

Formal Verification of Universal Numbers using Theorem Proving

Adnan Rashid*, Ayesha Gauhar, Osman Hasan*, Sa'ed Abed** and Imtiaz Ahmad**

*School of Electrical Engineering and Computer Science (SEECS)

National University of Sciences and Technology (NUST), Islamabad, Pakistan

{adnan.rashid,14mseegauhar,osman.hasan}@seecs.nust.edu.pk

**Computer Engineering Department

College of Engineering and Petroleum, Kuwait University, Kuwait

{s.abed,imtiaz.ahmad}@ku.edu.pk

Abstract

Universal number (Unum) is a number representation format that can reduce the memory contention issues in multicore processors and parallel computing systems by optimizing the bit storage in the arithmetic operations. Given the safety-critical nature of applications of Unum format, there is a dire need to rigorously assess the correctness of Unum based arithmetic operations. Unums are of three types, namely, Unum-I, Unum-II and Unum-III (commonly known as Posits). In this paper, we provide a higher-order-logic formalization of Unum-III (posits). In particular, we formally model a posit format (binary encoding of a posit), which is comprised of the sign, exponent, regime and fraction bits, using the HOL Light theorem prover. In order to prove the correctness of a posit format, we formally verify various properties regarding conversions of a real number to a posit and a posit to a real number and the scaling factors of the regime, exponent and fraction bits of a posit using HOL Light.

Keywords: Systems Biology, Molecular Pathways, Reaction Kinetics, Automated Reasoning, Theorem Proving

1 Introduction

Floating-point number format is widely used by the scientific community in application areas ranging from the aerospace, applied mathematics, physics to weather forecasting, for the representation of real numbers on a computer. Moreover, it is utilized for the execution of various arithmetic operations, i.e., addition, subtraction, division and multiplication, requiring an efficient hardware implementation. The IEEE-754 floating-point standard represents a real number as a signed fraction times an integer power of 2, i.e., $\pm(1 + f)2^e$, where f is a fraction and e is an exponent, and allows the representation of real numbers in computers using various bits. This includes the handling of the rounding and fraction bits, and various invalid results, such as Not-a-Number (NaN), which is returned as a result of an invalid arithmetic operation, such as $0/0$ or $\infty \times 0$. However, the IEEE-754 floating-point standard suffers from various limitations, such as limited numerical precision as a result of allocating a fixed number of exponent and mantissa bits, failure of the associative and distributive laws of real number arithmetic due to rounding and the hardware cost for handling the denormalized numbers.

John L. Gustafson, in 2015, proposed Universal Numbers (Unums) [1] that can overcome the above-mentioned limitations of the IEEE-754 floating point standard and provide a more precise representation of real numbers for performing computer arithmetics. There are three types of Unums, namely, Unum-I, Unum-II and Unum-III (commonly known as Posits). Unum-I [1] has a variable-length format as opposed to the fixed length floating-point number format and also provides a better numerical accuracy. However, its variable-length format makes it unexciting for hardware implementations. Unum-II [1] exhibits some interesting characteristics, such as calculating the exact reciprocal of a number and performing negation of a number simply. However, it requires pre-computed lookup tables to perform various arithmetic operations that makes it impractical for larger arithmetic word sizes. Unum-III or Posits [2, 3] are considered as the hardware-friendly version of Unums that provide an efficient utilization of fixed bit sizes, resulting in higher accuracy arithmetic for a given storage requirement, and are intended to be a drop-in replacement for the IEEE-754 format. Posit exhibits various features, such as simple rounding, larger dynamic range, better closure, no denormalized numbers to handle, and therefore simplifies the hardware and software implementations. Moreover, posit arithmetic provides identical answers on different computer systems, which is not possible using the IEEE-754 floating point arithmetic standard. Posits do not overflow to infinity or underflow to zero. Moreover, NaN provides an action rather than a bit pattern

as in floating-point numbers. Also, its processing unit takes less circuitry than the IEEE Floating-point Unit (FPU) [2]. All these features lead to an improved memory bandwidth and power efficiency. Moreover, **posits** have been implemented as an alternative to the floating-point number format in hardware and software. For example, the hardware architecture of Unum adder/subtractor and multiplier has been designed and implemented using Field-Programmable Gate-Arrays (FPGAs). Moreover, a Verilog Hardware Description Language (HDL) generator has been constructed for performing these arithmetic operations [4, 5, 6]. Software libraries for **posit**-based floating-point operations are also available for C# [7], C¹, C++² and Julia³ programming languages. Moreover, **posits** have outperformed the fixed point number system, in terms of accuracy and memory utilization, in various computational intensive applications, such as deep convolutional neural networks [8, 9].

The real number programs are widely used for analyzing the dynamics of the physical systems in various applications, such as aerospace, robotics and physics. They use the floating-point approximations resulting in the accumulation of floating-point inaccuracies that grow as the computation proceeds and thus introduce some unavoidable bugs that may lead to dire consequences. For example, an error in the Floating-point Division (FDIV) instruction of the Intel Pentium processors in 1994 resulted into a financial loss of \$475M⁴. Similarly, an uncaught floating-point exception resulted in the destruction of the Ariane 5 rocket shortly after its takeoff in 1996⁵. The cost of such errors in floating-point arithmetic is huge. The above-mentioned popular incidents due to such errors resulted into replacement of a large number of processors having FDIV instruction errors, leading to a huge financial loss of \$475M and destruction of the Ariane 5 rocket. Therefore, one can expect that similar kind of bugs today may cost tenfold of that loss without performing an exhaustive analysis of arithmetic based on **posits** [10]. Moreover, the conventional computer based simulation and numerical analysis techniques involve the unverified symbolic algorithms, discretization and numerical errors, and thus cannot ascertain exhaustive analysis of the safety-critical systems. Therefore, the formal verification of these number formats, performing various arithmetic, is a dire need.

Formal verification method [11] is a system analysis technique that mainly involves two steps; 1) developing a computer based mathematical model of the given system, 2) verifying that the system's model meets the rigorous specifications of the intended behaviour, based on deductive reasoning. Since deductive reasoning involves the use of the logical reasoning and evidence to reach a conclusion from one or more premises that are considered to be true. Therefore, the usage of this method increases the chances of catching the errors that are often ignored by the conventional simulation based and numerical analysis techniques. The idea of doing formal verification of a complex system is to identify its safety-critical components/parts that require an exhaustive analysis. For example, in the case of Ariane 5 rocket, for the identification the uncaught floating-point exception, it is sufficient to perform the formal analysis of a component providing the floating-point arithmetic. Therefore, it may not require formal verification of the whole system. Theorem proving [12] is one of the frequently used formal verification techniques that involves constructing a mathematical model of the given system based on logic and verifying its various properties by computer programs involving automated reasoning. Here, the automated reasoning refers to the computer-based deductive reasoning process that is based on the logical reasoning and evidence. Thus, it ensures the soundness of the theorem proving technique. Theorem proving can be interactive or automatic based on the choice of the underlying logic, which can be propositional, first and higher-order logic. Higher-order logic provides more expressiveness, which is important for analyzing the dynamics of physical systems. However, it requires user interaction for developing proofs within a theorem prover. Many theorem provers (automatic and interactive), such as HOL Light [13, 14], Coq [15, 16], ACL2 [17, 18] and PVS [19, 20, 21] have been used for the formal verification of the floating point numbers and their arithmetic. Moreover, there is a research group working on the verification of the different components/operations of the Intel processors, over the years. Some of the notable contributions are from Harrison [22, 23, 24, 25, 26], O'Leary [27], Narasimhan and Kaivola [28, 29], Slobodova [30] and Peter Tang [25] who have been working in the Intel research group. Similarly, Rockwell Collins Inc. and NASA have been successfully using formal methods for analyzing various aspects of avionics [31, 32, 33, 34]. However, none of these contributions cater for **posit**, which are intended as drop-in replacement for floating-point numbers in computer systems.

In this paper, we provide a formalization of **posits** (Type III Unums) using HOL Light. In particular, we formally model a **posit** format, which is composed of the sign, exponent, regime and fraction bits. Moreover, we formalize a conversion of a **posit** to its equivalent real number (decoding) and a real number to its equivalent **posit** representation (encoding), which mainly uses the notion of the fraction and exponential rounding. Finally, we formally verify various properties of the **posits** regarding these conversions and the scaling factors of the regime,

¹<https://github.com/libcg/bfp>

²<https://github.com/eruffaldi/cppPosit>

³<https://github.com/milankl/SoftPosit.jl>

⁴<https://www.intel.com/content/www/us/en/history/history-1994-annual-report.html>

⁵<https://www-users.math.umn.edu/~arnold/disasters/ariane.html>

exponential and fraction bits using **HOL Light**.

The novel contributions of the paper are:

- A higher-order logic formal model of a **posit** format, which includes the sign, exponent, regime and fraction bits, using the **HOL Light** theorem prover. **Posit** has not been formalized in any of the theorem prover before this paper.
- Higher-order logic formalization of the conversion of a **posit** to its equivalent real number and a real number to its equivalent **posit**.
- Formal verification of properties regarding conversions of a real number to a **posit** and a **posit** to a real number.
- Formal verification of properties regarding the scaling factors of the regime, exponential and fraction bits of a **posit** using **HOL Light**. These properties regarding the conversions and scaling factors ensure the correctness of our formalization of posit presented in Section 4.1 of the paper. Moreover, they would be useful for performing the arithmetic based on **posit**.

It is important to note here that our **HOL Light** code for the formal verification of **Unums** is publicly available for download at [35] and thus can be used by other researchers in the development of a formal library for **Unum** arithmetic.

2 Preliminaries

This section introduces the **HOL Light** theorem prover and posits.

2.1 **HOL Light** Theorem Prover

HOL Light [36] is a widely used interactive proof assistant for higher-order logic. The **HOL Light** is written in the strongly-typed functional programming language ML [37]. Theorems are formalized as axioms or inferred from the already verified theorems available in theories by inference rules. A theorem consists of a finite set Ω of Boolean terms (assumptions) and a Boolean term S as a conclusion. A new theorem is verified using any previously proved theorems and the primitive inference rules or applying existing axioms/inference rules in the **HOL Light** theorem proving environment that preserves the soundness of this approach. Many mathematical concepts have been formalized as **HOL Light** theories. A theory consists of a collection of valid **HOL Light** types, constants, axioms, definitions, and theorems. The **HOL Light** theorem proving system offers a wide range of theories, such as Boolean algebra, arithmetic, real numbers and list theories, which are extensively used in our formalization. Various automatic proof procedures [38] are also available in **HOL Light** to help and guide the user in conducting a proof effectively, efficiently and professionally. **HOL Light** has been used for the formal verification of floating-point numbers, the arithmetic involving these numbers and the associated algorithms. Some of the notable contributions are the formal verification of IA-64 division algorithms [23], square root algorithms [26], floating-point trigonometric functions [23] and floating-point exponential functions [13], development of a machine-checked theory of floating point arithmetic for the IA-64 architecture [14], parameterized floating-point formalization [39] and hierarchical verification of the IEEE-754 table-driven floating-point exponential function [40]. Similarly, some notable contributions in PVS related to arithmetic systems include the formalization of IEEE-854 floating-point standard [19], and the formal verification of IEEE rounding [41], IEEE compliant subtractive division algorithms [42], VAMP floating point unit [20] and IEEE floating point adder [43]. However, the **HOL Light** theorem prover supports automated reasoning of a larger set of computer arithmetic foundations that are widely used for analyzing the continuous dynamics of the engineering and physical systems, which is one of the motivations for choosing it for our proposed formalization of **posit**. This set of computer arithmetic foundational libraries will be extensively used in making a comparison between formal analysis of posits and floating-point numbers, which is one of our future directions. Moreover, the **HOL Light** theorem prover has the smallest trusted core (i.e., approximately 400 lines of Ocaml code) amongst all other higher-order-logic theorem prover and the underlying logic kernel has been verified in the CakeML project [44, 45].

Some standard symbols, their meanings and their **HOL Light** representations used in this paper are presented in Table 1.

In order to facilitate the understanding of the paper, we presented majority of the formalization of posits (Sections 4.1 and 4.2) in simple Math notation. However, for some of the **HOL Light** functions/symbols, we used the mathematical notations presented in Table 2. Some of these notations may not correlate with the traditional conventions. However, they have been considered only to facilitate the understanding of the paper.

Table 1: HOL Light Symbols

HOL Light Symbols	Standard Symbols	Meanings
\sim	not	Logical <i>negation</i>
$\langle = \rangle$	$=$	Equality in Boolean domain
num	\mathbb{N}	Natural numbers data type
real	\mathbb{R}	Real data type
SUC n	$(n + 1)$	Successor of natural number
&a	$\mathbb{N} \rightarrow \mathbb{R}$	Casting from a Natural number a to a Real number a
@f	Hilbert choice operator	Returns f if it exists
k DIV m	<i>quotient</i>	Returns the quotient of the division of two real numbers k and m
-- x	$-x$	Negative x
EL n l	<i>element</i>	Extracts n^{th} element of List l
LAST l	<i>last element</i>	Last element of List l
[a; b; c]	$[a, b, c]$	List having elements as a , b and c

2.2 Posits (Unum-III)

Posits (Unum-III) [2], utilize a fixed number of bits as opposed to Type I unums. The precise number may be chosen for a particular implementation, ranging from two bits up to many thousand bits. Posits may be simply implemented both in hardware and software. Moreover, they employ the similar type of low-level circuit building blocks that IEEE-754 floating-point numbers utilize for performing various arithmetic operations, such as integer addition and multiplication, and those also cover less chip area. Figure 1 presents a structure of an n -bit posit representation.

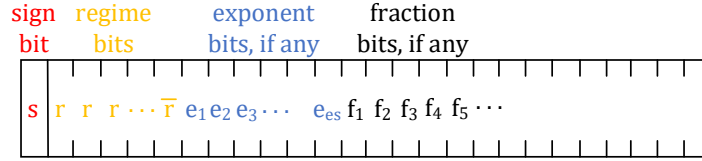


Figure 1: Generic Posit Format for Finite, Nonzero Values

Posit representation consists of sign, regime, exponent and fraction bits. It is to be highlighted that the only boundary is shown between the sign bit and the rest of the bits since the other boundaries are flexible and depend on number of the regime bits. The regime bits are a sequence of identical bits r (all 1s or 0s), which are terminated by the opposite bit \bar{r} for the case of non-zero exponent and fraction bits. In the case of zero exponent and fraction bits, identical bits r in a regime are terminated by the end of the posit. The sign bit serves the purpose of representing the positive and negative numbers, i.e., it is 0 for the positive numbers and 1 for the negative numbers. Moreover, we need to take the 2s complement for the negative numbers before decoding the regime, exponent, and fraction bits.

To capture the idea of regime, Figure 2 provides some binary strings and their corresponding interpretations as real numbers k determined by the run length of the regime bits. Here, the symbol x in a bit string models the *don't care* condition, i.e., the interpretation does not depend on the value of that bit.

Binary	0000	0001	001x	01xx	10xx	110x	1110	1111
Numerical meaning, k	-4	-3	-2	-1	0	1	2	3

Figure 2: Regime Bit Illustration

The leading bits in all bit strings (Figure 2) are known as the regime of the number. All binary strings start with some sequence of all 0 and all 1 bits in a row and terminate by either the complementary bit or the end of the posit. The identical bits r of the regime bits are color-coded in amber, whereas, the opposite bit \bar{r} that terminates the run, if any, is color-coded in brown. Assume m represents the number of identical bits in a run. If

Table 2: Conventions used for HOL Light Functions

HOL Light Functions	Mathematical Conventions	Descriptions
\wedge	\wedge	Logical <i>and</i>
\vee	\vee	Logical <i>or</i>
$\sim(a = b)$	$a \neq b$	a is not equal to b
$!x.t$	$\forall x.t$	For all $x : t$
$?x.t$	$\exists x.t$	There exists $x : t$
$\backslash x.t$	$\lambda x.t$	Function that maps x to $t(x)$
$==>$	\Rightarrow	Implication
$\&a$	\hat{a}	Casting from a Natural number a to a Real number a
<code>int_of_num a</code>	\hat{a}	Casting from a Natural number a to an Integer a
<code>num_of_int a</code>	\tilde{a}	Casting from an Integer a to a Natural number a
<code>z pow n</code>	z^n	z raise to power Natural number n
<code>x ipow y</code>	x^y	x raise to power integer y
<code>a EXP b</code>	a^b	a raise to power b , where a and b are the natural numbers
<code>nb_num</code>	nb_n	Casting from an Integer nb to a Natural number using <code>int_of_num</code>
<code>es_num</code>	es_n	Casting from an integer es to a Natural number using <code>int_of_num</code>
<code>TL l</code>	\underline{l}	Tail of List l
<code>CONS h t</code>	$h::t$	Concatenates head h of a List with its tail t
<code>HD l</code>	\overline{l}	Head of List l
<code>APPEND l1 l2</code>	$l_1 ++ l_2$	Append List l_1 with List l_2
<code>MEM m l</code>	$m \in l$	m is a member of List l
$\sim(\text{MEM } m \ l)$	$m \notin l$	m is not a member of List l
<code>NIL l</code>	$[\]$	List l is empty
<code>LENGTH l</code>	$ l $	Length of List l
<code>real_to_posit_check3</code>	$\text{posit}_{\text{real}}$	Conversion of a real number to its corresponding posit representation
<code>add_zero_real</code>	$\text{real}_{\text{posit}}$	Conversion of a posit representation to its corresponding real number
<code>exponential_rounding1</code>	round_e	Exponential rounding of a posit representation
<code>exponential_round_check1</code>	cond_e	Condition on the exponent bits in case of exponential rounding
<code>scale_factor_e</code>	scaling_e	Scaling factor of the exponent bits
<code>fraction_rounding1</code>	round_f	Fractional rounding of a posit representation
<code>fraction_residue_set1</code>	residue_f	Condition on the residue value in case of fractional rounding
<code>scale_factor_f</code>	scaling_f	Scaling factor of the fraction bits
<code>scale_factor_r</code>	scaling_r	Scaling factor of the regime bits

the identical bits in a regime bit are 1, then $k = m - 1$, otherwise, it is $k = -m$ as given in Figure 2. The regime provides a scale factor of $used^k$, where $used = 2^{2^{es}}$ with es representing maximum exponent bits.

The next bits in a posit structure are the exponent e bits that are color-coded in blue (Figure 1) and are considered as an unsigned integer. They model a scaling factor of 2^e . There can be a maximum of es exponent bits depending on the bits remaining on the right side of the regime.

Any bits left after the regime and the exponent bits in a posit model the fraction f and it is quite similar to the fraction $1.f$ in a floating-point number, with 1 as a hidden bit. Moreover, there are no subnormal numbers with a hidden bit of 0 as they are in floating-point numbers. The two exception values for posit are 0 and $\pm\infty$. When all bits of a posit are zero, it represents the number 0. Whereas, the first bit as 1 and the remaining bits as 0s represent the value $\pm\infty$.

Now, we illustrate the posit format (structure of a posit representation) described above, using an example of decoding a 16-bit posit 0000110111011101 as shown in Figure 3.

We have picked $es = 3$, which causes the value represented by the regime bits to provide a scaling factor between the negative and positive powers of $2^{2^3} = 256$. It is important to note here that the standard 16-bit posit consist of es of size 1. However, we have taken it as 3 for illustration purposes as shown in Figure 3. The

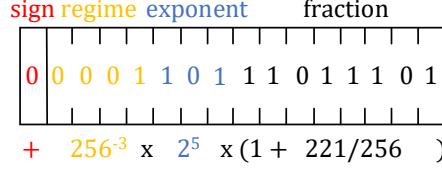


Figure 3: Decoding of a 16-bit Posit

sign bit of 0 asserts that it is a positive value/number. The regime bits consist of a run of three 0s that is terminated by a 1, making the power of useed equal to -3 . The regime bits present a scale factor of 256^{-3} . The exponent bits, 101, represent the decimal number 5 as an unsigned binary integer, and introduce another scale factor 2^5 . Finally, the fraction bits 1101101 represent 221 as an unsigned binary integer, so the fraction becomes $1 + \frac{221}{256}$. The overall value decoded by a 16-bit posit is given as follows:

$$256^{-3} \times 2^5 \times \left(1 + \frac{221}{256}\right) = 477 \times 2^{-27} \\ \approx 3.55393 \times 10^{-6}$$

3 Related Work

This section provides some related work regarding the formal verification of floating-point numbers, and hardware and software implementations of posits.

3.1 Formal Verification of Floating-point Numbers

Many theorem provers, such as HOL Light, Coq, ACL2 and PVS have been used for the formal verification of the floating point numbers and their arithmetic. Miner [19] employed the PVS theorem prover for a formalization of ANSI/IEEE-854 standard for Radix-Independent floating-point arithmetic. It mainly involves the mapping of floating-point numbers to reals, mapping of reals to floating-point numbers, rounding and various arithmetic operations, such as addition, subtraction, multiplication, division and square root operations. Similarly, Berg et al. [20] developed a formal library for IEEE rounding [41] in PVS while utilizing the formal definition of rounding provided by Miner. Moreover, the authors used it to formally verify the correctness of a fully IEEE compliant floating-point unit used in the VAMP processor. Some more notable contributions in PVS related to arithmetic systems include the formal verification of IEEE compliant subtractive division algorithms [42], VAMP floating point unit [20] and IEEE floating point adder [43].

Daumas et al. [15] provided a generic library to formally reason about the floating-point numbers using the Coq theorem prover. The proposed formal library for the floating-point arithmetic caters for an arbitrary floating-point format and an arbitrary base, i.e., it accommodates both bases 2 and 10 for the IEEE-754 standard. Similarly, Boldo et al. [16] proposed a framework for formally verifying the floating-point C programs. The authors extracted the verification conditions from C programs annotated at the source code level that are discharged using Coq.

Harrison [13] provided the formal verification of an algorithm for computation of the exponential function in IEEE-754 standard binary floating-point arithmetic using the HOL Light theorem prover. Later, Harrison [14] generalized the formal library of floating point arithmetic, by incorporating a wide variety of floating point formats. Moreover, the authors used their proposed formalization for the verification of the floating point arithmetic performed in Intel Itanium Architecture (IA)-64. Similarly, Harrison [23] provided a number of formally verified algorithms for the evaluation of the transcendental functions, such as sine and cosine, for Intel IA-64 using double-extended precision floating point arithmetic. Some more notable contributions in HOL Light are the formal verification of IA-64 division algorithms [23], square root algorithms [26], parameterized floating-point formalization [39] and hierarchical verification of the IEEE-754 table-driven floating-point exponential function [40].

O’Leary et al. [46] proposed a hybrid verification approach, based on theorem proving and model checking, for formally verifying the Intel’s FPU at the gate level. Akbarpour et al. [47] presented a formalization of fixed-point arithmetic using the HOL theorem prover. The authors formally modeled the fixed-point number system and provided specifications of various rounding modes, such as the directed and even rounding modes. Moreover,

they performed an error analysis for verifying rounding and various arithmetic operations, such as addition, subtraction, division and multiplication.

Moore et al. [17] used the automated theorem prover ACL2 for formally verifying the AMD-K5 floating-point division unit using the ACL2 theorem prover. Similarly, Rusinoff [18] formally verified the correctness of the floating point arithmetics, such as multiplication, division, and square root instructions of the AMD-K7 microprocessor using ACL2.

Intel has been applying formal verification after the incident of the infamous Intel processor bug. Indeed, there is a research group working on the verification of the different components/operations of the Intel processors, over the years. Some of the notable contributions are from Harrison [22, 23, 14, 24, 25, 26], O’Leary [27], Narasimhan and Kaivola [28, 29], Slobodova [30] and Peter Tang [25] who have been working in the Intel research group. Moreover, to the best of our knowledge, no bugs have been reported regarding Intel processors in literature after the infamous Pentium IV bug. The application of the formal verification could be one of the main reasons behind this as well. For example, Bentley [10] describes the steps for the identification of the bugs in the Pentium IV processor design prior to initial silicon. The author claims that he identified over 100 logic bugs and about 20 of them were *high quality* bugs that would not have been found using any other pre-silicon validation processes. Two out of those 20 bugs were classic floating point data space problems. In particular, the Floating ADD (FADD) instruction had a bug, where the 72-bit Floating Point (FP) adder was setting the carryout bit to 1 for a specific combination of source operands when there was no actual carryout. The author believes that if this error had not been caught, it may have resulted into a bug similar to the Floating Divide (FDIV) problem of the Pentium processor. More details about the validation of the Intel Pentium 4 microprocessor by the Intel research group can be found at [10]. Moreover, some experiences of O’Leary at Intel about the verification of the floating point arithmetic of the Intel Pentium 4 and Core i7 processors can be found at [48]. Similarly, many bugs have been found in different software over the past years. For example, Gesellensetter et al. [49] found a bug in the scheduler of the GNU Compiler Collection (GCC) compiler for a Very Long Instruction Word (VLIW) processor during its verification using the Isabelle/HOL theorem prover. Johnson [50] provides some natural history of bugs and discusses about the usage of the formal methods for analyzing the software issues in space related applications. Similarly, Fitzgerald et al. [51] discusses about the deployment of the formal verification method in industrial applications. Moreover, Zhang et al. [52] survey the successful deployment of formal methods in the industrial settings. Since posits are intended as drop-in replacement for floating-point numbers in computer systems, their formal verification is of utmost important to ensure the absence of any bugs before they are used in the processors.

Rockwell Collins Inc., a famous multinational company providing products and services regarding the aerospace applications, and NASA have been successfully using formal methods for analyzing various aspects of avionics. Whalen et al. [31] integrated formal methods with the model-based development tools, i.e., Simulink and SCADE Suite for the verification of the software during the design cycle for safety-critical avionics applications. The authors developed a set of tools that translate the Simulink models to formal models that can be used by the model checkers and theorem provers for the automatic analysis of these models. Moreover, they formally analyzed an Unmanned Aerial Vehicle (UAV) controller modeled in Simulink. During the analysis of the controller, the authors formally verified over 60 properties and they identified 10 modelling errors and 2 requirement errors in the relatively mature models of the system. Similarly, Miller et al. [32] performed the formal analysis of Flight Critical Software (FCS) 5000, a new family of flight control systems developed by Rockwell Collins Inc that is widely used in business and regional jet aircraft. Moreover, the research team at NASA has worked on the formal verification of a flight critical software [32], software safety analysis of a flight guidance system [33] and safety analysis of software intensive systems [34]. More details about their research contributions in this direction can be found at [53, 54].

Similarly, Barnat et al. [55] integrated the DIVINE model checker and HiLiTE, a tool for requirements-based verification of aerospace system components developed and used by Honeywell for the formal verification of the avionics Simulink models. The authors used their proposed framework for formally analyzing the Voter Core that is a sub-system of the common avionics triplex sensor voter. Cao et al. [56] presented a framework for formally verifying the airborne software based on DO-333 and thus providing guidance in the integration of formal methods in the development and analysis of the software. The authors used their proposed methodology for the formal verification of Air Data Computer (ADC) software. Moreover, the authors claim that they identified 16 errors during the verification of ADC software [56]. Some more notable contributions regarding application of formal methods in avionics and aerospace are [57, 58, 59, 60, 61]. However, none of these contributions cater for posit.

3.2 Hardware and Software Implementations of Posits

Posits have been implemented as an alternative to the floating-point number format in hardware and software. Lehoczky et al. [7] presented the software and hardware implementations of posits. The authors used C# programming language to implement posits on the .NET platform. Moreover, they used Hastlayer, which is a tool for converting .NET models to a language that can be implemented on FPGA, to develop a hardware based implementation of posit. Similarly, the hardware architecture of unum adder/subtractor and multiplier has been designed and implemented in FPGAs. Moreover, a Verilog HDL generator has been constructed for performing these arithmetic operations [4, 5, 6, 62]. The software implementations of posits are also available in C# [7], C⁶, C++⁷ and Julia⁸ programming languages. Moreover, it has been experimentally shown that posits perform better than the fixed point number system, in terms of accuracy and memory utilization, for both training and inferences of deep convolutional neural networks [8, 9, 63]. However, none of the above-mentioned works (presented in Sections 3.1 and 3.2) provide the verification of unums, which is the main scope of the paper.

4 Results

4.1 Formalization of Posits

This section provides a higher-order-logic based formalization of posits using the HOL Light theorem prover. It mainly involves a conversion from a posit representation having a bit pattern to its corresponding real number and vice versa. Moreover, a conversion from a posit representation to its corresponding real number is mainly based on extracting the regime, exponent and fraction bits.

A posit format representation comprises four components, namely, sign, exponent, regime and fraction bits. It is sufficient to define a computing environment for posits, i.e., a posit configuration, using the total number of bits (nb) having an integer value of greater than or equal to 2 and the number of exponent bits (es) with an integer value of greater than or equal to 0 [64]. For example, for the case of $nb = 2$ and $es = 0$, the two bits comprising the posit are the sign and the regime bits, respectively. We model a valid posit configuration as the following HOL Light function:

Definition 1: *Valid Posit*

$$\vdash_{def} \forall (nb:int) (es:int). \text{is_valid_posit} (nb,es) = (nb \geq \widehat{2}) \wedge (es \geq \widehat{0})$$

The function `is_valid_posit` accepts a pair of integers (nb, es), describing the total number of bits nb and the number of exponent bits es and returns a valid posit configuration providing the constraints on these number of bits. We model a posit configuration using *new type definition* feature of HOL Light as follows:

```
let posformat_tpbij = new_type_definition "posit" ("mk_posit", "dest_posit")
(prove('?(pst:int#int). is_valid_posit pst', REWRITE_TAC [PROOF_TYPE]));;
```

where `posit` models a new type by providing its name and bijection alongwith a theorem asserting that bijection. The function `mk_posit` projects a pair of integers to a posit type and `dest_posit` maps a posit to a pair of integers. Next, to model a valid length of a posit, we first extract the elements of the pair (nb, es) in HOL Light as follows:

Definition 2: *Extraction of the Elements of Pair (nb, es)*

$$\begin{aligned} \vdash_{def} \forall (P:posit). \text{nb } P &= \text{FST } (\text{dest_posit } P) \\ \vdash_{def} \forall (P:posit). \text{es } P &= \text{SND } (\text{dest_posit } P) \\ \vdash_{def} \forall (P:posit). \text{nb}_n P &= \widetilde{(\text{nb } P)} \\ \vdash_{def} \forall (P:posit). \text{es}_n P &= \widetilde{(\text{es } P)} \end{aligned}$$

The function `nb` accepts a posit configuration P and extracts the total numbers of bits (nb) of a posit as an integer. Similarly, the function `es` extracts the numbers of exponent bits (es) of a posit as an integer. The functions `nbn` and `esn` use `num_of_int` to cast the integers nb and es to natural numbers.

Now, we model a valid posit length as follows:

Definition 3: *Valid Posit Length*

$$\vdash_{def} \forall (P:posit) (L:bool \text{ list}). \text{is_valid_posit_length } P L = (|L| = \text{nb}_n P)$$

⁶<https://github.com/libcg/bfp>

⁷<https://github.com/eruffaldi/cppPosit>

⁸<https://github.com/milankl/SoftPosit.jl>

The function `is_valid_posit_length` accepts a posit configuration P and a list $L:\text{bool list}$, i.e., a posit representation, capturing the bit values of a posit format, and returns a valid length of a posit.

Next, we model the two exception values for posits [64] in HOL Light as follows:

Definition 4: *Exceptions (Zero and Infinity)*

$$\begin{aligned} \vdash_{\text{def}} \forall(L:\text{bool list}). \text{zero_exception } L &= T \notin L \\ \vdash_{\text{def}} \forall(L:\text{bool list}). \text{inf_exception } L &= \bar{L} \wedge (T \notin L) \end{aligned}$$

The functions `zero_exception` and `inf_exception` present the exception values *zero* and $\pm\infty$, respectively. If all bits of a posit representation are 0, it represents an exception value *zero*. Similarly, if the first bit is 1 and the rest of the bits are 0, it provides an exception value $\pm\infty$.

Next, we model the seed value for the posit P , described in Section 2.2, as the following HOL Light function:

Definition 5: *Seed Value*

$$\vdash_{\text{def}} \forall P. \text{useed } P = 2^{2^{(\text{es}_n P)}}$$

Definition 6: *Minimum and Maximum Positive Value of Posits*

$$\begin{aligned} \vdash_{\text{def}} \forall(P:\text{posit}). \text{maxpos } P &= (\text{useed } P)^{(\text{nb}_n P - 2)} \\ \vdash_{\text{def}} \forall(P:\text{posit}). \text{minpos } P &= \frac{1}{\text{maxpos } P} \end{aligned}$$

The functions `maxpos` and `minpos` model the largest and smallest positive real numbers (values) expressible as a posit P , respectively [2].

Definition 7: *Check Extreme Values of Posit*

$$\vdash_{\text{def}} \forall L. \text{checkmax } L = (F \notin L)$$

The HOL Light function `checkmax` accepts a posit representation L , containing all bit values of a posit format, and returns a Boolean value *true* (T) if all bits are equal to 1 except the first (leading) bit, which can be either 0 or 1. For the case of first bit, i.e., sign bit equal to 0, it captures the largest positive value, whereas, it models the largest negative value for a sign bit 1.

Next, to calculate the scale factor of the regime bits, i.e., useed^k , we first formalize k (power of the variable *useed*) in HOL Light as the following recursive function:

Definition 8: *Value of k for Scaling Factor of the Regime Bits*

$$\begin{aligned} \vdash_{\text{def}} \forall(h:\text{bool}) (t:\text{bool list}). \text{value_of_k } [] &= 0 \wedge \\ \text{value_of_k } (h:t) &= \text{if } \bar{t} \text{ then } (\text{if } [C_1](\bar{t} = \overline{(t)}) \wedge [C_2](1 < |t|) \\ &\quad \text{then } ((\text{value_of_k } t) + 1) \text{ else } 0) \\ &\quad \text{else } (\text{if } [C_3](\bar{t} = \overline{(t)}) \wedge [C_4](1 < |t|) \\ &\quad \text{then } ((\text{value_of_k } t) - 1) \text{ else } -1) \end{aligned}$$

The function `value_of_k` accepts a posit representation $L:\text{bool list}$ and returns k (power of the variable *useed*) for the scaling factor of the regime bits. If the identical bits in a regime for a run are 1, then `value_of_k` is equal to *run length* - 1, otherwise it is equal to the negation of *run length*. Here, *run length* corresponds to the variable m presented in Section 2.2 and models the length of the regime bits.

Now, the scaling factor of the regime bits is formalized in HOL Light as follows:

Definition 9: *Scaling Factor of the Regime Bits*

$$\vdash_{\text{def}} \forall(P:\text{posit}) (L:\text{bool list}). \text{scaling}_r P L = (\text{useed } P)^{(\text{value_of_k } L)}$$

To convert a posit to its equivalent real number, we require the scaling factors of the exponential and the fraction bits, which further need an extraction of these bits from a given posit representation. Moreover, for both these extractions, we need to pick elements (bit values) from a posit representation, which is formalized as the following HOL Light function:

Definition 10: *Pick Elements From a List*

$$\begin{aligned} \vdash_{\text{def}} \forall(L:\text{bool list}) (l:\text{num}) (u:\text{num}). \\ \text{pick_elements } L \mid u &= \text{pick_elements_simp } L \mid ((u - l) + 1) \end{aligned}$$

The function `pick_elements` accepts a list L , a lower index l and an upper index u and returns a list containing the elements of the input list from l to u indices. It uses a recursive function `pick_elements_simp` to extract the required elements from a given list.

Next, we extract the exponent bits of a posit representation as follows:

Definition 11: *Extracting Exponent Bits*

$\vdash_{def} \forall(P:\text{posit}) (L:\text{bool list}). \text{exp_bits } P \ L =$
 $\text{if } [C_1]((\text{regime_length } L) + 1) < (\text{nb}_n \ P) \wedge [C_2](1 \leq (\text{eb}_n \ P)) \text{ then}$
 $\quad (\text{if } [C_3]((\text{regime_length } L) + 1 + (\text{eb}_n \ P)) \leq |L| \text{ then}$
 $\quad \quad \text{pick_elements } L ((\text{regime_length } L) + 1) ((\text{regime_length } L) + (\text{eb}_n \ P))$
 $\quad \quad \text{else pick_elements } L ((\text{regime_length } L) + 1) (|L| - 1))$
 $\text{else } []$

The function `exp_bits` accepts a posit configuration P and a posit representation L and returns the exponent bits of the posit. Here, the function `regime_length` provides the length of the regime bits.

Now, the scaling factor of the exponent bits is formalized as the following HOL Light function:

Definition 12: *Scaling Factor of Exponent Bits*

$\vdash_{def} \forall(P:\text{posit}) (L:\text{bool list}). \text{scaling}_e \ P \ L = \dot{2}^{BV_n (\text{exp_bits } P \ L) * (2^{\text{eb}_n \ P} - \text{exp_length } P \ L)}$

where the function `BV_n` provides a natural number representation of the bit values. There can be a maximum of es exponent bits depending on the bit left on the right side of the regime in a posit representation. Therefore, the function `scalinge` provides a scale factor of the exponent bits by incorporating the scenario, where the exponent bits `exp_bits` e are less than es . Moreover, the exponent bits scales from 0 to 2^{es} .

Next, we extract the fraction bits of a posit representation as follows:

Definition 13: *Extracting Fraction Bits*

$\vdash_{def} \forall(P:\text{posit}) (L:\text{bool list}). \text{fraction_bits } P \ L =$
 $\text{if } [C]((\text{regime_length } L) + (\text{exp_length } P \ L) + 1 < (\text{nb}_n \ P)) \text{ then}$
 $\quad \text{pick_elements } L ((\text{regime_length } L) + (\text{exp_length } P \ L) + 1) ((\text{nb}_n \ P) - 1)$
 $\text{else } []$

The function `fraction_bits` accepts a posit configuration P and a posit representation L and returns the fraction bits of a posit representation, if any.

Now, we formalize the scaling factor of the fraction bits as the following HOL Light function:

Definition 14: *Scaling Factor of Fraction Bits*

$\vdash_{def} \forall(P:\text{posit}) (L:\text{bool list}). \text{scaling}_f \ P \ L = \dot{1} + \frac{(BV_n (\text{fraction_bits } P \ L))}{\dot{2}^{(\text{fraction_length } P \ L)}}$

where the function `fraction_length` provides the length of the fraction bits. The fraction bits serves the same functionality as they do in the floating-point numbers.

Finally, we formalize a conversion of a posit to its equivalent real number as follows:

Definition 15: *Posit to Real Number Conversion*

$\vdash_{def} \forall(P:\text{posit}) (L:\text{bool list}). \text{posit_to_signed_real } P \ L =$
 $\text{if } [C_1] \text{ zero_exception } L \text{ then } \dot{0}$
 $\text{else (if } [C_2] (\text{checkmax } L) \text{ then (if } [C_3] \text{ sign_bit } L \text{ then } -\frac{\&1}{\text{maxpos } P}$
 $\quad \quad \text{else maxpos } P)$
 $\quad \text{else (if } [C_4] (\text{sign_bit } L) \text{ then } -(\text{scaling}_r \ P \ L' * \text{scaling}_e \ P \ L' * \text{scaling}_f \ P \ L')$
 $\quad \quad \text{else } (\text{scaling}_r \ P \ L) * (\text{scaling}_e \ P \ L) * (\text{scaling}_f \ P \ L)))$

where $L' = (\text{sign_bit } L)::(\text{two_complement } L)$.

The function `posit_to_signed_real` accepts a posit configuration P and a posit representation L and returns a real value corresponding to the given posit. The first conditional statement of the function `posit_to_signed_real` (Condition C_1) checks all bits of a posit representation using a function `add_zero_real`. For the case of all bits equal to zero, it returns a real number/value 0. Otherwise, the second conditional statement (Condition C_2) uses the function `checkmax` (Definition 7) to confirm if the given posit representation provides a largest positive or a largest negative real value for the sign bit values of 0 and 1, respectively. For the scenario when a posit

representation does not capture any largest positive or negative values, it returns the corresponding real number, which can be any positive or negative value depending on the sign bit of the given posit representation. For example, for a sign bit 1, it uses the notion of 2s complement to represent a negative real number.

Now, we provide the formalization of the conversion function from a real number to posit, which is mainly based on the notion of the exponential and the fractional rounding. The approach for converting a real number to its corresponding posit representation is quite similar to the method used for transforming any real number to float in floating-point arithmetic. For the case of the floating-point numbers, the first step involves checking for the exception values, which are only 0 and $\pm\infty$ for the case of posits. If the number does not represent any extreme values, it is divided by 2 or multiplied by 2 until it is in the interval $[1, 2)$, and thus determining the fraction bits for the corresponding floating-point number. For the case of posits, the given real number is, first, repeatedly divided or multiplied by *used* until it is in the interval $[1, \textit{used})$. Then, the non-negative exponent for the posit is determined by repeatedly divided by 2 until it is in the interval $[1, 2)$. The fraction always consists of a leading 1 bit to the left of the binary point and does not require handling any subnormal exception values that have a 0 bit to the left of the binary point.

Definition 16: *Negative Real Number (Sign Bit)*

$\vdash_{\text{def}} \forall(x:\text{real}). \text{sign_real } x = (x < \dot{0})$

The function `sign_real` accepts a real number x and returns true if it is negative.

First, we formalize the regime bits (regime field) for a posit corresponding to a real number in HOL Light as follows:

Definition 17: *Regime Bits (Regime Field)*

$\vdash_{\text{def}} \forall(x:\text{real}) (P:\text{posit}). \text{regime_bits } x P =$

if $[C_1](1 \leq x)$ then (get_regime_ones $\times P ((nb_n P) - 2)$)
 else (get_regime_zeros $\times P ((nb_n P) - 2)$)

$\vdash_{\text{def}} \forall(x:\text{real}) (P:\text{posit}) (n:\text{num}).$

get_regime_zeros $\times P 0 =$ if $[C_2](x = \dot{0})$ then $[F]$ else $[T] \wedge$

get_regime_zeros $\times p (SUC n) =$ if $[C_3](1 \leq x)$ then $(T::[])$ else $F::(\text{get_regime_zeros } (x * (\textit{used } P))) P n$

$\vdash_{\text{def}} \forall(x:\text{real}) (P:\text{posit}) (n:\text{num}). \text{get_regime_ones } \times P 0 = [T] \wedge$

get_regime_ones $\times P (SUC n) =$ if $[C_4](1 \leq x < \textit{used } P)$ then $T::(F::[])$ else $T::(\text{get_regime_ones } \frac{x}{\textit{used } P} P n)$

The function `regime_bits` accepts a real number x and a posit configuration P and provides the regime bits of a posit representation corresponding to the given real number x . It mainly asserts a condition on the value of x , i.e., if $1 \leq x$, then it uses the function `get_regime_ones` to obtain identical regime bits 1 terminated with a 0. If the condition on x is false, it uses the function `get_regime_zeros` to generate a sequence of 0s in the regime field terminated by a 1.

Generally, three distinct cases arise during a conversion of a real number to its corresponding posit, i.e., 1) the resultant posit consists of the regime, exponent and fraction bits; 2) it has only the regime and the exponent bits (no fraction bits); 3) it has only regime bits (no exponent and fraction bits). These cases depend on the fact if the notion of rounding is involved in the conversion or not, i.e., if a real number is exactly expressible as a posit using the number of bits mentioned in a posit configuration or it requires more bits, where it is rounded to a nearest valid posit representation. The notion of fractional rounding is used for the case when we need more fraction bits than the number of bits left after the regime and the exponent bits in a posit representation, to express the given real value as a posit. Whereas, the exponential rounding captures the scenario, where the number of exponent bits e in a posit representation is less than the value es given in a posit configuration P . Moreover, in both types of rounding, i.e., the fractional and the exponential, if the input real value is at the tie-breaking point, it is rounded to the nearest even posit having the last bit equal to zero. The fractional rounding is quite similar to that of the floating-point numbers. In fractional rounding, the tie-breaking point is the arithmetic mean of the two choices (lower and upper bounds) for the rounding and the posit is rounded to the nearest fraction. However, in the case of the exponential rounding, the tie-breaking point is the geometric mean of the two choices (lower and upper bounds). Moreover, in the case of a real number not exactly expressible as a posit, i.e., the exponent bits are truncated and the real value is either rounded above or rounded down to a valid posit representation given in Equations (1) and (2) [65]. For example, for two posits 32 and 128, the tie-breaking point is 64. Therefore, any value greater than 64 maps to 128 and a value less than 64 is rounded to 32.

$$e^+ = \left(\left\lfloor \frac{e}{2^t} \right\rfloor + 1 \right) 2^t \quad (1)$$

$$e^- = \left\lfloor \frac{e}{2^t} \right\rfloor 2^t \quad (2)$$

Similarly, the fractional rounding is based on the residue left after the computation of the fraction bits and the tie-breaking point, which is $\frac{1}{2}$. If the residue is less than $\frac{1}{2}$, the corresponding posit is rounded down, otherwise it is rounded up.

Now, we formalize the exponential rounding in HOL Light as follows:

Definition 18: *Exponential Rounding*

$\vdash_{def} \forall(x:\text{real}) (P:\text{posit}). \text{round}_e x P =$
 if $[C_1] \sim(\text{exp_residue} \times P) \wedge [C_2] M = N$ then $(\text{exp_posit_tie} \times P)$
 else (if $[C_3] \text{exp_residue} \times P \wedge [C_4] (M = N)) \vee [C_5] (M > N)$ then $(\text{exp_posit_up} \times P)$
 else $(\text{exp_posit_down} \times P))$

where $M = \text{te_rounded_bits} \times P$ and $N = 2^{\text{te_bits} \times P - 1}$.

The function round_e accepts a real number x and a posit configuration P and returns the exponential rounding of the corresponding posit representation. It mainly asserts a condition on the exponential residue exp_residue and the exponent bits te_rounded_bits to be truncated. The function exp_residue (Entry 1, Table 3) checks if there is any value (residue) left after extracting the regime and exponent bits and returns true if the given real number x is not completely expressible using the exponent and regime bits. Similarly, the function te_rounded_bits (Entry 3, Table 3) returns the value of the exponent bits, to be truncated. Moreover, the function te_bits (Entry 8, Table 3) provides the number of truncated bits, which is the difference between the actual bits left for the exponent field and es defined in a posit configuration. Therefore, the condition involving the functions exp_residue and te_rounded_bits states that if the value of the truncated bits is equal to the tie breaking point, i.e., 2^{t-1} and the given real number is totally divisible by the calculated exponent, i.e., if there is an equal distance between the upper and the lower bounded posit, it is rounded to the even posit formalized as exp_posit_tie (Entry 2, Table 3) in HOL Light.

If the condition involving functions exp_residue and te_rounded_bits is not satisfied, the function round_e (Definition 18) checks if there is some residue or the value of the truncated bit is greater than the tie breaking point. If it returns the Boolean value true then the corresponding posit is rounded up (exp_posit_up , Entry 5, Table 3) otherwise it is rounded down (exp_posit_down , Entry 4, Table 3). The HOL Light functions exp_posit_up accepts a posit configuration P and a real number x , and returns a list containing the regime bits and the binary representation of the value of the exponent of the rounded up posit modeled by exp_up (Equation (1), Entry 7, Table 3). Similarly, exp_posit_down provides a list by appending the regime bits with the binary representation of the value of the exponent of the rounded down posit modeled by exp_down (Equation (2), Entry 6, Table 3). The HOL Light function DIV returns the quotient of the division of two real number and it is used as a floor function. Similarly, the function exp_list (Entry 10, Table 3) provides a list capturing the binary representation of the exponent computed using a recursive function exp_bit_value . The function exp_bits_value (Entry 11, Table 3) checks a given real number x for different conditions, i.e., if the value of x is in the interval $[1, 2)$, the value of exponent is zero. Otherwise, the function is called recursively with the real value divided by 2 and adds 1 to the returned value. The corresponding binary representation is obtained by using the HOL Light function num_BV_f [66], which accepts two inputs, i.e., a natural number and the length of the output list in which it is expressed, and returns the corresponding list.

Next, we modeled the fractional rounding as, if the residual value after computing the regime exponent and fraction bits is greater than $\frac{1}{2}$, then it is rounded up and it is rounded down for a residual value less than $\frac{1}{2}$. Moreover, if the residual value is equal to $\frac{1}{2}$, then it is rounded to the nearest even posit. We model it using the function round_f as follows:

Definition 19: *Fractional Rounding*

$\vdash_{def} \forall(x:\text{real}) (P:\text{posit}). \text{round}_f x P =$
 if $[C_1] \text{residue}_f \times P = 0$
 then $(\text{regime_bits} \times P) ++ ((\text{exp_list} \times P) ++ (\text{set_fraction_list} \times P))$
 else (if $[C_2] \text{residue}_f \times P = \frac{1}{2}$ then $(\text{frac_posit_tie} \times P)$
 else (if $[C_3] \text{residue}_f \times P > \frac{1}{2}$ then $(\text{frac_posit_up} \times P)$
 else $(\text{frac_posit_down} \times P)))$

The function round_f accepts a real number x and a posit configuration P and returns the fractional rounding of the corresponding posit representation. It checks a condition on the residue value (Condition C_1), if it is equal

Table 3: Definition 18 (Exponential Rounding)

$\vdash_{def} \forall (P:\text{posit}) (x:\text{real}). \text{exp_residue } P \ x = \frac{x}{(\text{scale_factor_r } P \ [\text{sign_real } x] ++ (\text{regime_bits } \times P)) * (\&2 \text{ pow } BV_n (\text{exp_list } \times P))} - i \neq 0$
$\vdash_{def} \forall (P:\text{posit}) (x:\text{real}). \text{exp_posit_tie } P \ x = \text{if } [C] \text{ LAST } (\text{exp_posit_down } \times P) \text{ then } (\text{exp_posit_up } \times P) \text{ else } (\text{exp_posit_down } \times P)$
$\vdash_{def} \forall (P:\text{posit}) (x:\text{real}). \text{te_rounded_bits } P \ x = BV_n \text{ pick_elements } (\text{exp_list } \times P) (nb_n P - (\text{regime_bits } \times P + 1)) (eb_n P - 1)$
$\vdash_{def} \forall (P:\text{posit}) (x:\text{real}). \text{exp_posit_down } P \ x = (\text{regime_bits } \times P) ++ (\text{num_BV_f } (nb_n P - (\text{regime_bits } \times P + 1)) (\text{exp_down } \times P))$
$\vdash_{def} \forall (P:\text{posit}) (x:\text{real}). \text{exp_posit_up } P \ x = \text{if } [C_1] (\text{exp_up } \times P) < (2 \text{ EXP } eb_n P) \text{ then } (\text{regime_bits } \times P) ++ (\text{num_BV_f } (nb_n P - (\text{regime_bits } \times P + 1)) (\text{exp_up } \times P)) \\ \text{else if } [C_2] \text{ HD } (\text{regime_bits } \times P) \text{ then } (T :: (\text{regime_bits } \times P)) ++ (\text{num_BV_f } ((nb_n P) - (\text{regime_bits } \times P + 1) - 1) 0) \\ \text{else } (TL (\text{regime_bits } \times P)) ++ (\text{num_BV_f } ((nb_n P) - (\text{regime_bits } \times P + 1) + 1) 0)$
$\vdash_{def} \forall (P:\text{posit}) (x:\text{real}). \text{exp_down } P \ x = BV_n (\text{exp_list } \times P) \text{ DIV } (2^{\text{te_bits } \times P}) * (2^{\text{te_bits } \times P})$
$\vdash_{def} \forall (P:\text{posit}) (x:\text{real}). \text{exp_up } P \ x = (BV_n (\text{exp_list } \times P) \text{ DIV } (2^{\text{te_bits } \times P}) + 1) * (2^{\text{te_bits } \times P})$
$\vdash_{def} \forall (P:\text{posit}) (x:\text{real}). \text{te_bits } P \ x = (eb_n P) - ((nb_n P) - (\text{regime_bits } \times P + 1))$
$\vdash_{def} \forall (P:\text{posit}) (x:\text{real}). \text{cond}_e P \ x = (nb_n P) - (\text{regime_bits } \times P + 1) < (eb_n P)$
$\vdash_{def} \forall (P:\text{posit}) (x:\text{real}). \text{exp_list } P \ x = \text{num_BV_f } (eb_n P) \left(\text{exp_bits_value } \left(\frac{x}{\text{scale_factor_r } P \ [\text{sign_real } x] ++ (\text{regime_bits } \times P)} \right) ((2^{eb_n P}) - 1) \right)$
$\vdash_{def} \forall (x:\text{real}) (n:\text{num}). \text{exp_bits_value } \times 0 = 0 \wedge \text{exp_bits_value } \times (\text{SUC } n) = \text{if } [C] (i \leq x < \dot{2}) \text{ then } 0 \text{ else } \left(\text{exp_bits_value } \frac{x}{2} n \right) + 1$

Table 4: Definition 19 (Fractional Rounding)

$\vdash_{def} \forall (P:\text{posit}) (x:\text{real}). \text{frac_posit_tie } P \ x =$ $\text{if } [C] \text{ LAST (frac_posit_down } x \ P) \text{ then (frac_posit_up } x \ P) \text{ else (frac_posit_down } x \ P)$
$\vdash_{def} \forall (P:\text{posit}) (x:\text{real}). \text{frac_posit_down } P \ x =$ $(\text{regime_bits } x \ P) ++ ((\text{exp_list } x \ P) ++ (\text{set_fraction_list } x \ P))$
$\vdash_{def} \forall (P:\text{posit}) (x:\text{real}). \text{frac_posit_up } P \ x = \text{num_BV_f } ((\text{nb}_n \ P) - 1)$ $(\text{BV_n } ((\text{regime_bits } x \ P) ++ ((\text{exp_list } x \ P) ++ (\text{set_fraction_list } x \ P))) + 1)$
$\vdash_{def} \forall (P:\text{posit}) (x:\text{real}). \text{residue}_f P \ x = \text{fraction_residue1}$ $\left(\frac{x}{(\text{scale_factor_r } P \ [\text{sign_real } x] ++ (\text{regime_bits } x \ P)) * (\&2 \text{ pow } \text{BV_n } \text{exp_list } x \ P) - i} - i \right)$ $((\text{nb}_n \ P) - \text{regime_bits } x \ P - (\text{eb}_n \ P) - 1)$
$\vdash_{def} \forall (x:\text{real})(n:\text{num}). \text{fraction_residue1 } x \ 0 = x \wedge \text{fraction_residue1 } x \ (\text{SUC } n) =$ $\text{if } [C] \ \dot{1} \leq x * \dot{2} \text{ then } \text{fraction_residue1 } (\dot{2} * x - 1) \ n \text{ else } \text{fraction_residue1 } (x * \dot{2}) \ n$
$\vdash_{def} \forall (P:\text{posit}) (x:\text{real}). \text{set_fraction_list } P \ x = \text{fraction_list}$ $\left(\frac{x}{(\text{scale_factor_r } P \ [\text{sign_real } x] ++ (\text{regime_bits } x \ P)) * (\&2 \text{ pow } \text{BV_n } \text{exp_list } x \ P) - i} - i \right)$ $((\text{nb}_n \ P) - \text{regime_bits } x \ P - (\text{eb}_n \ P) - 1)$
$\vdash_{def} \forall (x:\text{real})(n:\text{num}). \text{fraction_list } x \ 0 = [] \wedge \text{fraction_list } x \ (\text{SUC } n) =$ $\text{if } [C] \ \dot{1} \leq x * \dot{2} \text{ then } T::(\text{fraction_list } (\dot{2} * x - \dot{1}) \ n) \text{ else } F::(\text{fraction_list } (x * \dot{2}) \ n)$

to zero, i.e., in the case of no fractional rounding, it appends the fraction bits to the regime and the exponent bits. In case the above condition is not satisfied, it again checks the value of the residue (Condition C_2), and if it is equal to $\frac{1}{2}$ then the corresponding posit is rounded to the nearest even posit modeled by `frac_posit_tie` (Entry 1, Table 4). Similarly, for a residue value greater than $\frac{1}{2}$ (Condition C_3), it is rounded above using the function `frac_posit_up` (Entry 3, Table 4), otherwise it is rounded down using the function `frac_posit_down` (Entry 2, Table 4). The function `frac_posit_up` provides a rounded up posit by appending the regime bits with the exponent and the fraction bits. Similarly, the function `frac_posit_down` provides a rounded down posit by appending the binary representation obtained from `set_fraction_list` (Entry 6, Table 4) to the regime and the exponent bits. The function `frac_posit_tie` checks the last bit of both fractionally rounded up and rounded down posits and returns the nearest even posit.

Definition 20: *Minimum Positive Real Number in a Posit Representation*

$$\vdash_{\text{def}} \forall(P:\text{posit}). \text{minpos_posit } P = \text{num_BV_f } ((\text{nb}_n P) - 1) (1)$$

The function `minpos_posit` accepts a posit configuration P and returns a posit representation of a minimum positive real number expressible in a posit.

Similarly, the function `maxpos_posit` captures a posit representation of a maximum (largest) positive real number expressible in a posit.

Definition 21: *Maximum Positive Real Number in a Posit Representation*

$$\vdash_{\text{def}} \forall(P:\text{posit}). \text{maxpos_posit } P = \text{num_BV_f } ((\text{nb}_n P) - 1) (2^{(\text{nb}_n P) - 1} - 1)$$

Finally, we use Definitions 16 - 21 to formalize the conversion of a real number to its corresponding posit representation as the following HOL Light function:

Definition 22: *Real to Posit Conversion*

$$\begin{aligned} &\vdash_{\text{def}} \forall(x:\text{real}) (P:\text{posit}). \text{posit}_{\text{real}} x P = \\ &\text{if } ([C_1] \ x = 0 \vee [C_2] \ |x| \geq \text{maxpos } P \vee [C_3] \ |x| \leq \text{minpos } P) \text{ then} \\ &\quad (\text{if } [C_4] \ x = 0 \text{ then } [F]++(\text{regime_bits} \times P) \\ &\quad \text{else} \\ &\quad \quad (\text{if } [C_5] \ |x| \geq \text{maxpos } P \text{ then } (\text{sign_real } x)::(\text{if } [C_6] \ \text{sign_real } x \text{ then two_complement } A \text{ else } A) \\ &\quad \quad \text{else } (\text{sign_real } x)::(\text{if } [C_7] \ \text{sign_real } x \text{ then two_complement } B \text{ else } B))) \\ &\text{else} \\ &\quad (\text{if } [C_8] \ x > 0 \text{ then } (\text{if } [C_9] \ \text{cond}_e \times P \text{ then } [F]++H \text{ else } [F]++K) \\ &\quad \text{else } (\text{if } [C_{10}] \ (\text{cond}_e \ |x| \ P) \text{ then } [T]++(\text{two_complement } M) \text{ else } [T]++(\text{two_complement } N))) \end{aligned}$$

where $A = \text{maxpos_posit } P$, $B = \text{minpos_posit } P$, $H = \text{round}_e \times P$,
 $K = \text{round}_f \times P$, $M = \text{round}_e \ |x| \ P$, $N = \text{round}_f \ |x| \ P$.

The function `positreal` accepts a real number x and a posit configuration P and returns its equivalent posit representation. The satisfaction of the first conditional statement (Conditions $C_1 - C_7$) provides the posit representations for zero, minimum and maximum posits. Whereas, the satisfaction of the second conditional statement (Conditions $C_8 - C_{10}$) provides other posits corresponding to any positive or negative real numbers using the notions of the exponential and the fractional rounding. Moreover, for the case of the negative real numbers, it provides the 2s complement of the corresponding posit representation.

We use our higher-order-logic based formalization of posits, a conversion from a posit to its corresponding real number and a real number to its corresponding posit, presented in this section, to formally verify various properties providing the correctness of these conversions and the scaling factors of various bits, such as regime, exponential and the fraction bits in Section 4.2 of the paper.

4.2 Formal Verification of Posits

In this section, we present the formal verification of various properties of posits regarding the conversions and the scaling factors of the regime, exponential and the fraction bits using HOL Light. The verification of these properties not only ensures the correctness of our formal definitions presented in Section 4.1 but they are also quite vital for performing various arithmetic operations based on posits.

We formally verify the non-negativity of the value of k (power of the variable *used*) as the following HOL Light theorem:

Theorem 1: *Positive Value of k*

$$\vdash_{thm} \forall (L: \text{bool list}). [A] \overline{(L)} \Rightarrow (\text{value_of_k } L) \geq 0$$

The only assumption of the above theorem asserts a condition that the first regime bit of a posit representation is 1. The proof process of Theorem 1 is based on Definition 8 along with some arithmetic reasoning.

We verify a constraint on the minimum number of bits of a posit as follows:

Theorem 2: *Minimum Number of Bits of a Posit*

$$\vdash_{thm} \forall (P: \text{posit}). 2 \leq \text{nb}_n P$$

The verification of the above theorem is mainly based on Definitions 1 and 2. Next, we verify an upper bound on length of the exponent bits as the following HOL Light theorem:

Theorem 3: *Upper Bound on Length of the Exponent Bits*

$$\vdash_{thm} \forall (P: \text{posit}) (L: \text{bool list}). [A] \text{ is_valid_posit_length } P L \Rightarrow \text{exp_length } P L \leq \text{eb}_n P$$

The above theorem states that the length of exponent bits e in a posit representation is always less than or equal to es given in a posit configuration. Moreover, the assumption ensures valid length of the considered posit. The verification of Theorem 3 is based on the function `exp_length`, Definitions 2 and 3 along with some arithmetic reasoning.

Theorem 4: *Upper Bound on Scaling Factor of the Exponent Bits*

$$\vdash_{thm} \forall (P: \text{posit}) (L: \text{bool list}). [A] \text{ is_valid_posit_length } P L \Rightarrow \text{scaling}_e P L \leq (2^{(2^{\text{nb}_n P} - 1)})$$

The above theorem provides an upper bound of $2^{2^{es}-1}$ on the scaling factor of the exponent bits. Similarly, we verify the lower bound on the scaling factor of the exponent bits as the following HOL Light theorem:

Theorem 5: *Lower Bound on Scaling Factor of the Exponent Bits*

$$\vdash_{thm} \forall (P: \text{posit}) (L: \text{bool list}). [A] \text{ is_valid_posit_length } P L \Rightarrow 0 < \text{scaling}_e P L$$

The proof process of the above theorem is mainly based on Definition 12 along with some arithmetic reasoning. Next, we verify the upper bound of the scaling factor of the fraction bits as the following HOL Light theorem:

Theorem 6: *Upper Bound on Scaling Factor of the Fraction Bits*

$$\vdash_{thm} \forall (P: \text{posit}) (L: \text{bool list}). [A] \text{ is_valid_posit_length } P L \Rightarrow \text{scaling}_f P L < 2$$

The above theorem provides an upper bound of 2 on the scaling factor of the fraction bits, which is the same as that for the floating point numbers. The verification of Theorem 6 is mainly based on Definition 14 along with some arithmetic reasoning.

Next, we verify the scaling factor of the fraction bits in exponential rounding as follows:

Theorem 7: *Scaling Factor of the Fraction Bits in Exponential Rounding*

$$\vdash_{thm} \forall (P: \text{posit}) (L: \text{bool list}). [A_1] (\text{ is_valid_posit_length } P L) \wedge [A_2] \text{ exp_length } P L < \text{eb}_n P \\ \Rightarrow \text{scaling}_f P L = 1$$

Assumption A_1 ensures a valid length of the given posit. Similarly, Assumption A_2 asserts a condition that the number of exponent bits e in a posit representation is less than es bits defined in a posit configuration. Finally, the conclusion of Theorem 7 ensures that the scaling factor of the fraction bits is equal to 1, which indicates that no fraction bits are left after the exponential rounding.

Theorem 8: *Upper Bound on the Value of k*

$$\vdash_{thm} \forall (L: \text{bool list}). [A_1] \sim(\text{checkmax } L) \wedge [A_2] (2 \leq |L|) \Rightarrow \text{value_of_k } L < (|L| - 2)$$

The above theorem provides an upper bound of the value of k (power of the variable *used*). Assumption A_1 ensures that the given posit representation does not include any extreme (largest positive or negative) possible values. Similarly, Assumption A_2 provides a lower bound on the length of a posit representation. The verification of Theorem 8 is based mainly based on Definition 7 and Theorem 1.

Next, we verify an upper bound on the scaling factor of the regime bits as the following HOL Light theorem:

Theorem 9: *Upper Bound on Scaling Factor of the Regime Bits*

$$\vdash_{thm} \forall (P: \text{posit}) (L: \text{bool list}). [A_1] \text{ is_valid_posit_length } P L \wedge [A_2] \text{ is_valid_posit } (\text{dest_posit } P) \\ \Rightarrow \text{scaling}_r P L \leq \text{maxpos } P$$

Assumptions $A_1 - A_2$ provide a valid posit length and a valid posit configuration. The conclusion of Theorem 9 ensures that the scale factor is upper bounded by a maximum positive value of a posit. The verification of the above theorem is mainly based on Definitions 1, 3, 6 and 9 along with some arithmetic reasoning.

Theorem 10: *Lower Bound on Scaling Factor of the Regime Bits*

$$\vdash_{thm} \forall (P:\text{posit}) (L:\text{bool list}). [A] \text{ is_valid_posit_length } P L \Rightarrow \dot{0} < \text{scaling}_r P L$$

The above theorem provides a lower bound on the scaling factor of the regime bits.

Theorem 11: *Total Length of a Posit*

$$\vdash_{thm} \forall (P:\text{posit}) (L:\text{bool list}). [A] \text{ is_valid_posit_length } P L \\ \Rightarrow (\text{exp_length } P L) + (\text{regime_length } L) + (\text{fraction_length } P L) + 1 = \text{nb}_n P$$

The above theorem ensures that the total number of bits in a posit configuration is equal to the summation of the individual bits, such as regime, exponential, fraction and sign bits, extracted from a posit representation.

Next, we verify an upper bound on a real number/value obtained from its equivalent posit as the following HOL Light theorem:

Theorem 12: *Upper Bound on a Real Value Obtained from its Equivalent Posit*

$$\vdash_{thm} \forall (L:\text{bool list}) (P:\text{posit}). \\ [A_1] \text{ is_valid_posit } (\text{dest_posit } P) \wedge [A_2] \text{ is_valid_posit_length } P L \wedge [A_3] \sim (\text{zero_exception } L) \wedge \\ [A_4] \sim (\text{inf_exception } L) \\ \Rightarrow \text{posit_to_signed_real } P L \leq \text{maxpos } P$$

Assumptions $A_1 - A_2$ are same as that of Theorem 9. Assumptions $A_3 - A_4$ ensure that the given posit representation does not contain any exception values. Finally, the conclusion of the above theorem asserts that the real number obtained from a given posit is upper bounded by maxpos (Definition 6).

Next, we verify some important properties about the conversion of a real number to its equivalent posit (encoding) and a conversion of a posit to its equivalent real number (decoding) for the case of 0, largest positive and largest negative values. Given a real number $x = 0$, its encoding and decoding in subsequent steps is always equal to x , i.e., 0, which is verified as the following HOL Light theorem:

Theorem 13: *Encoding and Decoding of Zero*

$$\vdash_{thm} \forall (P:\text{posit}) (L:\text{bool list}) (x:\text{real}). \\ [A_1] \text{ is_valid_posit } (\text{dest_posit } P) \wedge [A_2] \text{ is_valid_posit_length } P (\text{posit_real } x P) \wedge [A_3] (x = \dot{0}) \\ \Rightarrow \text{real_posit } P (\text{posit_real } x P) = x$$

Assumption A_1 ensures that a posit obtained as a result of encoding x provides a valid posit configuration. Similarly, Assumption A_2 provides a valid length of a posit representation. Assumption A_3 asserts the condition $x = 0$. Finally, the conclusion of Theorem 13 ensures that the application of the encoding and decoding in subsequent steps to a real value 0 is always equal to the same number, i.e., 0. The verification of Theorem 13 is mainly based on the properties of the function `zero_exception` along with some real arithmetic reasoning.

Similarly, for a given real number capturing the largest possible value (maxpos) of a posit representation, the application of encoding and decoding in subsequent steps always results into maxpos and it is verified in HOL Light as the following theorem:

Theorem 14: *Encoding and Decoding of maxpos*

$$\vdash_{thm} \forall (P:\text{posit}) (L:\text{bool list}) (x:\text{real}). [A_1] \text{ is_valid_posit } (\text{dest_posit } P) \wedge \\ [A_2] \text{ is_valid_posit_length } P (\text{posit_real } x P) \wedge [A_3] \sim (\text{zero_exception } (\text{posit_real } x P)) \wedge \\ [A_4] x = \text{maxpos } P \Rightarrow \text{real_posit } P (\text{posit_real } x P) = x$$

Assumptions $A_1 - A_2$ are the same as that of Theorem 16. Assumption A_3 ensures that the given posit representation does not contain any zero exception value. Assumption A_4 provides the condition $x = \text{maxpos } P$. Finally, the conclusion ensures that the encoding and decoding of a real number, having value of maxpos , is always equal to maxpos .

Next, we verify an important property, which ensures that every real number greater than the largest negative value of its corresponding posit representation is mapped to $-\text{minpos}$.

Theorem 15: *For Values Greater than the Largest Negative Number*

$\vdash_{thm} \forall(P:\text{posit}) (L:\text{bool list}) (x:\text{real}).$

$[A_1] \text{ is_valid_posit } (\text{dest_posit } P) \wedge [A_2] \text{ is_valid_posit_length } P (\text{posit_real} \times P) \wedge$
 $[A_3] \sim(\text{zero_exception } (\text{posit_real} \times P)) \wedge [A_4] \sim(\text{inf_exception } (\text{posit_real} \times P)) \wedge$
 $[A_5] (x < \dot{0}) \wedge [A_6] (x \geq -\text{minpos } P)$
 $\Rightarrow \text{real_posit } P (\text{posit_real} \times P) = -\text{minpos } P$

Assumptions $A_1 - A_3$ are same as that of Theorem 14. Similarly, Assumption A_4 ensures that the given posit representation does not contain any infinity exception value. Assumptions $A_5 - A_6$ provide the conditions that the given real number is greater than the largest negative number of its corresponding posit representation. Finally, the conclusion ensures that the encoding and decoding of a real number, having value greater than $-\text{minpos}$, is always equal to $-\text{minpos}$.

Now, we verify the scaling factor of the exponent and the fraction bits during the conversion of a real number to its corresponding posit in a scenario of no fractional and exponential rounding. The scaling factor of the exponent bits is verified as the following HOL Light theorem:

Theorem 16: *Scaling factor of the Exponent Bits of the Encoded Posit*

$\vdash_{thm} \forall(P:\text{posit}) (L:\text{bool list}) (x:\text{real}).$

$[A_1] \text{ is_valid_posit } (\text{dest_posit } P) \wedge [A_2] \text{ is_valid_posit_length } P (\text{posit_real} \times P) \wedge$
 $[A_3] \sim(\text{zero_exception } (\text{posit_real} \times P)) \wedge [A_4] \sim(\text{inf_exception } (\text{posit_real} \times P)) \wedge$
 $[A_5] \sim(\text{checkmax } (\text{posit_real} \times P)) \wedge [A_6] (x > \dot{0}) \wedge [A_7] (\text{minpos } P < |x| < \text{maxpos } P) \wedge$
 $[A_8] \sim(\text{cond}_e \times P) \wedge [A_9] \text{residue}_f \times P = \dot{0} \wedge$
 $[A_{10}] \text{regime_length } ([F]++((\text{regime_bits} \times P)++((\text{exp_list} \times P)++$
 $(\text{set_fraction_list } \times P)))) = |\text{regime_bits} \times P|$
 $\Rightarrow \text{scaling}_e P (\text{posit_real} \times P) = 2^{\text{BV}_n (\text{exp_list} \times P)}$

Assumptions $A_1 - A_4$ are same as that of Theorem 15. Assumption A_5 ensures that a conversion of a real number to its posit does not contain any extreme (largest positive and negative) values. Assumptions $A_6 - A_7$ provide the bounds on the real number x . Similarly, Assumptions $A_8 - A_9$ assert no exponential and fractional rounding conditions. Assumption A_{10} provides an equality of the length of the regime bits extracted after appending all fields of a posit and the length computed over regime bits only. Finally, the conclusion provides the scaling factor of the exponent bits. The proof process of Theorem 16 is based on the corresponding definitions, properties of lists along with some arithmetic reasoning.

Next, we verify the scaling factor of the fraction bits as follows:

Theorem 17: *Scaling factor of the Fraction Bits of the Encoded Posit*

$\vdash_{thm} \forall(P:\text{posit}) (L:\text{bool list}) (x:\text{real}).$

$[A_1] \text{ is_valid_posit } (\text{dest_posit } P) \wedge [A_2] \text{ is_valid_posit_length } P (\text{posit_real} \times P) \wedge$
 $[A_3] \sim(\text{zero_exception } (\text{posit_real} \times P)) \wedge [A_4] \sim(\text{inf_exception } (\text{posit_real} \times P)) \wedge$
 $[A_5] \sim(\text{checkmax } (\text{posit_real} \times P)) \wedge [A_6] (x > \dot{0}) \wedge$
 $[A_7] (\text{minpos } P < |x| < \text{maxpos } P) \wedge [A_8] \sim(\text{cond}_e \times P) \wedge [A_9] \text{residue}_f \times P = \dot{0} \wedge$
 $[A_{10}] \text{regime_length } ([F]++((\text{regime_bits} \times P)++((\text{exp_list} \times P)++$
 $(\text{set_fraction_list } \times P)))) = |\text{regime_bits} \times P|$
 $\Rightarrow \text{scaling}_f P (\text{posit_real} \times P) = \dot{1} + \frac{\text{BV}_n (\text{set_fraction_list} \times P)}{(\dot{2} - \text{set_fraction_list} \times P)}$

Assumptions $A_1 - A_{10}$ are same as that of Theorem 16. Moreover, the conclusion provides the scaling factor of the fraction bits. The verification of Theorem 17 is based on the properties regarding the length of the fraction bits, the function `pick_element` along with some real arithmetic reasoning.

Now, we verify the scaling factors of the exponent and the fraction bits during the conversion of a real number to its corresponding posit in a fractional rounding scenario, i.e., the posit is rounded below. The scaling factor of the exponent bits is verified as follows:

Theorem 18: *Scaling Factor of the Exponent Bits*

$\vdash_{thm} \forall(P:\text{posit}) (L:\text{bool list}) (x:\text{real}).$

$[A_1] (\text{ is_valid_posit } (\text{dest_posit } P)) \wedge [A_2] \text{ is_valid_posit_length } P (\text{posit_real} \times P) \wedge$
 $[A_3] \sim(\text{zero_exception } (\text{posit_real} \times P)) \wedge [A_4] \sim(\text{inf_exception } (\text{posit_real} \times P)) \wedge$
 $[A_5] (x > \dot{0}) \wedge [A_6] (\text{minpos } P < |x| < \text{maxpos } P) \wedge \wedge$
 $[A_7] \sim(\text{checkmax } (\text{posit_real} \times P)) \wedge [A_8] \sim(\text{cond}_e \times P) \wedge$

$$\begin{aligned}
& [A_9] \text{residue}_f \times P \neq \dot{0} \wedge [A_{10}] \text{residue}_f \times P < \frac{1}{2} \wedge \\
& [A_{11}] \text{regime_length } ([F]++((\text{regime_bits} \times P)++((\text{exp_list} \times P)++ \\
& \quad (\text{set_fraction_list} \times P)))) = |\text{regime_bits} \times P| \\
& \Rightarrow \text{scaling}_e P ([F]++(\text{frac_posit_down} \times P)) = \dot{2}^{\text{BV}_n (\text{exp_list} \times P)}
\end{aligned}$$

Assumptions $A_1 - A_8$ are same as that of Theorem 17. Assumption A_9 asserts the existence condition of the fractional rounding. Similarly, Assumption A_{10} ensures that the residue value after extracting the fraction bits is less than $\frac{1}{2}$. Assumption A_{11} provides an equality of the length of the regime bits extracted after appending all fields of a posit and the length computed over regime bits only. Finally, the conclusion provides the scaling factor of the exponent bits in the case of the fractional rounding. The proof process of Theorem 21 is based on the properties regarding the length of the exponent and the fraction bits and function residue_f along with some arithmetic reasoning.

Next, we verify the scaling factor of the fraction bits as the following HOL Light theorem:

Theorem 19: *Scaling Factor of the Fraction Bits*

$$\begin{aligned}
& \vdash_{thm} \forall (P:\text{posit}) (L:\text{bool list}) (x:\text{real}). \\
& [A_1] \text{is_valid_posit } (\text{dest_posit } P) \wedge [A_2] \text{is_valid_posit_length } P (\text{posit}_{\text{real}} \times P) \wedge \\
& [A_3] \sim(\text{zero_exception } (\text{posit}_{\text{real}} \times P)) \wedge [A_4] \sim(\text{inf_exception } (\text{posit}_{\text{real}} \times P)) \wedge \\
& [A_5] (x > \dot{0}) \wedge [A_6] (\text{minpos } P < |x| < \text{maxpos } P) \wedge \\
& [A_7] \sim(\text{checkmax } (\text{posit}_{\text{real}} \times P)) \wedge [A_8] \sim(\text{cond}_e \times P) \wedge \\
& [A_9] \text{residue}_f \times P \neq \dot{0} \wedge [A_{10}] \text{residue}_f \times P < \frac{1}{2} \wedge \\
& [A_{11}] \text{regime_length } ([F]++((\text{regime_bits} \times P)++((\text{exp_list} \times P)++ \\
& \quad (\text{set_fraction_list} \times P)))) = |\text{regime_bits} \times P| \\
& \Rightarrow \text{scaling}_f P ([F]++(\text{frac_posit_down} \times P)) = \dot{1} + \frac{\text{BV}_n (\text{set_fraction_list} \times P)}{(\dot{2}^{|\text{set_fraction_list} \times P|})}
\end{aligned}$$

Assumptions $A_1 - A_{11}$ are same as that of Theorem 18. Moreover, the conclusion provides the scaling factor of the fraction bits in the case of the fractional rounding. The proof process of Theorem 19 is based on the properties regarding the length of the exponent and the fraction bits and function residue_f along with some arithmetic reasoning.

Theorem 20: *Completely Encoded fraction*

$$\begin{aligned}
& \vdash_{thm} \forall (n:\text{num}) (x:\text{real}). [A_1] \text{fraction_residue1} \times n = \&0 \\
& \Rightarrow x = \frac{(\text{BV}_n (\text{fraction_list} \times n))}{(\dot{2}^{|\text{fraction_list} \times n|})}
\end{aligned}$$

Theorem 20 verifies the value of the real number in the case of no fractional rounding.

Next, we verify some important properties about the conversion of a real number to its equivalent posit (encoding) and a conversion of a posit to its equivalent real number (decoding) for the scenario of no exponential or fractional rounding during the process of decoding, i.e., the decoding of an exactly representable real number. Given an exactly representable positive real number x , its encoding and decoding in subsequent steps is always equal to x , which is verified as the following HOL Light theorem:

Theorem 21: *Decoding of Encoded Positive Real Numbers*

$$\begin{aligned}
& \vdash_{thm} \forall (P:\text{posit}) (L:\text{bool list}) (x:\text{real}). \\
& [A_1] \text{is_valid_posit } (\text{dest_posit } P) \wedge [A_2] \text{is_valid_posit_length } P (\text{posit}_{\text{real}} \times P) \wedge \\
& [A_3] \sim(\text{zero_exception } (\text{posit}_{\text{real}} \times P)) \wedge [A_4] \sim(\text{inf_exception } (\text{posit}_{\text{real}} \times P)) \wedge \\
& [A_5] (x > \dot{0}) \wedge [A_6] (\text{minpos } P < |x| < \text{maxpos } P) \wedge \\
& [A_7] \sim(\text{checkmax } (\text{posit}_{\text{real}} \times P)) \wedge [A_8] \sim(\text{cond}_e \times P) \wedge \\
& [A_9] \text{residue}_f \times P = \dot{0} \wedge \\
& [A_{10}] \text{regime_length } ([F]++((\text{regime_bits} \times P)++((\text{exp_list} \times P)++ \\
& \quad (\text{set_fraction_list} \times P)))) = |\text{regime_bits} \times P| \wedge \\
& [A_{11}] \text{value_of_k } (\text{posit}_{\text{real}} \times P) = \text{value_of_k } ([\text{sign_real } x]++(\text{regime_bits} \times P)) \\
& \Rightarrow \text{real}_{\text{posit}} P (\text{posit}_{\text{real}} \times P) = x
\end{aligned}$$

Assumptions $A_1 - A_{10}$ are same as that of Theorem 17. Assumption A_{11} ensures that the value of k calculated over a posit representation is equal to its value calculated over the regime bits only. Moreover, the conclusion ensures that the application of the encoding and decoding on a positive real number, in subsequent steps, results into a same number. The verification of Theorem 21 is based on the properties regarding the length and bounds of the exponent, fraction and regime bits along with some real arithmetic reasoning.

Similarly, we verify the encoding and decoding of a negative real number, with no exponential or fractional rounding involved in decoding, i.e., an exactly representable negative real number, as the following HOL Light theorem:

Theorem 22: *Decoding of Encoded Negative Real Numbers*

$$\begin{aligned} & \vdash_{thm} \forall (P:\text{posit}) (L:\text{bool list}) (x:\text{real}). \\ & [A_1] \text{ is_valid_posit } (\text{dest_posit } P) \wedge [A_2] \text{ is_valid_posit_length } P (\text{posit_real } \times P) \wedge \\ & [A_3] \sim(\text{zero_exception } (\text{posit_real } \times P)) \wedge [A_4] \sim(\text{inf_exception } (\text{posit_real } \times P)) \wedge \\ & [A_5] x \neq 0 \wedge [A_6] (\text{minpos } P < |x| < \text{maxpos } P) \wedge \\ & [A_7] \sim(\text{checkmax } (\text{posit_real } \times P)) \wedge [A_8] \sim(\text{cond}_e |x| P) \wedge \\ & [A_9] (\text{residue}_f |x| P = 0) \wedge [A_{10}] (\forall A. \text{two_complement } (\text{two_complement } A) = A) \wedge \\ & [A_{11}] (\text{regime_length } ((\text{sign_real } x) :: ((\text{regime_bits } |x| P) ++ ((\text{exp_list } |x| P) ++ \\ & \quad (\text{set_fraction_list } (|x| P)))) = |\text{regime_bits } (|x| P)|) \wedge \\ & [A_{12}] (\text{value_of_k } ((\text{sign_real } x) :: ((\text{regime_bits } (|x| P) ++ ((\text{exp_list } (|x| P) ++ \\ & \quad (\text{set_fraction_list } (|x| P)))) = \\ & \quad \text{value_of_k } ([\text{sign_real } (|x|)] ++ (\text{regime_bits } (|x| P)))) \\ & \Rightarrow \text{real_posit } P (\text{posit_real } \times P) = x \end{aligned}$$

Table 5: Verification Details for Each Theorem

Formalized Theorems	Proof Lines	Man-Hours	Complexity of Proofs
Theorem 1 (Positive Value of k)	8	1	Easy
Theorem 2 (Minimum Number of Bits of a Posit)	20	1	Easy
Theorem 3 (Upper Bound on Length of the Exponent Bits)	15	1	Easy
Theorem 4 (Upper Bound on Scaling Factor of the Exponent Bits)	40	7	Easy
Theorem 5 (Lower Bound on Scaling Factor of the Exponent Bits)	3	0.5	Easy
Theorem 6 (Upper Bound on Scaling Factor of the Fraction Bits)	20	2	Easy
Theorem 7 (Scaling Factor of the Fraction Bits in Exponential Rounding)	26	1	Easy
Theorem 8 (Upper Bound on the Value of k)	69	19	Medium
Theorem 9 (Upper Bound on Scaling Factor of the Regime Bits)	64	17	Medium
Theorem 10 (Lower Bound on Scaling Factor of the Regime Bits)	7	0.5	Easy
Theorem 11 (Total Length of a Posit)	39	2	Easy
Theorem 12 (Upper Bound on a Real Value Obtained from its Equivalent Posit)	300	84	Hard
Theorem 13 (Encoding and Decoding of Zero)	10	2	Easy
Theorem 14 (Encoding and Decoding of maxpos)	25	2	Easy
Theorem 15 (For Values Greater than the Largest Negative Number)	96	27	Hard
Theorem 16 Scaling factor of the Exponent Bits of the Encoded Posit)	70	19	Medium
Theorem 17 (Scaling factor of the Fraction Bits of the Encoded Posit)	82	21	Medium
Theorem 18 (Scaling Factor of the Exponent Bits of the Encoded Rounded Posit)	100	40	Hard
Theorem 19 (Scaling Factor of the Fraction Bits of the Encoded Rounded Posit)	126	40	Hard
Theorem 20 (Completely Encoded fraction)	45	36	Hard
Theorem 21 (Decoding of Encoded Positive Real Numbers)	80	46	Hard
Theorem 22 (Decoding of Encoded Negative Real Numbers)	450	70	Hard

Assumptions $A_1 - A_9$ are same as that of Theorem 21. Assumption A_{10} ensures that taking the 2s complement of a number is same as the given number. Assumption A_{11} is same as that of Theorem 18. Assumption A_{12} ensures that the value of k extracted over a whole posit representation is equal to the one extracted over the regime bits only. Finally, the conclusion ensures that the application of encoding and decoding, on a negative

real number, in subsequent steps results into a same number. The proof process of Theorem 22 is very similar to the verification of Theorem 21.

The formal verification of the above theorems ensures the correctness of our formalization of posits, presented in Section 4.1, i.e., the formal model of posits and conversion from a real number to posit and vice versa, and its various parameters, such as the scaling factors of the regime, exponential and the fraction bits using HOL Light. Moreover, these formalization results, presented in Sections 4.1 and 4.2, can be further used for the verification of various arithmetic operations, such as addition, subtraction, multiplication and division operators.

5 Discussion

The distinguished feature of the proposed formalization is that all the proved theorems are of generic nature, i.e., all the functions and variables are universally quantified and hence, can be specialized based on the requirement of the unum arithmetics, like the encoding or decoding of any particular unums. Moreover, the inherent correctness of the theorem proving approach ensures that all the necessary assumptions are explicitly present with the respective theorem. The effort spent in verification of each theorem is represented in the form of proof lines and the man-hours as shown in Table 5. The man-hours are calculated based on two factors. The first factor includes the number of lines of HOL Light code per hour by a person with an average expertise and the second factor is the complexity of the proof. Moreover, there is no direct method to access the complexity of the proof. We often consider three major factors to estimate it. 1) The complexity of the mathematical results that are used in the analysis or proof of a theorem. For example, a proof involving integrals will be more complex than that involving matrices/vectors that are easy to handle.; 2) The expertise of a researcher regarding a particular proof goal.; 3) How many lemmas that are directly used in a proof of a theorem and are not available in a library. More lemmas to prove, make a formal proof of a theorem more complex and vice versa. Therefore, lines number of the proof script do not have a direct relationship with the man-hours. For instance, the man-hours for the verification of Theorems 18 and 19 are identical, while the proof lines for the former are less than that for the later. Moreover, there are few inherent limitations of our proposed higher-order-logic theorem proving approach. 1) Our proposed approach involves a lot of human interaction due to the undecidable nature of higher-order-logic, i.e., the user is involved in the process of formal verification along with the machine.; 2). Sometimes there is a significant gap between the traditional mathematical proof and its formal proof. Therefore, we need to identify the additional steps at our own that are required for developing a complete formal proof. ; 3) We have identified all formally verified properties, presented in Section 4.2, at our own to ensure the correctness of the formalization of posit provided in Section 4.1 of the paper. Moreover, to the best of our knowledge, these properties are not mentioned in the literature.

6 Conclusions

The unum format is a number representation format that provides an improved memory bandwidth and the power efficiency as compared to the floating-point numbers [1]. As a first step towards the verification of the unum arithmetic, this paper provides a formalization of posit, which is a Type III unums. In particular, we provide a conversion of a real number to its corresponding posit representation and a posit representation to its corresponding real number. We also verify some important properties regarding scaling factors of its regime, exponential and fraction bits using HOL Light that are widely used to perform various arithmetic operations involving posits. In future, we plan to verify different arithmetic operations [65], such as addition, subtraction, multiplication, exponential and division for posits. We also plan to utilize our proposed formalization to formally verify the computations of the transcendental functions, such as sine, cosine and exponential functions [1]. Another future direction is to make a comparison of the formal libraries of the floating-point number and posits.

References

- [1] John L Gustafson. *The End of Error: Unum Computing*. CRC Press, 2017.
- [2] John L Gustafson and Isaac T Yonemoto. Beating floating point at its own game: Posit Arithmetic. *Supercomputing Frontiers and Innovations*, 4(2):71–86, 2017.
- [3] Anwar A Esmael, Sa’ed Abed, Bassam J Mohd, Abbas A Fairouz, et al. POSIT vs. Floating Point in Implementing IIR Notch Filter by Enhancing Radix-4 Modified Booth Multiplier. *Electronics*, 11(1):163, 2022.

- [4] Manish Kumar Jaiswal and Hayden K-H So. Architecture Generator for Type-3 Unum Posit Adder/Subtractor. In *Circuits and Systems*, pages 1–5. IEEE, 2018.
- [5] Artur Podobas and Satoshi Matsuoka. Hardware Implementation of POSITs and their Application in FPGAs. In *Parallel and Distributed Processing*, pages 138–145. IEEE, 2018.
- [6] Manish K Jaiswal and Hayden K-H So. Universal Number Posit Arithmetic Generator on FPGA. In *Design, Automation & Test in Europe*, pages 1159–1162. IEEE, 2018.
- [7] Zoltán Lehóczky, András Retzler, Richárd Tóth, Álmos Szabó, Benedek Farkas, and Krisztián Somogyi. High-level. NET Software Implementations of Unum Type I and Posit with Simultaneous FPGA Implementation using Hastlayer. In *Next Generation Arithmetic*, pages 1–7, 2018.
- [8] Hamed F Langroudi, Vedant Karia, John L Gustafson, and Dhireesha Kudithipudi. Adaptive Posit: Parameter Aware Numerical Format for Deep Learning Inference on the Edge. In *Computer Vision and Pattern Recognition*, pages 726–727, 2020.
- [9] Raul Murillo, Alberto A Del Barrio, and Guillermo Botella. Deep PeNSieve: A Deep Learning Framework based on the Posit Number System. *Digital Signal Processing*, page 102762, 2020.
- [10] Bob Bentley. Validating the Intel Pentium 4 Microprocessor. In *Design Automation*, pages 244–248, 2001.
- [11] Edmund M Clarke and Jeannette M Wing. Formal Methods: State of the Art and Future Directions. *ACM Computing Surveys*, 28(4):626–643, 1996.
- [12] John Harrison. *Handbook of Practical Logic and Automated Reasoning*. Cambridge University Press, 2009.
- [13] John Harrison. Floating Point Verification in HOL Light: the Exponential Function. In *Algebraic Methodology and Software Technology*, pages 246–260. Springer, 1997.
- [14] John Harrison. A Machine-checked Theory of Floating Point Arithmetic. In *Theorem Proving in Higher Order Logics*, pages 113–130. Springer, 1999.
- [15] Marc Daumas, Laurence Rideau, and Laurent Théry. A Generic Library for Floating-point Numbers and its Application to Exact Computing. In *Theorem Proving in Higher Order Logics*, volume 2152 of *LNCS*, pages 169–184. Springer, 2001.
- [16] Sylvie Boldo and Jean-Christophe Filliâtre. Formal verification of floating-point programs. In *Symposium on Computer Arithmetic*, pages 187–194. IEEE, 2007.
- [17] J Strother Moore, Thomas W Lynch, and Matt Kaufmann. A Mechanically Checked Proof of the AMD5K86TM Floating-point Division Program. *IEEE Transactions on Computers*, (9):913–926, 1998.
- [18] David Russinoff. A Mechanically Checked Proof of IEEE Compliance of a Register-transfer-level Specification of the AMD-K7 Floating-point Multiplication, Division, and Square Root Instructions. *LMS Journal of Computation and Mathematics*, 1:148–200, 1998.
- [19] Paul S Miner. Defining the IEEE-854 Floating-Point Standard in PVS. 1995.
- [20] Christoph Berg and Christian Jacobi. Formal Verification of the VAMP Floating-point Unit. In *Correct Hardware Design and Verification Methods*, volume 2144 of *LNCS*, pages 325–339. Springer, 2001.
- [21] Christian Jacobi. Formal Verification of a Fully IEEE Compliant Floating Point Unit. 2002.
- [22] Marius Cornea, John Harrison, Cristina Anderson, Ping Tak Peter Tang, Eric Schneider, and Evgeny Gvozdev. A Software Implementation of the IEEE 754R Decimal Floating-point Arithmetic using the Binary Encoding Format. *Transactions on Computers*, 58(2):148–162, 2008.
- [23] John Harrison. Formal Verification of Floating Point Trigonometric Functions. In *Formal Methods in Computer-aided Design*, pages 254–270. Springer, 2000.
- [24] John Harrison. Floating-point Verification using Theorem Proving. In *Formal Methods for the Design of Computer, Communication and Software Systems*, pages 211–242. Springer, 2006.
- [25] John Harrison, Ted Kubaska, Shane Story, et al. The Computation of Transcendental Functions on the IA-64 Architecture. In *Intel Technology Journal*. Citeseer, 1999.

- [26] John Harrison. Formal Verification of Square Root Algorithms. *Formal Methods in System Design*, 22(2):143–153, 2003.
- [27] Robert B Jones, John W O’Leary, C-JH Seger, Mark D Aagaard, and Thomas F Melham. Practical Formal Verification in Microprocessor Design. *Design & Test of Computers*, 18(4):16–25, 2001.
- [28] N. Narasimhan and R. Kaivola. Formal Verification of the Pentium® 4 Floating-Point Multiplier. In *Design, Automation & Test in Europe*, pages 1–8. IEEE, 2002.
- [29] Roope Kaivola. Intel Core™ i7 Processor Execution Engine Validation in a Functional Language Based Formal Framework. In *Practical Aspects of Declarative Languages*, pages 414–429. Springer, 2011.
- [30] Anna Slobodová. Formal Verification of Hardware Support for Advanced Encryption Standard. In *Formal Methods in Computer-Aided Design*, pages 1–4. IEEE, 2008.
- [31] Michael Whalen, Darren Cofer, Steven Miller, Bruce H Krogh, and Walter Storm. Integration of Formal Analysis into a Model-based Software Development Process. In *Formal Methods for Industrial Critical Systems*, pages 68–84. Springer, 2007.
- [32] Steven Miller, Elise Anderson, Lucas Wagner, Michael Whalen, and Matts Heimdahl. Formal Verification of Flight Critical Software. In *AIAA Guidance, Navigation, and Control Conference and Exhibit*, page 6431, 2005.
- [33] Alan C Tribble, DL Lempia, and Steven P Miller. Software Safety Analysis of a Flight Guidance System. In *Digital Avionics Systems Conference*, volume 2, pages 13C1–13C1. IEEE, 2002.
- [34] AC Tribble and SP Miller. Safety Analysis of Software Intensive Systems. *IEEE Aerospace and Electronic Systems*, 19(10):21–26, 2004.
- [35] https://github.com/adrashid/posits_verification.
- [36] J. Harrison. HOL Light: A Tutorial Introduction. In Mandayam Srivas and Albert Camilleri, editors, *Proceedings of the First International Conference on Formal Methods in Computer-Aided Design (FMCAD’96)*, volume 1166 of *Lecture Notes in Computer Science*, pages 265–269. Springer-Verlag, 1996.
- [37] L.C. Paulson. *ML for the Working Programmer*. Cambridge University Press, 1996.
- [38] J. Harrison. Formalized Mathematics. Technical Report 36, Turku Centre for Computer Science, Finland, 1996.
- [39] Charles Jacobsen, Alexey Solovyev, and Ganesh Gopalakrishnan. A Parameterized Floating-point Formalization in HOL Light. *Electronic Notes in Theoretical Computer Science*, 317:101–107, 2015.
- [40] Amr Talaat Abdel-Hamid. *A hierarchical Verification of the IEEE-754 Table-driven Floating-point Exponential Function using HOL*. PhD thesis, Concordia University, 2001.
- [41] Silvia M Müller and Wolfgang J Paul. *Computer Architecture: Complexity and Correctness*. Springer Science & Business Media, 2013.
- [42] Paul S Miner and James F Leathrum. Verification of IEEE Compliant Subtractive Division Algorithms. In *Formal Methods in Computer-Aided Design*, volume 1166 of *LNCS*, pages 64–78. Springer, 1996.
- [43] Christoph Berg. Formal Verification of an IEEE Floating Point Adder. *Master’s Thesis, Saarland University, Germany*, 2001.
- [44] John Harrison. Towards self-verification of HOL Light. In *International Joint Conference on Automated Reasoning*, pages 177–191. Springer, 2006.
- [45] Ramana Kumar. Self-compilation and Self-verification. Technical report, University of Cambridge, Computer Laboratory, 2016.
- [46] John O’Leary, Xudong Zhao, Rob Gerth, and Carl-Johan H Seger. Formally Verifying IEEE Compliance of Floating-point Hardware. *Intel Technology Journal*, 3(1):1–14, 1999.
- [47] Behzad Akbarpour, Abdelkader Dekdouk, and Sofiène Tahar. Formalization of Cadence SPW Fixed-Point Arithmetic in HOL. In *Integrated Formal Methods*, volume 2335 of *LNCS*, pages 185–204. Springer, 2002.

- [48] John O’Leary. Theorem Proving in Intel Hardware Design. 2009.
- [49] Lars Gesellensetter, Sabine Glesner, and Elke Salecker. Formal Verification with Isabelle/HOL in Practice: Finding a Bug in the GCC Scheduler. In *Formal Methods for Industrial Critical Systems*, pages 85–100. Springer, 2007.
- [50] Chris W Johnson. The Natural History of Bugs: Using Formal Methods to Analyse Software Related Failures in Space Missions. In *Formal Methods*, pages 9–25. Springer, 2005.
- [51] John Fitzgerald, Juan Bicarregui, Peter Gorm Larsen, and Jim Woodcock. Industrial Deployment of Formal Methods: Trends and Challenges. In *Industrial Deployment of System Engineering Methods*, pages 123–143. Springer, 2013.
- [52] Feng Zhang, Wensheng Niu, et al. A Survey on Formal Specification and Verification of System-level Achievements in Industrial Circles. *Academic Journal of Computing & Information Science*, 2(1), 2019.
- [53] <https://shemesh.larc.nasa.gov/fm/fm-main-research.html>.
- [54] <https://shemesh.larc.nasa.gov/fm/fm-collins-intro.html>.
- [55] Jiri Barnat, Jan Beran, Lubos Brim, Tomas Kratochvíla, and Petr Ročkal. Tool Chain to Support Automated Formal Verification of Avionics Simulink Designs. In *Formal Methods for Industrial Critical Systems*, pages 78–92. Springer, 2012.
- [56] Zongyu Cao, Wanyou Lv, Yanhong Huang, Jianqi Shi, and Qin Li. Formal Analysis and Verification of Airborne Software Based on DO-333. *Electronics*, 9(2):327, 2020.
- [57] Johanna Nellen, Thomas Rambow, Md Tawhid Bin Waez, Erika Ábrahám, and Joost-Pieter Katoen. Formal Verification of Automotive Simulink Controller Models: Empirical Technical Challenges, Evaluation and Recommendations. In *Formal Methods*, pages 382–398. Springer, 2018.
- [58] Haiyang Xu and Ping Wang. Real-time Reliability Verification for UAV Flight Control System Supporting Airworthiness Certification. *PloS ONE*, 11(12):e0167168, 2016.
- [59] Darren Cofer. Formal Methods in the Aerospace Industry: Follow the Money. In *Formal Engineering Methods*, pages 2–3. Springer, 2012.
- [60] Lennon Chaves, Iury V Bessa, Hussama Ismail, Adriano Bruno dos Santos Frutuoso, Lucas Cordeiro, and Ed-die Batista de Lima Filho. DSVerifier-aided Verification Applied to Attitude Control Software in Unmanned Aerial Vehicles. *Transactions on Reliability*, 67(4):1420–1441, 2018.
- [61] Virginie Wiels, Robert Delmas, David Doose, Pierre-Loïc Garoche, Jacques Cazin, and Guy Durrieu. Formal Verification of Critical Aerospace Software. *AerospaceLab*, (4):p-1, 2012.
- [62] Manish K Jaiswal and Hayden K-H So. Pacogen: A Hardware Posit Arithmetic Core Generator. *ACCESS*, 7:74586–74601, 2019.
- [63] Seyed Hamed Fatemi Langroudi, Tej Pandit, and Dhireesha Kudithipudi. Deep Learning Inference on Embedded Devices: Fixed-point Vs Posit. In *Energy Efficient Machine Learning and Cognitive Computing for Embedded Applications*, pages 19–23. IEEE, 2018.
- [64] John L Gustafson. Posit Arithmetic. *Mathematica Notebook Describing the Posit Number System*, 30, 2017.
- [65] Shin Yee Chung. Provably Correct Posit Arithmetic with Fixed-point Big Integer. In *Next Generation Arithmetic*, pages 1–10, 2018.
- [66] Sumayya Shiraz and Osman Hasan. A Library for Combinational Circuit Verification using the HOL Theorem Prover. *Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 37(2):512–516, 2017.