

# COMPILING A HIGHER-ORDER SMART CONTRACT LANGUAGE TO LLVM

Vaivaswatha Nagaraj, Jacob Johannsen, Anton Trunov, George Pîrlea, Amrit Kumar and Ilya Sergey

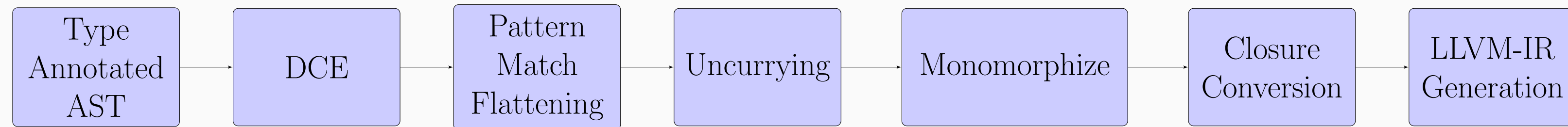
Zilliqa Research

## Scilla - Smart Contract Intermediate Level Language

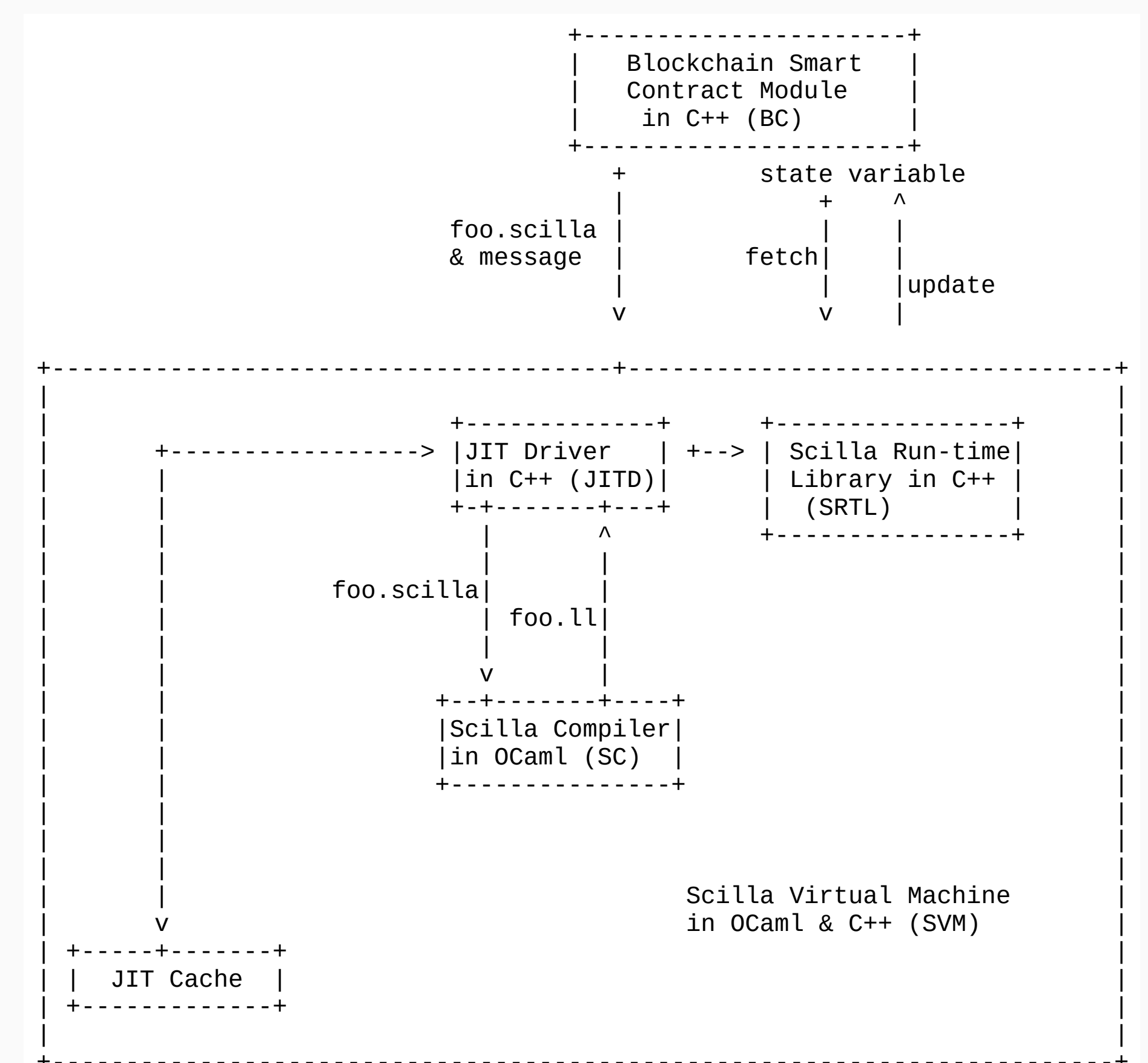
- DSL for blockchains: Used to write smart-contracts.
- Typed, higher-order, polymorphic language in the ML family.
- Currently interpreted in production. This work explores its compilation.

## The Compiler and Runtime

- The compiler, written in OCaml (using LLVM's OCaml bindings), shares the parser and type-checker with the existing interpreted implementation.
- The runtime implements language built-in operations and provides an interface b/w the JIT'ed code and the blockchain.



## System Architecture



## Experimental Results

In comparison to the reference interpreter:

- Synthetic tests like Ackermann function computation and Church numerals show improvements from 47x to more than 100x.
- Production benchmarks show improvements ranging from 5x to 10x.

## Part I: Front-End Transformations

### Pattern Match Flattening

- Nested patterns are flattened into nested pattern matches.
- Enables direct translation to efficient LLVM switch statements.

Before

```
1 fun (p : List (Option Int32)) =>
2   match p with
3   | Nil => z
4   | Cons (Some x) xs => x
5   | Cons _ _ => z
6   end
```

After

```
1 fun (p : List (Option Int32)) =>
2   match p with
3   | Nil => z
4   | Cons a xs =>
5     match a with
6     | Some x => x
7     | _ => z
8   end
9   end
```

### Uncurrying

- Scilla uses curried semantics. Transform the AST to be uncurried.
- Transforms nested sequence of functions to a single function with multiple arguments.

Before

```
1 fun (x : Int32) =>
2   fun (y : Int32) =>
3     builtin add x y
```

After

```
1 fun (x : Int32, y : Int32) =>
2   builtin add x y
```

The syntax used in these examples is not always legal Scilla syntax. We don't have a textual syntax for every intermediate form.

### Monomorphization

- Parametric polymorphism is implemented via monomorphization (as against type erasure).
- A higher-order type-flow analysis (conservatively) computes the possible ground types that may flow into a type variable.
- Polymorphic expressions are specialized with all types computed by the analysis.

Before

```
1 let option_to_bool =
2   tfun 'A =>
3     fun (a : Option 'A) =>
4       match a with
5       | Some _ => True
6       | None => False
7   in
8   let otb_i32 = @option_to_bool Int32 in
9   let otb_string = @option_to_bool String in ...
```

After

```
1 let option_to_bool =
2   [ Int32 ->
3     fun (a : Option Int32) =>
4       match a with
5       | Some _ => True
6       | None => False
7   ; String ->
8     fun (a : Option String) =>
9       match a with
10      | Some _ => True
11      | None => False ]
12 in
13 let otb_i32 = option_to_bool[Int32] in
14 let otb_string = option_to_bool[String] in ...
```

### Closure Conversion

- In Scilla, functions can be nested: inner functions have access to the **free** variables that are defined outside their bodies.
- Closure conversion eliminates free variables by collecting them in **environments** taken as a function's additional argument.
- The modified function with no free variable accesses is lifted to global level (thus translatable to an LLVM function).
- Every reference to the function is now a pair, of the function's pointer and its environment.

Before

```
1 let f1 =
2   let y = Int32 1 in
3   fun (x : Int32) =>
4     builtin add x y
5   in ...
```

After

```
1 let f2 =
2   fun (e : {y : Int32}, x : Int32) =>
3     let y = e.y in
4     builtin add x y
5   in
6   let f1 =
7     let y = Int32 1 in
8     let e = new {y : Int32} in
9     e.y ← y;
10    {f2, e}
11   in ...
```

## Part II: Mapping Scilla to LLVM-IR

### Type Descriptors

- Type Descriptors are data structures defined in the runtime library, objects of which are created in the JIT'ed code.
- A global array of type descriptors, with an entry for every type that occurs in the Scilla program, with each type identified by its index in the array, is inserted by the compiler.
- Some functions in the runtime library (such as the JSON (de)serializer) that can operate on different Scilla types, take in a type descriptor argument.

### Primitive Types

- Integer types are wrapped with an LLVM struct type to differentiate signedness. For example the Scilla type `Int32` becomes the LLVM type `%Int32 = type { i32 }`.
- `String` and `ByStr` (bytes of arbitrary length) types are translated to LLVM types `%String = { i8*, i32 }` and `%Bystr = { i8*, i32 }`.
- Fixed sized byte sequences `ByStrX` (i.e.,  $X$  is known statically) are expressed as LLVM array types `[ X × i8 ]`.

### Maps

- Map values are boxed, handled through an opaque pointer.
- The runtime library creates and operates on Map values.
- Internally, the runtime library uses `unordered_map<string, any>`. `any` is used to enable representing values of any Scilla type, including nested maps and user-defined types.

### Algebraic Data Types

- Algebraic data types, aka variants or tagged unions, are composite types that have one or more "constructors" (sum type), each defining a tuple (product type).
- A canonical example of `List` defined in Scilla, with two constructors `Cons` and `Nil`:
  - `%tname = type { i8, ... }`
  - `%nil = type <{ i8 }`
  - `%cons = type <{ i8, Int32, %tname* }`

```
1 type MyList =
2 | Nil
3 | Cons of Int32 MyList
```

- ADT objects are boxed, i.e. represented using a pointer.

The `i8` field (called the tag) in `%tname` determines the constructor used. No object of type `%tname` is built. Only pointers to it are type-casted to the right constructor before use. Packed LLVM structs are used to simplify access in the runtime library.

### Closures

- A closure is represented by an anonymous struct `type { fundef_sig*, void* }`.
- `fundef_sig` is the signature of the LLVM function definition. The `void*` represents the environment pointer.
- All Scilla functions are represented as closures, with their first argument being the environment pointer.
- If the function's return type cannot be "by value", then the second argument will be a stack pointer ("stret") where the return value must be stored.
- To avoid ABI complexities, generated LLVM functions and the hand written functions in the runtime library that they may call, all follow the simple rule that if the value size is larger than two eightbytes, we define that parameter (or return value) to be passed by reference (i.e. via a stack pointer).