\frac{1}{2}\right)^{k+1} X_k) \quad (3)$$

This lambda abstraction function accepts a natural number n and generates an n -bit Standard Uniform random variable using the computation principle of Equation (2). The continuous Standard Uniform random variable can now be represented as a special case of Equation (3) when n tends to infinity

$$\lim_{n \rightarrow \infty} (\lambda n. \sum_{k=0}^{n-1} \left(\frac{1}{2}\right)^{k+1} X_k) \quad (4)$$

The advantage of expressing the sampling algorithm of Equation (2) in these two steps is that now it can be specified in HOL. The mathematical relationship of Equation (3) can be specified in HOL by a recursive function using Hurd’s methodology as it consumes a finite number of random bits, i.e., n . Then, the formalization of the mathematical concept of limit of a real sequence [10] in HOL can be used to specify the mathematical relation of Equation (4).

Next, we present the HOL formalization of the above steps. We first formalize a discrete Standard Uniform random variable that produces any one of the equally spaced 2^n dyadic rationals in the interval $[0, 1 - (\frac{1}{2})^n]$ with the same probability $(\frac{1}{2})^n$. This random variable is specified in HOL as a recursive function using Hurd’s formalization framework of Section 2.2.

Definition 1:

$$\begin{aligned} &\vdash (\text{std_unif_disc } 0 = \text{unit } 0) \wedge \\ &\quad \forall n. (\text{std_unif_disc } (\text{SUC } n) = \\ &\quad \text{bind } (\text{std_unif_disc } n) \\ &\quad (\lambda m. \text{bind } \text{sdest} \\ &\quad (\lambda b. \text{unit } (\text{if } b \text{ then } ((\frac{1}{2})^{n+1} + m) \text{ else } m)))) \end{aligned}$$

The function, *std_unif_disc*, models an n -bit discrete Standard Uniform random variable based on the principle of Equation (2) by simply converting the first n random bits $B_0, B_1, B_2, \dots, B_{n-1}$ of the infinite Boolean sequence to their equivalent real number with the binary representation $0.B_0B_1B_2 \dots B_{n-1}$. It returns a pair such that the first component is the n -bit discrete Standard Uniform random variable and the second component is the unused portion of the infinite Boolean sequence. The following properties for the discrete Standard Uniform random variable may be proved by induction on its argument n

Lemma 1:

$$\vdash \forall n, s. 0 \leq \text{fst}(\text{std_unif_disc } n \text{ } s) \leq 1 - (\frac{1}{2})^n$$

Lemma 2:

$$\begin{aligned} &\vdash \forall m, n, x. \mathbb{P}\{s \mid \text{fst}(\text{std_unif_disc } n \text{ } s) \leq x\} = \\ &\quad \text{if } (x < 0) \text{ then } 0 \text{ else} \\ &\quad \text{if } (x \geq 1) \text{ then } 1 \text{ else} \\ &\quad \text{if } (x = \frac{m}{2^n}) \text{ then } \frac{\text{SUC } m}{2^n} \\ &\quad \text{else } 0 \end{aligned}$$

Lemma 3:

$$\begin{aligned} &\vdash \forall m, n, x. \mathbb{P}\{s \mid \text{fst}(\text{std_unif_disc } n \text{ } s) = x\} = \\ &\quad \text{if } (x < 0) \text{ then } 0 \text{ else} \\ &\quad \text{if } (x \geq 1) \text{ then } 0 \text{ else} \\ &\quad \text{if } (x = \frac{m}{2^n}) \text{ then } \frac{1}{2^n} \\ &\quad \text{else } 0 \end{aligned}$$

The term $\frac{m}{2^n}$ in Lemmas 2 and 3 represents all the dyadic rationals in the interval $[0, 1 - (\frac{1}{2})^n]$ since the variable m belongs to the HOL datatype *num* for natural numbers $\{0, 1, 2, \dots\}$. Collectively Lemmas 1, 2 and 3, illustrated in Figure 1, formally prove that the first component of the function *std_unif_disc* is a discrete uniform random variable.

The function *std_unif_disc* can now be used to model the real sequence of Equation (3). We proved in HOL that this sequence is convergent, i.e., it approaches a unique value when n tends to infinity.

Lemma 4:

$$\vdash \forall s. \text{convergent } (\lambda n. \text{fst}(\text{std_unif_disc } n \text{ } s))$$

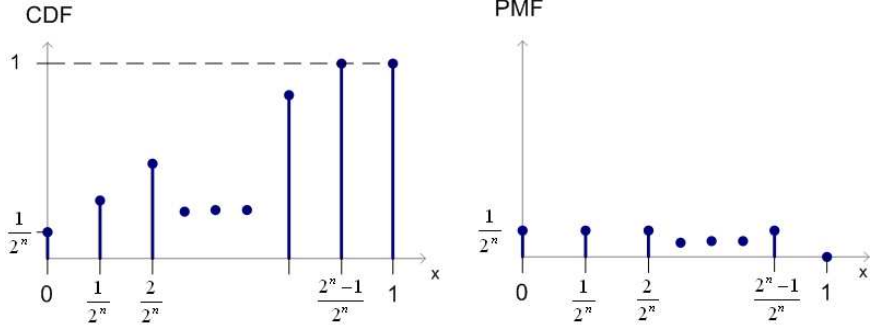


Fig. 1. CDF and PMF for *std_unif_disc*

where, *convergent* represents the HOL function of a convergent real sequence [10]. Based on Lemma 4, we are able to formally specify the Standard Uniform random variable in HOL according to Equation (4).

Definition 2:

$$\vdash \forall s. \text{std_unif_cont } s = \text{lim } (\lambda n. \text{fst}(\text{std_unif_disc } n \ s))$$

where, *lim* M represents the HOL formalization of the limit of a real sequence M (i.e., $\lim_{n \rightarrow \infty} M(n) = \text{lim } M$) [10]. The following properties may be proved using the real analysis theorems [10] and the function definition for *std_unif_disc*.

Lemma 5:

$$\vdash \forall s. 0 \leq \text{std_unif_cont } s \leq 1$$

Lemma 6:

$$\vdash \forall s, n. \text{fst}(\text{std_unif_disc } n \ s) \leq \text{std_unif_cont } s \leq \text{fst}(\text{std_unif_disc } n \ s) + \left(\frac{1}{2}\right)^n$$

Lemma 5 formally shows that the value for the function *std_unif_cont* always lies in the real interval $[0,1]$. The minimum and maximum values of 0 and 1 correspond to the cases when all the elements of the infinite Boolean sequence s are False or True respectively. Lemma 6 highlights the relationship between the values of the first component of the function *std_unif_disc* and the function *std_unif_cont*, i.e., if the value for the former is a , then the value of the later lies in the interval $[a, a + (\frac{1}{2})^n]$.

$$\begin{aligned} (\text{fst}(\text{std_unif_disc } n \ s) = a) &\Rightarrow \\ (a \leq (\text{std_unif_cont } s) \leq a + \left(\frac{1}{2}\right)^n) &\quad (5) \end{aligned}$$

4 Verification of Standard Uniform Random Variable in HOL

In this section, we verify the correctness of the proposed specification of the Standard Uniform random variable using its probability distribution properties. We present the major lemmas and theorems regarding this verification in this paper and more details can be found in [11].

The CDF of a random variable represents the probability that its value is less than or equal to some real number x and is a distinguishing characteristic for all random variables. We formally verify in HOL that the function *std_unif_cont* correctly models the Standard Uniform random variable by proving its CDF to be equal to the theoretical CDF of a Standard Uniform random variable.

$$\mathbb{P}\{s \mid \text{std_unif_cont } s \leq x\} = \begin{cases} 0 & \text{if } x < 0; \\ x & \text{if } 0 \leq x < 1; \\ 1 & \text{if } 1 \leq x. \end{cases} \quad (6)$$

The proof for the cases ($x < 0$) and ($1 \leq x$) is a straightforward implication of Lemma 5 which states that for all infinite Boolean sequences the value of the function *std_unif_cont* lies in the interval $[0,1]$. The probability that the function *std_unif_cont* acquires a value less than 0 is 0 as there is no infinite Boolean sequence that satisfies this condition. Similarly, the probability that the function *std_unif_cont* acquires a value less than or equal to 1 is 1 since all infinite Boolean sequences fulfill this condition.

Lemma 7:

$$\vdash \forall x. (x < 0) \Rightarrow \mathbb{P}\{s \mid \text{fst}(\text{std_unif_cont } s) \leq x\} = 0$$

Lemma 8:

$$\vdash \forall x. (1 \leq x) \Rightarrow \mathbb{P}\{s \mid \text{fst}(\text{std_unif_cont } s) \leq x\} = 1$$

Evaluating the probability of Equation (6) for the interval $[0,1)$ is a surprisingly difficult problem in the HOL theorem prover. However, given that we have evaluated the CDF for the first component of the function *std_unif_disc*, which represents a discrete Standard Uniform random variable, a reasonable approach is to find a discrete approximation to the CDF of the function *std_unif_cont*, which represents the Standard Uniform random variable. The key to this approach is to be able to express the CDF of the function *std_unif_cont* in terms of the CDF of the first component of the function *std_unif_disc*. In order to do this, we need to identify the closest values of the first component of the function *std_unif_disc* corresponding to any given value of the function *std_unif_cont*. We know from Lemmas 1, 2 and 3 that the

first component of the function *std_unif_disc* generates dyadic rationals with denominator 2^n in the interval $[0, 1 - (\frac{1}{2})^n]$. On the other hand, the function *std_unif_cont* can attain any real value in the interval $[0, 1]$. Therefore, the two values of the first component of the function *std_unif_disc* that are closest to any given value, say y , of the function *std_unif_cont* are the two consecutive dyadic rationals (with denominators 2^n) such that the smaller dyadic rational is less than y and the greater dyadic rational is greater than or equal to y . The mathematical concept of *ceiling*, that represents the smallest integer number greater than or equal to a real number, can be used in identifying these dyadic rationals. We proved in HOL that for any positive real number y the above mentioned dyadic rationals are $\frac{\lceil 2^n y \rceil - 1}{2^n}$ and $\frac{\lceil 2^n y \rceil}{2^n}$

Lemma 9:

$$\vdash \forall n, y. (0 \leq y) \Rightarrow \frac{\lceil 2^n y \rceil - 1}{2^n} < y \leq \frac{\lceil 2^n y \rceil}{2^n}$$

where, $\lceil z \rceil$ denotes our HOL definition for the ceiling function that returns the smallest *num* value greater than or equal to its real argument z .

Now we will show how to express the CDF of the function *std_unif_cont* in terms of the CDF of the first component of the function *std_unif_disc* using the above dyadic rationals. It is important to note that the set $\{s \mid \text{fst}(\text{std_unif_disc } n \ s) \leq \frac{m}{2^n}\}$ contains all the infinite Boolean sequences for which the value of the first n bits based on the algorithm implemented by the function *std_unif_disc* is less than or equal to the dyadic rational $\frac{m}{2^n}$. Using Equation (5), we can say that the value produced by the algorithm implemented by the function *std_unif_cont*, for any infinite Boolean sequence that is present in the set $\{s \mid \text{std_unif_disc } n \ s \leq \frac{m}{2^n}\}$, must be less than or equal to $\frac{m+1}{2^n}$. We used this useful reasoning along with Lemma 9 to prove the following

Lemma 10:

$$\begin{aligned} \vdash \forall x, n. (0 \leq x) \Rightarrow \\ \{s \mid \text{fst}(\text{std_unif_disc } n \ s) \leq \frac{\lceil 2^n x \rceil - 2}{2^n}\} \subseteq \\ \{s \mid \text{std_unif_cont } s \leq x\} \subseteq \\ \{s \mid \text{fst}(\text{std_unif_disc } n \ s) \leq \frac{\lceil 2^n x \rceil}{2^n}\} \end{aligned}$$

The first set $\{s \mid \text{fst}(\text{std_unif_disc } n \ s) \leq \frac{\lceil 2^n x \rceil - 2}{2^n}\}$, based on the above reasoning, contains all the infinite Boolean sequences for which the value produced by the algorithm implemented by the function *std_unif_cont* lies in the interval $[0, \frac{\lceil 2^n x \rceil - 1}{2^n}]$. This set is a subset of the set $\{s \mid \text{std_unif_cont } s \leq x\}$, which contains all the infinite Boolean sequences for which the value produced by the algorithm implemented by the function *std_unif_cont* lies in the interval $[0, x]$, as $\frac{\lceil 2^n x \rceil - 1}{2^n}$ is always less than x according to Lemma 9. Similarly, the set $\{s \mid \text{std_unif_cont } s \leq x\}$ is a subset of the set $\{s \mid \text{fst}(\text{std_unif_disc } n \ s) \leq \frac{\lceil 2^n x \rceil}{2^n}\}$, which contains all the infinite Boolean sequences for which the value produced by the algorithm implemented by the function *std_unif_cont* lies in

the interval $[0, \frac{\lceil 2^n x \rceil + 1}{2^n}]$, as x is always less than or equal to $\frac{\lceil 2^n x \rceil + 1}{2^n}$ according to Lemma 9.

Lemma 10 and the monotonic property of the probability function \mathbb{P} , formalized in [12], which states that the probability of a measurable set is always less than or equal to the probability of its measurable superset

$$\vdash \forall s t. \text{measurable}(s) \wedge \text{measurable}(t) \wedge s \subseteq t \Rightarrow \mathbb{P} s \leq \mathbb{P} t \quad (7)$$

can be used to obtain the required relationship between the CDFs of the function *std_unif_cont* and the first component of the function *std_unif_disc*. But, in order to use Equation (7), we must prove all the sets in Lemma 10 to be measurable, i.e., they are in \mathcal{E} . It has been shown in [12], that if a function accesses the infinite boolean sequence using only the `unit`, `bind` and `sdest` primitives then the function is guaranteed to be measurable and thus leads to measurable sets. The function *std_unif_disc* satisfies this condition and thus Hurd's formalization framework can be used to prove

Lemma 11:

$$\vdash \forall x, n. \text{measurable} \{s \mid \text{fst}(\text{std_unif_disc } n \ s) \leq x\}$$

On the other hand, the definition of the function *std_unif_cont* involves the *lim* function and thus the corresponding sets can not be proved to be measurable in a very straight forward manner. Therefore, in order to prove this, we leveraged the fact that each set in the sequence of sets $(\lambda n. \{s \mid \text{fst}(\text{std_unif_disc } n \ s) \leq x\})$ is a subset of the set before it, in other words, this sequence of sets is a monotonically decreasing sequence. Thus, the countable intersection of all sets in this sequence can be proved to be equal to the set $\{s \mid \text{std_unif_cont } s \leq x\}$

Lemma 12:

$$\vdash \forall x. \{s \mid \text{std_unif_cont } s \leq x\} = \bigcap_n (\lambda n. \{s \mid \text{fst}(\text{std_unif_disc } n \ s) \leq x\})$$

Now the set $\{s \mid \text{std_unif_cont } s \leq x\}$ can be proved to be measurable since measurable sets are closed under countable intersections [12] and all the sets in the sequence $(\lambda n. \{s \mid \text{fst}(\text{std_unif_disc } n \ s) \leq x\})$ are measurable according to Lemma 11.

Lemma 13:

$$\vdash \forall x. \text{measurable} \{s \mid \text{std_unif_cont } s \leq x\}$$

Lemmas 10, 11 and 13 along with Equation (7) can be used to obtain the desired relationship between the CDFs. The result can be further simplified

by using the CDF relation for the first component of the function *std_unif_disc* given in Lemma 2.

Lemma 14:

$$\vdash \forall x, n. (0 \leq x) \wedge (x < 1) \Rightarrow \frac{\lfloor 2^n x \rfloor - 1}{2^n} \leq \mathbb{P}\{s \mid \text{std_unif_cont } s \leq x\} \leq \frac{\lfloor 2^n x \rfloor + 1}{2^n}$$

As n approaches infinity both the fractions in Lemma 11 approach x . This fact led us to prove the CDF relation of Equation (6) for the interval $[0,1)$ in the HOL theorem prover. Now, Lemmas 7, 8 and 14 can be used to prove the desired CDF property for the function *std_unif_cont*

Theorem 1:

$$\vdash \forall x. \mathbb{P}\{s \mid \text{std_unif_cont } s \leq x\} = \begin{cases} 0 & \text{if } (x < 0) \\ x & \text{if } (x < 1) \\ 1 & \text{else} \end{cases}$$

Theorem 1 proves that the CDF of the function *std_unif_cont* is the same as the theoretical value of the CDF for a Standard Uniform random variable given in Equation (1) and thus is a formal argument to support the correctness of the fact that the function *std_unif_cont* models a Standard Uniform random variable.

Using similar reasoning as above, we also proved in the HOL theorem prover that the PMF of the function *std_unif_cont* is equal to 0.

Theorem 2:

$$\vdash \forall x. \mathbb{P}\{s \mid \text{std_unif_cont } s = x\} = 0$$

It follows from Theorem 2 that every outcome of the function *std_unif_cont* has a probability 0; which is a unique characteristic of all continuous random variables. Thus, Theorem 2 can be used to formally regard the function *std_unif_cont* as a continuous random variable.

Theorems 1 and 2 along with the measurability proofs of the corresponding sets allow us to formally reason about interesting probabilistic properties of the Standard Uniform random variable within the HOL theorem prover. The measurability of the sets $\{s \mid \text{std_unif_cont } s = x\}$ and $\{s \mid \text{std_unif_cont } s \leq x\}$ can be used to prove that any set that involves a relational property with the random variable *std_unif_cont*, e.g. $\{s \mid \text{std_unif_cont } s < x\}$ and $\{s \mid \text{std_unif_cont } s \geq x\}$, is measurable because of the closed under complements and countable unions property of \mathcal{E} . The CDF and the PMF property proofs given in Theorems 1 and 2 along with the proofs of some basic probability laws [12] can then be used to determine the corresponding probabilistic quantities.

5 Formalization of Continuous Uniform Distribution

The Standard Uniform random variable defined in the previous section can be used to formalize other continuous probability distributions in the HOL theorem prover using Non-uniform random number generation techniques [5]. In this section, we illustrate this by formalizing the Continuous Uniform random variable in the HOL theorem prover.

Continuous Uniform random variable is a random variable for which the probability that it will belong to a subinterval of $[a,b]$ is proportional to the length of that subinterval. It can be characterized by the CDF as follows

$$\mathbb{P}(X \leq x) = \begin{cases} 0 & \text{if } x \leq a; \\ \frac{x-a}{b-a} & \text{if } a < x \leq b; \\ 1 & \text{if } b < x. \end{cases} \quad (8)$$

The process of obtaining random variates with arbitrary distributions using a uniform *random number generator* (RNG) is termed as Non-uniform random number generation. All computer based RNGs generate uniformly distributed numbers [14] in the interval $[0,1]$ and Non-uniform random generation methods are quite commonly used in applications which call for other kinds of distributions. One such commonly used method is the *Inverse Transform method* (ITM) [5] that is used to generate random variates for any continuous probability distribution for which the inverse of the CDF can be expressed in a closed mathematical form. The ITM is based on the following proposition:

Let U be a Standard Uniform random variable. For any continuous distribution F , the random variable X defined by $X = F^{-1}(U)$ has distribution F , where $F^{-1}(u)$ is defined to be that value of x such that $F(x) = u$.

A formal proof for the above proposition may be found in [5].

The Continuous Uniform $[a,b]$ random variable can be formally specified using the ITM since the inverse of its CDF exists in a closed mathematical form.

Definition 3:

$\vdash \forall a,b,s.$

`uniform_cont a b s = (b - a) * (std_unif_cont s) + a`

The probabilistic function *uniform_cont* accepts two real valued parameters a , b and the infinite Boolean sequence s and returns a real number in the interval $[a,b]$. It can be verified that the function *uniform_cont* correctly models a Continuous Uniform random variable by proving its CDF

$(\mathbb{P}\{s | \text{uniform_cont } a \ b \ s \leq x\})$ to be equal to the theoretical value given in Equation (8).

The first step is to express the set $\{s | \text{uniform_cont } a \ b \ s \leq x\}$ in such a way that $(\text{std_unif_cont } s)$ is the only term that remains on the left side of the inequality. This can be done by using the definition of the function uniform_cont and some real number theorems in HOL

Lemma 15:

$$\vdash \forall a \ b \ x. \ \{s \mid \text{uniform_cont } a \ b \ s \leq x\} = \{s \mid \text{std_unif_cont } s \leq \frac{x-a}{b-a}\}$$

The measurability of the set $\{s | \text{uniform_cont } a \ b \ s \leq x\}$ can be proved using Lemmas 13 and 15.

Lemma 16:

$$\vdash \forall a \ b \ x. \ \text{measurable } \{s \mid \text{uniform_cont } a \ b \ s \leq x\}$$

Similarly, the CDF of the random variable uniform_cont can now be proved using the PMF and CDF properties of the Standard Uniform random variable proved in Theorem 1 and 2.

Theorem 3:

$$\vdash \forall a, b, x. \ (a < b) \Rightarrow \mathbb{P}\{s \mid \text{uniform_cont } a \ b \ s \leq x\} = \begin{cases} 0 & \text{if } (x \leq a) \\ \frac{x-a}{b-a} & \text{if } (x \leq b) \\ 1 & \text{else} \end{cases}$$

Other continuous distributions, for which the inverse of the CDF can be expressed in a closed mathematical form such as $\text{Rayleigh}(\sigma)$, $\text{Exponential}(\lambda)$ and $\text{Cauchy}(a, b)$, can also be formally specified in the HOL theorem prover in terms of the formalized Standard Uniform random variable according to the ITM [5]. Their corresponding CDF properties can then be proved using the formal proofs of PMF and CDF properties of the Standard Uniform random variable as was demonstrated in the case of Continuous Uniform random variable in this section.

6 Illustrative Example

In this section, we present a simplified probabilistic analysis example of round-off error in a digital processor to illustrate how the proposed higher-order-logic theorem proving based approach can be used to perform precise quantitative analysis of probabilistic systems.

Consider the problem of determining the probability of the event when the

roundoff error in a digital processor fluctuates the actual value by $\pm 4 \times 10^{-12}$. Assume that the roundoff error for this particular processor is uniformly distributed over the interval $[-5 \times 10^{-12}, 5 \times 10^{-12}]$. The formalized Continuous Uniform random variable of Section 5, represented by the function *uniform_cont*, can be used to determine this probability within the environment of the HOL theorem prover. We use the measurability of the set $\{s \mid \text{uniform_cont } a \ b \ s \leq x\}$, proved in Lemma 16, along with the formalization of the additive law of probability, given in Section 2.2, and the set theory in HOL to prove

Lemma 17:

$$\begin{aligned} \vdash \forall a, b, x, y. (a < b) \Rightarrow \\ \mathbb{P}\{s \mid (x \leq \text{uniform_cont } a \ b \ s) \wedge \\ (\text{uniform_cont } a \ b \ s \leq y)\} = \\ \mathbb{P}\{s \mid \text{uniform_cont } a \ b \ s \leq y\} - \\ \mathbb{P}\{s \mid \text{uniform_cont } a \ b \ s \leq x\} \end{aligned}$$

Now, the CDF relation for the function *uniform_cont*, that we proved in Theorem 3, can be utilized to evaluate the probability under consideration by specializing Lemma 17 for the case when a , b , x and y are equal to -5×10^{-12} , 5×10^{-12} , -4×10^{-12} and 4×10^{-12} , respectively.

$$\begin{aligned} \vdash \mathbb{P}\{s \mid (-4 \times 10^{-12} \leq \text{uniform_cont } -5 \times 10^{-12} \ 5 \times 10^{-12} \ s) \wedge \\ (\text{uniform_cont } -5 \times 10^{-12} \ 5 \times 10^{-12} \ s \leq 4 \times 10^{-12})\} = 0.8 \end{aligned}$$

Thus, we are able to determine the unknown probability with 100% precision; a novelty which is not available in the simulation based approaches. This added benefit comes at the cost of a significant amount of time and effort spent, while formalizing the system behavior, by the user.

7 Related Work

Due to the vast application domain of continuous random variables, many researchers around the world are trying to improve the modeling techniques for continuous probability distributions in computer based environments. The ultimate goal is to come up with a probabilistic analysis framework that includes robust and accurate analysis methods, has the ability to perform analysis for large-scale problems and is easy to use. In this section, we provide a brief account of the state-of-the-art and some related work in this field.

Hurd's PhD thesis, reviewed in Section 2.2 of this paper, presents a methodology to formalize probabilistic algorithms in higher-order-logic. Hurd's main contributions include the formalization of the mathematical measure and probability theories along with the development of a framework to formalize prob-

abilistic algorithms in the HOL theorem prover. Hurd also presented the formalization of several discrete random variables but did not touch the topic of formalizing continuous random variables, which has been the main topic of discussion in the current paper. We have proposed to build upon Hurd’s PhD thesis to formalize the continuous random variables so that the HOL theorem prover can be used as a probabilistic analysis platform for systems which exhibit continuous random behaviour.

A number of *probabilistic languages*, e.g., Probabilistic cc [9], λ_o [20] and IBAL [22], have been proposed that are capable of modeling random variables. The probabilistic languages allow programmers to perform probabilistic computations at the level of probability distributions by treating probability distributions as primitive data types. These probabilistic languages are quite expressive and have been successfully used to perform probabilistic analysis of systems which exhibit continuous random behavior but they have their own limitations. For example, either they require a special treatment such as the lazy list evaluation strategy in IBAL and the limiting process in Probabilistic cc or they do not support precise reasoning as in the case of λ_o . The proposed theorem proving based approach, on the other hand, is not only capable of formally expressing most continuous probability distributions but also to precisely reason about them.

It is interesting to note that one of the probabilistic languages, λ_o , proposed by Park *et. al* in [20] is based on sampling functions. A sampling function is defined as a mapping from the unit interval $[0,1]$ to a probability domain \mathfrak{D} . Given a random number drawn from a Standard Uniform distribution, it returns a sample in \mathfrak{D} , and thus specifies a unique probability distribution. Thus, this approach is very similar to what we have proposed in this paper, as it also utilizes the Standard Uniform random variable to obtain other continuous random variables. [20] contains sampling algorithms for various continuous random variables which can be utilized to formalize the respective random variables in the HOL theorem prover using our formalized Standard Uniform random variable.

Another alternative for formal probabilistic verification is to use probabilistic model checking techniques, e.g., [1], [23]. Like the traditional model checking, it involves the construction of a precise mathematical model of the probabilistic system which is then subjected to exhaustive analysis to verify if it satisfies a set of formal properties. This approach is capable of providing precise solutions in an automated way; however it is limited for systems that can only be expressed as a probabilistic finite state machine. Our proposed theorem proving based approach, in contrast, is capable of handling all kinds of probabilistic systems including the *unbounded* ones. Another major limitation of the probabilistic model checking approach is the state space explosion [4], which is not an issue with our approach.

Knuth and Yao [15] presented procedures that minimize the average number of bits required to generate random numbers from arbitrary probability distributions in arbitrary systems of notation. This classical paper about Non-uniform random number generation also covers the main ideas behind the formalization that we have presented. An interesting topic for further research is to utilize the formalized Standard Uniform random variable to formalize Knuth and Yao's results in the HOL theorem prover.

8 Conclusions

We have presented a formalization methodology for the Standard Uniform random variable in higher-order logic. We also demonstrated that the formalized Standard Uniform random variable can be used in conjunction with Non-uniform random number generation methods to formalize other continuous random variables. To the best of our knowledge, this is the first time that a methodology for the formal reasoning of continuous probability distributions in a mechanical theorem prover has been proposed. Therefore, the presented work can be considered a significant step towards the development of a formal reasoning framework for systems involving continuous probability distributions.

The proposed formalization opens the doors for a new era in the fields of computer science, operations research and statistics. Computer scientists may use the formalized continuous random variables in program verification and comparisons of algorithms. In operations research, the formalized probabilistic analysis may be used to complement large scale simulations. Statisticians may use the formalized continuous random variables to formally test and compare estimators before using them in real life.

A distinguishing characteristic of the proposed probabilistic analysis approach is the ability to perform precise quantitative analysis of probabilistic systems. Though, this approach cannot be automated and thus involves considerable user interaction. On the other hand, simulation is capable of handling analytically complex probabilistic analysis problems in an automated way but the solutions provided are approximated. Therefore, we consider simulation and higher-order theorem proving as complementary techniques, i.e., the methods have to play together for a successful probabilistic analysis framework. For example, the proposed theorem proving based approach can be used for the safety critical parts of the design and simulation based approaches can handle the rest.

As a next step towards a complete formalization framework for continuous random variables, we are formalizing the ITM itself in HOL. This will con-

siderably ease the formalization process for the corresponding continuous random variables. Similarly, the formalization of the Standard Uniform random variable can also be transformed to formalize other continuous probability distributions, for which the inverse CDF is not available in a closed mathematical form, by exploring the formalization of other Non-uniform random number generation techniques such as Box-Muller and acceptance/rejection [5]. In order to have a complete formal probabilistic analysis framework, it is also essential to formalize the theory of expectation and the basic statistical quantities of mean, moment, and variance. In order to have a complete formal probabilistic analysis framework, it is also essential to formalize the theory of expectation and the basic statistical quantities of mean, moment and variance.

Acknowledgements The authors would like to thank Joe Hurd from Oxford University, UK, for his precious advice and feedback on our methodology. We would also like to thank the anonymous referees for their comments on how to improve the paper.

References

- [1] C. Baier, B. Haverkort, H. Hermanns, and J. P. Katoen. Model Checking Algorithms for Continuous time Markov chains. *IEEE Transactions on Software Engineering*, 29(4):524–541, 2003.
- [2] P. Billingsley. *Probability and Measure*. John Wiley, 1995.
- [3] A. Church. A Formulation of the Simple Theory of Types. *Journal of Symbolic Logic*, 5:56–68, 1940.
- [4] E. M. Clarke, O. Grumberg, and D. A. Peled. *Model Checking*. The MIT Press, 2000.
- [5] L. Devroye. *Non-Uniform Random Variate Generation*. Springer-Verlag, 1986.
- [6] W. Feller. *An Introduction to Probability Theory and Its Applications*, volume 2. Wiley, 1971.
- [7] M. J. C. Gordon. Mechanizing Programming Logics in Higher-Order Logic. In *Current Trends in Hardware Verification and Automated Theorem Proving*, pages 387–439. Springer-Verlag, 1989.
- [8] M. J. C. Gordon and T.F. Melham. *Introduction to HOL: A Theorem Proving Environment for Higher-Order Logic*. Cambridge University Press, 1993.
- [9] V. T. Gupta, R. Jagadeesan, and P. Panangaden. Stochastic Processes as Concurrent Constraint Programs. In *Principles of Programming Languages*, pages 189–202. ACM Press, 1999.
- [10] J. Harrison. *Theorem Proving with the Real Numbers*. Springer-Verlag, 1998.

- [11] O. Hasan and S. Tahar. Standard Uniform Distribution Theory in HOL-4. Technical Report, Concordia University, Montreal, Canada, 2007.
- [12] J. Hurd. *Formal Verification of Probabilistic Algorithms*. PhD thesis, University of Cambridge, Cambridge, UK, 2002.
- [13] W. J. Kennedy and J. E. Gentle. *Statistical Computing*. Marcel-Dekker, 1980.
- [14] D. E. Knuth. *The Art of Computer Programming*, volume 2. Addison-Wesley Professional, 1998.
- [15] D. E. Knuth and A. C. Yao. The Complexity of Nonuniform Random Number Generation. *Algorithms and Complexity: New Directions and Recent Results*, pages 357–428, 1976.
- [16] D. J. C. MacKay. Introduction to Monte Carlo methods. In *Learning in Graphical Models, NATO Science Series*, pages 175–204. Kluwer Academic Press, 1998.
- [17] B. D. McCullough. Assessing the Reliability of Statistical Software: Part I. *The American Statistician*, 52(4):358–366, 1998.
- [18] B. D. McCullough. Assessing the Reliability of Statistical Software: Part II. *The American Statistician*, 53(2):149–159, 1999.
- [19] R. Milner. A Theory of Type Polymorphism in Programming. *Journal of Computer and System Sciences*, 17:348–375, 1978.
- [20] S. Park, F. Pfenning, and S. Thrun. A Probabilistic Language based upon Sampling Functions. In *Principles of Programming Languages*, pages 171–182. ACM Press, 2005.
- [21] L. C. Paulson. *ML for the Working Programmer*. Cambridge University Press, 1996.
- [22] A. Pfeffer. IBAL: A Probabilistic Rational Programming Language. In *International Joint Conferences on Artificial Intelligence*, pages 733–740. Morgan Kaufmann Publishers, 2001.
- [23] J. Rutten, M. Kwaiatkowska, G. Normal, and D. Parker. Mathematical Techniques for Analyzing Concurrent and Probabilistic Systems. *CRM Monograph*, 23, 2004.
- [24] K. S. Tridevi. *Probability and Statistics with Reliability, Queuing and Computer Science Applications*. Wiley-Interscience, 2002.
- [25] B. Widrow. Statistical Analysis of Amplitude-quantized Sampled Data Systems. *AIEE Trans. (Applications and Industry)*, 81:555–568, January 1961.