

Data Type Representation

We define LLVM representations for a variety of classical and quantum data types. In each case, we specify an LLVM type alias for the underlying type. Compilers should always emit the type alias so that the semantic type can be determined by LLVM-level compiler passes.

Unit

The `Unit` data type is not represented as an LLVM type because its primary use is as a return type, and all callables use a tuple as the return type. If there is a need to represent the `Unit` type for some purpose, a tuple type with no contained elements should be used.

The one possible value of `Unit`, `()`, is represented as a null tuple pointer.

Simple Types

The simple types are those whose values are fixed-size and do not contain pointers. They are represented as follows:

Type	LLVM Representation	Comments
<code>Int</code>	<code>%Int = type i64</code>	
<code>Double</code>	<code>%Double = double</code>	
<code>Bool</code>	<code>%Bool = i8</code>	0 is false, 1 is true.
<code>Result</code>	<code>%Result = i8</code>	0 is Zero, 1 is One.
<code>Pauli</code>	<code>%Pauli = i8</code>	0 is PauliI, 1 is PauliX, 3 is PauliY, and 2 is PauliZ.
<code>Qubit</code>	<code>%Qubit = i8 addrspace(2) *</code>	A qubit is logically a pointer into an alternate address space. The only meaningful operation on this type, other than passing as an argument to a quantum instruction, is equality comparison.
<code>Range</code>	<code>%Range = {%Int, %Int, %Int}</code>	In order, these are the start, step, and end of the range. When passed as a function argument or return value, ranges should be passed by value.

The following global constants are defined for use with the `%Result` and `%Pauli` types:

```
@Result.One = %Result 1
@Result.Zero = %Result 0

@Pauli.I = %Pauli 0
@Pauli.X = %Pauli 1
@Pauli.Y = %Pauli 3
@Pauli.Z = %Pauli 2
```

Variably-Sized Types

The variably-sized types are those whose values vary in size but do not contain pointers. Values of these types are represented as pointers to LLVM structures. The structures may be allocated on the heap or be located in static memory for compile-time constants.

Each structure starts with an `i32` reference count for memory management use. Next there is an `i32` size field, followed by the actual data. The reference count field is managed using the library functions described below. For static (compile-time constant) values of these types, the reference count should be set to -1. The `reference` and `unreference` functions recognize this value and will act appropriately, not updating the reference count or attempting to release the value.

| Type | LLVM Representation | Comments | |-----|-----|-----| | String |
`%String = type { i32, i32, [0 x i8] }` | The byte array is the UTF-8 representation of the string. The size field is the size in bytes, not characters. | | BigInt | `%BigInt = type { i32, i32, [0 x i64] }` | The size field is the size in i64s. |

Strings

The following utility functions are provided by the classical runtime to support strings:

Function	Signature	Description
<code>quantum.rt.string_init</code>	<code>void(%String*)</code>	Initializes a heap-allocated string by setting the reference count to 1.
<code>quantum.rt.string_reference</code>	<code>void(%String*)</code>	Increments the reference count of a string.
<code>quantum.rt.string_unreference</code>	<code>void(%String*)</code>	Decrements the string's reference count and releases the string if the count is now zero.
<code>quantum.rt.string_concatenate</code>	<code>%String* (%String*, %String*)</code>	Creates a new heap-allocated string that is the concatenation of the two argument strings.
<code>quantum.rt.string_equal</code>	<code>i1(%String*, %String*)</code>	Returns true if the two strings are equal, false otherwise.

The following utility functions support converting values of other types to strings. In every case, the returned string is allocated on the heap; the string can't be allocated by the caller because the length of the string depends on the actual value.

Function	Signature	Description
<code>quantum.rt.int_to_string</code>	<code>%String*(%Int)</code>	Returns a string representation of the integer.
<code>quantum.rt.double_to_string</code>	<code>%String*(%Double)</code>	Returns a string representation of the double.
<code>quantum.rt.bool_to_string</code>	<code>%String*(%Bool)</code>	Returns a string representation of the Boolean.
<code>quantum.rt.result_to_string</code>	<code>%String*(%Result)</code>	Returns a string representation of the result.

Function	Signature	Description
quantum.rt.pauli_to_string	<code>%String*(%Pauli)</code>	Returns a string representation of the Pauli.
quantum.rt.qubit_to_string	<code>%String*(%Qubit)</code>	Returns a string representation of the qubit.
quantum.rt.range_to_string	<code>%String*(%Range)</code>	Returns a string representation of the range.
quantum.rt.bigint_to_string	<code>%String*(%BigInt*)</code>	Returns a string representation of the big integer.

Big Integers

The following utility functions are provided by the classical runtime to support big integers. Note that all returned big integers will be allocated on the heap.

Function	Signature	Description
quantum.rt.bigint_init	<code>void(%BigInt*)</code>	Initializes a heap-allocated big integer by setting the reference count to 1.
quantum.rt.bigint_reference	<code>void(%BigInt*)</code>	Increments the reference count of a big integer.
quantum.rt.bigint_unreference	<code>void(%BigInt*)</code>	Decrements the big integer's reference count and releases the big integer if the count is now zero.
quantum.rt.bigint_negate	<code>%BigInt* (%BigInt*)</code>	Returns the negative of the big integer.
quantum.rt.bigint_add	<code>%BigInt* (%BigInt*, %BigInt*)</code>	Adds two big integers and returns their sum.
quantum.rt.bigint_subtract	<code>%BigInt* (%BigInt*, %BigInt*)</code>	Subtracts the second big integer from the first and returns their difference.
quantum.rt.bigint_multiply	<code>%BigInt* (%BigInt*, %BigInt*)</code>	Multiplies two big integers and returns their product.
quantum.rt.bigint_divide	<code>%BigInt* (%BigInt*, %BigInt*)</code>	Divides the first big integer by the second and returns their quotient.
quantum.rt.bigint_modulus	<code>%BigInt* (%BigInt*, %BigInt*)</code>	Returns the first big integer modulo the second.
quantum.rt.bigint_power	<code>%BigInt* (%BigInt*, i64)</code>	Returns the big integer raised to the integer power.

Function	Signature	Description
quantum.rt.bigint_bitand	<code>%BigInt* (%BigInt*, %BigInt*)</code>	Returns the bitwise-AND of two big integers.
quantum.rt.bigint_bitor	<code>%BigInt* (%BigInt*, %BigInt*)</code>	Returns the bitwise-OR of two big integers.
quantum.rt.bigint_bitxor	<code>%BigInt* (%BigInt*, %BigInt*)</code>	Returns the bitwise-XOR of two big integers.
quantum.rt.bigint_bitnot	<code>%BigInt* (%BigInt*)</code>	Returns the bitwise complement of the big integer.
quantum.rt.bigint_shiftright	<code>%BigInt* (%BigInt*, i64)</code>	Returns the big integer arithmetically shifted left by the integer amount of bits.
quantum.rt.bigint_shiftright	<code>%BigInt* (%BigInt*, i64)</code>	Returns the big integer arithmetically shifted right by the integer amount of bits.
quantum.rt.bigint_equal	<code>i1(%BigInt*, %BigInt*)</code>	Returns true if the two big integers are equal, false otherwise.
quantum.rt.bigint_greater	<code>i1(%BigInt*, %BigInt*)</code>	Returns true if the first big integer is greater than the second, false otherwise.
quantum.rt.bigint_greater_eq	<code>i1(%BigInt*, %BigInt*)</code>	Returns true if the first big integer is greater than or equal to the second, false otherwise.

Container Types

Container types are those whose values contain other values. For instance, a tuple may contain several values of various types. The container types defined here are tuples and arrays.

Values of these types are represented as handles to LLVM structures, where the handles may be simple pointers or may contain additional information. The structures may be allocated on the heap or be located in static memory for compile-time constants. The handles should be treated as scalars and stored directly in LLVM variables or as direct fields of other data structures.

Each structure has an `i32` reference count for memory management use. The reference count field is managed using the library functions described in each section below. For static (compile-time constant) values of these types, the reference count should be set to -1. The `reference` and `unreference` functions recognize this value and will act appropriately, not updating the reference count or attempting to release the value.

Structures also contain a pointer to a destructor function; that is, to a function that is used to release values that are pointed to by this value. This function is generated by the compiler based on the specifics of the complex type. The precise details of the destructor function varies depending on the specific value type,

although the signature of the destructor function is always `void(i8*)`, where the argument passed to the function is a pointer to the value to be destroyed. Values that the compiler knows don't require destructor functions should have a null pointer in this field.

We define an LLVM type for a pointer to a destructor:

```
%Destructor = type void(i8*)*
```

Tuples and User-Defined Types

Tuple data, including values of user-defined types, is represented as an LLVM structure type containing the reference count and destructor fields described above as a standard header, followed by the tuple fields. We define an LLVM type for the tuple header:

```
%TupleHeader = type {i32, %Destructor}
```

For instance, a tuple containing two integers, `(Int, Int)`, would be represented in LLVM as `type {%TupleHeader, i64, i64}`.

Tuple handles are simple pointers to the underlying data structure. To satisfy LLVM's strong typing, tuple pointers are passed as pointers to the initial tuple header field, and then cast to pointers to the correct data structures by the receiver:

```
%TuplePointer = type %TupleHeader*
```

The destructor for a tuple will get called with a pointer to the tuple structure as argument when the tuple is being released. The destructor should recursively unreferencce any variably-sized or container components of the tuple. If the tuple only contains simple types, then the destructor pointer should be null.

Many languages provide immutable tuples, along with operators that allow a modified copy of an existing tuple to be created. In QIR, this is implemented by creating a new copy of the existing tuple on the heap, and then modifying the newly-created tuple in place. In some cases, if the compiler knows that the existing tuple is not used after the creation of the modified copy, it is possible to avoid the copy and modify the existing tuple as long as there are no other references to the tuple. The `tuple_is_writable` utility function is used to test if the tuple pointer is the only reference.

The following utility functions are provided by the classical runtime to support tuples and user-defined types:

Function	Signature	Description
----------	-----------	-------------

Function	Signature	Description
quantum.rt.tuple_init_stack	<code>void(%TuplePointer)</code>	Initializes a stack-allocated tuple by setting the reference count to -1 and the destructor to null. Any code required to recursively unreference contained components should be generated as in-line clean-up code.
quantum.rt.tuple_init_heap	<code>void(%TuplePointer, %Destructor)</code>	Initializes a heap-allocated tuple by setting the reference count to 1 and filling in the destructor.
quantum.rt.tuple_reference	<code>void(%TuplePointer)</code>	Increments the reference count of the tuple.
quantum.rt.tuple_unreference	<code>void(%TuplePointer)</code>	Decrements the tuple's reference count; calls its destructor and releases the tuple if the reference count is now zero.
quantum.rt.tuple_is_writable	<code>i1(%TuplePointer)</code>	Returns true if it is safe to modify the contents of the tuple; that is, it returns true if the reference count of the tuple is equal to 1.

Arrays

Array data is stored as an LLVM structure with the actual element data stored in a variably-sized byte array:

```
%ArrayData = type {i32, %Destructor, i32, i64, [0 x i8]}
```

The first `i32` is the reference count, and the `%Destructor` is the destructor for the array. The second `i32` is the size of each array element in bytes. The `i64` is the total number of elements in the array. Finally, the byte array holds the actual elements for the array. Pointers to elements will get cast to the actual element type by code that accesses the array.

The handle for an array is more complicated than a simple pointer. It holds the information that allows a list of array indices to be translated into a specific element offset:

```
%ArrayPointer = type {%ArrayData*, i32, [0 x i64]}
```

The `i32` is the number of dimensions in the array. The list of `i64`s provides information on the slice represented by this handle. Specifically, the first element of the list is an offset into the array data, and the next `d` elements (where `d` is the number of dimensions in the array) are the strides for each array. Finally, the last `d` elements are the lengths of the corresponding dimensions; these are used to validate that a set of dimension indices are within the array's bounds. Thus, the list has `2d+1` elements.

To access an element of the underlying array, the element index into the `ArrayData` is calculated by multiplying each dimension's index by the corresponding stride, adding the results together, and then adding

the offset. For example, for a 3-dimensional array with offset 7 and strides 24, 8, and 1, the element at `[3, 2, 5]` would be at element $7 + 3 \times 24 + 2 \times 8 + 5 \times 1 = 100$ in the underlying `%ArrayData`.

Array destructors are invoked slightly differently than tuple destructors. Because every element of an array is the same and requires the same processing, the `unreference` routine for arrays invokes the destructor on each element of the array in turn, rather than on the entire array. This allows the "boilerplate" loop over elements to exist in one place, rather than being replicated in every destructor.

Many languages provide immutable arrays, along with operators that allow a modified copy of an existing array to be created. In QIR, this is implemented by creating a new copy of the existing arrays on the heap, and then modifying the newly-created arrays in place. In some cases, if the compiler knows that the existing arrays is not used after the creation of the modified copy, it is possible to avoid the copy and modify the existing arrays as long as there are no other references to the arrays. The `arrays_is_writable` utility function is used to test if the array handle is the only reference.

There are two special operations on arrays:

- An array *slice* is specified by providing a dimension to slice on and a `%Range` to slice with. This is equivalent to the current Q# array slicing feature, generalized to multi-dimensional arrays.
- An array *projection* is specified by providing a dimension to project on and an `i64` index value to project to. The resulting array has one fewer dimension than the original array, and is the segment of the original array with the given dimension fixed to the given index value.

Both slicing and projecting are implemented by creating a new `%ArrayPointer%` with updated dimension count, offset, stride, and length information.

The following utility functions are provided by the classical runtime to support arrays:

Function	Signature	Description
<code>quantum.rt.array_create</code>	<code>void(%ArrayPointer, i32, %Destructor, i32, [0 x i64])</code>	Creates a heap-allocated array by setting the reference count to 1 and filling in the element size and the destructor. The second <code>i32</code> is the dimension count. The list of <code>i64</code> s contains the length of each dimension. This routine allocates and initializes an <code>%ArrayData</code> structure and fills in the supplied <code>%ArrayPointer</code> structure.
<code>quantum.rt.array_copy</code>	<code>void(%ArrayPointer, %ArrayPointer)</code>	Initializes the first <code>%ArrayPointer</code> to be the same as the second, and increments the reference count of the underlying <code>%ArrayData</code> .
<code>quantum.rt.array_get_dim</code>	<code>i32(%ArrayPointer)</code>	Returns the number of dimensions in the array.
<code>quantum.rt.array_get_length</code>	<code>%Int(%ArrayPointer, i32)</code>	Returns the length of a dimension of the array. The <code>i32</code> is the zero-based dimension to return the length of.

Function	Signature	Description
quantum.rt.array_validate	<code>i1(%ArrayPointer, [0 x i64])</code>	Returns true if the list of indices is valid for the array and false otherwise.
quantum.rt.array_get_element	<code>i8*(%ArrayPointer, [0 x i64])</code>	Returns a pointer to the indicated element of the array.
quantum.rt.array_slice	<code>void(%ArrayPointer, \$ArrayPointer, i32, %Range)</code>	Initializes the first <code>%ArrayPointer</code> to be a slice of the second. The <code>i32</code> indicates which dimension the slice is on, and the <code>%Range</code> specifies the slice. The reference count of the underlying <code>%ArrayData</code> is incremented.
quantum.rt.array_project	<code>void(%ArrayPointer, \$ArrayPointer, i32, i64)</code>	Initializes the first <code>%ArrayPointer</code> to be a projection of the second. The <code>i32</code> indicates which dimension the projection is on, and the <code>i64</code> specifies the specific index value to project. The new <code>ArrayPointer</code> will have one fewer dimension than the original. The reference count of the underlying <code>%ArrayData</code> is incremented.
quantum.rt.array_reference	<code>void(%ArrayPointer)</code>	Increments the reference count of the array. This is almost never needed because new <code>%ArrayPointers</code> are typically created with one of the above four routines, all of which update the reference count properly.
quantum.rt.array_unreference	<code>void(%ArrayPointer)</code>	Decrements the array's reference count; calls its destructor and releases the array if the reference count is now zero.
quantum.rt.array_is_writable	<code>i1(%ArrayPointer)</code>	Returns true if it is safe to modify the contents of the tuple; that is, it returns true if the reference count of the tuple is equal to 1.

Callables

We use the term *callable* to mean a subroutine in the source language. Different source languages use different names for this concept.

Code

Callables may have up to four different implementations to handle different combinations of functors. Each implementation is represented as an LLVM function that takes three tuple pointers as input and returns no output; that is, as an LLVM `void(%TuplePointer, %TuplePointer, %TuplePointer)`. The first input is the capture tuple, which will be null for top-level callables. The second input is the argument tuple. The third input points to the result tuple, which will be allocated by the caller.

We define the LLVM type `%CallableImpl = type void(%TuplePointer, %TuplePointer, %TuplePointer)` for convenience.

We use a caller-allocates strategy for the result tuple because this allows us to avoid a heap allocation in many cases. If the callee allocates space for the result tuple, that space has to be on the heap because a stack-based allocation would be released when the callee returns. The caller can usually allocate the result tuple on the stack, or reuse the result tuple pointer it received for tail calls.

The names of the implementation should be the namespace-qualified name of the callable with periods, followed by a hyphen ('-'), followed by `body` for the default implementation, `adj` for the adjoint implementation, `ctl` for the controlled implementation, and `ctladj` for the controlled adjoint implementation.

For instance, for a callable named `Some.Namespace.Symbol` with all four implementations, the compiler should generate the following in LLVM:

```
define void Some.Namespace.Symbol-body (%TuplePointer capture, %TuplePointer args,
%TuplePointer result)
{
    ; code goes here
}

define void Some.Namespace.Symbol-adj (%TuplePointer capture, %TuplePointer args,
%TuplePointer result)
{
    ; code goes here
}

define void Some.Namespace.Symbol-ctl (%TuplePointer capture, %TuplePointer args,
%TuplePointer result)
{
    ; code goes here
}

define void Some.Namespace.Symbol-ctladj (%TuplePointer capture, %TuplePointer
args, %TuplePointer result)
{
    ; code goes here
}
```

Each implementation should start with a prologue that decomposes the argument and capture tuples. Depending on the result type, it should end with an epilogue that fills in the result tuple. If the result tuple contains another container such as an array, it may be simpler and more efficient to fill in the sub-container in-line in the implementation, rather than creating a separate array and then copying it to the result tuple.

Implementation Table

For each defined callable, a four-entry table is created with pointers to these four labels; implementations that don't exist for a specific callable have a null pointer in that place. The table is defined as a global constant with

the namespace-qualified name, with periods, of the callable.

For the example above, the following would be generated:

```
@Some.Namespace.Symbol = constant [4 x %CallableImpl*]
[
  %CallableImpl* @Some.Namespace.Symbol-body,
  %CallableImpl* @Some.Namespace.Symbol-adj,
  %CallableImpl* @Some.Namespace.Symbol-ctl,
  %CallableImpl* @Some.Namespace.Symbol-ctladj
]
```

We define the LLVM type `%CallableImplTable = type [4 x %CallableImpl*]` for convenience.

External Callables

Callables may be specified as external; that is, they are declared in the quantum source, but are defined in an external component that is statically or dynamically linked with the compiled quantum code. For such callables, the compiler should generate a normal callable in LLVM, and generate a function body that unwraps the argument tuple, calls the external function, and then wraps the return value into a result tuple.

Generating the proper linkage is the responsibility of the target-specific compilation phase.

Generics

QIR does not provide support for generic or type-parameterized callables. It relies on the language-specific compiler to generate a new, uniquely-named callable for each combination of concrete type parameters. The LLVM representation treats these generated callables as the actual callables to generate code for; the original callables with open type parameters are not represented in LLVM.

Callable Values

Callable values are represented by a pointer to a structure with the following format:

Field	Data Type	Description
Specialization index	i8	Identifies the specialization to use when invoking this callable.
Implementation table	%CallableImplTable*	Points to the implementation table described above.
Controlled count	i8	Counts how many times the Controlled functor has been applied to this callable.
Capture tuple	%TuplePointer	A pointer to the tuple of captured values. This is null if the callable has no captured values.

The specialization index and controlled count are used and maintained by the callable utility functions to implement functors.

For convenience, the LLVM Callable type is defined:

```
%Callable = type { i8, %CallableImplTable*, i8, %TuplePointer }
```

To create a callable value, allocate space for the value on the stack or heap and use the `callable_init` utility function to initialize it from an implementation table and capture tuple or the `callable_copy` utility function to initialize it as a copy of another callable value.

Invoking a Callable

An explicit call to a top-level callable may be translated into a direct LLVM call instruction to the callable's body implementation:

```
call void @Symbol-body (%TuplePointer null, %TuplePointer %arg-tuple,
%TuplePointer %res-tuple)
```

The capture tuple pointer is null because top-level named callables never have a capture tuple.

Similarly, an explicit call to the adjoint of a named callable may be translated as:

```
call void @Symbol-adj (%TuplePointer null, %TuplePointer %arg-tuple, %TuplePointer
%res-tuple)
```

Calls to the controlled or controlled adjoint specializations of a named callable follow the same pattern; in these cases, the argument tuple should point to a two-tuple whose first element is the array of control qubits and whose second element is a tuple of the remaining arguments.

To invoke a callable value represented as a `%Callable*`, the `callable_invoke` utility function should be used. This function uses the information in the callable value to invoke the proper implementation with the appropriate parameters.

Implementing Functors

The Adjoint and Controlled functors are important for expressing quantum algorithms. They are implemented by the `callable_adjoint` and `callable_control` utility functions, which update a `%Callable` in place by applying the Adjoint or Controlled functors, respectively.

To support cases where the original, unmodified `%Callable` is still needed after functor application, the `callable_copy` routine may be used to create a new copy of the original `%Callable`; the functor may then be applied to the new `%Callable`. For instance, to implement the following:

```
let f = someOp;
let g = Adjoint f;
// ... code that uses both f and g ...
```

The following snippet of LLVM code could be generated:

```
%f = alloca %Callable
call %quantum.rt.callable_init(%Callable* %f, %CallableImplTable* @someOp,
%TuplePointer null)
%g = alloca %Callable
call %quantum.rt.callable_copy(%g, %f)
call %quantum.rt.callable_adjoint(%g)
```

Implementing Lambdas

The language-specific compiler should generate a new top-level callable of the appropriate type with implementation provided by the anonymous body of the lambda; this is known as "lifting" the lambda. A unique name should be generated for this callable. Lifted callables can support functors, just like any other callable. The language-specific compiler is responsible for determining the set of functors that should be supported by the lifted callable and generating code for them accordingly.

At the point where a lambda is created as a value in the code, a new callable data structure should be created with the appropriate contents. Any values referenced inside the lambda that are defined in a scope external to the lambda should be added to the lambda's capture tuple.

For instance, in the following Q# code, assuming a possible syntax for lambda expressions:

```
let x = 1;
let y = 2;
let f = lambda (z) { return x + y + z; };
```

The capture tuple for the callable value stored in `f` will consist of a tuple of two integers, `(1, 2)`. The code generated for the `let f =` line would do the following:

1. Build a tuple on the stack or heap from the current values of `x` and `y`. The tuple must be allocated on the heap if there is a possibility of the resulting callable being returned from this code.
2. Allocate a `%Callable` structure on the stack or heap. This structure should be allocated in the same place as the capture tuple.
3. Call the `callable_init` library routine with the callable pointer, the implementation table for the compiler-generated callable, and the capture pointer.
4. Assign the address of the `%Callable` structure to the local variable `%f`.

`callable_init` initializes the new callable structure with the provided parameters, a specialization index of 0, and a control count of 0.

Implementing Partial Application

A partial application is a form of closure value; that is, it is a lambda value, although the source syntax is different.

Partial applications should be rewritten into lambdas by the language-specific compiler. The lambda body may include additional code that performs argument tuple construction before calling the underlying callable. For instance, the following Q# partial application:

```
f(x, _)
```

could be rewritten as:

```
lambda (y) {
  let z = (x, y);
  return f(z);
}
```

Utility Functions

The following utility functions are provided by the classical runtime to support callables:

Function	Signature	Description
quantum.rt.callable_init	void(%Callable*, %CallableImplTable*, %TuplePointer)	Initializes the callable with the provided function table and capture tuple.
quantum.rt.callable_copy	void(%Callable*, %Callable*)	Initializes the first callable to be the same as the second callable.
quantum.rt.callable_invoke	void(%Callable*, %TuplePointer, %TuplePointer)	Invokes the callable with the provided argument tuple and fills in the result tuple.
quantum.rt.callable_adjoint	void(%Callable*)	Updates the callable bu applying the Adjoint functor.
quantum.rt.callable_control	void(%Callable*)	Updates the callable by applying the Controlled functor.