

The Intention Language Reference

Anna Bukowska, Marek Kaput

May 28, 2018

Contents

Preface	iii
I The Intention Language	1
1 Introduction	2
1.1 Notational Conventions	2
1.2 Compile-time and Run-time	4
1.3 Program Structure	4
1.4 Values, Types, Terms and Results	5
1.5 Namespaces	5
2 Lexical structure	7
2.1 Input format	7
2.2 Special Lexical Productions	8
2.3 Identifiers, Keywords and Operators	9
2.4 Comments	9
2.5 Literals	10

<i>CONTENTS</i>	ii
2.5.1 Numeric Literals	10
2.5.2 String Literals	10
3 Expressions	14
4 Items	16
5 Modules and Assemblies	17
II The Intentio Standard Library	18
6 Introduction	19
A Influences	20
Bibliography	21
Index	22

Preface

We live in times of rapid emerging of new, modern programming languages. Some of them, like Rust[8], Go[7] or Swift[1], have proved that programming language styles have not settled down and there is still room for new ideas, especially for merging existing paradigms. Fundamentally, one can observe a shift from imperative programming to functional programming.

Despite all these changes, not all ideas get a chance to shine. Some of them are becoming forgotten and treated as esoteric. One of these is goal-oriented evaluation, with Icon[11] being one of its most *iconic* implementers. Authors of this document believe that Icon exposes some very interesting ideas and they made an attempt to recreate them in a new programming language: *Intentio*, named after Latin for *intention*.

This document is the primary reference for the Intentio programming language. It consists of three parts:

- Chapters that *semi-formally*¹ describe each language construct and their use.
- Chapters that *semi-formally* describe runtime services and standard library that are core part of the language.
- Various appendix chapters.

This document does not serve as a beginner-friendly introduction to the language.

¹This document tries to maintain reasonably formal description of all items, but there are no guarantees all cases are described. As a fallback, the *intentioc*[2] compiler can be used as secondary reference.

Goals

The primary goals when designing Intentio was for language to satisfy following constraints:

1. It should feature core concepts of Icon language: goal-oriented evaluation and generators
2. It should support Unicode character set
3. It should be fast and easy to build prototype applications, the language should be *ergonomic*² from developer perspective **and** *IDE friendly*³
4. It should feature rich capabilities in processing textual data

Our main goals were **not**:

1. The language should be general purpose language
2. Compilation time should be short
3. Memory usage of Intentio programs should be low
4. It should be easy to integrate with other languages

Acknowledges

The structure and parts of content of this document are inspired by two existing language specifications which we believe are good examples to follow: The Rust Reference[9] and Haskell 2010 Language Report[5].

²Source code of Intentio programs should be concise, pleasant to write and easy to reason about.

³Source code of Intentio programs should be easily processable by text editors and analysis tools.

Part I

The Intentio Language

Chapter 1

Introduction

Intentio is domain specific, imperative programming language oriented for processing textual data. Intentio provides goal-oriented execution, generators, strong dynamic typing with optional type annotations, and a rich set of primitive data types, including Unicode strings, lists, arrays, maps, sets, arbitrary and fixed precision integers, and floating-point numbers. Intentio tries to incorporate ideas of the Icon[11] programming language into modern programming patterns.

This part defines the syntax of Intentio language and informal abstract semantics for the meaning of such programs. We leave as implementation detail how Intentio programs are manipulated, interpreted, compiled, etc. This includes all steps from source code to running program, programming environments and error messages. This also means that this document do not describe the reference compiler for Intentio language - *intentioc*[2].

1.1 Notational Conventions

Grammar

Throughout this document a BNF-style notational syntax is used to describe lexical structure and grammar:

$$nonterminal \rightarrow terminal \mid alternative$$

Following conventions are used for presenting productions syntax:

token	terminal symbol (in fixed-width font)
<i>rule</i>	nonterminal symbol (in italic font)
(pat)	grouping
$[pat]$	optional (0 or 1 times)
$\langle pat \rangle$	repetition (1 to n times)
$\{pat\}$	optional repetition (0 to n times)
$pat_1 pat_2$	concatenation
$pat_1 \mid pat_2$	alternation
$pat_{\langle pat' \rangle}$	difference (symbols generated by pat except those generated by pat')

Parameterized productions

Some productions in Intention's grammar (like raw string literals) cannot be expressed using context-free grammar with finite number of productions. In order to reduce need of falling back to natural language, a concept of *parameterized productions* is used throughout this document.

A production $A(x_1, x_2, \dots, x_n) \rightarrow B$, parameterized over arguments x_1, x_2, \dots, x_n , defines different production for each combination of its arguments.

Unicode productions

A few productions in Intention's grammar permit Unicode^[10] code points outside the ASCII range. These productions are defined in terms of character properties specified in the Unicode standard, rather than in terms of ASCII-range code points. Intention compilers are expected to make use of new versions of Unicode as they are made available.

Source code listings

Examples of Intention program fragments are given in fixed-width font:


```

fun main() {
  x, y := 4, 3;
  writeln(f"sum = ${x + y}");
}

```

In some situations, there are *placeholders* in program fragments representing arbitrary pieces of Intentio code are written in italics. By convention *e* will mean expressions, *d* - item declarations, *t* - types, etc.:

```

if e1 { e2 } else { e3 }

```

1.2 Compile-time and Run-time

Intentio's semantics obey a *phase distinction* between compile-time and run-time¹. Semantic rules that have a static interpretation govern the success or failure of compilation, while semantic rules that have a dynamic interpretation govern the behaviour of the program at run-time.

1.3 Program Structure

An Intentio program is structured syntactically and semantically into five abstract levels:

1. At the topmost of each Intentio program or library is an *assembly*. In compiled environments assembly is an unit of compilation, while in interpreted environments assembly is a whole set of loaded modules.
2. At the topmost of each assembly is a set of *modules*. Modules provide a way to control namespaces and to re-use software in larger programs. A particular source code file of Intentio program consists of one module. Module structure is flat, there is no concept of submodule.
3. The top level of each module is a set of *item declarations*. An item is a component of module, such as a function, type definition or constant variable.

¹In interpreter environments, compile-time would consist of syntactic analysis and linting.

4. Items which contain the real, executable code are built of *expressions*. Expression denotes how to evaluate a *term* and evaluating expression returns a *result*.
5. At the bottom level is Intentio's *lexical structure*. It describes how to build tokens - the most basic blocks of program's source code from sequences of characters in source file.

1.4 Values, Types, Terms and Results

A *value* is a representation of some entity which can be manipulated by a program. A *type* is a tag that defines the interpretation of value representation. Values and types are not mixed in Intentio. Values by itself are untyped, value's type is required to perform any kind of operation on value.

A (*value, type*) pair is called a *term*. Terms represent data yielded from evaluating expressions. Intentio is *strongly typed* so implicit type conversions do not exist in the language, but it is not prohibited to include casts in expressions semantics (thus `5 + 4.0` runs successfully).

Evaluating expressions may either succeed or fail. A tagged union of successfully evaluated result or failure with information describing what failed is called a *result*. Terms and results are the basic blocks of representing information in Intentio.

Following Haskell-style code listing describes relationships between these concepts:

```
data Value = ...
data Type = ...

newtype Term = (Value, Type)

data Result = Succ Term
           | Fail Term
```

1.5 Namespaces

There are three distinct namespaces in Intentio:

Item namespace Consists item and variable names.

Module namespace Consists of module names and import renames.

Type namespace Consists of type names.

There are no constraints on names belonging to particular namespace, therefore it is possible for name `Int` to simultaneously denote an item/variable, module and type.

Chapter 2

Lexical structure

This chapter describes the lexical structure of Intentio. Most of the details may be skipped in a first reading of the reference.

In this chapter all white space is expressed explicitly in syntax descriptions, there is no implicit space between juxtaposed symbols. Terminal characters represent real characters in program source code.

2.1 Input format

Intentio program source is interpreted as a sequence of Unicode code points encoded in UTF-8, though most grammar rules are defined in terms of printable ASCII code points.

Intentio is *case sensitive* language and each code point is distinct; for instance, upper and lower case letters are different characters.

The NUL character (U+0000) may be not allowed in whole program source text.

If an UTF-8-encoded byte order mark (U+FEFF) is the first Unicode code point in program source text, it may be ignored. Byte order mark may be not allowed anywhere else in program source text.

2.2 Special Lexical Productions

Following productions define Unicode character sets which are used to define non pure ASCII productions. These productions do not have any semantic meaning themselves.

- *xidstart* and *xidcont* are sets of characters that have properties *XID_start* and *XID_continue* as in Unicode Standard Annex #31[3], these productions define valid identifier characters
- *any* is any single Unicode character with all implementation-specific character restrictions applied
- *eol* matches either line feed character `\n` (U+000A, Unix-style newline marker) or carriage return and then line feed characters `\r\n` (U+000D U+000A, Windows-style newline marker)

Additionally:

noneol \rightarrow *any*_{*eol*}

decdig \rightarrow 0 | 1 | ... | 9

bindig \rightarrow 0 | 1

octdig \rightarrow 0 | 1 | ... | 7

hexdig \rightarrow 0 | 1 | ... | 9 | A | B | ... | F | a | b | ... | f

2.3 Identifiers, Keywords and Operators

$$\begin{aligned}
 id &\rightarrow (\textcolor{blue}{xidstart} \{ \textcolor{blue}{xidcont} \mid ' \})_{\langle keyword \rangle} \\
 qid &\rightarrow [\textcolor{blue}{id} :] id \\
 keyword &\rightarrow \text{abstract} \mid \text{and} \mid \text{break} \mid \text{case} \mid \text{const} \mid \text{continue} \\
 &\quad \mid \text{do} \mid \text{else} \mid \text{enum} \mid \text{fail} \mid \text{fun} \mid \text{if} \mid \text{impl} \mid \text{import} \\
 &\quad \mid \text{in} \mid \text{is} \mid \text{loop} \mid \text{mod} \mid \text{not} \mid \text{or} \mid \text{return} \mid \text{struct} \\
 &\quad \mid \text{type} \mid \text{where} \mid \text{while} \mid \text{yield} \mid _ \\
 operator &\rightarrow + \mid - \mid * \mid / \\
 &\quad \mid (\mid) \mid [\mid] \mid \{ \mid \} \mid : \\
 &\quad \mid == \mid < \mid <= \mid > \mid >= \\
 &\quad \mid := \mid <- \mid \$ \mid \%
 \end{aligned}$$

An *identifier* consists of a "letter" or underscore followed by zero or more letters, digits, underscores, and single quotes. Simple, unqualified identifiers (*id*) are always resolved within current module and scope. In order to be able to specify which module identifier belongs to, in most places identifier may be prefixed with module name and ":" character to form a *qualified identifier* (*qid*).

Keywords are identifier-like tokens which have special meaning in the grammar, all of them are excluded from the *id* rule. *Operators* are another special tokens, these ones are formed from symbol characters. *keyword* and *operator* productions have no use in Intentio grammar definition, instead particular tokens are used.

Implementations that offer lints or warnings for unused parameters/variables/items are encouraged to suppress such warnings for identifiers beginning with underscore. This allows programmers to use `_arg` for a parameter that they expect to be unused.

2.4 Comments

$$comment \rightarrow \# \{ \textcolor{blue}{noneol} \} \textcolor{blue}{eol}$$

Comments are valid white space. Comments in Intentio are only line-based, there is no concept of block comment.

2.5 Literals

literal \rightarrow *integer* | *float* | *string*

2.5.1 Numeric Literals

decimal \rightarrow *decdig* { *decdig* | *_* }
binary \rightarrow (0b | 0B) *bindig* { *bindig* | *_* }
octal \rightarrow (0o | 0O) *octdig* { *octdig* | *_* }
hexadecimal \rightarrow (0x | 0X) *hexdig* { *hexdig* | *_* }

integer \rightarrow *binary* | *octal* | *decimal* | *hexadecimal*

exponent \rightarrow (e | E) [+ | -] { *_* } *decimal*

float \rightarrow *decimal* . *decimal* [*exponent*]
| *decimal* *exponent*

A *numeric literal* is either an *integer literal* or *floating-point literal*. Integer literals may be given in decimal (the default), binary (prefixed by 0b or 0B), octal (prefixed by 0o or 0O) or hexadecimal notation (prefixed by 0x or 0X). Floating-point literals are always decimal. A floating literal must contain digits both before and after the decimal point. A "_" character is allowed inside numeric literals for visual separation of digit groups, for instance 476_981__109_528_ is equal to 476981109528. Negative numeric literals are described grammatically, not lexically.

2.5.2 String Literals

string \rightarrow [*stringmod*] (*string'* | *charstring* | *rawstring* | *regexstring*)

Intentionio features very flexible string literals syntax. String literals can be written in four forms (regular, character, raw and regular expression) that make the literal

compile to one of three value types (**String**, **Char** or **Regex**). String literals can be also prefixed with few modifiers that alter literal value before compiling it to value.

Regular String Literals

$$string' \rightarrow " \{ \textcolor{blue}{any}_{\langle " | \backslash \rangle} \mid \textcolor{blue}{escseq} \} "$$

A *string literal* is a sequence of Unicode characters, typed either directly or via escape sequence. String literals are compiled to **String** terms.

Character Literals

$$charstring \rightarrow \text{c} " (\textcolor{blue}{any}_{\langle , | \backslash \rangle} \mid \textcolor{blue}{escseq}) "$$

A *character literal* is a single Unicode character string, typed either directly or via escape sequence. Character literals are compiled to **Char** terms. The implementation is required to verify that character literal makes only one character.

Raw String Literals

$$\begin{aligned} rawstring &\rightarrow \text{r } \textcolor{blue}{rawstring}'(0) \\ rawstring'(n) &\rightarrow " \textcolor{blue}{rawstring}''(n) " \\ &\mid \# \textcolor{blue}{rawstring}'(n+1) \# \\ rawstring''(n) &\rightarrow \text{An } \{ \textcolor{blue}{any} \} \text{ that does not contain} \\ &\quad " \text{ followed by } \# \text{ repeated } n \text{ times.} \end{aligned}$$

A *raw string literal* does not process any escape sequences. It starts with letter **r**, followed by zero or more repetitions of hash symbol (**#**) and double quote (**"**). The raw string body can contain any sequence of Unicode characters and is terminated only by another double quote (**"**) followed by the same number of hashes (**#**) that preceded the opening quote.

Examples of raw string literals:


```

"foo"      == r"foo"      # foo
"\\"foo\\" == r#"foo"#    # "foo"
"x #"y"    == r##"x #"y"  # x #" y

```

Regular Expression Literals

$$regexstring \rightarrow x (\textcolor{blue}{string}' \mid rawstring)$$

A *regular expression literal* is a string literal that represents a regular expression and is compiled to **Regex** term. The exact syntactic and semantic details of regular expressions in Intention are implementation dependent.

Modifiers

$$stringmod \rightarrow t \mid u$$

A *string literal modifier* is a special flag that, when enabled, adds a step of processing of the literal value before compiling it to Intention term. The order of modifiers is respected, they are processed from left-most modifier to right-most one.

Because modifiers alter literal contents before its compiling, they can fundamentally change their meaning, for instance the literal `tc" x "` should successfully compile and evaluate to single character `"x"`.

Intention provides following modifiers:

- **t** - *trim*: The trim modifier removes all white space characters from both sides of the string literal value.
- **u** - *unindent*: The unindent modifier gets the white-space-only prefix of the string literal value and then removes it from each line of the value.

Escape Sequences

$charescseq \rightarrow \backslash ' \mid \backslash " \mid \backslash n \mid \backslash r \mid \backslash t \mid \backslash \backslash \mid \backslash 0$
 $asciiescseq \rightarrow \backslash x \text{ *hexdig hexdig*}$
 $unicodeescseq \rightarrow \backslash u \{ \text{ *hexdig* } \mid _ \}$

 $escseq \rightarrow \text{ *charescseq* } \mid \text{ *asciiescseq* } \mid \text{ *unicodeescseq*}$

Some *escape sequences* are available in non-raw string literals. An escape starts with a backslash character (\) and continues with one of the following forms:

- An *8-bit code point escape sequence* starts with letter x and is followed by exactly two hex digits with value up to $0x7f$. It denotes an ASCII character with value equal to provided hex value. Higher values are not permitted because it is ambiguous whether they mean Unicode code points or byte values, though the implementation should accept them on lexical level for better user experience.
- A *24-bit code point escape sequence* starts with letter u and is followed by up to six hex digits surrounded by braces "{" and "}". It denotes the Unicode code point equal to the provided hex value. The implementation should accept more digits on lexical level for better user experience.
- The *character escape sequences* are convenience shortcuts for 8-bit code point escape sequences. Following table describes exact translations:

Escape	Character	Unicode
$\backslash '$	'	U+0027
$\backslash "$	"	U+0022
$\backslash n$	$\backslash n$	U+000A
$\backslash r$	$\backslash r$	U+000D
$\backslash t$	$\backslash t$	U+0009
$\backslash \backslash$	\backslash	U+005C
$\backslash 0$	NUL	U+0000

Chapter 3

Expressions

This chapter describes syntax and semantics of Intentionio *expressions*. Intentionio is an expression language. This means that all forms of result-producing or effect-causing evaluation fall into uniform syntax category of expressions. Usually each kind of expression can nest within each other kind of expression, and rules for evaluation of expressions involve specifying both the result produced by the expression and the order in which its sub-expressions are themselves evaluated.

Intentionio does not have a concept of *statement* known from other programming languages.

$expr$	\rightarrow	<code>qid</code>	(variable or item)
		<code>literal</code>	
		<code>blockexpr</code>	(block expression)
		<code>unopexpr</code>	(unary operator expression)
		<code>binopexpr</code>	(binary operator expression)
		<code>(expr)</code>	(parenthesized expression)
		<code>callepr</code>	(function call)
		<code>loopexpr</code>	(loops)
		<code>ifexpr</code>	(conditionals)
		<code>returnexpr</code>	(return expression)

An *expression* is a syntactic construct that *evaluates* to a term. It may either *succeed* or *fail* resulting in a result. During the evaluation, expression may perform *side effects*, for example it can mutate some state or perform execution jump. The

meaning of each kind of expression dictates several things:

- Whether or not to evaluate the sub-expressions when evaluating the expression
- The order in which to evaluate the sub-expressions
- How to combine the sub-expressions' results to obtain the result of the expression

Chapter 4

Items

Chapter 5

Modules and Assemblies

Modules and items are entirely determined at compile-time, remain fixed during execution, and may reside in read-only memory. This limitation does not apply to assemblies, it is allowed to provide mechanisms to dynamically compile, load and unload assemblies at run-time.

Part II

The Intentio Standard Library

Chapter 6

Introduction

This part defines the Intentio Standard Library (shortly *stdlib*), its contents, semantics and the *prelude* which is automatically imported in each Intentio program. This library provides essentials for building proper Intentio programs, some of which should be used by the implementation through compiler intrinsics. For developer convenience it also provides common utilities which ease application development and make a standard for code interoperation. The standard library must be distributed with each implementation of the Intentio language.

Appendix A

Influences

Intentio is not particularly original language, having the language Icon as the main source of inspiration, but also borrowing design element from wide range of other sources. Some of these are listed below:

- Icon[11]: goal-oriented execution, generators
- Rust[8]: syntax, string literals
- Erlang[4]: syntax, modules
- Python[6]: syntax

Bibliography

- [1] Apple Inc. The Swift Programming Language. <https://swift.org/>.
- [2] A. Bukowska and M. Kaput. *intentioc* - The reference Intentio compiler. <https://github.com/intentio-lang/intentio>.
- [3] M. Davis. Unicode Standard Annex #31: Unicode Identifier and Pattern Syntax. Technical Report Version 9.0.0.
- [4] Ericsson AB. Erlang programming language. <http://www.erlang.org/>.
- [5] S. Marlow. Haskell 2010 language report. <https://www.haskell.org/definition/haskell2010.pdf>.
- [6] Python Software Foundation. Python language reference, version 3.6. <https://docs.python.org/3.6/reference/>.
- [7] The Go Authors. The Go Programming Language Specification. <https://golang.org/ref/spec>.
- [8] The Rust Project Developers. The Rust Programming Language. <https://doc.rust-lang.org/book/>.
- [9] The Rust Project Developers. The Rust Reference. <https://doc.rust-lang.org/reference/>.
- [10] The Unicode Consortium. The Unicode Standard. Technical Report Version 9.0.0, Unicode Consortium.
- [11] University of Arizona. The Icon Programming Language. <https://www2.cs.arizona.edu/icon/>.

Index

character literal, [11](#)
comment, [9](#)
compile-time, [4](#)

escape sequence, [13](#)
expression, [14](#)

floating-point literal, [10](#)

identifier, [9](#)
integer literal, [10](#)
item namespace, [6](#)

keyword, [9](#)

module namespace, [6](#)

numeric literal, [10](#)

operator, [9](#)

qualified identifier, [9](#)

raw string literal, [11](#)
regular expression literal, [12](#)
result, [5](#)
run-time, [4](#)

string literal, [11](#)
string literal modifier, [12](#)

term, [5](#)
type, [5](#)
type namespace, [6](#)

value, [5](#)