# Formal Specification of the Plutus Core Language

Plutus Team

18th January 2023

## **DRAFT**

## Abstract

This is intended to be a reference guide for developers who want to utilise the Plutus Core infrastructure. We lay out the grammar and syntax of untyped Plutus Core terms, and their semantics and evaluation rules. We also describe the built-in types and functions. The Appendices include a list of supported builtins in each era and some aspects of Plutus Core which have been mechanically formalised.

This document only describes untyped Plutus Core: a subsequent version will also include the syntax and semantics of Typed Plutus Core and describe its relation to untyped Plutus Core.

## Contents

Co	Contents						
1	Introduction	4					
2	Some Basic Notation	4					
_	2.1 Sets	4					
	2.2 Lists	5					
2	The Comment of Distance Comment	_					
3	The Grammar of Plutus Core 3.1 Lexical grammar	<b>5</b> 5					
		6					
	3.3 Notes	6					
4	Interpretation of Built-in Types and Functions	7					
	4.1 Built-in types	7					
	4.1.1 Type variables	7					
	4.1.2 Polymorphic types	8					
	4.1.3 Type assignments	9					
	4.2 Built-in functions	9					
	4.2.1 Inputs to built-in functions	9					
	4.2.2 Signatures and denotations of built-in functions	10					
	4.2.3 Denotations of built-in functions	11					
	4.2.4 Results of built-in functions	12					
	4.2.5 Parametricity for *-polymorphic arguments	12					
	4.3 Evaluation of built-in functions	13					
	4.3.1 Compatibility of inputs and signature entries	13					
	4.3.2 Evaluation	13					
	4.5.2 Evaluation	13					
5	Term Reduction	14					
	5.1 Values in Plutus Core	14					
	5.2 Term reduction	16					
6	The CEK Machine	18					
	6.1 Converting CEK evaluation results into Plutus Core terms	20					
7	Cost Accounting for Untyped Plutus Core	21					
8	Typed Plutus Core	21					
<b>A</b>	<b>Built-in Types and Functions Supported in the Alonzo Release</b>	22					
A	A.1 Built-in types and type operators	22					
	A.2 Alonzo built-in functions						
	A.2 Alonzo bulit-ili functions	24					
В	<b>Built-in Types and Functions Supported in the Vasil Release</b>	28					
	B.1 Built-in types and type operators	28					
	B.2 Built-in functions	28					
C	Formally Verified Behaviours	29					

D	Seria	alising data Objects Using the CBOR Format	29
	D.1	Introduction	29
	D.2	Notation	29
	D.3	The CBOR format	30
	D.4	Encoding and decoding the heads of CBOR items	30
	D.5	Encoding and decoding bytestrings	31
	D.6	Encoding and decoding integers	32
	D.7	Encoding and decoding data	33
E	Seria	alising Plutus Core Terms and Programs Using the flat Format	35
	E.1	Encoding and decoding	35
		E.1.1 Padding	36
	E.2	Basic flat encodings	36
		E.2.1 Fixed-width natural numbers	36
		E.2.2 Lists	37
		E.2.3 Natural numbers	37
		E.2.4 Integers	37
		E.2.5 Bytestrings	38
		E.2.6 Strings	39
	E.3	Encoding and decoding Plutus Core	39
		E.3.1 Programs	39
		E.3.2 Terms	40
		E.3.3 Built-in types	41
		E.3.4 Constants	42
		E.3.5 Built-in functions	43
		E.3.6 Variable names	44
	E.4	Cardano-specific serialisation issues	45
		E.4.1 Scope checking	45
		E.4.2 CBOR wrapping	45
	E.5	Example	45
Re	feren	ces	47
In	dex of	f Notation	49

## 1 Introduction

Plutus Core (more correctly, Untyped Plutus Core) is an eagerly-evaluated version of the untyped lambda calculus extended with some "built-in" types and functions; it is intended for the implementation of validation scripts on the Cardano blockchain. This document presents the syntax and semantics of Plutus Core, a specification of an efficient evaluator, a description of the built-in types and functions available in the Alonzo release of Cardano, and a specification of the binary serialisation format used by Plutus Core.

Since Plutus Core is intended for use in an environment where computation is potentially expensive and excessively long computations can be problematic we have also developed a costing infrastructure for Plutus Core programs. A description of this will be added in a later version of this document.

We also have a typed version of Plutus Core which provides extra robustness when untyped Plutus Core is used as a compilation target, and we will eventually provide a specification of the type system and semantics of Typed Plutus Core here as well, together with its relationship to Untyped Plutus Core.

## 2 Some Basic Notation

We begin with some notation which will be used throughout the document.

## **2.1** Sets

- Given a set X,  $X^*$  denotes the set of finite sequences of elements of X:

$$X^* = \left\{ + \right\} \{ X^n : n \in \mathbb{N} \}.$$

- $\mathbb{N} = \{0, 1, 2, 3, \dots\}.$
- $\mathbb{N}^+ = \{1, 2, 3, \dots\}.$
- $\bullet \ \mathbb{N}_{[a,b]} = \{n \in \mathbb{N} : a \le n \le b\}.$
- $\mathbb{B} = \mathbb{N}_{[0.255]}$ , the set of 8-bit bytes.
- $\mathbb{B}^*$  is the set of all bytestrings.
- $\mathbb{Z} = \{\ldots, -2, -1, 0, 1, 2, \ldots\}.$
- U denotes the set of Unicode scalar values, as defined in [24, Definition D76].
- U\* is the set of all Unicode strings.
- We assume that there is a special symbol × which does not appear in any other set we mention. The symbol × is used to indicate that some sort of error condition has occurred, and we will often need to consider situations in which a value is either × or a member of some set S. For brevity, if S is a set then we define

$$S_{\times} := S \uplus \{\times\}.$$

## 2.2 Lists

- The symbol [] denotes an empty list.
- The notation  $[x_m, \dots, x_n]$  denotes a list containing the elements  $x_m, \dots, x_n$ . If m > n then the list is empty.
- The length of a list L is denoted by  $\ell(L)$ .
- Given two lists  $L = [x_1, \dots, x_m]$  and  $L' = [y_1, \dots, y_n]$ ,  $L \cdot L'$  denotes their concatenation  $[x_1, \dots, x_m, y_1, \dots, y_n]$ .
- Given an object x and a list  $L = [x_1, ..., x_n]$ , we denote the list  $[x, x_1, ..., x_n]$  by  $x \cdot L$ .
- Given a list  $L = [x_1, \dots, x_n]$  and an object x, we denote the list  $[x_1, \dots, x_n, x]$  by  $L \cdot x$ .
- Given a syntactic category V, the symbol  $\overline{V}$  denotes a possibly empty list  $[V_1, \dots, V_n]$  of elements  $V_i \in V$ .

## 3 The Grammar of Plutus Core

This section presents the grammar of Plutus Core in a Lisp-like form. This is intended as a specification of the abstract syntax of the language; it may also by used by tools as a concrete syntax for working with Plutus Core programs, but this is a secondary use and we do not make any guarantees of its completeness when used in this way. The primary concrete form of Plutus Core programs is the binary format described in Appendix E.

## 3.1 Lexical grammar

```
Name n ::= [a-zA-Z] [a-zA-Z0-9_']* name Var x ::= n term variable BuiltinName bn ::= n built-in function name Version v ::= [0-9]*.[0-9]* version Constant c ::= \langle literal constant \rangle
```

Figure 1: Lexical grammar of Plutus Core

## 3.2 Grammar

```
L, M, N ::= x
Term
                                                                                             variable
                               (con Tc)
                                                                                             constant
                               (builtin b)
                                                                                              builtin
                               (lam x M)
                                                                                        \lambda abstraction
                               [M N]
                                                                                function application
                               (delay M)
                                                                           delay execution of a term
                               (force M)
                                                                           force execution of a term
                               (\operatorname{constr} i \ M_0 \ldots M_{m-1})
                                                          constructor with tag i and m arguments
                               (\operatorname{case} M \ N_0 \dots N_{m-1})
                                                                   case analysis with m alternatives
                               (error)
                                                                                                error
Program
                        ::=
                               (program v M)
                                                                                  versioned program
```

Figure 2: Grammar of untyped Plutus Core

## 3.3 Notes

Scoping. For simplicity, we assume throughout that the body of a Plutus Core program is a closed term, ie, that it contains no free variables. Thus (program 1.0.0 (lam x x)) is a valid program but (program 1.0.0 (lam x y)) is not, since the variable y is free. This condition should be checked before execution of any program commences, and the program should be rejected if its body is not closed. The assumption implies that any variable x occurring in the body of a program must be bound by an occurrence of lam in some enclosing term; in this case, we always assume that x refers to the most recent (ie, innermost) such binding.

**Iterated applications.** An application of a term M to a term N is represented by  $[M\ N]$ . We may occasionally write  $[M\ N_1\ ...\ N_k]$  or  $[M\ \overline{N}]$  as an abbreviation for an iterated application  $[\ ...\ [[M\ N_1]\ N_2]\ ...\ N_k]$ , and tools may also use this as concrete syntax.

**Built-in types and functions.** The language is parameterised by a set  $\mathcal{U}$  of *built-in types* (we sometimes refer to  $\mathcal{U}$  as the *universe*) and a set  $\mathcal{B}$  of *built-in functions* (*builtins* for short), both of which are sets of Names. Briefly, the built-in types represent sets of constants such as integers or strings; constant expressions (con Tc) represent values of the built-in types (the integer 123 or the string "string", for example), and built-in functions are functions operating on these values, and possibly also general Plutus Core terms. Precise details are given in Section 4. Plutus Core comes with a default universe and a default set of builtins, which are described in Appendix A.

**De Bruijn indices.** The grammar defines names to be textual strings, but occasionally (specifically in Appendix E) we want to use de Bruijn indices ([11], [4, C.3]), and for this we redefine names to be natural numbers. In de Bruijn terms,  $\lambda$ -expressions do not need to bind a variable, but in order to re-use our existing syntax we arbitrarily use 0 for the bound variable, so that all  $\lambda$ -expresssions are of the form (lam 0 M); other variables (ie, those not appearing immediately after a lam binder) are represented by natural number greater than zero.

**Lists in constructor and case terms** The grammar defines constructor and case terms to have a variable number of subterms written in sequence with no delimiters. This corresponds to the concrete syntax, e.g.

we write (constr  $0 t_1 t_2 t_3$ ). However, in the rest of the specification we will abuse notation and treat these terms as having *lists* of subterms.

## 4 Interpretation of Built-in Types and Functions

As mentioned above, Plutus Core is generic over a universe  $\mathcal U$  of types and a set  $\mathcal B$  of built-in functions. As the terminology suggests, built-in functions are interpreted as functions over terms and elements of the built-in types: in this section we make this interpretation precise by giving a specification of built-in types and functions in a set-theoretic denotational style. We require a considerable amount of extra notation in order to do this, and we emphasise that nothing in this section is part of the syntax of Plutus Core: it is meta-notation introduced purely for specification purposes.

## 4.1 Built-in types

We require some extra syntactic notation for built-in types: see Figure 3.

```
egin{array}{lll} at & ::= & n & & & \text{Atomic type} \\ op & ::= & n & & \text{Type operator} \\ T & ::= & at \mid op(T,T,...,T) & & \text{Built-in type} \\ \end{array}
```

Figure 3: Type names and operators

We assume that we have a set  $\mathcal{U}_0$  of atomic type names and a set  $\mathcal{O}$  of type operator names. Each type operator name  $op \in \mathcal{O}$  has an argument count  $|op| \in \mathbb{N}^+$ , and a type name  $op(T_1, \dots, T_n)$  is well-formed if and only if n = |op|. We define the universe  $\mathcal{U}$  to be the closure of  $\mathcal{U}_0$  under repeated applications of operators in  $\mathcal{O}$ :

$$\mathcal{U}_{i+1} = \mathcal{U}_i \cup \{op(T_1, \dots, T_{|op|}) : op \in \mathcal{O}, T_1, \dots, T_{|op|} \in \mathcal{U}_i\}$$
$$\mathcal{U} = \bigcup \{\mathcal{U}_i : i \in \mathbb{N}^+\}$$

The universe  $\mathcal U$  consists entirely of *names*, and the semantics of these names are given by *denotations*. Each built-in type  $T \in \mathcal U$  is associated with some mathematical set  $[\![T]\!]$ , the *denotation* of T. For example, we might have  $[\![boolean]\!] = \{\text{true}, \text{false}\}$  and  $[\![integer]\!] = \mathbb Z$  and  $[\![pair(a,b)]\!] = [\![a]\!] \times [\![b]\!]$ . See Appendix A for a description of the built-in types and type operators available in the Alonzo release of Plutus Core.

For non-atomic type names  $T = op(T_1, ..., T_r)$  we would generally expect the denotation of T to be obtained in some uniform way (ie, parametrically) from the denotations of  $T_1, ..., T_r$ ; we do not insist on this though.

## 4.1.1 Type variables

Built-in functions can be polymorphic, and to deal with this we need *type variables*. An argument of a polymorphic function can be either restricted to built-in types or can be an arbitrary term, and we define two different kinds of type variables to cover these two situations. See Figure 4.

TypeVariable tv ::=  $n_*$  fully polymorphic type variable  $n_\#$  built-in-polymorphic type variable

Figure 4: Type variables

We denote the set of all possible type variables by  $\mathcal{V}$ , the set of all fully-polymorphic type variables by  $\mathcal{V}_*$ , and the set of all built-in-polymorphic type variables  $v_\#$  by  $\mathcal{V}_\#$ . Note that  $\mathcal{V} \cap \mathcal{U} = \emptyset$  since the symbols  $_*$  and  $_\#$  do not occur in names in  $\mathcal{U}$ . The two kinds of type variable are required because we have two different types of polymorphism. Later on we will see that built-in functions can take arguments which can be of a type which is unknown but must be in  $\mathcal{U}$ , whereas other arguments can range over a larger set of values such as the set of all Plutus Core terms. Type variables in  $\mathcal{V}_\#$  are used in the former situation and  $\mathcal{V}_*$  in the latter.

Given a variable  $v \in \mathcal{V}$  we sometimes write

v::# if  $v\in\mathcal{V}_\#$ 

and

$$v :: * \text{ if } v \in \mathcal{V}_*.$$

## 4.1.2 Polymorphic types

We also need to talk about polymorphic types, and to do this we define an extended universe of polymorphic types  $\mathcal{U}_{\#}$  by adjoining  $\mathcal{V}_{\#}$  to  $\mathcal{U}_{0}$  and closing under type operators as before:

$$\begin{split} \mathcal{U}_{\#,0} &= \mathcal{U}_0 \cup \mathcal{V}_\# \\ \mathcal{U}_{\#,i+1} &= \mathcal{U}_{\#,i} \cup \{op(T_1,\dots,T_{|op|}) : op \in \mathcal{O}, T_1,\dots,T_{|op|} \in \mathcal{U}_{\#,i}\} \\ \mathcal{U}_\# &= \bigcup \big\{\mathcal{U}_{\#,i} \ : \ i \in \mathbb{N}^+ \big\}. \end{split}$$

We will denote a typical element of  $\mathcal{U}_{\#}$  by the symbol P (possibly subscripted). We define the set of *free* #-variables of an element of  $\mathcal{U}_{\#}$  by

$$\mathsf{FV}_\#(P) = \emptyset \text{ if } P \in \mathcal{U}_0$$
 
$$\mathsf{FV}_\#(v_\#) = \{v_\#\}$$
 
$$\mathsf{FV}_\#(op(P_1,\dots,P_k)) = \mathsf{FV}_\#(P_1) \cup \mathsf{FV}_\#(P_2) \cup \dots \cup \mathsf{FV}_\#(P_r).$$

Thus  $\mathsf{FV}_\#(P) \subseteq \mathcal{V}_\#$  for all  $P \in \mathcal{U}$ . We say that a type name  $P \in \mathcal{U}_\#$  is *monomorphic* if  $\mathsf{FV}_\#(P) = \emptyset$  (in which case we actually have  $P \in \mathcal{U}$ ); otherwise P is *polymorphic*. The fact that type variables in  $\mathcal{U}_\#$  are only allowed to come from  $\mathcal{V}_\#$  will ensure that values of polymorphic types such as lists and pairs can only contain values of built-in types: in particular, we will not be able to construct types representing things such as lists of Plutus Core terms.

## 4.1.3 Type assignments

A type assignment is a function  $S: D \to \mathcal{U}$  where D is some subset of  $\mathcal{V}_{\#}$ . As usual we say that D is the domain of S and denote it by dom S.

We can extend a type assignment S to a map  $\hat{S}: \mathcal{U}_{\#} \uplus \mathcal{V}_{*} \to \mathcal{U}_{\#} \uplus \mathcal{V}_{*}$  by defining

$$\begin{split} \hat{S}(v_\#) &= S(v_\#) \quad \text{if } v_\# \in \text{dom } S \\ \hat{S}(v_\#) &= v_\# \quad \text{if } v_\# \in \mathcal{V}_\# \backslash \text{dom } S \\ \hat{S}(T) &= T \quad \text{if } T \in \mathcal{U}_0 \\ \hat{S}(op(P_1, \dots, P_n)) &= op(\hat{S}(P_1), \dots, \hat{S}(P_n)) \\ \hat{S}(v_*) &= v_* \quad \text{if } v_* \in \mathcal{V}_*. \end{split}$$

If  $P \in \mathcal{U}_{\#}$  and S is a type assignment with  $\mathsf{FV}_{\#}(P) \subseteq \mathsf{dom}\, S$  then in fact  $\hat{S}(P) \in \mathcal{U}$ ; in this case we say that  $\hat{S}(P)$  is an *instance* or a *monomorphisation* of P ( $via\, S$ ). If T is an instance of P then there is a unique smallest S (with  $\mathsf{FV}_{\#}(P) = \mathsf{dom}\, S$ ) such that  $T = \hat{S}(P)$ : we write  $T \leq_S P$  to indicate that T is an instance of P via S and S is minimal.

Constructing type assignments. We say that a collection  $\{S_i : 1 \le i \le n\}$  of type assignments is consistent if  $S_i|_{D_{ij}} = S_j|_{D_{ij}}$  for all i and j, where | denotes function restriction and  $D_{ij} = \text{dom } S_i \cap \text{dom } S_j$ . If this is the case then (viewing functions as sets of pairs in the usual way)  $S_1 \cup \cdots \cup S_n$  is also a well-formed type assignment (each variable in its domain is associated with exactly one type).

Given  $T \in \mathcal{U}$  and  $P \in \mathcal{U}_{\#}$  it can be shown that  $T \leq_S P$  if and only if one of the following holds:

- T = P and  $S = \emptyset$ .
- $P \in \mathcal{V}_{\#}$  and  $S = \{(v_{\#}, T)\}.$
- $$\begin{split} \bullet & \quad \ T = op(T_1, \dots, T_n) \text{ with each } T_i \in \mathcal{U}. \\ & \quad \ P = op(P_1, \dots, P_n) \text{ with each } P_i \in \mathcal{U}_\#. \\ & \quad \ T_i \preceq_{S_i} P_i \text{ for } 1 \leq i \leq n. \\ & \quad \ \{S_1, \dots, S_n\} \text{ is consistent.} \\ & \quad \ S = S_1 \cup \dots \cup S_n. \end{split}$$

This allows us to decide whether  $T \in \mathcal{U}$  is an instance of  $P \in \mathcal{U}_{\#}$  and, if so, to construct an S with  $T \leq_S P$ .

## 4.2 Built-in functions

## 4.2.1 Inputs to built-in functions

To treat the typed and untyped versions of Plutus Core uniformly it is necessary to make the machinery of built-in functions generic over a set  $\mathcal{I}$  of *inputs* which are taken as arguments by built-in functions. In practice  $\mathcal{I}$  will be the set of Plutus Core values or something very closely related.

We require  $\ensuremath{\mathfrak{I}}$  to have the following two properties:

•  $\Im$  is disjoint from  $\llbracket T \rrbracket$  for all  $T \in \mathcal{U}$ 

• There should be disjoint subsets  $\mathcal{C}_T \subseteq \mathcal{I}$  (where  $T \in \mathcal{U}$ ) of constants of type T and maps  $\llbracket \cdot \rrbracket_T : \mathcal{C}_T \to \llbracket T \rrbracket$  (denotation) and  $\llbracket \cdot \rrbracket_T : \llbracket T \rrbracket \to \mathcal{C}_T$  (reification) such that  $\llbracket \llbracket c \rrbracket_T \rrbracket_T = c$  for all  $c \in \mathcal{C}_T$ . We do not require these maps to be bijective (for example, there may be multiple inputs with the same denotation), but the condition implies that  $\llbracket \cdot \rrbracket_T$  is surjective and  $\llbracket \cdot \rrbracket_T$  is injective.

It is also convenient to let  $[\![ \mathbb{J} ]\!] = \mathbb{J}$  and define both  $[\![ \cdot ]\!]_{\mathbb{J}}$  and  $\{\![ \cdot ]\!]_{\mathbb{J}}$  to be the identity function.

For example, we could take  $\mathcal{I}$  to be the set of all Plutus Core values (see Section 5.1),  $\mathcal{C}_T$  to be the set of all terms of the form (con Tc), and  $[\cdot]_T$  to be the function which maps (con Tc) to c. For simplicity we are assuming that mathematical entities occurring as members of type denotations [T] are embedded directly as values c in Plutus Core constant terms. In reality, tools which work with Plutus Core will need some concrete syntactic representation of constants; we do not specify this here, but see Section A.1 for suggested syntax for the built-in types included in the Alonzo release.

## 4.2.2 Signatures and denotations of built-in functions

We will consistently use the symbol  $\tau$  and subscripted versions of it to denote members of  $\mathcal{U}_{\#} \uplus \mathcal{V}_{*}$  in the rest of the document; these indicate the types of values consumed and returned by built-in functions.

We also define a class of *quantifications* which are used to introduce type variables: a quantification is a symbol of the form  $\forall v$  with  $v \in \mathcal{V}$ ; the set of of all possible quantifications is denoted by  $\mathcal{Q}$ .

**Signatures.** Every built-in function  $b \in \mathcal{B}$  has a *signature*  $\sigma(b)$  which describes the types of its arguments and its return value: a signature is of the form

$$[\iota_1,\ldots,\iota_n] \to \tau$$

with

- $\iota_i \in \mathcal{U}_\# \uplus \mathcal{V}_* \uplus \mathcal{Q}$  for all j
- $\tau \in \mathcal{U}_{\#} \uplus \mathcal{V}_{*}$
- $|\{j : \iota_i \notin \mathcal{Q}\}| \ge 1 \text{ (so } n \ge 1)$
- If  $\iota_j$  involves  $v \in \mathcal{V}$  then  $\iota_k = \forall v$  for some k < j, and similarly for  $\tau$ ; in other words, any type variable v must be introduced by a quantification before it is used. (Here  $\iota$  involves v if either  $\iota = T \in \mathcal{U}_\#$  and  $v \in \mathsf{FV}_\#(T)$  or  $\iota = v$  and  $v \in \mathcal{V}_*$ .)
- If  $\tau$  involves  $v \in \mathcal{V}$  then some  $\iota_i$  must involve v; this implies that  $\mathsf{FV}_\#(\tau) \subseteq \bigcup \{\mathsf{FV}_\#(\iota_i) : \iota_i \in \mathcal{U}_\#\}$ .
- If  $j \neq k$  and  $\iota_j, \iota_k \in \Omega$  then  $\iota_j \neq \iota_k$ ; ie, no quantification appears more than once.
- If  $i_i = \forall v \in \Omega$  then some  $i_j \notin \Omega$  with j > i must involve v (signatures are not allowed to contain phantom type variables).

For example, in our default set of built-in functions we have the functions mkCons with signature  $[\forall a_\#, a_\#, a_\#] \rightarrow \text{list}(a_\#)$  and ifThenElse with signature  $[\forall a_*, \text{boolean}, a_*, a_*] \rightarrow a_*$ . When we use mkCons its arguments must be of built-in types, but the two final arguments of ifThenElse can be any Plutus Core values.

If b has signature  $[\iota_1, \dots, \iota_n] \to \tau$  then we define the arity of b to be

$$\alpha(b) = [\iota_1, \ldots, \iota_n].$$

We also define

$$\chi(b) = n$$
.

We may abuse notation slightly by using the symbol  $\sigma$  to denote a specific signature as well as the function which maps built-in function names to signatures, and similarly with the symbol  $\alpha$ .

Given a signature  $\sigma = [\iota_1, \dots, \iota_n] \to \tau$ , we define the *reduced signature*  $\bar{\sigma}$  to be

$$\bar{\sigma} = [\iota_j \, : \, \iota_j \not \in \mathfrak{Q}] \to \tau$$

Here we have extended the usual set comprehension notation to lists in the obvious way, so  $\bar{\sigma}$  just denotes the signature  $\sigma$  with all quantifications omitted. We will often write a reduced signature in the form  $[\tau_1, \dots, \tau_m] \to \tau$  to emphasise that the entries are *types*, and  $\forall$  does not appear.

Also, given an arity =  $[\iota_1, ..., \iota_n]$ , the *reduced arity* is

$$\bar{\alpha} = [\iota_i : \iota_i \notin \Omega].$$

**Commentary.** What is the intended meaning of the notation introduced above? In Typed Plutus Core we have to instantiate polymorphic functions (both built-in functions and polymorphic lambda terms) at concrete types before they can be applied, and in Untyped Plutus Core instantiation is replaced by an application of force. When we are applying a built-in function we supply its arguments one by one, and we can also apply force (or perform type instantiation in the typed case) to a partially-applied builtin "between" arguments (and also after the final argument); no computation occurs until all arguments have been supplied and all forces have been applied. The arity (read from left to right) specifies what types of arguments are expected and how they should be interleaved with applications of force, and  $\chi(b)$  tells you the total number of arguments and applications of force that a built-in function b requires. A fully-polymorphic type variable  $a_*$  indicates that an arbitrary value from  $\mathfrak I$  can be provided, whereas a type from  $\mathfrak U_\#$  indicates that a value of the specified built-in type is expected. Occurrences of quantifications indicate that force is to be applied to a partially-applied builtin; we allow this purely so that partially-applied builtins can be treated in the same way as delayed lambda-abstractions: force has no effect unless it is the very last item in the signature. In Plutus Core, partially-applied builtins are values which can be treated like any others (for example, by being passed as an argument to a lam-expression): see Section 5.1.

## 4.2.3 Denotations of built-in functions

The basic idea is that a built-in function b should represent some mathematical function on the denotations of the types of its inputs. However, this is complicated by the presence of polymorphism and we have to require that there is such a function for every possible monomorphisation of b.

More precisely, suppose that we have a builtin b with reduced signature  $[\tau_1, \dots \tau_n] \to \tau$ . For every type assignment S with dom  $S = \mathsf{FV}_\#(\tau_1) \cup \dots \cup \mathsf{FV}_\#(\tau_n)$  (which contains  $\mathsf{FV}_\#(\tau)$  by the conditions on signatures in Section 4.2.2) we require a *denotation of b at S*, a function

$$[\![b]\!]_S: [\![\hat{S}(\tau_1)]\!] \times \cdots \times [\![\hat{S}(\tau_n)]\!] \to [\![\hat{S}(\tau)]\!]_{\times}.$$

where

$$\llbracket v_* \rrbracket = \Im \text{ for } v_* \in \mathcal{V}_*.$$

This makes sense because  $\hat{S}(\tau_i) \in \mathcal{U} \uplus \mathcal{I}$  for all i, so  $[\hat{S}(\tau_i)]$  is always defined, and similarly for  $\tau$ .

If  $FV_{\#}(\bar{\sigma}(b)) = \emptyset$  (in which case we say that *b* is *monomorphic*) then the only relevant type assignment will be the empty one; in this case we have a single denotation

$$\llbracket b \rrbracket_{\varnothing} : \llbracket \tau_1 \rrbracket \times \cdots \times \llbracket \tau_n \rrbracket \to \llbracket \tau \rrbracket_{\times}.$$

Denotations of builtins are mathematical functions which terminate on every possible input; the symbol x can be returned by a function to indicate that something has gone wrong, for example if an argument is out of range.

In practice we expect most builtins to be *parametrically polymorphic* [25, 22], so that the denotation  $[\![b]\!]_S$  will be the "same" for all type assignments S; we do not insist on this though.

## 4.2.4 Results of built-in functions.

If *r* is the result of the evaluation of some built-in function there are thus three possibilities:

- 1.  $r \in [T]$  for some  $T \in \mathcal{U}$ .
- $2. r \in \mathcal{I}.$
- 3. r = x.

In other words,

$$r\in\mathcal{R}:=\left\{+\right\}\{[\![T]\!]\,:\,T\in\mathcal{U}\}\uplus\mathcal{I}\uplus\{\times\}.$$

Our assumptions on the set  $\Im$  (Section 4.2.1) allow us define a function

$$\{\![\cdot]\!]: \mathcal{R} \to \mathcal{I}_{\mathsf{x}}$$

which converts results of built-in functions back into inputs (or the x symbol):

- 1. If  $r \in [T]$ , then  $\{r\} = \{r\}_T \in \mathcal{C}_T \subseteq \mathcal{I}$ .
- 2. If  $r \in \mathcal{I}$  then  $\{ |r| \} = r$ .
- 3.  $\{ |x| \} = x$ .

## 4.2.5 Parametricity for \*-polymorphic arguments

A built-in function b can only inspect arguments which are values of built-in types; other arguments (occurring as  $a_*$  in  $\bar{\sigma}(b)$ ) are treated opaquely, and can be discarded or returned as (part of) a result, but cannot be altered or examined (in particular, they cannot be compared for equality): b is parametrically polymorphic in such arguments. This implies that if a builtin returns a value  $v \in \mathcal{I}$ , then v must have been an argument of the function.

## 4.3 Evaluation of built-in functions

## 4.3.1 Compatibility of inputs and signature entries

The previous section describes how a built-in function is interpreted as a mathematical function. When a Plutus Core built-in function b is applied to a sequence of arguments, the arguments must have types which are compatible with the signature of b; for example, if b has signature  $[\forall a_{\#}, \forall b_{\#}, a_{\#}, b_{\#}, a_{\#}, c_{*}, c_{*}] \rightarrow c_{*}$  and b is applied to a sequence of inputs  $V_{1}, V_{2}, V_{3}, V_{4}, V_{5}$  then  $V_{1}, V_{2}$ , and  $V_{3}$  must all be constants of some monomorphic built-in types and the types of  $V_{1}$  and  $V_{3}$  must be the same;  $V_{4}$  and  $V_{5}$  can be arbitrary inputs. This section describes the conditions for type compatibility.

In detail, given a reduced arity  $\bar{\alpha} = [\tau_1, \dots, \tau_n]$ , a sequence  $\bar{V} = [V_1, \dots, V_m]$ , and a type assignment S we say that  $\bar{V}$  is *compatible with*  $\bar{\alpha}$  (*via* S) if and only if n = m and, letting  $I = \{i : 1 \le i \le n, \tau_i \in \mathcal{U}_\#\}$  (so  $\tau_j \in \mathcal{V}_*$  if  $j \notin I$ ), there exist type assignments  $S_i$   $(1 \le i \le n)$  such that all of the following are satisfied

- For all  $i \in I$  there exists  $T_i \in \mathcal{U}$  such that  $V_i \in \mathcal{C}_{T_i}$  and  $T_i \leq_{S_i} \tau_i$ .
- $\{S_i : i \in I\}$  is consistent (see Section 4.1.3).
- $S = \{ | \{S_i : i \in I\} \}$ .

If these conditions are all satisfied then we can find suitable  $S_i$  using the procedure described in Section 4.1.3 and this allows us to construct S explicitly since the  $S_i$  are consistent. Note that in this case dom  $S = \text{dom } S_1 \cup \ldots \cup \text{dom } S_n = \text{FV}_\#(\tau_1) \cup \cdots \cup \text{FV}_\#(\tau_n) = \text{FV}_\#(\alpha)$ , so S is minimal in the sense that no S' with dom S' strictly smaller than dom S is sufficient to monomorphise all of the  $\tau_i$  simultaneously. We write

$$[V_1, \ldots, V_m] \approx_S [\tau_1, \ldots, \tau_n]$$

in this case. If  $\bar{V}$  is not compatible with  $\bar{\alpha}$  then we write  $\bar{V} \not\approx \bar{\alpha}$ .

#### 4.3.2 Evaluation

For later use we define a function Eval which attempts to evaluate an application of a built-in function b to a sequence of inputs  $[V_1, \ldots, V_m]$ . This fails if the number of inputs is incorrect or if the inputs are not compatible with  $\bar{\alpha}(b)$ :

$$\text{Eval}(b, [V_1, \dots, V_n]) = \times \text{ if } [V_1, \dots, V_n] \not\approx \bar{\alpha}(b).$$

Otherwise, the conditions for the existence of a denotation of b are met and we can apply that denotation to the denotations of the inputs and then reify the result. If  $[V_1, \ldots, V_n] \approx_S \bar{\alpha}(b) = [\tau_1, \ldots, \tau_n]$ , let  $T_i = \hat{S}(\tau_i)$  for  $1 \le i \le n$ ; then we define

$$\mathsf{Eval}(b, [V_1, \dots, V_n]) = \{\![\![b]\!]_S ([\![V_1]\!]_{T_1}, \dots, [\![V_n]\!]_{T_n}) \}\!\}.$$

It can be checked that the compatibility condition guarantees that this makes sense according to the definition of  $[\![b]\!]_S$  in Section 4.2.3.

## Notes.

- All of the machinery which we have defined for built-in functions is parametric over the set  $\mathcal{I}$  of inputs and the sets  $\mathcal{C}_T \subseteq \mathcal{I}$  of constants. This also applies to the Eval function, so its meaning is not fully defined until we have given concrete definitions of the sets of inputs and constants.
- The error value × can occur in two different ways: either because the arguments are not compatible with the signature, or because the builtin itself returns × to signal some error condition.

• The symbol x is not part of Plutus Core; when we define reduction rules and evaluators for Plutus Core later some extra translation will be required to convert the result of Eval into something appropriate to the context.

## 5 Term Reduction

This section defines the semantics of (untyped) Plutus Core.

## **5.1** Values in Plutus Core

The semantics of built-in functions in Plutus Core are obtained by instantiating the sets  $\mathcal{C}_T$  of constants of type T (see Section 4.2.1) to be the expressions of the form (con T c) and the set  $\mathfrak{I}$  to be the set of Plutus Core values, terms which cannot immediately undergo any further reduction, such as lambda terms and delayed terms. Values also include partial applications of built-in functions such as [(builtin modInteger) (con integer 5)], which cannot perform any computation until a second integer argument is supplied. However, partial applications must also be well-formed, in the sense that applications of force must be correctly interleaved with genuine arguments, and the arguments must themselves be values.

We define syntactic classes V of Plutus Core values and A of partial builtin applications simultaneously:

```
Value V ::= (\operatorname{con} T c) (\operatorname{delay} M) (\operatorname{lam} x M) (\operatorname{constr} i \overline{V})
```

Figure 5: Values in Plutus Core

Here *A* is the class of well-formed partial applications, and to define this we first define a class of possibly ill-formed iterated applications *B* for each built-in function  $b \in \mathcal{B}$ :

```
B ::= (builtin b)
[B V]
(force B)
```

Figure 6: Partial built-in function application

We let B denote the set of terms generated by the grammar in Figure 6 and we define a function  $\beta$  which extracts the name of the built-in function occurring in a term in B:

```
\begin{array}{ll} \beta((\mbox{builtin}\,b)) &= b \\ \beta([B\,V]) &= \beta(B) \\ \beta((\mbox{force}\,B)) &= \beta(B) \end{array}
```

We also define a function  $\|\cdot\|$  which measures the size of a term  $B \in B$ :

$$\begin{array}{ll} \| (\mbox{builtin } b) \| &= 0 \\ \| [B\ V] \| &= 1 + \| B \| \\ \| (\mbox{force } B) \| &= 1 + \| B \| \end{array}$$

**Well-formed partial applications.** A term  $B \in B$  is an application of  $b = \beta(B)$  to a number of values in S, interleaved with applications of force. We now define what it means for B to be a *well-formed partial application*. Suppose that  $\alpha(b) = [\iota_1, \ldots, \iota_n]$ . Firstly we require that ||B|| < n, so that b is not fully applied; in this case we put  $\iota = \iota_{||B||}$ , the element of b's signature which describes what kind of "argument" b currently expects. The definition is completed by induction on the structure of B:

- 1. B = (builtin b) is always well-formed.
- 2.  $B = [B' \ V]$  is well-formed if B' is well-formed and  $\iota \in \mathcal{U}_{\#}$  or  $\iota \in \mathcal{V}_{*}$  (equivalently,  $\iota \notin \mathcal{Q}$ ).
- 3. B = (force B') is well-formed if B' is well-formed and  $\iota \in \Omega$ .

The definition of values in Figure 5 is now completed by defining *A* to be the syntactic class of well-formed *partial* built-in function applications:

```
A = \{B \in B : B \text{ is a well-formed partial application}\}.
```

Note that this definition does not impose any requirements of type correctness. For example, with the types and functions defined in Appendix A the term X = [builtin modInteger) (con string "blue")] is a valid value which could be treated like any other, for instance by being passed as an argument to a lam expression. However, the evaluation rules described in the next section require that when a built-in function b becomes fully applied the types of the arguments are checked against the signature of b using the relation a and the function Eval defined in Sections 4.3.1 and 4.3.2, so an error would arise if the term a were ever applied to another argument.

**More notation.** Suppose that A is a well-formed partial application with  $\alpha(\beta(A)) = [\iota_1, \dots, \iota_n]$ . We define a function next which extracts the next argument (or force) expected by A:

$$\mathsf{next}(A) = \iota_{\|A\|+1}.$$

This makes sense because in a well-formed partial application A we have ||A|| < n.

We also define a function args which extracts the arguments which b has received so far in A:

```
\begin{array}{ll} \operatorname{args}((\operatorname{builtin} b)) &= [] \\ \operatorname{args}([A \ V]) &= \operatorname{args}(A) \cdot V \\ \operatorname{args}((\operatorname{force} A)) &= \operatorname{args}(A). \end{array}
```

## **5.2** Term reduction

We define the semantics of Plutus Core using contextual semantics (or reduction semantics): see [15] or [13] or [16, 5.3], for example. We use A to denote a partial application of a built-in function as in Section 5.1 above. For builtin evaluation, we instantiate the set  $\Im$  of Section 4.2.1 to be the set of Plutus Core values. Thus all builtins take values as arguments and return a value or  $\times$ . Since values are terms here, we can take  $\{V\} = V$ .

The notation [V/x]M below denotes substitution of the value V for the variable x in M. This is *capture-avoiding* in that substitution is not performed on occurrences of x inside subterms of M of the form  $(lam\ x\ N)$ .

Frame 
$$f$$
 ::= [\_M] left application [V \_] right application (force \_) force (constr  $i \overline{V} - \overline{M}$ ) constructor argument (case \_  $\overline{M}$ ) case scrutinee

(a) Grammar of reduction frames for Plutus Core

$$M \rightarrow M'$$

Term M reduces in one step to term M'.

(b) Reduction via Contextual Semantics

$$\mathsf{Eval'}(b,[V_1,\dots,V_n]) = \begin{cases} (\texttt{error}) & \text{if } \mathsf{Eval}(b,[V_1,\dots,V_n]) = \mathsf{X} \\ \mathsf{Eval}(b,[V_1,\dots,V_n]) & \text{otherwise} \end{cases}$$

(c) Built-in function application

Figure 7: Term reduction for Plutus Core

It can be shown that any closed Plutus Core term whose evaluation terminates yields either (error) or a value. Recall from Section 3.3 that we require the body of every Plutus Core program to be closed.

## **6** The CEK Machine

This section contains a description of an abstract machine for efficiently executing Plutus Core. This is based on the CEK machine of Felleisen and Friedman [14].

The machine alternates between two main phases: the *compute* phase ( $\triangleright$ ), where it recurses down the AST looking for values, saving surrounding contexts as frames (or *reduction contexts*) on a stack as it goes; and the *return* phase ( $\triangleleft$ ), where it has obtained a value and pops a frame off the stack to tell it how to proceed next. In addition there is an error state  $\spadesuit$  which halts execution with an error, and a halting state  $\square$  which halts execution and returns a value to the outside world.

To evaluate a program (program vM), we first check that the version number v is valid, then start the machine in the state  $[]; [] \triangleright M$ . It can be proved that the transitions in Figure 10 always preserve validity of states, so that the machine can never enter a state such as  $[] \triangleleft M$  or s, (force  $\_) \triangleleft (lam x A M)$  which isn't covered by the rules. If such a situation were to occur in an implementation then it would indicate that the machine was incorrectly implemented or that it was attempting to evaluate an ill-formed program (for example, one which attempts to apply a variable to some other term).

```
State \Sigma ::= s; \rho \rhd M \mid s \lhd V \mid \blacklozenge \mid \Box V Stack s ::= f^* CEK value V ::= \langle \operatorname{con} T c \rangle \mid \langle \operatorname{delay} M \rho \rangle \mid \langle \operatorname{lam} x M \rho \rangle \mid \langle \operatorname{constr} i \overline{V} \rangle \mid \langle \operatorname{builtin} b \overline{V} \eta \rangle Environment \rho ::= [] \mid \rho[x \mapsto V] Expected builtin arguments \eta ::= [i] \mid \iota \cdot \eta
```

Figure 8: Grammar of CEK machine states for Plutus Core

Figure 9: Grammar of CEK stack frames

Figures 8 and 9 define some notation for *states* of the CEK machine: these involve a modified type of value adapted to the CEK machine, environments which bind names to values, and a stack which stores partially evaluated terms whose evaluation cannot proceed until some more computation has been performed (for example, since Plutus Core is a strict language function arguments have to be reduced to values before application takes place, and because of this a lambda term may have to be stored on the stack while its

argument is being reduced to a value). Environments are lists of the form  $\rho = [x_1 \mapsto V_1, \dots, x_n \mapsto V_n]$  which grow by having new entries appended on the right; we say that x is bound in the environment  $\rho$  if  $\rho$  contains an entry of the form  $x \mapsto V$ , and in that case we denote by  $\rho[x]$  the value V in the rightmost (ie, most recent) such entry.\*

To make the CEK machine fit into the built-in evaluation mechanism defined in Section 4 we define  $\mathcal{I} = V$  and  $\mathcal{C}_T = \{\langle \text{con } T c \rangle : T \in \mathcal{U}, c \in [T]\}$ .

The rules in Figure 10 show the transitions of the machine; if any situation arises which is not included in these transitions (for example, if a frame  $\lceil \langle con \ T \ c \rangle \rceil$  is encountered or if an attempt is made to apply force to a partial builtin application which is expecting a term argument), then the machine stops immediately in an error state.

<sup>\*</sup>The description of environments we use here is more general than necessary in that it permits a given variable to have multiple bindings; however, in what follows we never actually retrieve bindings other than the most recent one and we never remove bindings to expose earlier ones. The list-based definition has the merit of simplicity and suffices for specification purposes but in an implementation it would be safe to use some data structure where existing bindings of a given variable are discarded when a new binding is added.

 $\Sigma \mapsto \Sigma'$ 

Machine takes one step from state  $\Sigma$  to state  $\Sigma'$ 

```
s: \rho \triangleright x
                                                                                                                             \mapsto s \triangleleft \rho[x] if x is bound in \rho
                                                                                                                            \mapsto s \triangleleft \langle con T c \rangle
                                                          s; \rho \triangleright (\operatorname{con} T c)
                                                          s; \rho \triangleright (\operatorname{lam} x M)
                                                                                                                            \mapsto s \triangleleft \langle lam \ x \ M \ \rho \rangle
                                                          s; \rho \triangleright (\text{delay } M)
                                                                                                                           \mapsto s \triangleleft \langle \text{delay } M \rho \rangle
                                                                                                                            \mapsto (force _) \cdot s; \rho \triangleright M
                                                          s; \rho \triangleright (force M)
                                                          s; \rho \triangleright [M \ N]
                                                                                                                            \mapsto [ (N, \rho)] \cdot s; \rho \triangleright M
                                                                                                                            \mapsto (constr i (\overline{M}, \rho)) \cdot s; \rho \triangleright M
                                                          s; \rho \triangleright (\text{constr } i M \cdot \overline{M})
                                                          s; \rho \triangleright (\text{constr } i \mid )
                                                                                                                            \mapsto s \triangleleft \langle constri\rangle
                                                          s; \rho \triangleright (case N \overline{M})
                                                                                                                           \mapsto (case (\overline{M}, \rho)) \cdot s; \rho \triangleright N
                                                                                                                            \mapsto s \triangleleft \langle \text{builtin } b \mid \alpha(b) \rangle
                                                          s; \rho \triangleright (\text{builtin } b)
                                                          s; \rho \triangleright (error)
                                                                                                                            \mapsto \sqcap V
                                                              [(M, \rho)] \cdot s \triangleleft V
                                                                                                                            \mapsto [V ] · s; \rho \triangleright M
                    [\langle lam \ x \ M \ \rho \rangle \ \_] \cdot s \triangleleft V
                                                                                                                           \mapsto s; \rho[x \mapsto V] \triangleright M
                                    (force ) \cdot s \triangleleft \langle \text{delay } M \rho \rangle
                                                                                                                           \mapsto s; \rho \triangleright M
                                     (force _) \cdot s \triangleleft (builtin b \overline{V} (\iota \cdot \eta)) \mapsto s \triangleleft (builtin b \overline{V} \eta) if \iota \in \Omega
                                    (\texttt{force}\_) \cdot s \lhd (\texttt{builtin}\ b\ \overline{V}\ [\iota]) \quad \mapsto \ \mathsf{Eval}_{\mathsf{CFK}}\ (s,b,\overline{V}) \ \text{if}\ \iota \in \mathfrak{Q}
(\operatorname{constr} i \overline{V} \ (M \cdot \overline{M}, \rho)) \cdot s \triangleleft V
                                                                                                                          \mapsto (constr i \overline{V} \cdot V \quad (\overline{M}, \rho) \cdot s : \rho \triangleright M
                                                                                                                           \mapsto s \triangleleft \langle \text{constr } i \, \overline{V} \cdot V \rangle
          (\text{constr } i \ \overline{V} \ \_([], \rho)) \cdot s \triangleleft V
     (\mathsf{case}\,\_(M_0 \ldots M_n, \rho)) \cdot s \lhd \langle \mathsf{constr}\, i\, V_0 \ldots V_m \rangle \ \mapsto \ [\_V_m] \cdot \ldots \cdot [\_V_0] \cdot s; \rho \rhd M_i \ \text{if} \ i \in [0, n]
   [\langle \text{builtin } b \ \overline{V} (\iota \cdot \eta) \rangle] \cdot s \triangleleft V
                                                                                                                       \mapsto s \triangleleft \langle \text{builtin } b (\overline{V} \cdot V) \eta \rangle \text{ if } \iota \in \mathcal{U}_{\#} \cup \mathcal{V}_{\#}
        [\langle \mathtt{builtin}\; b\; \overline{V}\; [\imath] \rangle\; \_] \cdot s \lhd V
                                                                                                                           \mapsto \operatorname{Eval}_{CFK}(s, b, \overline{V} \cdot V) \text{ if } \iota \in \mathcal{U}_{\#} \cup \mathcal{V}_{*}
                                                                    (a) CEK machine transitions for Plutus Core
```

$$\mathsf{Eval}_{\mathsf{CEK}}(s,b,[V_1,\dots,V_n]) = \begin{cases} \blacklozenge & \text{if } \mathsf{Eval}\,(b,[V_1,\dots,V_n]) = \mathsf{x} \\ s \lhd \mathsf{Eval}\,(b,[V_1,\dots,V_n]) & \text{otherwise} \end{cases}$$

(b) Evaluation of built-in functions

Figure 10: A CEK machine for Plutus Core

## 6.1 Converting CEK evaluation results into Plutus Core terms

The purpose of the CEK machine is to evaluate Plutus Core terms, but in the definition in Figure 10 it does not return a Plutus Core term; instead the machine can halt in two different ways:

• The machine can halt in the state  $\square V$  for some CEK value V.

• The machine can halt in the state •.

To get a complete evaluation strategy for Plutus Core we must convert these states into Plutus Core terms. The term corresponding to  $\spadesuit$  is (error), and to obtain a term from  $\square V$  we perform a process which we refer to as *discharging* the CEK value V (also known as *unloading*: see [21, pp. 129–130], [12, pp. 71ff]). This process substitutes bindings in environments for variables occurring in the value V to obtain a term  $\mathcal{U}(V)$ : see Figure 11a. Since environments contain bindings  $x \mapsto W$  of variables to further CEK values, we have to recursively discharge those bindings first before substituting: see Figure 11b, which defines an operation  $@\rho$  which does this. As before [N/x]M denotes the usual (capture-avoiding) process of substituting the term N for all unbound occurrences of the variable x in the term M. Note that in Figure 11b we substitute the rightmost (ie, the most recent) bindings in the environment first.

```
\mathcal{U}(\langle\operatorname{con}Tc\rangle) = (\operatorname{con}Tc) \mathcal{U}(\langle\operatorname{delay}M\rho\rangle) = (\operatorname{delay}M)@\rho \mathcal{U}(\langle\operatorname{lam}xM\rho\rangle) = (\operatorname{lam}xM)@\rho \mathcal{U}(\langle\operatorname{constr}i\overline{V}\rangle) = (\operatorname{constr}i\overline{\mathcal{U}(V)}) \mathcal{U}(\langle\operatorname{builtin}bV_1V_2\dots V_k\eta\rangle) = [\dots [[(\operatorname{builtin}b)(\mathcal{U}(V_1))](\mathcal{U}(V_2))]\dots(\mathcal{U}(V_k))] (a) Discharging CEK values M@\rho = [(\mathcal{U}(V_1))/x_1] \cdots [(\mathcal{U}(V_n))/x_n]M \quad \text{if } \rho = [x_1 \mapsto V_1, \dots, x_n \mapsto V_n] (b) Iterated substitution/discharging
```

Figure 11: Discharging CEK values to obtain Plutus Core terms

We can prove that if we evaluate a closed Plutus Core term in the CEK machine and then convert the result back to a term using the above procedure then we get the result that we should get according to the semantics in Figure 7.

## 7 Cost Accounting for Untyped Plutus Core

To follow.

## **8 Typed Plutus Core**

To follow.

## Appendix A Built-in Types and Functions Supported in the Alonzo Release

## A.1 Built-in types and type operators

The Alonzo release of the Cardano blockchain (September 2021) supports a default set of built-in types and type operators defined in Tables 1 and 2. We also include concrete syntax for these; the concrete syntax is not strictly part of the language, but may be useful for tools working with Plutus Core.

Type	Denotation	Concrete Syntax
integer	Z	-?[0-9]*
bytestring	$\mathbb{B}^*$ , the set of sequences of bytes or 8-bit	#([0-9A-Fa-f][0-9A-Fa-f])*
	characters.	
string	$\mathbb{U}^*$ , the set of sequences of Unicode char-	See note below.
	acters.	
bool	{true, false}	True   False
unit	{()}	()
data	See below.	Not yet supported.

Table 1: Atomic Types

Operator op	op	Denotation	Concrete Syntax
list	1	$\llbracket \mathtt{list}(t) \rrbracket = \llbracket t \rrbracket^*$	Not yet supported
pair	2	$[\![\operatorname{pair}(t_1,t_2)]\!] = [\![t_1]\!] \times [\![t_2]\!]$	Not yet supported

Table 2: Type Operators

**Concrete syntax for strings.** Strings are represented as sequences of Unicode characters enclosed in double quotes, and may include standard escape sequences.

Concrete syntax for higher-order types. Types such as list (integer) and pair (bool, string)) are represented by application at the type level, thus: [(con list) (con integer)] and [(con pair) (con bool) (con string)]. Each higher-order type will need further syntax for representing constants of those types. For example, we might use [] for list values and (,) for pairs, so the list [11, 22, 33] might be written as

Note however that this syntax is not currently supported by most Plutus Core tools at the time of writing.

**The** data **type.** We provide a built-in type data which permits the encoding of simple data structures for use as arguments to Plutus Core scripts. This type is defined in Haskell as

```
data Data =
   Constr Integer [Data]
   | Map [(Data, Data)]
   | List [Data]
   | I Integer
   | B ByteString
```

In set-theoretic terms the denotation of data is defined to be the least fixed point of the endofunctor F on the category of sets given by  $F(X) = ([[integer]] \times X^*) \uplus (X \times X)^* \uplus X^* \uplus [[integer]] \uplus [[bytestring]]$ , so that

```
[data] = ([integer] \times [data]^*) \uplus ([data] \times [data])^* \uplus [data]^* \uplus [integer] \uplus [bytestring].
```

We have injections

```
\begin{split} & \operatorname{inj}_C : \llbracket \operatorname{integer} \rrbracket \times \llbracket \operatorname{data} \rrbracket^* \to \llbracket \operatorname{data} \rrbracket \\ & \operatorname{inj}_M : \llbracket \operatorname{data} \rrbracket \times \llbracket \operatorname{data} \rrbracket^* \to \llbracket \operatorname{data} \rrbracket \\ & \operatorname{inj}_L : \llbracket \operatorname{data} \rrbracket^* \to \llbracket \operatorname{data} \rrbracket \\ & \operatorname{inj}_I : \llbracket \operatorname{integer} \rrbracket \to \llbracket \operatorname{data} \rrbracket \\ & \operatorname{inj}_R : \llbracket \operatorname{bytestring} \rrbracket \to \llbracket \operatorname{data} \rrbracket \end{split}
```

and projections

```
\begin{split} &\operatorname{proj}_C: \llbracket \operatorname{data} \rrbracket \to (\llbracket \operatorname{integer} \rrbracket \times \llbracket \operatorname{data} \rrbracket^*)_{\mathsf{X}} \\ &\operatorname{proj}_M: \llbracket \operatorname{data} \rrbracket \to (\llbracket \operatorname{data} \rrbracket \times \llbracket \operatorname{data} \rrbracket^*)_{\mathsf{X}} \\ &\operatorname{proj}_L: \llbracket \operatorname{data} \rrbracket \to \llbracket \operatorname{data} \rrbracket^*_{\mathsf{X}} \\ &\operatorname{proj}_I: \llbracket \operatorname{data} \rrbracket \to \llbracket \operatorname{integer} \rrbracket_{\mathsf{X}} \\ &\operatorname{proj}_B: \llbracket \operatorname{data} \rrbracket \to \llbracket \operatorname{bytestring} \rrbracket_{\mathsf{X}} \end{split}
```

which extract an object of the relevant type from a data object D, returning  $\times$  if D does not lie in the expected component of the disjoint union; also there are functions

$$is_C, is_M, is_I, is_I, is_B : [data] \rightarrow [bool]$$

which determine whether a data value lies in the relevant component.

**Note:** Constr tag values. The Constr constructor of the data type is intended to represent values from algebraic data types (also known as sum types and discriminated unions, among other things; data itself is an example of such a type), where Constr i [ $d_1, \ldots, d_n$ ] represents a tuple of data items together with a tag i indicating which of a number of alternatives the data belongs to. The definition above allows tags to be any integer value, but because of restrictions in the serialisation format for data (see Section D.7) we recommend that in practice **only tags** i **with**  $0 \le i \le 2^{64} - 1$  **should be used**: descrialisation will fail for data items (and programs which include such items) involving tags outside this range.

## A.2 Alonzo built-in functions

The default set of built-in functions for the Alonzo release is shown in Table 3. The table indicates which functions can fail during execution, and conditions causing failure are specified either in the denotation given in the table or in a relevant note. Recall also that a built-in function will fail if it is given an argument of the wrong type: this is checked in conditions involving the  $\sim$  relation and the Eval function in Figures 7 and 10. Note also the some of the functions are #-polymorphic. According to Section 4.2.3 we require a denotation for every possible monomorphisation of these; however all of these functions are parametrically polymorphic so to simplify notation we have given a single denotation for each of them with an implicit assumption that it applies at each possible monomorphisation in an obvious way.

Function	Signature	Denotation	Can Fail?	Note
addInteger	[integer, integer] → integer	+		
subtractInteger	[integer, integer] → integer	_		
multiplyInteger	[integer, integer] → integer	×		
${ t divideInteger}$	[integer, integer] → integer	div	Yes	1
modInteger	[integer, integer] → integer	mod	Yes	1
${\tt quotientInteger}$	[integer, integer] → integer	quot	Yes	1
${\tt remainderInteger}$	$[integer, integer] \rightarrow integer$	rem	Yes	1
equalsInteger	[integer, integer] $\rightarrow$ bool	=		
${\tt lessThanInteger}$	[integer, integer] $\rightarrow$ bool	<		
${\tt lessThanEqualsInteger}$	[integer, integer] $\rightarrow$ bool	≤		
appendByteString	[bytestring, bytestring]	$([c_1, \dots, c_m], [d_1, \dots, d_n])$		
	ightarrow bytestring	$\mapsto [c_1, \dots, c_m, d_1, \dots, d_n]$		
consByteString	[integer, bytestring]	$(c,[c_1,\ldots,c_n])$		
	ightarrow bytestring	$\mapsto [\operatorname{mod}(c, 256), c_1, \dots, c_n]$		
sliceByteString	[integer, integer, bytestring]	$(s,k,[c_0,\ldots,c_n])$		2
	ightarrow bytestring	$\mapsto [c_{\max(s,0)}, \dots, c_{\min(s+k-1,n-1)}]$		
lengthOfByteString	$[bytestring] \rightarrow integer$	$[] \mapsto 0, [c_1, \dots, c_n] \mapsto n$		
${\tt indexByteString}$	[bytestring, integer]	$([c_0,\ldots,c_{n-1}],j)$	Yes	
	$\rightarrow$ integer	$\int c_i$ if $0 < j < n-1$		
		$\mapsto \begin{cases} c_i & \text{if } 0 \le j \le n-1 \\ \times & \text{otherwise} \end{cases}$		
agual aBrot aCt min m	[butostring butostring]			3
equalsByteString	[bytestring, bytestring]  → bool	=		3
lessThanByteString				3
ressinancy testing	[bytestring, bytestring]  → bool	<		3
lessThanEqualsByteString	bool   [bytestring, bytestring]	<b>≤</b>		3
ressinanequarsbytestring	→ bool			3
appendString	$[string, string] \rightarrow string$	$([u_1,\ldots,u_m],[v_1,\ldots,v_n])$		
		$\mapsto [u_1, \dots, u_m, v_1, \dots, v_n]$		
equalsString	$[string, string] \rightarrow bool$	=		
encodeUtf8	$[string] \rightarrow bytestring$	utf8		4
decodeUtf8	[bytestring] → string	utf8 <sup>-1</sup>	Yes	4
sha2_256	[bytestring] → bytestring	Hash a bytestring using SHA256.		
sha3_256	[bytestring] → bytestring	Hash a bytestring using SHA3-		
		256.		

Table 3: Built-in Functions

blake2b_256  verifyEd25519Signature  ifThenElse  chooseUnit	$[bytestring]  o bytestring$ $[bytestring, bytestring, bytestring]  o bool [\forall a_*, bool, a_*, a_*]  o a_* [\forall a_*, unit, a_*]  o a_*$	Hash a bytestring using Blake2B256. Verify an Ed25519 digital signature. $(\text{true}, t_1, t_2) \mapsto t_1$	Fail? Yes	
verifyEd25519Signature	$ \begin{array}{c} \text{[bytestring, bytestring,} \\ \text{bytestring]} \rightarrow \text{bool} \\ [\forall a_*, \text{bool}, a_*, a_*] \rightarrow a_* \end{array} $	Blake2B256. Verify an Ed25519 digital signature.	Yes	
ifThenElse	$\begin{array}{c} \text{bytestring}] \rightarrow \text{bool} \\ [\forall a_*, \text{bool}, a_*, a_*] \rightarrow a_* \end{array}$	ture.	Yes	·
ifThenElse	$[\forall a_*, bool, a_*, a_*] \to a_*$			5, 6
		$(true, t_1, t_2) \mapsto t_1$		l
choosellnit	$[\forall a \text{ unit. } a ] \rightarrow a$			l
choosellnit	$[\forall a \text{ unit. } a  1 \rightarrow a$	$(false, t_1, t_2) \mapsto t_2$		l
CHOOPEONIT	[ • a <sub>*</sub> , a <sub>111</sub> o, a <sub>*</sub> ]	$((),t)\mapsto t$		l
trace	$[\forall a_*, \mathtt{string}, a_*] \rightarrow a_*$	$(s,t)\mapsto t$		7
fstPair	$[\forall a_\#, \forall b_\#, \mathtt{pair}(a_\#, b_\#)]  o a_\#$	$(x,y)\mapsto x$		l
sndPair	$[\forall a_{\#}, \forall b_{\#}, \mathtt{pair}(a_{\#}, b_{\#})] \rightarrow b_{\#}$	$(x,y)\mapsto y$		l
chooseList	$[\forall a_\#, \forall b_*, \mathtt{list}(a_\#), b_*, b_*]  o b_*$	$([],t_1,t_2)\mapsto t_1,$		l
		$([x_1, \dots, x_n], t_1, t_2) \mapsto t_2 \ (n \ge 1).$		l
mkCons	$[\forall a_\#, a_\#, \mathtt{list}(a_\#)] \to \mathtt{list}(a_\#)$	$(x, [x_1, \dots, x_n]) \mapsto [x, x_1, \dots, x_n]$		l
headList	$[\forall a_\#, \mathtt{list}(a_\#)]  o a_\#$	$[] \mapsto x, [x_1, x_2, \dots, x_n] \mapsto x_1$	Yes	l
tailList	$[\forall a_{\#}, \mathtt{list}(a_{\#})] \rightarrow \mathtt{list}(a_{\#})$	[] → <b>x</b> ,	Yes	l
		$[x_1, x_2, \dots, x_n] \mapsto [x_2, \dots, x_n]$		l
nullList	$[\forall a_{\#}, \mathtt{list}(a_{\#})] \rightarrow \mathtt{bool}$	$[] \mapsto true, [x_1, \dots, x_n] \mapsto false$		l
chooseData	$[\forall a_*, \mathtt{data}, a_*, a_*, a_*, a_*, a_*]  ightarrow a_*$	$(d, t_C, t_M, t_L, t_I, t_B)$		l
		$\int t = ific (d)$		l
		$t_C$ if is $t_C(d)$		l
		$\int_{t}^{t_{M}} \inf_{i \in S_{M}(d)} S_{M}(d)$		l
		$\mapsto \begin{cases} t_M & \text{if is}_M(d) \\ t_L & \text{if is}_L(d) \\ t_I & \text{if is}_I(d) \end{cases}$		l
		$ \begin{vmatrix} t_I & \text{if } \text{is}_I(u) \\ t_B & \text{if } \text{is}_B(d) \end{vmatrix} $		l
		$(B \cap B \cap B \cap B)$		
constrData	[integer, list (data)] → data	$\inf_{C}$		
mapData	[list(pair(data,data))  → data	$\operatorname{inj}_M$		
listData	[list(data)] → data	$ \operatorname{inj}_L $		l
iData	[integer] → data	$\inf_{I}$		l
bData	[bytestring] → data	$\inf_{B}$		l
unConstrData	[data]	$ \operatorname{proj}_C $	Yes	l
	$\rightarrow$ pair (integer, list (data))			l
unMapData	[data]	$\operatorname{proj}_M$	Yes	l
-	→ list(pair(data,data))	1 VIII		l
unListData	[data] → list(data)	$\operatorname{proj}_L$	Yes	l
unIData	[data] → integer	$  proj_I $	Yes	l
unBData	[data] → bytestring	$ \operatorname{proj}_{B} $	Yes	l
equalsData	[data, data] → bool			l
mkPairData	[data, data]	$(x,y)\mapsto (x,y)$		l
	→ pair (data, data)			l
mkNilData	$[unit] \rightarrow list(data)$	() → []		l
mkNilPairData	[unit]	() → []		l
	→ list(pair(data,data))			l

Table 3: Built-in Functions (continued)

Note 1. Integer division functions. We provide four integer division functions: divideInteger,

modInteger, quotientInteger, and remainderInteger, whose denotations are mathematical functions div, mod, quot, and rem which are modelled on the corresponding Haskell operations. Each of these takes two arguments and will fail (returning  $\times$ ) if the second one is zero. For all  $a, b \in \mathbb{Z}$  with  $b \neq 0$  we have

$$\operatorname{div}(a, b) \times b + \operatorname{mod}(a, b) = a$$

$$|\operatorname{mod}(a, b)| < |b|$$

and

$$quot(a, b) \times b + rem(a, b) = a$$

$$|\operatorname{rem}(a, b)| < |b|$$
.

The div and mod functions form a pair, as do quot and rem; div should not be used in combination with mod, not should quot be used with mod.

For positive divisors b, div truncates downwards and mod always returns a non-negative result  $(0 \le \text{mod}(a, b) \le b - 1)$ . The quot function truncates towards zero. Table 4 shows how the signs of the outputs of the division functions depend on the signs of the inputs; + means  $\ge 0$  and - means  $\le 0$ , but recall that for b = 0 all of these functions return the error value  $\times$ .

a	b	div	mod	quot	rem
+	+	+	+	+	+
-	+	_	+	–	_
+	_	_	+	+	+
-	_	+	_	+	_

Table 4: Behaviour of integer division functions

**Note 2. The** sliceByteString **function.** The application [[(builtin sliceByteString) (con integer s)] (con integer k)] (con bytestring b)] returns the substring of b of length k starting at position s; indexing is zero-based, so a call with s=0 returns a substring starting with the first element of b, s=1 returns a substring starting with the second, and so on. This function always succeeds, even if the arguments are out of range: if  $b=[c_0,\ldots,c_{n-1}]$  then the application above returns the substring  $[c_i,\ldots,c_i]$  where  $i=\max(s,0)$  and  $j=\min(s+k-1,n-1)$ ; if j< i then the empty string is returned.

Note 3. Comparisons of bytestrings. Bytestrings are ordered lexicographically in the usual way. If we have  $a = [a_1, ..., a_m]$  and  $b = [b_1, ..., b_n]$  then (recalling that if m = 0 then a = [], and similarly for b),

- a = b if and only if m = n and  $a_i = b_i$  for  $1 \le i \le m$ .
- $a \le b$  if and only if one of the following holds:
  - -a=[]
  - m, n > 0 and  $a_1 < b_1$
  - -m, n > 0 and  $a_1 = b_1$  and  $[a_2, \dots, a_m] \le [b_2, \dots, b_n]$ .
- a < b if and only if  $a \le b$  and  $a \ne b$ .

For example, #23456789 < #24 and #2345 < #234500. The empty bytestring is equal only to itself and is strictly less than all other bytestrings.

**Note 4. Encoding and decoding bytestrings.** The encodeUtf8 and decodeUtf8 functions convert between the string type and the bytestring type. We have defined [string] to consist of sequences of Unicode characters without specifying any particular character representation, whereas [bytestring] consists of sequences of 8-bit bytes. We define the denotation of encodeUtf8 to be the function

utf8: 
$$\mathbb{U}^* \to \mathbb{B}^*$$

which converts sequences of Unicode characters to sequences of bytes using the well-known UTF-8 character encoding [24, Definition D92]. The denotation of decodeUtf8 is the partial inverse function

$$utf8^{-1}: \mathbb{B}^* \to \mathbb{U}^*_{\checkmark}.$$

UTF-8 encodes Unicode characters encoded using between one and four bytes: thus in general neither function will preserve the length of an object. Moreover, not all sequences of bytes are valid representations of Unicode characters, and decodeUtf8 will fail if it receives an invalid input (but encodeUtf8 will always succeed).

**Note 5. Digital signature verification functions.** We use a uniform interface for digital signature verification algorithms. A digital signature verification function takes three bytestring arguments (in the given order):

- a public key vk (in this context vk is also known as a verification key)
- a message m
- a signature s.

A signature verification function may require one or more arguments to be well-formed in some sense (in particular an argument may need to be of a specified length), and in this case the function will fail (returning  $\times$ ) if any argument is malformed. If all of the arguments are well-formed then the verification function returns true if the private key corresponding to vk was used to sign the message m to produce s, otherwise it returns false.

Note 6. Ed25519 signature verification. The verifyEd25519Signature function<sup>†</sup> performs cryptographic signature verification using the Ed25519 scheme [5, 18], and conforms to the interface described in Note 5. The arguments must have the following sizes:

- vk: 32 bytes
- m: unrestricted
- s: 64 bytes.

**Note 7. The** trace **function.** An application [(builtin trace) s v] (s a string, v any Plutus Core value) returns v. We do not specify the semantics any further. An implementation may choose to discard s or to perform some side-effect such as writing it to a terminal or log file.

<sup>†</sup>verifyEd25519Signature was formerly called verifySignature but was renamed to avoid ambiguity when further signature verification functions were introduced in the Vasil release (see Section B.2).

## **Appendix B Built-in Types and Functions Supported in the Vasil Release**

The Vasil release of Cardano (June 2022) extends the set of built-in functions slightly.

## **B.1** Built-in types and type operators

The built-in types and type operators remain unchanged from the Alonzo release (Appendix A.1).

## **B.2** Built-in functions

The Vasil release continues to support the Alonzo built-in functions (Table 3) and adds three new ones: these are described in Table 5.

Function	Signature	Denotation	Can	Note
			Fail?	
serialiseData	$[data] \rightarrow bytestring$	$\mathcal{E}_{ ext{data}}$		1
verifyEcdsaSecp256k1Signature	[bytestring, bytestring,	Verify an SECP-256k1	Yes	2
	$bytestring] \rightarrow bool$	ECDSA signature		
verifySchnorrSecp256k1Signature	[bytestring, bytestring,	Verify an SECP-256k1	Yes	3
	bytestring] → bool	Schnorr signature		

Table 5: Built-in Functions

Note 1. Serialising data objects. The serialiseData function takes a data object and converts it into a bytestring using a CBOR encoding. A full specification of the encoding (including the definition of  $\mathcal{E}_{\text{data}}$ ) is provided in Appendix D.

**Note 2.** Secp256k1 ECDSA Signature verification. The verifyEcdsaSecp256k1Signature function performs elliptic curve digital signature verification [1, 2, 17] over the secp256k1 curve [9, §2.4.1] and conforms to the interface described in Note 5 of Section A.2. The arguments must have the following sizes:

- vk: 33 bytes
- m: 32 bytes
- s: 64 bytes.

The public key vk is expected to be in the 33-byte compressed form described in [6]. Moreover, the ECDSA scheme admits two distinct valid signatures for a given message and private key, and we follow the restriction imposed by Bitcoin (see [20], LOW\_S) and **only accept the smaller signature**; verifyEcdsa-Secp256k1Signature will return false if the larger one is supplied.

**Note 3.** Secp256k1 Schnorr Signature verification. The verifySchnorrSecp256k1Signature function performs verification of Schnorr signatures [23, 19] over the secp256k1 curve and conforms to the interface described in Note 5 of Section A.2. The arguments are expected to be of the forms specified in BIP-340 [19] and thus should have the following sizes:

• vk: 32 bytes

- m: unrestricted
- s: 64 bytes.

## **Appendix C** Formally Verified Behaviours

To follow.

## Appendix D Serialising data Objects Using the CBOR Format

## **D.1** Introduction

In this section we define a CBOR encoding for the data type introduced in Section A.1. For ease of reference we reproduce the definition of the Haskell Data type, which we may regard as the definition of the Plutus data type. Other representations are of course possible, but this is useful for the present discussion.

```
data Data =
   Constr Integer [Data]
   | Map [(Data, Data)]
   | List [Data]
   | I Integer
   | B ByteString
```

The CBOR encoding defined here uses basic CBOR encodings as defined in the CBOR standard [8], but with some refinements. Specifically

- We use a restricted encoding for bytestrings which requires that bytestrings are serialised as sequences of blocks, each block being at most 64 bytes long. Any encoding of a bytestring using our scheme is valid according to the CBOR specification, but the CBOR specification permits some encodings which we do not accept. The purpose of the size restriction is to prevent arbitrary data from bring stored on the blockchain.
- Large integers (less than −2<sup>64</sup> or greater than 2<sup>64</sup> − 1) are encoded via the restricted bytestring encoding; other integers are encoded as normal. Again, our restricted encodings are compatible with the CBOR specification.
- The Constr case of the data type is encoded using a scheme which is an early version of a proposed extension of the CBOR specification to include encodings for discriminated unions. See [10] and [7, Section 9.1].

## **D.2** Notation

We introduce some extra notation for use here and in Appendix E.

The notation  $f: X \to Y$  indicates that f is a partial map from X to Y. We denote the empty bytestring by  $\epsilon$  and (as in 2.2) use  $\ell(s)$  to denote the length of a bytestring s and  $\cdot$  to denote the concatenation of two bytestrings, and also the operation of prepending or appending a byte to a bytestring. We will also make use of the div and mod functions described in Note 1 in Appendix A.

**Encoders and decoders.** Recall that  $\mathbb{B} = \mathbb{N}_{[0,255]}$ , the set of integral values that can be represented in a single byte, and that we identify bytestrings with elements of  $\mathbb{B}^*$ . We will describe the CBOR encoding of the data type by defining families of encoding functions (or *encoders*)

$$\mathcal{E}_X:X\to\mathbb{B}^*$$

and decoding functions (or decoders)

$$\mathcal{D}_{X}:\mathbb{B}^{*} \to \mathbb{B}^{*} \times X$$

for various sets X, such as the set  $\mathbb{Z}$  of integers and the set of all data items. The encoding function  $\mathcal{E}_X$  takes an element  $x \in X$  and converts it to a bytestring, and the decoding function  $\mathcal{D}_X$  takes a bytestring s, decodes some initial prefix of s to a value  $x \in X$ , and returns the remainder of s together with s. Decoders for complex types will often be built up from decoders for simpler types. Decoders are *partial* functions because they can fail, for instance, if there is insufficient input, or if the input is not well formed, or if a decoded value is outside some specified range.

Many of the decoders which we define below involve a number of cases for different forms of input, and we implicitly assume that the decoder fails if none of the cases applies. We also assume that if a decoder fails then so does any other decoder which invokes it, so any failure when attempting to decode a particular data item in a bytestring will cause the entire decoding process to fail (immediately).

## **D.3** The CBOR format

A CBOR-encoded item consists of a bytestring beginning with a *head* which occupies 1,2,3,5, or 9 bytes. Depending on the contents of the head, some sequence of bytes following it may also contribute to the encoded item. The first three bits of the head are interpreted as a natural number between 0 and 7 (the *major type*) which gives basic information about the type of the following data. The remainder of the head is called the *argument* of the head and is used to encode further information, such as the value of an encoded integer or the size of a list of encoded items. Encodings of complex objects may occupy the bytes following the head, and these will typically contain further encoded items.

## D.4 Encoding and decoding the heads of CBOR items

For  $i \in \mathbb{N}$  we define a function  $b_i : \mathbb{N} \to \mathbb{B}$  which returns the *i*-th byte of an integer, with the 0-th byte being the least significant:

$$b_i(n) = \text{mod}(\text{div}(n, 256^i), 256).$$

We use this to define for each  $k \ge 1$  a partial function  $e_k : \mathbb{N} \to \mathbb{B}^*$  which converts a sufficiently small integer to a bytestring of length k (possibly with leading zeros):

$$e_k(n) = [b_{k-1}(n), \dots, b_0(n)]$$
 if  $n \le 256^k - 1$ .

This function fails if the input is too large to fit into a *k*-byte bytestring.

We also define inverse functions  $d_k : \mathbb{B}^* \to \mathbb{N}$  which decode a k-byte natural number from the start of a bytestring, failing if there is insufficient input:

$$d_k(s) = (s', \sum_{i=0}^{k-1} 256^i b_i)$$
 if  $s = [b_{k-1}, \dots, b_0] \cdot s'$ .

We now define an encoder  $\mathcal{E}_{\text{head}}: \mathbb{N}_{[0,7]} \times \mathbb{N}_{[0,2^{64}-1]} \to \mathbb{B}^*$  which takes a major type and a natural number and encodes them as a CBOR head using the standard encoding:

$$\mathcal{E}_{\mathsf{head}}(m,n) = \begin{cases} [32m+n] & \text{if } n \leq 23 \\ (32m+24) \cdot \mathsf{e}_1(n) & \text{if } 24 \leq n \leq 255 \\ (32m+25) \cdot \mathsf{e}_2(n) & \text{if } 256 \leq n \leq 256^2 - 1 \\ (32m+26) \cdot \mathsf{e}_4(n) & \text{if } 256^2 \leq n \leq 256^4 - 1 \\ (32m+27) \cdot \mathsf{e}_8(n) & \text{if } 256^4 \leq n \leq 256^8 - 1. \end{cases}$$

The corresponding decoder  $\mathcal{D}_{\mathsf{head}}:\mathbb{B}^* \to \mathbb{B}^* \times \mathbb{N}_{[0,7]} \times \mathbb{N}_{[0,2^{64}-1]}$  is given by

$$\mathcal{D}_{\mathsf{head}}(n \cdot s) = \begin{cases} (s, \mathsf{div}(n, 32), \mathsf{mod}(n, 32)) & \text{if } \mathsf{mod}(n, 32) \leq 23 \\ (s', \mathsf{div}(n, 32), k) & \text{if } \mathsf{mod}(n, 32) = 24 \text{ and } \mathsf{d}_1(s) = (s', k) \\ (s', \mathsf{div}(n, 32), k) & \text{if } \mathsf{mod}(n, 32) = 25 \text{ and } \mathsf{d}_2(s) = (s', k) \\ (s', \mathsf{div}(n, 32), k) & \text{if } \mathsf{mod}(n, 32) = 26 \text{ and } \mathsf{d}_4(s) = (s', k) \\ (s', \mathsf{div}(n, 32), k) & \text{if } \mathsf{mod}(n, 32) = 27 \text{ and } \mathsf{d}_8(s) = (s', k). \end{cases}$$

This function is undefined if the input is the empty bytestring  $\epsilon$ , if the input is too short, or if its initial byte is not of the expected form.

**Heads for indefinite-length items.** The functions  $\mathcal{E}_{head}$  and  $\mathcal{D}_{head}$  defined above are used for a number of purposes. One use is to encode integers less than 64 bits, where the argument of the head is the relevant integer. Another use is for "definite-length" encodings of items such as bytestrings and lists, where the head contains the length n of the object and is followed by some encoding of the object itself (for example a sequence of n bytes for a bytestring or a sequence of n encoded objects for the elements of a list). It is also possible to have "indefinite-length" encodings of objects such as lists and arrays, which do not specify the length of an object in advance: instead a special head with argument 31 is emitted, followed by the encodings of the individual items; the end of the sequence is marked by a "break" byte with value 255. We define an encoder  $\mathcal{E}_{indef}$ :  $\mathbb{N}_{[2,5]} \to \mathbb{B}^*$  and a decoder  $\mathcal{D}_{indef}$ :  $\mathbb{B}^* \to \mathbb{B}^* \times \mathbb{N}_{[2,5]}$  which deal with indefinite heads for a given major type:

$$\mathcal{E}_{\mathsf{indef}}(m) = [32m + 31]$$
 
$$\mathcal{D}_{\mathsf{indef}}(n \cdot s) = (s, m) \quad \text{if } n = 32m + 31.$$

Note that  $\mathcal{E}_{indef}$  and  $\mathcal{D}_{indef}$  are only defined for  $m \in \{2, 3, 4, 5\}$  (and we shall only use them in these cases). The case m = 31 corresponds to the break byte and for  $m \in \{0, 1, 6\}$  the value is not well formed: see [8, 3.2.4].

## D.5 Encoding and decoding bytestrings

The standard CBOR encoding of bytestrings encodes a bytestring as either a definite-length sequence of bytes (the length being given in the head) or as an indefinite-length sequence of definite-length "chunks" (see [8, §§3.1 and 3.4.2]). We use a similar scheme, but only allow chunks of length up to 64. To this end, suppose that  $a = [a_1, \ldots, a_{64k+r}] \in \mathbb{B}^* \setminus \{\epsilon\}$  where  $k \ge 0$  and  $0 \le r \le 63$ . We define the *canonical 64-byte decomposition*  $\bar{a}$  of a to be

$$\bar{a} = [[a_1, \dots, a_{64}], [a_{65}, \dots, a_{128}], \dots, [a_{64(k-1)+1}, \dots, a_{64k}]] \in (\mathbb{B}^*)^*$$

if r = 0 and

$$\bar{a} = [[a_1, \dots, a_{64}], [a_{65}, \dots, a_{128}], \dots, [a_{64(k-1)+1}, \dots, a_{64k}], [a_{64k+1}, \dots, a_{64k+r}]] \in (\mathbb{B}^*)^*$$

if r > 0. The canonical decomposition of the empty list is  $\bar{\epsilon} = []$ .

We define the encoder  $\mathcal{E}_{\mathbb{B}^*}$ :  $\mathbb{B}^* \to \mathbb{B}^*$  for bytestrings by encoding bytestrings of size up to 64 using the standard CBOR encoding and encoding larger bytestrings by breaking them up into 64-byte chunks (with the final chunk possibly being less than 64 bytes long) and encoding them as an indefinite-length list (major type 2 indicates a bytestring):

$$\mathcal{E}_{\mathbb{B}^*}(s) = \begin{cases} \mathcal{E}_{\mathsf{head}}(2, \ell(s)) \cdot s & \text{if } \ell(s) \leq 64 \\ \mathcal{E}_{\mathsf{indef}}(2) \cdot \mathcal{E}_{\mathsf{head}}(2, \ell(c_1)) \cdot c_1 \cdot \mathcal{E}_{\mathsf{head}}(2, \ell(c_2)) \cdot \cdots \\ \cdots \cdot c_{n-1} \cdot \mathcal{E}_{\mathsf{head}}(2, \ell(c_n)) \cdot c_n \cdot 255 & \text{if } \ell(s) > 64 \text{ and } \bar{s} = [c_1, \dots, c_n]. \end{cases}$$

The decoder is slightly more complicated. Firstly, for every  $n \ge 0$  we define a decoder  $\mathcal{D}_{\text{bytes}}^{(n)}$ :  $\mathbb{B}^* \times \mathbb{B}^*$  which extracts an n-byte prefix from its input (failing in the case of insufficient input):

$$\mathcal{D}_{\mathsf{bytes}}^{(n)}(s) = \begin{cases} (s, \epsilon) & \text{if } n = 0\\ (s'', b \cdot t) & \text{if } s = b \cdot s' \text{ and } \mathcal{D}_{\mathsf{bytes}}^{(n-1)}(s') = (s'', t). \end{cases}$$

Secondly, we define a decoder  $\mathcal{D}_{block}$ :  $\mathbb{B}^* \to \mathbb{B}^* \times \mathbb{B}^*$  which attempts to extract a bytestring of length at most 64 from its input;  $\mathcal{D}_{block}$  (and any other function which calls it) will fail if it encounters a bytestring which is greater than 64 bytes.

$$\mathcal{D}_{\text{block}}(s) = \mathcal{D}_{\text{bytes}}^{(n)}(s')$$
 if  $\mathcal{D}_{\text{head}}(s) = (s', 2, n)$  and  $n \le 64$ .

Thirdly, we define a decoder  $\mathcal{D}_{blocks}$ :  $\mathbb{B}^* \to \mathbb{B}^* \times \mathbb{B}^*$  which decodes a sequence of blocks and returns their concatenation.

$$\mathcal{D}_{\mathsf{blocks}}(s) = \begin{cases} (s', \epsilon) & \text{if } s = 255 \cdot s' \\ (s'', t \cdot t') & \text{if } \mathcal{D}_{\mathsf{block}}(s) = (s', t) \text{ and } \mathcal{D}_{\mathsf{blocks}}(s') = (s'', t'). \end{cases}$$

Finally we define the decoder  $\mathcal{D}_{\mathbb{R}^*}: \mathbb{B}^* \to \mathbb{B}^* \times \mathbb{B}^*$  for bytestrings by

$$\mathcal{D}_{\mathbb{B}^*}(s) = \begin{cases} (s', t) & \text{if } \mathcal{D}_{\mathsf{block}}(s) = (s', t) \\ \mathcal{D}_{\mathsf{blocks}}(s') & \text{if } \mathcal{D}_{\mathsf{indef}}(s) = (s', 2). \end{cases}$$

This looks for either a single block or an indefinite-length list of blocks, in the latter case returning their concatenation. It will accept the output of  $\mathcal{E}_{\mathbb{B}^*}$  but will reject bytestring encodings containing any blocks greater than 64 bytes long, even if they are valid bytestring encodings according to the CBOR specification.

## D.6 Encoding and decoding integers

As with bytestrings we use a specialised encoding scheme for integers which prohibits encodings with overly-long sequences of arbitrary data. We encode integers in  $\mathbb{N}_{[-2^{64},2^{64}-1]}$  as normal (see [8, §3.1]: the major type is 0 for positive integers and 1 for negative ones) and larger ones by emitting a CBOR tag (major type 6; argument 2 for positive numbers and 3 for negative numbers) to indicate the sign, then converting the integer to a bytestring and emitting that using the encoder defined above. This encoding scheme is the same as the standard one except for the size limitations.

We firstly define conversion functions itos:  $\mathbb{N} \to \mathbb{B}^*$  and stoi:  $\mathbb{B}^* \to \mathbb{N}$  by

$$itos(n) = \begin{cases} \epsilon & \text{if } n = 0\\ itos(div(n, 256)) \cdot mod(n, 256) & \text{if } n > 0. \end{cases}$$

and

$$\mathsf{stoi}(l) = \begin{cases} 0 & \text{if } l = \epsilon \\ 256 \times \mathsf{stoi}(l') + n & \text{if } l = l' \cdot n \text{ with } n \in \mathbb{B}. \end{cases}$$

The encoder  $\mathcal{E}_{\mathbb{Z}}:\mathbb{Z}\to\mathbb{B}^*$  for integers is now defined by

$$\mathcal{E}_{\mathbb{Z}}(n) = \begin{cases} \mathcal{E}_{\mathsf{head}}(0,n) & \text{if } 0 \leq n \leq 2^{64} - 1 \\ \mathcal{E}_{\mathsf{head}}(6,2) \cdot \mathcal{E}_{\mathbb{B}^*}(\mathsf{itos}(n)) & \text{if } n \geq 2^{64} \\ \mathcal{E}_{\mathsf{head}}(1,-n-1) & \text{if } -2^{64} \leq n \leq -1 \\ \mathcal{E}_{\mathsf{head}}(6,3) \cdot \mathcal{E}_{\mathbb{B}^*}(\mathsf{itos}(-n-1)) & \text{if } n \leq -2^{64} - 1. \end{cases}$$

The decoder  $\mathcal{D}_{\mathbb{Z}}: \mathbb{B}^* \to \mathbb{B}^* \times \mathbb{Z}$  inverts this process. The decoder is in fact slightly more permissive than the encoder because it also accepts small integers encoded using the scheme for larger ones. However, the CBOR standard permits integer encodings which contain bytestrings longer than 64 bytes and it will not accept those.

$$\mathcal{D}_{\mathbb{Z}}(s) = \begin{cases} (s', n) & \text{if } \mathcal{D}_{\mathsf{head}}(s) = (s', 0, n) \\ (s', -n - 1) & \text{if } \mathcal{D}_{\mathsf{head}}(s) = (s', 1, n) \\ (s'', \mathsf{stoi}(b)) & \text{if } \mathcal{D}_{\mathsf{head}}(s) = (s', 6, 2) \text{ and } \mathcal{D}_{\mathbb{B}^*}(s') = (s'', b) \\ (s'', -\mathsf{stoi}(b) - 1) & \text{if } \mathcal{D}_{\mathsf{head}}(s) = (s', 6, 3) \text{ and } \mathcal{D}_{\mathbb{B}^*}(s') = (s'', b). \end{cases}$$

## D.7 Encoding and decoding data

It is now quite straightforward to encode most data values. The main complication is in the encoding of constructor tags (the number i in Constr i l).

**The encoder.** The encoder is given by

$$\begin{split} \mathcal{E}_{\text{data}}(\text{Map }l) &= \mathcal{E}_{\text{head}}(5, \mathcal{E}(l)) \cdot \mathcal{E}_{(\text{data}^2)^*}(l) \\ \mathcal{E}_{\text{data}}(\text{List }l) &= \mathcal{E}_{\text{indef}}(4) \cdot \mathcal{E}_{\text{data}^*}(l) \cdot 255 \\ \mathcal{E}_{\text{data}}(\text{Constr }i\,l) &= \mathcal{E}_{\text{ctag}}(i) \cdot \mathcal{E}_{\text{indef}}(4) \cdot \mathcal{E}_{\text{data}^*}(l) \cdot 255 \\ \mathcal{E}_{\text{data}}(\text{I}\,n) &= \mathcal{E}_{\mathbb{Z}}(n) \\ \mathcal{E}_{\text{data}}(\text{B}\,s) &= \mathcal{E}_{\mathbb{B}^*}(s). \end{split}$$

This definition uses encoders for lists of data items, lists of pairs of data items, and constructor tags as follows:

$$\mathcal{E}_{\text{data}^*}([d_1, \dots, d_n]) = \mathcal{E}_{\text{data}}(d_1) \cdot \dots \cdot \mathcal{E}_{\text{data}}(d_n)$$

$$\mathcal{E}_{(\text{data}^2)^*}([(k_1, d_1), \dots, (k_n, d_n)]) = \mathcal{E}_{\text{data}}(k_1) \cdot \mathcal{E}_{\text{data}}(d_1) \cdot \dots \cdot \mathcal{E}_{\text{data}}(k_n) \cdot \mathcal{E}_{\text{data}}(d_n)$$

$$\mathcal{E}_{\text{ctag}}(i) = \begin{cases} \mathcal{E}_{\text{head}}(6, 121 + i) & \text{if } 0 \leq i \leq 6 \\ \mathcal{E}_{\text{head}}(6, 1280 + (i - 7)) & \text{if } 7 \leq i \leq 127 \\ \mathcal{E}_{\text{head}}(6, 102) \cdot \mathcal{E}_{\text{head}}(4, 2) \cdot \mathcal{E}_{\mathbb{Z}}(i) & \text{otherwise.} \end{cases}$$

In the final case of  $\mathcal{E}_{\text{ctag}}$  we emit a head with major type 4 and argument 2. This indicates that an encoding of a list of length 2 will follow: the first element of the list is the constructor number and the second is the argument list of the constructor, which is actually encoded in  $\mathcal{E}_{\text{data}}$ . It might be conceptually more accurate to have a single encoder which would encode both the constructor tag and the argument list, but this would increase the complexity of the notation even further. Similar remarks apply to  $\mathcal{D}_{\text{ctag}}$  below.

**The decoder.** The decoder is given by

$$\mathcal{D}_{\text{data}}(s) = \begin{cases} (s'', \text{Map } l) & \text{if } \mathcal{D}_{\text{head}}(s) = (s', 5, n) \text{ and } \mathcal{D}_{(\text{data}^2)^*}^{(n)}(s') = (s'', l) \\ (s', \text{List } l) & \text{if } \mathcal{D}_{\text{data}^*}(s) = (s', l) \\ (s'', \text{Constr } i \, l) & \text{if } \mathcal{D}_{\text{ctag}}(s) = (s', i) \text{ and } \mathcal{D}_{\text{data}^*}(s') = (s'', l) \\ (s', \text{I } n) & \text{if } \mathcal{D}_{\mathbb{Z}}(s) = (s', n) \\ (s', \text{B } b) & \text{if } \mathcal{D}_{\mathbb{B}^*}(s) = (s', b) \end{cases}$$

where

$$\mathcal{D}_{\text{data*}}(s) = \begin{cases} \mathcal{D}_{\text{data*}}^{(n)}(s') & \text{if } \mathcal{D}_{\text{head}}(s) = (s', 4, n) \\ \mathcal{D}_{\text{data*}}^{\text{indef}}(s') & \text{if } \mathcal{D}_{\text{indef}}(s) = (s', 4) \end{cases}$$

$$\mathcal{D}_{\text{data*}}^{(n)}(s) = \begin{cases} (s, \epsilon) & \text{if } n = 0 \\ (s'', d \cdot l) & \text{if } \mathcal{D}_{\text{data}}(s) = (s', d) \text{ and } \mathcal{D}_{\text{data*}}^{(n-1)}(s') = (s'', l) \end{cases}$$

$$\mathcal{D}_{\text{data*}}^{\text{indef}}(s) = \begin{cases} (s', \epsilon) & \text{if } s = 255 \cdot s' \\ (s'', d \cdot l) & \text{if } \mathcal{D}_{\text{data}}(s) = (s', d) \text{ and } \mathcal{D}_{\text{data*}}^{\text{indef}}(s') = (s'', l) \end{cases}$$

$$\begin{cases} (s, \epsilon) & \text{if } n = 0 \end{cases}$$

$$\mathcal{D}_{(\mathtt{data}^2)^*}^{(n)}(s) = \begin{cases} (s, \epsilon) & \text{if } n = 0 \\ (s''', (k, d) \cdot l) & \text{if } n > 0 \\ \text{and } \mathcal{D}_{\mathtt{data}}(s) = (s', k) \\ \text{and } \mathcal{D}_{\mathtt{data}}(s') = (s'', d) \\ \text{and } \mathcal{D}_{(\mathtt{data}^2)^*}^{(n-1)}(s'') = (s''', l) \end{cases}$$

$$\mathcal{D}_{\text{ctag}}(s) = \begin{cases} (s', i - 121) & \text{if } \mathcal{D}_{\text{head}}(s) = (s', 6, i) \text{ and } 121 \le i \le 127 \\ (s', (i - 1280) + 7) & \text{if } \mathcal{D}_{\text{head}}(s) = (s', 6, i) \text{ and } 1280 \le i \le 1400 \\ & \text{if } \mathcal{D}_{\text{head}}(s) = (s', 6, 102) \\ & \text{and } \mathcal{D}_{\text{head}}(s') = (s'', 4, 2) \\ & \text{and } \mathcal{D}_{\mathbb{Z}}(s'') = (s''', i) \\ & \text{and } 0 \le i \le 2^{64} - 1. \end{cases}$$

Note that the decoders for List and Constr accept both definite-length and indefinite-length lists of encoded data values, but the decoder for Map only accepts definite-length lists (and the length is the number of *pairs* in the map). This is consistent with CBOR's standard encoding of arrays and lists (major type 4) and maps (major type 5).

Note also that the encoder  $\mathcal{E}_{\text{ctag}}$  accepts arbitrary integer values for Constr tags, but (for compatibility with [10]) the decoder  $\mathcal{D}_{\text{ctag}}$  only accepts tags in  $\mathbb{N}_{[0,2^{64}-1]}$ . This means that some valid Plutus Core programs can be serialised but not descrialised, and is the reason for the recommendation in Section A.1 that only constructor tags between 0 and  $2^{64}-1$  should be used.

## **Appendix E Serialising Plutus Core Terms and Programs Using the** flat **Format**

We use the flat format [3] to serialise Plutus Core terms, and we regard this format as being the definitive concrete representation of Plutus Core programs. For compactness we generally (and *always* for scripts on the blockchain) replace names with de Bruijn indices (see Section 3.3) in serialised programs.

We use bytestrings for serialisation, but it is convenient to define the serialisation and deserialisation process in terms of strings of bits. Some extra bits of padding are added at the end of the encoding of a program to ensure that the number of bits in the output is a multiple of 8, and this allows us to regard serialised programs as bytestrings in the obvious way.

See Section E.4 for some restrictions on serialisation specific to the Cardano blockchain.

**Note:** flat **versus CBOR.** Much of the Cardano codebase uses the CBOR format for serialisation; however, it is important that serialised scripts not be too large. CBOR pays a price for being a self-describing format. The size of the serialised terms is consistently larger than a format that is not self-describing: benchmarks show that flat encodings of Plutus Core scripts are smaller than CBOR encodings by about 35% (without using compression).

## E.1 Encoding and decoding

Let  $S = \{0, 1\}^*$ , the set of all finite sequences of bits. For brevity we write a sequence of bits in the form  $b_{n-1} \cdots b_0$  instead of  $[b_{n-1}, \dots, b_0]$ : thus 011001 instead of [0, 1, 1, 0, 0, 1]). We denote the empty sequence by  $\epsilon$ , and use  $\ell(s)$  to denote the length of a sequence of bits, and  $\cdot$  to denote concatenation (or prepending or appending a single bit to a sequence of bits).

Similarly to the CBOR encoding for data described in Appendix D, we will describe the flat encoding by defining families of encoding functions (or *encoders*)

$$\mathsf{E}_X: \mathbb{S} \times X \to \mathbb{S}$$

and (partial) decoding functions (or decoders)

$$D_X : \mathbb{S} \rightharpoonup \mathbb{S} \times X$$

for various sets X, such as the set  $\mathbb{Z}$  of integers and the set of all Plutus Core terms. The encoding function  $\mathsf{E}_X$  takes a sequence  $s \in \mathbb{S}$  and an element  $x \in X$  and produces a new sequence of bits by appending the encoding of x to s, and the decoding function  $\mathsf{D}_X$  takes a sequence of bits, decodes some initial prefix of s to a value  $x \in X$ , and returns the remainder of s together with s.

Encoding functions basically operate by decomposing an object into subobjects and concatenating the encodings of the subobject; however it is sometimes necessary to add some padding between subobjects

in order to make sure that parts of the output are aligned on byte boundaries, and for this reason (unlike the CBOR encoding for data) all of our encoding functions have a first argument containing all of the previous output, so that it can be examined to determine how much alignment is required.

As in the case of CBOR, decoding functions are partial: they can fail if, for instance, there is insufficient input, or if a decoded value is outside some specified range. To simplify notation we will mention any preconditions separately, with the assumption that the decoder will fail if the preconditions are not met; we also make a blanket assumption that all decoders fail if there is not enough input for them to proceed. Many of the definitions of decoders construct objects by calling other decoders to obtain subobjects which are then composed, and these are often introduced by a condition of the form "if  $D_X(s) = x$ ". Conditions like this should be read as implicitly saying that if the decoder  $D_X$  fails then the whole decoding process fails.

## E.1.1 Padding

The encoding functions mentioned above produce sequences of *bits*, but we sometimes need sequences of *bytes*. To this end we introduce a functions pad:  $\mathbb{S} \to \mathbb{S}$  which adds a sequence of 0s followed by a 1 to a sequence *s* to get a sequence whose length is a multiple of 8; if *s* is a sequence such that  $\ell(s)$  is already a multiple of 8 then pad still adds an extra byte of padding; pad is used both for internal alignment (for example, to make sure that the contents of a bytestring are aligned on byte boundaries) and at the end of a complete encoding of a Plutus Core program to to make the length a multiple of 8 bits. Symbolically,

$$pad(s) = s \cdot p_k$$
 if  $\ell(s) = 8n + k$  with  $n, k \in \mathbb{N}$  and  $0 \le k \le 7$ 

where

 $p_0 = 00000001$ 

 $p_1 = 0000001$ 

 $p_2 = 000001$ 

 $p_3 = 00001$ 

 $p_4 = 0001$ 

 $p_5 = 001$ 

 $p_6 = 01$ 

 $p_7 = 1$ .

We also define a (partial) inverse function unpad :  $\mathbb{S} \to \mathbb{S}$  which discards padding:

unpad
$$(q \cdot s) = s$$
 if  $q = p_i$  for some  $i \in \{0, 1, 2, 3, 4, 5, 6, 7\}$ .

This can fail if the padding is not of the expected form or if the input is the empty sequence  $\epsilon$ .

## E.2 Basic flat encodings

#### E.2.1 Fixed-width natural numbers

We often wish to encode and decode natural numbers which fit into some fixed number of bits, and we do this simply by encoding them as their binary expansion (most significant bit first), adding leading zeros if necessary. More precisely for  $n \ge 1$  we define an encoder

$$\mathsf{E}_n: \mathbb{S} \times \mathbb{N}_{[0,2^{n-1}-1]} \to \mathbb{S}$$

by

$$\mathsf{E}_n(s, \sum_{i=0}^{n-1} b_i 2^i) = s \cdot b_{n-1} \cdots b_0 \quad (b_i \in \{0, 1\})$$

and a decoder

$$\mathsf{D}_n: \mathbb{S} \to \mathbb{S} \times \mathbb{N}_{[0,2^{n-1}-1]}$$

by

$$D_n(b_{n-1} \cdots b_0 \cdot s) = (s, \sum_{i=0}^{n-1} b_i 2^i).$$

As in Appendix D,  $\mathbb{N}_{[a,b]}$  denotes the closed interval of integers  $\{n \in \mathbb{Z} : a \le n \le b\}$ . Note that n here is a variable (not a fixed label) so we are defining whole families of encoders  $E_1, E_2, E_3, \ldots$  and and decoders  $D_1, D_2, D_3, \ldots$ 

#### E.2.2 Lists

Suppose that we have a set X for which we have defined an encoder  $\mathsf{E}_X$  and a decoder  $\mathsf{D}_X$ ; we define an encoder  $\mathsf{E}_X$  which encodes lists of elements of X by emitting the encodings of the elements of the list, each preceded by a 1 bit, then emitting a 0 bit to mark the end of the list.

$$\vec{\mathsf{E}}_X(s,[]) = s \cdot 0$$

$$\vec{\mathsf{E}}_X(s,[x_1,\ldots,x_n]) = \vec{\mathsf{E}}_X(s \cdot 1 \cdot \mathsf{E}_X(x_1),[x_2,\ldots,x_n]).$$

The corresponding decoder is given by

$$\overrightarrow{\mathsf{D}}_X(0 \cdot s) = (s, [])$$

$$\overrightarrow{\mathsf{D}}_Y(1 \cdot s) = (s'', x \cdot l) \quad \text{if } D_Y(s) = (s', x) \text{ and } \overrightarrow{\mathsf{D}}_Y(s') = (s'', l).$$

#### **E.2.3** Natural numbers

We encode natural numbers by splitting their binary representations into sequences of 7-bit blocks, then emitting these as a list with the **least significant block first**:

$$\mathsf{E}_{\mathbb{N}}(s, \sum_{i=0}^{n-1} k_i 2^{7i}) = \vec{\mathsf{E}}_7(s, [k_0, \dots, k_{n-1}])$$

(where  $k_i \in \mathbb{Z}$  and  $0 \le k_i \le 127$ ). The decoder is

$$\mathsf{D}_{\mathbb{N}}(s) = (s', \sum_{i=0}^{n-1} k_i 2^{7i}) \quad \text{if } \overrightarrow{\mathsf{D}}_7(s) = (s', [k_0, \dots, k_{n-1}]).$$

#### E.2.4 Integers

Signed integers are encoded by converting them to natural numbers using the zigzag encoding  $(0 \mapsto 0, -1 \mapsto 1, 1 \mapsto 2, -2 \mapsto 3, 2 \mapsto 4, ...)$  and then encoding the result using  $E_{\mathbb{N}}$ :

$$\mathsf{E}_{\mathbb{Z}}(s,n) = \begin{cases} \mathsf{E}_{\mathbb{N}}(s,2n) & \text{if } n \geq 0 \\ \mathsf{E}_{\mathbb{N}}(s,-2n-1) & \text{if } n < 0. \end{cases}$$

The decoder is

$$\mathsf{D}_{\mathbb{Z}}(s) = \begin{cases} (s', \frac{n}{2}) & \text{if } n \equiv 0 \pmod{2} \\ (s', -\frac{n+1}{2}) & \text{if } n \equiv 1 \pmod{2} \end{cases} \quad \text{if } \mathsf{D}_{\mathbb{N}}(s) = (s', n).$$

#### E.2.5 Bytestrings

Bytestrings are encoded by dividing them into nonempty blocks of up to 255 bytes and emitting each block in sequence. Each block is preceded by a single unsigned byte containing its length, and the end of the encoding is marked by a zero-length block (so the empty bytestring is encoded just as a zero-length block). Before emitting a bytestring, the preceding output is padded so that its length (in bits) is a multiple of 8; if this is already the case a single padding byte is still added; this ensures that contents of the bytestring are aligned to byte boundaries in the output.

Recall that  $\mathbb B$  denotes the set of 8-bit bytes,  $\{0,1,\ldots,255\}$ . For specification purposes we may identify the set of bytestrings with the set  $\mathbb B^*$  of (possibly empty) lists of elements of  $\mathbb B$ . We denote by C the set of bytestring chunks of **nonempty** bytestrings of length at most 255:  $C = \{[b_1,\ldots,b_n]: b_i \in \mathbb B, 1 \leq n \leq 255\}$ , and define a function  $E_C: C \to \mathbb S$  by

$$E_C([b_1,\ldots,b_n]) = \mathsf{E}_8(n) \cdot \mathsf{E}_8(b_1) \cdot \cdots \cdot \mathsf{E}_8(b_n).$$

We define an encoder  $E_{C^*}$  for lists of chunks by

$$\mathsf{E}_{C^*}(s,[c_1,\ldots,c_n]) = s \cdot E_C(c_1) \cdot \cdots \cdot E_C(c_n) \cdot 00000000.$$

Note that each  $c_i$  is required to be nonempty but that we allow the case n = 0, so that an empty list of chunks encodes as 00000000.

To encode a bytestring we decompose it into a list L of chunks and then apply  $\mathsf{E}_{C^*}$  to L. However, there will usually be many ways to decompose a given bytestring a into chunks. For definiteness we recommend (but do not demand) that a is decomposed into a sequence of chunks of length 255 possibly followed by a smaller chunk. Formally, suppose that  $a = [a_1, \ldots, a_{255k+r}] \in \mathbb{B}^* \setminus \{\epsilon\}$  where  $k \ge 0$  and  $0 \le r \le 254$ . We define the *canonical 256-byte decomposition*  $\tilde{a}$  of a to be

$$\tilde{a} = [[a_1, \dots, a_{255}], [a_{256}, \dots, a_{510}], \dots [a_{255(k-1)+1}, \dots, a_{255k}]] \in C^*$$

if r = 0 and

$$\tilde{a} = [[a_1, \dots, a_{255}], [a_{256}, \dots, a_{510}], \dots [a_{255(k-1)+1}, \dots, a_{255k}], [a_{255k+1}, \dots, a_{255k+r}]] \in C^*$$

if r > 0.

For the empty bytestring we define

$$\tilde{\epsilon} = [].$$

Given all of the above, we define the canonical encoding function  $E_{\mathbb{R}^*}$  for bytestrings to be

$$\mathsf{E}_{\mathbb{B}^*}(s,a) = E_{C^*}(\mathsf{pad}(s),\tilde{a}).$$

Non-canonical encodings can be obtained by replacing  $\tilde{a}$  with any other decomposition of a into nonempty chunks, and the decoder below will accept these as well.

To define a decoder for bytestrings we first define a decoder  $D_C$  for bytestring chunks:

$$D_C(s) = D_C^{(n)}(s', [])$$
 if  $D_8(s) = (s', n)$ 

where

$$\mathsf{D}_C^{(n)}(s,l) = \begin{cases} (s,l) & \text{if } n = 0\\ \mathsf{D}_C^{(n-1)}(s',l \cdot x) & \text{if } n > 0 \text{ and } \mathsf{D}_8(s) = (s',x). \end{cases}$$

Now we define

$$\mathsf{D}_{C^*}(s) = \begin{cases} (s', []) & \text{if } D_C(s) = (s', []) \\ (s'', x \cdot l) & \text{if } \mathsf{D}_C(s) = (s', x) \text{ with } x \neq [] \text{ and } \mathsf{D}_{C^*}(s') = (s'', l). \end{cases}$$

The notation is slightly misleading here:  $D_{C^*}$  does not decode to a list of bytestring chunks, but to a single bytestring. We could alternatively decode to a list of bytestrings and then concatenate them later, but this would have the same overall effect.

Finally, we define the decoder for bytestrings by

$$\mathsf{D}_{\mathbb{R}^*}(s) = \mathsf{D}_{C^*}(\mathsf{unpad}(s)).$$

## E.2.6 Strings

We have defined values of the string type to be sequences of Unicode characters. As mentioned earlier we do not specify any particular internal representation of Unicode characters, but for serialisation we use the UTF-8 representation to convert between strings and bytestrings and then use the bytestring encoder and decoder:

$$\mathsf{E}_{\mathbb{H}^*}(s,u) = \mathsf{E}_{\mathbb{H}^*}(s,\mathsf{utf8}(u))$$

$$D_{\mathbb{H}^*}(s) = (s', \text{utf } 8^{-1}(a)) \text{ if } D_{\mathbb{H}^*}(s) = (s', a)$$

where utf8 and utf8<sup>-1</sup> are the UTF8 encoding and decoding functions mentioned in Appendix A. Recall that utf8<sup>-1</sup> is partial (not all bytestrings represent valid Unicode sequences), so  $D_{\mathbb{U}^*}$  may fail if the input is invalid.

## E.3 Encoding and decoding Plutus Core

#### E.3.1 Programs

A program is encoded by encoding the three components of the version number in sequence then encoding the body, and possibly adding some padding to ensure that the total number of bits in the output is a multiple of 8 (and hence the output can be viewed as a bytestring).

$$\mathsf{E}_{\mathsf{program}}((\mathsf{program}\ a.b.c\ t)) = \mathsf{pad}(\mathsf{E}_{\mathsf{term}}(\mathsf{E}_{\mathbb{N}}(\mathsf{E}_{\mathbb{N}}(\mathsf{E}_{\mathbb{N}}(\varepsilon,a),b),c),t)).$$

The decoding process is the inverse of the encoding process: three natural numbers are read to obtain the version number and then the body is decoded. After this we discard any padding in the remaining input and check that all of the input has been consumed.

$$\mathsf{D}_{\mathsf{program}}(s) = (\mathsf{program}\ a.b.c\ t) \quad \begin{cases} \text{if} & \mathsf{D}_{\mathbb{N}}(s) = (s',a) \\ \text{and} & \mathsf{D}_{\mathbb{N}}(s') = (s'',b) \\ \text{and} & \mathsf{D}_{\mathbb{N}}(s'') = (s''',c) \\ \text{and} & \mathsf{D}_{\mathsf{term}}(s''') = (r,t) \\ \text{and} & \mathsf{unpad}(r) = \epsilon. \end{cases}$$

#### E.3.2 Terms

Plutus Core terms are encoded by emitting a 4-bit tag identifying the type of the term (see Table 6; recall that [] denotes application) then emitting the encodings for any subterms. We currently only use eight of the sixteen available tags: the remainder are reserved for potential future expansion.

Term type	Binary	Decimal
Variable	0000	0
delay	0001	1
lam	0010	2
[]	0011	3
const	0100	4
force	0101	5
error	0110	6
builtin	0111	7
constr	1000	8
case	1001	9

Table 6: Term tags

The encoder for terms is given below: it refers to other encoders (for names, types, and constants) which will be defined later.

```
\begin{split} & \mathsf{E}_{\mathsf{term}}(s,x) &= \mathsf{E}_{\mathsf{name}}(s \cdot 0000,x) \\ & \mathsf{E}_{\mathsf{term}}(s,(\mathsf{delay}\,t)) &= \mathsf{E}_{\mathsf{term}}(s \cdot 0001,t) \\ & \mathsf{E}_{\mathsf{term}}(s,(\mathsf{lam}\,x\,t)) &= \mathsf{E}_{\mathsf{term}}(\mathsf{E}_{\lambda\mathsf{var}}(s \cdot 0010,x),t) \\ & \mathsf{E}_{\mathsf{term}}(s,[t_1\,t_2]) &= \mathsf{E}_{\mathsf{term}}(\mathsf{E}_{\mathsf{term}}(s \cdot 0011,t_1),t_2) \\ & \mathsf{E}_{\mathsf{term}}(s,(\mathsf{const}\,tn\,c)) &= \mathsf{E}_{\mathsf{term}}^{tn}(\mathsf{E}_{\mathsf{type}}(s \cdot 0100,T),c) \\ & \mathsf{E}_{\mathsf{term}}(s,(\mathsf{force}\,t)) &= \mathsf{E}_{\mathsf{term}}(s \cdot 0101,t) \\ & \mathsf{E}_{\mathsf{term}}(s,(\mathsf{error})) &= s \cdot 0110 \\ & \mathsf{E}_{\mathsf{term}}(s,(\mathsf{builtin}\,b)) &= \mathsf{E}_{\mathsf{builtin}}(s \cdot 0111,b) \\ & \mathsf{E}_{\mathsf{term}}(s,(\mathsf{constr}\,i\,\bar{t})) &= \vec{\mathsf{E}}_{\mathsf{term}}(\mathsf{E}_{\mathsf{N}}(s \cdot 1000,i),\bar{t}) \\ & \mathsf{E}_{\mathsf{term}}(s,(\mathsf{case}\,u\,\bar{t})) &= \vec{\mathsf{E}}_{\mathsf{term}}(\mathsf{E}_{\mathsf{term}}(s \cdot 1001,u),\bar{t}) \end{split}
```

The decoder for terms is given below. To simplify the definition we use some pattern-matching syntax for inputs to decoders: for example the argument  $0101 \cdot s$  indicates that when the input is a string beginning with 0101 the definition after the = sign should be used (and the remainder of the input is available in s there). If the input is not long enough to permit the indicated decomposition then the decoder fails. The decoder also fails if the input begins with a prefix which is not listed; that does not happen here, but does in some later decoders.

```
D_{\text{term}}(0000 \cdot s) = (s', x)
                                                                       if D_{name}(s) = (s', x)
D_{\text{term}}(0001 \cdot s) = (s', (\text{delay } t))
                                                                       if D_{term}(s) = (s', t)
D_{term}(0010 \cdot s) = (s'', (lam x t))
                                                                       if D_{\lambda var}(s) = (s', x) and D_{term}(s') = (s'', t)
D_{\text{term}}(0011 \cdot s) = (s'', [t_1 \ t_2])
                                                                       if D_{\text{term}}(s) = (s', t_1) and D_{\text{term}}(s') = (s'', t_2)
\mathsf{D}_{\mathsf{term}}(\mathsf{0100} \cdot s) = (s'', (\mathsf{const} \ tn \ c)) \quad \text{if } \mathsf{D}_{\mathsf{type}}(s) = (s', T) \quad \text{and } \mathsf{D}_{\mathsf{constant}}^T(s') = (s'', c)
D_{term}(0101 \cdot s) = (s', (force t))
                                                                       if D_{term}(s) = (s', t)
D_{\text{term}}(0110 \cdot s) = (s, (\text{error}))
D_{\text{term}}(0111 \cdot s) = (s', b)
                                                                       if D_{\text{builtin}}(s) = (s', b)
\mathsf{D}_{\mathsf{term}}(1000 \cdot s) = (s', (\mathsf{constr}\ i\ \bar{t})) \quad \text{if } \mathsf{D}_{\mathbb{N}}(s) = (s', i) \quad \text{and } \overrightarrow{\mathsf{D}}_{\mathsf{term}}(s') = (s'', \bar{t})
D_{\text{term}}(1001 \cdot s) = (s', (case \, u \, \overline{t}))
                                                                       if D_{\text{term}}(s) = (s', u) and \overrightarrow{D}_{\text{term}}(s') = (s'', \overline{t})
```

## E.3.3 Built-in types

Constants from built-in types are essentially encoded by emitting a sequence of 4-bit tags representing the constant's type and then emitting the encoding of the constant itself. However the encoding of types is somewhat complex because it has to be able to deal with type operators such as list and pair. The tags are given in Table 7: they include tags for the basic types together with a tag for a type application operator.

Туре	Binary	Decimal	
integer	0000	0	
bytestring	0001	1	
string	0010	2	
unit	0011	3	
bool	0100	4	
list	0101	5	
pair	0110	6	
(type application)	0111	7	
data	1000	8	

Table 7: Type tags

We define auxiliary functions  $e_{type}:\mathcal{U}\to\mathbb{N}^*$  and  $d_{type}:\mathbb{N}^*\to\mathbb{N}^*\times\mathcal{U}$  ( $d_{type}$  is partial and  $\mathcal{U}$  denotes the universe of types used in Alonzo and Vasil) by

$$\begin{split} \mathbf{e}_{\mathsf{type}}(\mathsf{integer}) &= [0] \\ \mathbf{e}_{\mathsf{type}}(\mathsf{bytestring}) &= [1] \\ \mathbf{e}_{\mathsf{type}}(\mathsf{string}) &= [2] \\ \mathbf{e}_{\mathsf{type}}(\mathsf{unit}) &= [3] \\ \mathbf{e}_{\mathsf{type}}(\mathsf{bool}) &= [4] \\ \mathbf{e}_{\mathsf{type}}(\mathsf{list}(t)) &= [7,5] \cdot \mathbf{e}_{\mathsf{type}}(t) \\ \mathbf{e}_{\mathsf{type}}(\mathsf{pair}(t_1,t_2)) &= [7,7,6] \cdot \mathbf{e}_{\mathsf{type}}(t_1) \cdot \mathbf{e}_{\mathsf{type}}(t_2) \\ \mathbf{e}_{\mathsf{type}}(\mathsf{data}) &= [8]. \end{split}$$

$$\begin{split} & \mathsf{d}_{\mathsf{type}}(0 \cdot l) &= (l, \mathsf{integer}) \\ & \mathsf{d}_{\mathsf{type}}(1 \cdot l) &= (l, \mathsf{bytestring}) \\ & \mathsf{d}_{\mathsf{type}}(2 \cdot l) &= (l, \mathsf{string})) \\ & \mathsf{d}_{\mathsf{type}}(3 \cdot l) &= (l, \mathsf{unit}) \\ & \mathsf{d}_{\mathsf{type}}(4 \cdot l) &= (l, \mathsf{bool}) \\ & \mathsf{d}_{\mathsf{type}}([7, 5] \cdot l) &= (l', \mathsf{list}(t)) & \text{if } \mathsf{d}_{\mathsf{type}}(l) = (l', t) \\ & \mathsf{d}_{\mathsf{type}}([7, 7, 6] \cdot l) = (l'', \mathsf{pair}(t_1, t_2)) & \begin{cases} \mathsf{if} & \mathsf{d}_{\mathsf{type}}(l) = (l', t_1) \\ \mathsf{and} & \mathsf{d}_{\mathsf{type}}(l') = (l'', t_2) \end{cases} \\ & \mathsf{d}_{\mathsf{type}}(8 \cdot l) &= (l, \mathsf{data}). \end{split}$$

The encoder and decoder for types is obtained by combining  $e_{type}$  and  $d_{type}$  with  $\vec{E}_4$  and  $\vec{D}_4$ , the encoder and decoder for lists of four-bit integers (see Section E.2).

$$\mathsf{E}_{\mathsf{type}}(s,t) = \vec{\mathsf{E}}_{\mathsf{4}}(s,\mathsf{e}_{\mathsf{type}}(t))$$
 
$$\mathsf{D}_{\mathsf{type}}(s) = (s',t) \quad \text{if } \vec{\mathsf{D}}_{\mathsf{4}}(s) = (s',l) \text{ and } \mathsf{d}_{\mathsf{type}}(l) = ([],t).$$

### E.3.4 Constants

Values of built-in types can mostly be encoded quite simply by using encoders already defined. Note that the unit value (con unit ()) does not have an explicit encoding: the type has only one possible value, so there is no need to use any space to serialise it.

The data type is encoded by converting to a bytestring using the CBOR encoder  $\mathcal{E}_{\text{data}}$  described in Appendix D and then using  $E_{\mathbb{B}^*}$ . The decoding process is the opposite of this: a bytestring is obtained using  $D_{\mathbb{B}^*}$  and this is then decoded from CBOR using  $\mathcal{D}_{\text{data}}$  to obtain a data object.

$$\begin{split} & \mathsf{E}_{\mathsf{constant}}^{\mathsf{integer}}(s,n) &= \mathsf{E}_{\mathbb{Z}}(s,n) \\ & \mathsf{E}_{\mathsf{constant}}^{\mathsf{bytestring}}(s,a) &= \mathsf{E}_{\mathbb{B}^*}(s,a) \\ & \mathsf{E}_{\mathsf{constant}}^{\mathsf{string}}(s,t) &= \mathsf{E}_{\mathbb{U}^*}(s,t) \\ & \mathsf{E}_{\mathsf{constant}}^{\mathsf{unit}}(s,c) &= s \\ & \mathsf{E}_{\mathsf{constant}}^{\mathsf{bool}}(s,\mathsf{False}) &= s \cdot 0 \\ & \mathsf{E}_{\mathsf{constant}}^{\mathsf{bool}}(s,\mathsf{True}) &= s \cdot 1 \\ & \mathsf{E}_{\mathsf{constant}}^{\mathsf{bool}}(s,l) &= \vec{\mathsf{E}}_{\mathsf{constant}}^T(s,l) \\ & \mathsf{E}_{\mathsf{constant}}^{\mathsf{alist}(T)}(s,l) &= \vec{\mathsf{E}}_{\mathsf{constant}}^T(s,l) \\ & \mathsf{E}_{\mathsf{constant}}^{\mathsf{pair}(T_1,T_2)}(s,(c_1,c_2)) &= \mathsf{E}_{\mathsf{constant}}^{T_2}(\mathsf{E}_{\mathsf{constant}}^{T_1}(s,c_1),c_2) \\ & \mathsf{E}_{\mathsf{constant}}^{\mathsf{data}}(s,d) &= \mathsf{E}_{\mathbb{B}^*}(s,\mathcal{E}_{\mathsf{data}}(d)). \end{split}$$

```
\begin{array}{lll} \mathsf{D}_{\mathsf{constant}}^{\mathsf{integer}}(s) &= \mathsf{D}_{\mathbb{Z}}(s) \\ \mathsf{D}_{\mathsf{constant}}^{\mathsf{bytestring}}(s) &= \mathsf{D}_{\mathbb{B}^*}(s) \\ \mathsf{D}_{\mathsf{constant}}^{\mathsf{string}}(s) &= \mathsf{D}_{\mathbb{U}^*}(s) \\ \mathsf{D}_{\mathsf{constant}}^{\mathsf{unit}}(s) &= s \\ \mathsf{D}_{\mathsf{constant}}^{\mathsf{bool}}(0 \cdot s) &= (s, \mathsf{False}) \\ \mathsf{D}_{\mathsf{constant}}^{\mathsf{bool}}(1 \cdot s) &= (s, \mathsf{True}) \\ \mathsf{D}_{\mathsf{constant}}^{\mathsf{bool}}(s) &= \overrightarrow{\mathsf{D}}_{\mathsf{constant}}^T(s, l) \\ \mathsf{D}_{\mathsf{constant}}^{\mathsf{pair}\,(T_1, T_2)}(s) &= (s'', (c_1, c_2)) \begin{cases} \mathsf{if} & \mathsf{D}_{\mathsf{constant}}^{T_1}(s) = (s', c_1) \\ \mathsf{and} & \mathsf{D}_{\mathsf{constant}}^{T_2}(s') = (s'', c_2) \end{cases} \\ \mathsf{D}_{\mathsf{constant}}^{\mathsf{data}}(s) &= (s', d) & \mathsf{if} \; \mathsf{D}_{\mathbb{B}*}(s) = (s', t) \; \mathsf{and} \; \mathcal{D}_{\mathsf{data}}(t) = (t', d) \; \mathsf{for} \; \mathsf{some} \; t'. \end{cases}
```

#### **E.3.5** Built-in functions

Built-in functions are represented by seven-bit integer tags and encoded and decoded using  $E_7$  and  $D_7$ . The tags are specified in Tables 8 and 9. We assume that there are (partial) functions tag and tag<sup>-1</sup> which convert back and forth between builtin names and their tags.

$$\mathsf{E}_{\mathsf{builtin}}(s,b) = \mathsf{E}_7(s,\mathsf{tag}(b))$$
 
$$\mathsf{D}_{\mathsf{builtin}}(s) = (s',\mathsf{tag}^{-1}(n)) \quad \text{if } \mathsf{D}_7(s) = (s',n).$$

Builtin	Binary	Decimal	Builtin	Binary	Decimal
addInteger	0000000	0	ifThenElse	0011010	26
subtractInteger	0000001	1	chooseUnit	0011011	27
multiplyInteger	0000010	2	trace	0011100	28
divideInteger	0000011	3	fstPair	0011101	29
quotientInteger	0000100	4	sndPair	0011110	30
remainderInteger	0000101	5	chooseList	0011111	31
modInteger	0000110	6	mkCons	0100000	32
equalsInteger	0000111	7	headList	0100001	33
lessThanInteger	0001000	8	tailList	0100010	34
lessThanEqualsInteger	0001001	9	nullList	0100011	35
appendByteString	0001010	10	chooseData	0100100	36
consByteString	0001011	11	constrData	0100101	37
sliceByteString	0001100	12	mapData	0100110	38
lengthOfByteString	0001101	13	listData	0100111	39
indexByteString	0001110	14	iData	0101000	40
equalsByteString	0001111	15	bData	0101001	41
lessThanByteString	0010000	16	unConstrData	0101010	42
lessThanEqualsByteString	0010001	17	unMapData	0101011	43
sha2_256	0010010	18	unListData	0101100	44
sha3_256	0010011	19	unIData	0101101	45
blake2b_256	0010100	20	unBData	0101110	46
verifyEd25519Signature	0010101	21	equalsData	0101111	47
appendString	0010110	22	mkPairData	0110000	48
equalsString	0010111	23	mkNilData	0110001	49
encodeUtf8	0011000	24	mkNilPairData	0110010	50
decodeUtf8	0011001	25			

Table 8: Tags for Alonzo builtins

Builtin	Binary	Decimal
serialiseData	0110011	51
verifyEcdsaSecp256k1Signature	0110100	52
verifySchnorrSecp256k1Signature	0110101	53

Table 9: Extra tags for Vasil builtins

# E.3.6 Variable names

Variable names are encoded and decoded using the  $E_{name}$  and  $D_{name}$  functions, and variables bound in lam expressions are encoded and decoded by the  $E_{\lambda var}$  and  $D_{\lambda var}$  functions.

**De Bruijn indices.** We use serialised de Bruijn-indexed terms for script transmission because this makes serialised scripts significantly smaller. Recall from Section 3.3 that when we want to use our syntax with de Bruijn indices we replace names with natural numbers and the bound variable in a lam expression with 0. During serialisation the zero is ignored, and during deserialisation no input is consumed and the index 0 is always returned:

$$\mathsf{E}_{\lambda\mathsf{var}}(s,n) = s$$

$$D_{\lambda var}(s) = 0.$$

For variables we always use indices which are greater than zero, and our encoder and decoder for names are given by

$$\mathsf{E}_{\mathsf{name}} = \mathsf{E}_{\mathbb{N}}$$

and

$$\mathsf{D}_{\mathsf{name}}(s) = (s', n) \quad \text{if } \mathsf{D}_{\mathbb{N}} = (s', n) \text{ and } n > 0.$$

Other types of name. One can serialise code involving other types of name by providing suitable encoders and decoders for name. For example, for textual names one could use  $E_{\lambda var} = E_{name} = E_{U^*}$  and  $D_{\lambda var} = D_{name} = D_{U^*}$ . Depending on the method used to represent variable names it may also be necessary to check during descrialisation the more general requirement that variables are well-scoped, but this problem will not arise if de Bruijn indices are used.

## E.4 Cardano-specific serialisation issues

## E.4.1 Scope checking

To execute a Plutus Core program on the blockchain it will be necessary to deserialise it to some inmemory representation, and during or immediately after deserialisation it should be checked that the body of the program is a closed term (see the requirement in Section 3.3); if this is not the case then evaluation should fail immediately.

#### E.4.2 CBOR wrapping

Plutus Core programs are not stored on the Cardano chain directly as flat bytestrings; for consistency with other objects used on the chain, the flat bytestrings are in fact wrapped in a CBOR encoding. This should not concern most users, but we mention it here to avoid possible confusion.

#### E.5 Example

```
Consider the program
```

```
(program 5.0.2
  [
    [(builtin indexByteString)(con bytestring #1a5f783625ee8c)]
  (con integer 54321)
])
```

Suppose this is stored in index.uplc. We can convert it to flat by running

```
$ cabal run exec uplc convert -- -i index.uplc --of flat -o index.flat
```

The serialised program looks like this:

Figure 12 shows how this encodes the original program. Sequences of bits are followed by explanatory comments and lines beginning with # provide further commentary on preceding bit sequences.

```
00000101 : Final integer chunk: 0000101 \rightarrow 5
00000000 : Final integer chunk: 0000000 \rightarrow 0
00000010 : Final integer chunk: 0000000 \rightarrow 2
           # Version: 5.0.2
          : Term tag 3: apply
0011
          : Term tag 3: apply
0011
          : Term tag 7: builtin
0111
0001110 : Builtin tag 14
           # builtin indexByteString
0100
           : Term tag 4: constant
           : Start of type tag list
1
0001
           : Type tag 1
           : End of list
0
           # Type tags: [1] \rightarrow \text{bytestring}
           : Padding before bytestring
00000111: Bytestring chunk size: 7
00011010 : 0x1a
01011111 : 0x5f
01111000 : 0x78
00110110 : 0x36
00100101 : 0x25
11101110 : 0xee
10001100 : 0x8c
00000000 : Bytestring chunk size: 0 (end of list of chunks)
           # con bytestring #1a5f783625ee8c
0100
           : Term tag 4: constant
           : Start of type tag list
0000
           : Type tag 0
           : End of list
           # Type tags: [0] \rightarrow integer
11100010: Integer chunk 1100010 (least significant)
11010000 : Integer chunk 1010000
00000110 : Final integer chunk 0000110 (most significant)
           # 0000110 · 1010000 · 1100010 \rightarrow 108642 decimal
           # Zigzag encoding: 108642/2 \rightarrow +54321
           # con integer 54321
000001
           : Padding
```

Figure 12: flat encoding of index.uplc

# References

- [1] ANSI. X9.62: Public Key Cryptography for the Financial Services Industry: the Elliptic Curve Digital Signature Algorithm (ECDSA), 2005.
- [2] ANSI. X9.142: Public Key Cryptography for the Financial Services Industry: the Elliptic Curve Digital Signature Algorithm (ECDSA), 2020.
- [3] Pasqualino 'Titto' Assini. Flat format specification. http://quid2.org/docs/Flat.pdf.
- [4] Hendrik Pieter Barendregt. *The Lambda Calculus its Syntax and Semantics*, volume 103 of *Studies in Logic and the Foundations of Mathematics*. North-Holland, 1985.
- [5] Daniel J. Bernstein, Niels Duif, Tanja Lange, Peter Schwabe, and Bo-Yin Yang. High-speed high-security signatures. In *CHES*, volume 6917 of *Lecture Notes in Computer Science*, pages 124–142. Springer, 2011.
- [6] Bitcoin Wiki. Elliptic Curve Digital Signature Algorithm, 2022.
- [7] Carsten Bormann. Notable CBOR Tags. https://www.ietf.org/archive/id/draft-bormann-cbor-notable-tags-06.html.
- [8] Carsten Bormann and Paul E. Hoffman. RFC 8949: Concise Binary Object Representation (CBOR). https://www.rfc-editor.org/info/rfc8949, December 2020.
- [9] Certicom Research. Standards for Efficient Cryptography 2 (SEC 2). https://www.secg.org/SEC2-Ver-2.0.pdf, 2010.
- [10] Duncan Coutts, Michael Peyton Jones, and Carsten Bormann. CBOR Tags for Discriminated Unions. https://github.com/cabo/cbor-discriminated-unions/.
- [11] N.G de Bruijn. Lambda calculus notation with nameless dummies, a tool for automatic formula manipulation, with application to the Church-Rosser theorem. *Indagationes Mathematicae (Proceedings)*, 75(5):381–392, 1972.
- [12] Matthias Felleisen. Programming languages and lambda calculi, 2007.
- [13] Matthias Felleisen, Robert Bruce Findler, and Matthew Flatt. Semantics Engineering with PLT Redex. MIT Press, 2009.
- [14] Matthias Felleisen and Daniel P. Friedman. Control operators, the SECD-machine, and the lambda-calculus. In *3rd Working Conference on the Formal Description of Programming Concepts*, August 1986.
- [15] Matthias Felleisen and Robert Hieb. The revised report on the syntactic theories of sequential control and state. *Theor. Comput. Sci.*, 103(2):235–271, September 1992.
- [16] Robert Harper. *Practical Foundations for Programming Languages*. Cambridge University Press, New York, NY, USA, 2012.
- [17] Don Johnson, Alfred Menezes, and Scott A. Vanstone. The elliptic curve digital signature algorithm (ECDSA). *Int. J. Inf. Sec.*, 1(1):36–63, 2001.
- [18] Simon Josefsson and Ilari Liusvaara. RFC 8032: Edwards-Curve Digital Signature Algorithm (Ed-DSA). https://www.rfc-editor.org/info/rfc8032, January 2017.

- [19] Johnson Lau, Jonas Nick, and Tim Ruffing. Bitcoin Improvement Proposal 340: Schnorr Signatures for secp256k1. https://github.com/bitcoin/bips/blob/master/bip-0340.mediawiki, 2020.
- [20] Johnson Lau and Pieter Wuilie. Bitcoin Improvement Proposal 146: Dealing with signature encoding malleability. https://github.com/bitcoin/bips/blob/master/bip-0146.mediawiki, 2016.
- [21] Gordon D. Plotkin. Call-by-name, call-by-value and the lambda-calculus. *Theor. Comput. Sci.*, 1(2):125–159, 1975.
- [22] John C. Reynolds. Types, abstraction and parametric polymorphism. In R. E. A. Mason, editor, *Information Processing 83, Proceedings of the IFIP 9th World Computer Congress, Paris, France, September 19-23, 1983*, pages 513–523. North-Holland/IFIP, 1983.
- [23] Claus-Peter Schnorr. Efficient identification and signatures for smart cards. In Gilles Brassard, editor, *CRYPTO*, volume 435 of *Lecture Notes in Computer Science*, pages 239–252. Springer, 1989.
- [24] The Unicode Consortium. The Unicode Standard. https://www.unicode.org/versions/latest/.
- [25] Philip Wadler. Theorems for free! In FPCA '89: Proceedings of the fourth international conference on Functional programming languages and computer architecture, pages 347–359, New York, NY, USA, 1989. ACM.

# **Index of Notation**

```
Sets
\mathbb{B}
                   \{n \in \mathbb{Z} : 0 \le n \le 255\}, 4
\mathbb{B}^*
                   The set of all bytestrings, 4
\mathbb{N}
                   \{0, 1, 2, 3, \dots\}, 4
\mathbb{N}^+
                   \{1, 2, 3, \dots\}, 4
                   \{n \in \mathbb{N} : a \le n \le b\}, 4
\mathbb{N}_{[a,b]}
                   The set of strings of bits, 35
\mathbb{U}
                   The set of Unicode scalar values, 4
\mathbb{U}^*
                   The set of Unicode strings, 4
\mathbb{Z}
                   \{\ldots, -2, -1, 0, 1, 2, \ldots\}, 4
                   S \uplus \{x\} (S \text{ a set}), 4
S_{\mathsf{x}}
                   Disjoint union of sets, 4
X^*
                   The set of all finite sequences of elements of a set X, 4
Lists
[]
                   The empty list, 5
L \cdot L'
                   Concatenation of lists L and L', 5
x \cdot L
                   [x] \cdot L, 5
L \cdot x
                   L \cdot [x], 5
\overline{V}
                   A sequence [V_1, \ldots, V_n], 5
\ell(\cdot)
                   Length of a list or bytestring, 5
                   The empty bytestring, 29
Plutus Core grammar
bn, b
                   The name of a built-in function, 5
                   A literal constant, 5
c
                   An integer, 5
i
L, M, N
                   A term, 5
                   A name, 5
n
P
                   A Plutus Core program, 5
                   Plutus Core version, 5
υ
                   A variable name, 5
x
Built-in types
                   Constants of built-in type T, 10
\mathcal{C}_T
O
                   The set of built-in type operator names, 7
\mathcal{U}
                   The universe of built-in types, 7
\mathcal{U}_0
                   The set of atomic type names, 7
                   The set of builtin-polymorphic types, 8
\mathcal{U}_{\#}
\nu
                   The set of all type variables, \mathcal{V}_{\#} \cup \mathcal{V}_{*}, 8
\mathcal{V}_*
                   The set of fully-polymorphic type variables, 8
                   The set of builtin-polymorphic type variables, 8
\mathcal{V}_{\#}
                   A built-in type operator (an element of 0), 7
op
```

```
A fully-polymorphic type variable, 8
v_*
                 A builtin-polymorphic type variable, 8
v_{\#}
P
                 A builtin-polymorphic type, 8
T
                 A built-in type (an element of \mathcal{U}), 7
T \leq_S P
                 T is an instance of P via S and dom S = FV_{\#}(P), 9
Ŝ
                 The extension of a type assignment S to \mathcal{U}_{\#} \cup \mathcal{V}_{*}, 9
S
                 A type assignment, 9
FV_{\#}(P)
                 Free #-variables of a polymorphic type P \in \mathcal{U}_{\#}, 8
                 Denotation of a type T \in \mathcal{U}, 7
\llbracket T 
rbracket
\llbracket \cdot 
rbracket_T
                 Denotation of constants of type T, 10
                 Reification of constants of type T, 10
\{\![\,\cdot\,]\!\}_T
\{\![\,\cdot\,]\!\}
                 Reification of the result of a built-in function application, 12
Built-in functions
\mathcal{B}
                 The set of built-in functions, 6
J
                 The set of inputs to built-in functions, 9
Ω
                 The set of all type quantifications, 10
                 Arity of built-in function, 11
α
                 Reduced arity, 11
\bar{\alpha}
                 Number of arguments of built-in function, 11
χ
                 Signature item, 10
ı
                 Signature of built-in function, 10
\sigma
                 Reduced signature, 11
\bar{\sigma}
                 A member of \mathcal{U}_{\#} \uplus \mathcal{V}_{*}, ie a type or type variable, 10
                 Compatibility of built-in function arguments with function arity via S, 13
\approx_S
\forall v
                 Type quantification, 10
                 The denotation of the built-in function b at the type assignment S, 11
[\![b]\!]_{\mathcal{S}}
Eval
                 Evaluation of built-in functions, 13
Term reduction
\boldsymbol{A}
                 Well-formed partial built-in function application, 14
В
                 Partial built-in function application (possibly ill-formed), 14
В
                 Set of all partial built-in function applications, 14
V
                 Plutus Core value, 14
\beta(B)
                 Function in partial builtin application B, 14
||B||
                 Size of partial builtin application B, 15
next(A)
                 Next argument type (or force) required by a partial builtin application A, 15
args(A)
                 Term arguments received so far by partial builtin application A, 15
[V/x]M
                 Capture-avoiding substitution of value V for variable x in term M, 16
                 Reduction frame for contextual semantics: [\_M], [V\_], (force \_), 17
CEK machine
\triangleright
                 CEK compute phase, 18
⊲
                 CEK return phase, 18
                 CEK error state, 18
CEK halting state, 18
```

```
Σ
                   CEK machine state, 18
                   CEK machine stack, 18
                   \text{CEK stack frame: (force\_), [\_(M,\rho)], [V\_], (constriM_0 \dots M_{i-1}\_V_{i+1} \dots V_n), (case\_M_0 \dots M_n),}
f
                   CEK value: \langle \text{con } T \, c \rangle, \langle \text{delay } M \, \rho \rangle, \langle \text{lam } x \, M \, \rho \rangle, \langle \text{builtin } b \, \overline{V} \, \eta \rangle, 18
V
                   CEK environment, 18
ρ
                   Value bound to variable x in environment \rho, 18
\rho[x]
                   Arguments expected by partial builtin application, 18
\mathcal{U}(V)
                   Discharge a CEK value V to obtain a Plutus Core term, 21
M@\rho
                   Discharge all variables bound by \rho in the term M, 21
Serialisation and deserialisation
                   CBOR decoder for data, 30
```

 $\mathcal{D}_X$  $\mathcal{E}_X$ CBOR encoder for data, 30

 $\mathsf{D}_X$ Flat decoder, 35 Flat encoder, 35  $\mathsf{E}_X$