

# The Welkin Standard

Oscar Bender-Stone

December 2023

## LICENSE

The Welkin Standard 0.1 is licensed under a Creative Commons Attribution 4.0 International License. A copy of this license is provided at <http://creativecommons.org/licenses/by-sa/4.0/>.

This document describes the Welkin information language, a programming language aimed at preserving, analyzing, and extending information.

This edition of the standard, in English, is the basis for all other translations. Only the grammars will be given in an English (specifically ASCII) and should be copied identically. However, these grammars can be built upon via Section ??, and any other terms in this document may be translated or changed as necessary.

Throughout this document, every instance of “Welkin grammar” means “standard Welkin grammar.” Every instance of “Welkin interpreter” means “conformant Welkin interpreter.” For a definition of conformance, refer to Section ??

## 1 Preliminaries

### 1.1 Mathematical Background

This standard requires a cursory background in discrete mathematics, parsing, and Backus-Naur Form (BNF). A reading of [1] and [2] suffices to understand this document. We clarify on our mathematical notation below.

Let  $A, B, C$  be sets. We denote by  $A^*$  the set of finite sequences of  $A$  (or  $\mathbf{A}^*$  in grammars), including the empty sequence  $\varepsilon := \emptyset$ , and we call each  $w \in A^*$  a **word**.

Let  $f : A \rightarrow B$ ,  $g : B \rightarrow C$  be functions. We denote the composite of  $f$  and  $g$  by  $g \circ f$ , and by  $f^* : A \rightarrow B^*$  the extension of  $f$  to finite sequences of  $A$  componentwise to finite sequences in  $B$ .

### 1.2 Character Encodings

We define text, character encodings, and character decodings as abstract notions. The discussion here may be carried out in terms of bytes and specific data formats, but these concepts are beyond the scope of this standard.

Let Char be a finite set. An **encoding** is an injective mapping  $\mathcal{E} : \mathbb{N} \rightarrow \text{Char}$ . The associated **decoding** is the left-inverse  $\mathcal{D} : \text{Char} \rightarrow \mathbb{N}$  of  $\mathcal{E}$ . We denote  $\text{CHAR} = \mathcal{D}(\text{Char})$ .

Character encodings may be given as finite tables, matching natural numbers with characters. Several major encodings are defined in the following sources.

- US-ASCII [1]. We will refer to this simply as ASCII<sup>1</sup>, but there are subtle variations across specific nationalities and applications (see [1]).
- UTF-8, UTF-16 [2]. The Unicode Standard defines encodings across numerous human languages and unique characters.

Although ASCII is a subset of UTF-8, this standard will prioritize ASCII as much as possible. Every BNF grammar (see Section 2.2) for this standard will be written in ASCII as a unifying encoding, but users may create grammars using UTF-8, UTF-16, or their own encodings (see Section 2.2).

**Definition 1.1.** A BNF **encoding** under Char is an encoding  $\mathcal{E} : \mathbb{N} \rightarrow \text{Char}$  along with four subsets  $\text{CHAR}^*$ ,

- WHITE\_SPACES
- NUMBERS
- DELIMITERS
- STRING\_ESCAPES

such that WHITE\_SPACES is pair-wise disjoint with every other subset.

Moreover, let STRING be the set of strings over DELIMITERS, where a string is a word  $d_1 w d_2$  such that  $w \in (\text{CHAR} \setminus \text{DELIMITERS} \cup \text{STRING\_ESCAPES}(d_1, d_2))^*$ . Notice that DELIMITERS, STRING\_ESCAPES can be defined in terms of STRING.

In our BNF encoding, we assume all sets above are non-empty. The first is an optional addition, adding support for defining machine representable numbers; the second distinguishes words from one another; and the last two are used for arbitrary strings. Any of these sets, however, can be made empty with a different grammar (see Section 2.2).

Every Welkin grammar, written in a BNF metasyntax, uses the ASCII encoding with

- NUMBERS =  $\emptyset$ ,
- WHITE\_SPACES = ASCII whitespace characters,
- STRING = all ASCII characters enclosed by single or double quotes. We escape single (double) quotes via  $\backslash$  (").

A **text** is a subset of  $\text{CHAR}^*$ . We will not consider streaming issues, i.e., we will assume every Welkin file is present at one time.

---

<sup>1</sup>American Standard Code for Information Exchange

### 1.3 BNF Variant

Our variant of BNF uses the notation shown in Table ?? and in Definition ?.?. Our notation, as well as every standard Welkin grammar, is available in ASCII.

Each BNF has an associated subset  $\text{RESERVED} \subseteq \text{CHAR}^*$  for any literals that appear in the grammar. We will explicitly state these for Welkin variants in the next section.

Concept	Notation	Example
Rule Assignment	=	term = atom
Empty Word (In ASCII)	empty	term = empty
Alternation		boolean = true   false
Concatenation (No white spaces inbetween rules)	Separate with a period (.)	function = name.“(”.number.“)”
Concatenation (Zero or more white spaces allowed)	Separate with white space	data = date name
Groupings	Parantheses ( )	boolean = (true   false)
Literals	“word”	boolean = “true”   “false”
Choice Names	terms → rule	boolean = “true” → true   “false” → false equivalent to boolean = true   false true = “true” false = “false”
Rule Substitution (Definition ?.?)	rule1[rule2 → rule3]	rule1 = rule2 , equivalent to s

We introduce one new notion that is definable with BNFs.

**Definition 1.2.** (Rule Substitution) Let  $rule_1, rule_2$  be rules with  $rule_2$  appearing on the right hand side of  $rule_1$ . Suppose  $rule_3$  appears on the lefthand side of  $rule_2$ . Then  $rule_1[rule_2 \rightarrow rule_3]$  is  $rule_1$  with every instance of  $rule_2$  replaced by  $rule_3$ .

Some particular cases:

- $term[empty]$  means no rule should be applied.

An conformant parser for all three grammars should include the ability to compose and override rules. This ensures that any updates to grammars, if necessary, are isolated. See more details in Section ??.

## 2 The Welkin Language

There are three fundamental variants of Welkin that define the foundation for the language:

- Base Welkin, mirroring the key properties of the core data structure.
- Attribute Welkin, extending Base Welkin with attributes. Attributes are a limited type of directive that can customize how the interpreter accepts or presents data.
- Binder Welkin, enabling arbitrary evaluation of Welkin files and access to the user’s operating system. This is equivalent to Attribute Welkin with three new directives: `@eval`, `@exec`, and `@bind`.

Each of these variants can be parsed with LALR parsers and fundamentally have the same semantics. However, in Binder Welkin, `@eval` makes the interpreter Turing complete (see Section ??), and using `@exec` can run external programs that impact the user’s system. For this reason, Binder Welkin is a separate, optional component, as detailed in the Section ??

### Syntax

Define each RESERVED set as follows.

- $\text{RESERVED}_{\text{base}} = \{\{, \}, ., -, ->, <\},$
- $\text{RESERVED}_{\text{attribute}} = \{\{, \}, (, ), [, ], ., , , -, ->, <-, @\},$
- $\text{RESERVED}_{\text{binder}} = \text{RESERVED}_{\text{attribute}}.$

Each grammar is provided in Table 3.?

### Semantics

We break down our semantics first by terms. Directives are handled separately in the next section.

**Definition 2.1.** Equality of terms.

- **Basis.** Two units are equal if they are the same kind and obey one of the following.
  - **ident** terms are equal if their corresponding characters are equal,
  - **string** terms are equal if their corresponding contents are equal. Thus, if `"`, `'` are delimiters, `"A"` coincides with `'A'`,

Variant	Grammar
Base Welkin	<pre> terms      = term* term       = (graph   connections   member   unit) graph      = unit? "{" terms "}" connections = term (connector term)+ connector  = "&lt;" term "&lt;" → edge               "&gt;" term "&gt;" → left_arrow               "&lt;-" term "&lt;'" → right_arrow member     = unit? (".(ident   string)?   "#".num )+ unit       = ident   string   num ident      = CHAR* string     = STRING num        = NUMBER </pre>
Attribute Welkin	<pre> %import grammars/base.txt %override term term      = "@".(directive   graph[directive])             construct   graph   connection             member   unit directive = attributes attributes = "import".tuple → import             "self".(member?) → self             "alias".graph[empty] → alias             "resource".graph[unit] → resources             "metadata".graph[unit] → metadata             "input".graph → input             "parse".(graph   unit) → parse             "validate".tuple → validate             "record".term → record             "output".graph → output             graph[unit] → custom construct = operation   tuple   list   series   all_terms operation = term.tuple   term unit term tuple     = "(" series ")" list      = "[" series "]" </pre>
Binder Welkin	<pre> %import grammars/attribute.txt %override directives directives = attributes   binders binders    = "eval".tuple[unit] → eval             "exec".tuple[string] → exec             "bind".graph[empty] → bind </pre>

- **num** terms are equal if they represent the same value. Thus, **1** coincides with **10E**.

- **Recursion.**

- Two members are the same if they contain the same list of units.
- Two connectors are equal if they are equal as terms.
- Two connections are equal if they connect the same terms and have equal connectors.
- Two graphs are equal if they contain the same terms.

A **scope** is recursively defined and intuitively is a level of terms.

**Definition 2.2.** (Scope) Let  $t$  be a term.

- If  $t$  is not contained in a graph, then  $\text{scope}(\text{term}) = 0$ ,
- If  $G' \in G$  are both graphs and  $\text{scope}(G') = n$ , then for all  $t \in G'$ ,  $\text{scope}(t) = \text{scope}(G) + 1$ .

A **valid** base Welkin file consists satisfies a unique naming rule: in every scope, there are no name collisons. In particular, every graph must **only be defined once**. Note that, by the way equality was defined between two numbers, there can only be one representation of a given number in a scope. For example, using **1** and **10E-1** in the same scope would produce a name collison. We first form an Abstract Syntax Tree (AST), from which we form the final stored data in a **Welkin Information Graph**.

**Definition 2.3.** Base Welkin is parsed into the following AST  $\mathcal{A}$ .

- Every term is a new subtree with its contents as children.
- Every graph is an ordered pair of its name and list of children.
- Every connection is an ordered pair:
  - Left arrows  $u \xrightarrow{e} v$  correspond to a triple  $(u, e, v)$ ;
  - Right arrows  $u \xrightarrow{e} v$  correspond to the triple  $(v, e, u)$ ;
  - Edges correspond to both a left and right arrow.
- Every unit is encoded via  $\mathcal{E}^*$ . Numbers are further transformed into a machine representable form, which is dependent on the implementation.

**Definition 2.4.** A **Welkin Information Graph**  $G = (V, C, L, l)$  consists of

- set  $V$  of **vertices**, with  $\varepsilon \notin V$ ,
- a set  $C \subseteq V \times (V \cup \{\epsilon\}) \times V$  **connections**,
- a set  $L$  of **labels**, and

- an injective function  $l : L \rightarrow V$ , called the **labeling function**,

such that  $C$  is reflexive: for all  $A \in V$ ,  $(A, A, A) \in C$ . The last condition implies that the function  $i : V \rightarrow C$ , given by  $i(A) = (A, A, A)$ , is well-defined. From this map, we abuse notation and write  $A$  and  $(A, A, A)$  synonymously. (The use of this notation is clear from context).

**Definition 2.5.** Let  $G = (V_G, C_G, L_G, l_G)$ ,  $H = (V_H, C_H, L_H, l_H)$  be WIGS. A **morphism**  $f : G \rightarrow H$  is a pair of functions  $f_V : V_G \rightarrow V_H$ ,  $f_C : C_G \rightarrow C_H$  with  $f_C(A, e, B) = (f_V(A), f_V(e), f_V(B))$  for  $A, e, B \in V_G$ . Equivalently,  $(A, e, B) \in C_G$  implies  $(f_V(A), f_V(e), f_V(B)) \in C_H$ .

- $f$  *preserves connections*:  $(A, e, B) \in C_G$  implies  $(f(A), f(e), f(B)) \in C_H$ ,
- $f$  *is 3-injective*  $f(A) = f(e) = f(B)$  and  $(f(A), f(e), f(B)) \in C_H$  implies  $A = e = B$  and, subsequently,  $(A, e, B) \in C_k$ .

**Definition 2.6.** Let  $G$  be a WIG. The **Welkin Canonical Form (WCF)** of  $G$  is the WIG  $\text{Can}(G) = (V', E', C', L, l')$ , where

- $\text{Can}(V) = \{0 \leq k \leq |G|\}$ ,
- $(\text{Can}(A), \text{Can}(B)) \in \text{Can}(E)$  if  $(A, B) \in E$
- $(\text{Can}(A), \text{Can}(E), \text{Can}(B)) \in \text{Can}(C)$  if  $(A, E, B) \in C$ .
- (Condition on labeling)

The next lemma defines the basis for the recorder. Every WIG must be put into a universal form that can be labeled in completely different ways. The underlying structure of the WIG is always preserved, no matter what labels are added.

**Lemma 2.1.** For any WIG  $G$ ,  $G \cong \text{Can}(G)$ .

In fact, the structure of a WIG can only be appended. Any new additions to the graph do not remove existing structure.

## 3 General Application Behavior

Note that all apparent structures may be adjusted under cryptomorphism.

### 3.0.1 Bootstrap

The welkin/bootstrap file facilitates the user API to Attribute Welkin. It is currently located at <https://github.com/AstralBearStudios/welkin> and is essential for creating a stable Welkin interpreter.

Every Welkin interpreter must implement Base Welkin. To implement Attribute Welkin, the parser must be equipped with a parser for this variant, create

a WIG for **attribute** (with all imported units in bootstrap), and follow the General Application Guidelines for its implementation.

Finally, to support Base Welkin, the interpreter must be written with a programming language that is Turing complete and can execute system resources. The final command, **@bind**, is defined in terms of these notions.

## Directives

Each directive relies on the following components.

- Input: processes an input or set of inputs.
- Parser: takes in a Welkin file and generates an AST,
- Validator: ensures that the AST is valid, raising an error that directly points to a violation,
- From here, an AST may be processed by three different means:
  - Recorder: takes the AST, converts it into a WIG in WCF, serializes the data,
  - Output: displays some information provided in a Welkin file to the user.
  - Evaluator: evaluates, executes, or binds commands into Welkin units.

All official implementations may implement

- Base Welkin alone,
- Attribute and Base Welkin, or
- All three variants.

In each case, these variants must be implemented according to the welkin/bootstrap file. This file bootstraps the essential information from this standard. (See Section ?? on a specification of this bootstrap)

## 3.1 Customization

All Welkin files are infinitely customizable via the welkin config file, which is written in Attribute Welkin. Any attribute can be used, and other Welkin files can be imported. A base config file is required to customize a Welkin grammar. From there, configs can be arbitrarily nested to create and connect any desired (context-free) grammar, validator, recorder, and displayer.

# 4 Core Algorithms

## 4.1 Graph Encoding