

# Dead Code Detection for Scala Programs

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Running Example</b>	<b>2</b>
<b>3</b>	<b>AST Approach</b>	<b>3</b>
3.1	Represent Data dependences with AST . . . . .	3
3.2	Obtain Data Dependences at Runtime from AST . . . . .	4
3.3	Analyzing Data Dependences Extracted from AST . . . . .	6
<b>4</b>	<b>ICode Approach</b>	<b>6</b>
4.1	Represent Data Dependences with ICodes . . . . .	6
4.2	Obtain Data Dependences at Runtime from ICodes . . . . .	7
4.3	Analyze Data Dependences Extracted from ICodes . . . . .	8
4.3.1	Build Data Dependence Graph . . . . .	8
4.3.2	Search for Dead Codes . . . . .	9
4.3.3	Get the Result . . . . .	9
<b>5</b>	<b>Experiments and Results</b>	<b>9</b>

# 1 Introduction

Scala is a programming language uniting both functional and objective-oriented styles, which gains popularity due gradually these days. To guarantee the efficiency of Scala programs, it's better to use an automatic code reviewer or coverage tool. We develop a tool set to do runtime data dependence analysis and detect dead codes in Scala programs.

Compared with existing code coverage tool, our tool can dig into data dependences among variables and values. Scoverage only analyzes usage of source code at runtime, like how many lines of codes are executed during one execution. Our tool does even more. Given a target, which can be either a variable or value in the program, it can find out all codes that contribute to the final value of this target, thus detects useless data flow at run time.

The basic idea of our tool is to change the program with a special compiler or compiler plugin, and make it provide data dependence messages for every operation in the program. Then a small tool will be used to analyze these messages at real time when the program is running.

First we need to consider two essential problems about data dependence analysis

1. How to represent data dependences?
2. How to obtain data dependences at runtime?
3. How to analyze the dependences at runtime?

Each of these questions will be answered respectively for AST and ICode approaches in following sections.

# 2 Running Example

Through out this report we'll use a simple running example to show how a data dependences can be obtained from the AST or ICode of a program and how to analyze these dependences. The example is a Scala class defines an `add` function, which simply increase the parameter by 1 and return. Also it has another function `callAdd` inside which the `add` function is called and the result is stored into `addResult` value.

```
1  class Add {  
2      def add(x: Int): Int = {  
3          x + 1  
4      }  
5  
6      def callAdd() {  
7          val addResult = add(8)  
8      }  
9  }
```

## 3 AST Approach

### 3.1 Represent Data dependences with AST

To represent data dependences, we can make use of information provided by Scala compiler. Since we use the compiler to change the program’s behavior, it is very convenient to use informations from AST or ICode. ICode is a intermediate representation in Scala compiler between AST and byte code, which is translated into JVM byte code in the end. This section we focus on representing data dependences using information from ASTs. In section 4.1 we consider how to do it with ICode.

Abstract Syntax Tree (AST) is a tree data structure used by compilers, which representing the parsed structure of a program. First let’s see how the function definition `def add(...){...}` is parsed into AST

//Graph for the AST

Now we want to represent the dependences related to value `addResult`. First, notice see that `addResult` is an identifier, corresponding to `Ident` AST in Scala. And the function call `add(8)` is an `Apply` AST in Scala. Using `>>>` to represent dependence, then

`Apply >>> Ident`

which means the result of function call `add(8)` determines the value of `addResult`. However, the `Apply` and `Ident` should be given more specific names for analysis. For example, use `addResult` as the name of this `Ident` tree and `add` as the name of this `Apply` tree. Thus the information we provide to the analysis unit should be

`add >>> addResult`

Similarly, other tree structures should have their specific names. For tree structures without a inherited suitable name, we will create a synthesis name, for example, a `Block` tree may use the name `Block124`, where 124 is the ID of this `Block` provided by the compiler. Here is the table about how names are given to various ASTs.

AST	Name
<code>Apply(args, fun)</code>	Name of fun
<code>DefDef(mods, name, tparams, vparamss, tpt, rhs)</code>	name
<code>Block(stats, expr)</code>	“Block” + AST ID
<code>Return(expr)</code>	“Return” + AST ID
<code>ValDef(mods, name, tpt, rhs)</code>	name
<code>If(cond, thenp, elsep)</code>	“If” + AST ID
<code>Select(qualifier, selector)</code>	selector
<code>Match(selector, cases)</code>	“Match” + AST ID
<code>New(tpt)</code>	“New” + AST ID
<code>Try(block, catches, finalizer)</code>	“Try” + AST ID

In addition, we need another name for each function definition to represent its parameters. We use the function name appended by “\$” to represent the pa-

rameters of this function. The usage of this special name will be clear in the following example. Here's the full example showing how the value of `addResult` is determined. The “>>>” means “determines”.

```
8 >>> add$
add$ >>> x
x 1 >>> add
add >>> addResult
```

The first two lines `8 >>> add$` and `add$ >>> x` say 8 is passed as the argument of function `add`. Here `add$` acts as the bridge that connects 8 and `x`. Note that we don't use `8 >>> x` directly because number 8 is known when the function is called, while name of parameter `x` is only accessible inside the function definition.

## 3.2 Obtain Data Dependences at Runtime from AST

We've shown how to use information from AST to represent data dependences. Next problem is how can these information be passed to our analyzing module.

The basic idea is, we can change the Scala program, using compiler plugins, to insert some “log” operations inside the program as parts of the program itself. A “log” operation is simply a function call to the analyzing module, taking the data dependence information as a string argument. So that when the program is executed, data dependence information is passed to the analyzing module simultaneously.

Here is how we want to insert the “log” operations into the program.

```
1  class Add {
2      def add(x: Int): Int = {
3          ScalaTrace.logger.log("add$ >>> x")
4          ScalaTrace.logger.log("x 1 >>> add")
5          x + 1
6      }
7
8      def callAdd() {
9          ScalaTrace.logger.log("8 >>> add$")
10         val addResult = add(8)
11         ScalaTrace.logger.log("add >>> addResult")
12     }
13 }
```

Once the original program is modified into program above, data dependences `add$ >>> x`, `x 1 >>> add`, `8 >>> add$`, `add >>> addResult` will be passed to analyzing module sequentially when it is executed. It is convenient to insert the “log”'s after the program is parsed into AST.

To explain how to insert the “log”'s into ASTs, first we should know basically how Scala compiler works. Scala compiler has over 20 compilation phases. Almost all phases take an AST as input, modify the AST and pass it to the next phase.

Our goal is to modify the AST by inserting our “log” operations. Scala compiler plugin is a convenient tool to do this. A compiler plugin is a program written by user that acts as an additional compilation phase. The plugin can be inserted between almost any two standard compilation phases. Just like

standard phases, a plugin also takes an AST as input, modify it and produce a new AST, which is passed to next phase.

To insert “log” operations, we should create ASTs for the function calls of `ScalaTrace.logger.log`, then insert them in the appropriate place of the original AST. Following shows a code snippet in our compiler plugin that creates such ASTs.

```

15 def genLog(dataDependence: String, positionInSourceFile: Position): Tree = {
16   def extractEssentialPath(path: String): String = {
17     if(path.contains("NoPosition"))
18       "NoPosition"
19     else
20       path.substring(path.indexOf("/src/") + 4)
21   }
22
23   val liter0 = Literal(Constant(extractEssentialPath(positionInSourceFile.focus.toString)))
24   liter0.setType(typeOf[String])
25
26   val liter1 = Literal(Constant(dataDependence))
27   liter1.setType(typeOf[String])
28
29   val scalatrace = rootMirror.staticModule("ScalaTrace")
30   val logger = definitions.getMember(scalatrace, newTermName("logger"))
31   val callLog = global.gen.mkMethodCall(logger, newTermName("log"), List(liter0, liter1))
32   val Apply(fun, _) = callLog
33   fun.setSymbol(logger.info.decl(newTermName("log")))
34   fun.setType(logger.info.decl(newTermName("log")).tpe)
35   val Select(qual, _) = fun
36   qual.setType(logger.tpe)
37   callLog.setType(typeOf[Unit])
38   callLog
39 }

```

Line 23 and 26 creates ASTs for two string arguments of `log` method. Line 31 creates the function call of `log`, which is an `Apply` AST. The compiler plugin should runs after standard phase `typer`. Because after `typer`, the names and types of ASTs are calculated, so that these can be used in the data dependence information.

However note that AST for `log` is created by us in the plugin, which runs after `typer`. So we have to set the types and symbols for it and its subtrees by ourselves so that it can be a legal tree for the following phases. That’s why there are a lot of `setType`, `setSymbol` in the snippet above.

Then our plugin traverses the original AST, calls `genLog` to create the “log” ASTs and insert them at appropriate places. Following shows the ASTs before and after our plugin phase.

```

class Add extends scala.AnyRef {
  ...
  def add(x: Int): Int = x.+(1);
  def callAdd(): Unit = {
    val addResult: Int = Add.this.add(8);
    ()
  }
}

```

```

class Add extends scala.AnyRef {
  ...
  def add(x: Int): Int = {
    val newvalue = x.+({
      ScalaTrace.this.logger.log("Add.add$ >>> Add.x", "/print/hello.scala,line-10,offset=160");
      1
    });
    ScalaTrace.this.logger.log("Add.x 1 >>> Add.add", "/print/hello.scala,line-11,offset=191");
  }
  def callAdd(): Unit = {
    ...
    val addResult: Int = Add.this.add({
      ScalaTrace.this.logger.log("8 >>> Add.add$", "/print/hello.scala,line-15,offset=253");
      8
    });
    ScalaTrace.this.logger.log("Add.add >>> Add.addResult", "/print/hello.scala,line-15,offset=238");
    ...
  }
}

```

### 3.3 Analyzing Data Dependences Extracted from AST

## 4 ICode Approach

### 4.1 Represent Data Dependences with ICodes

ICode imitates push/pop operations of a stack machine. Each ICode instruction can first take some values as input from the stack, and then put its result onto the stack. Here's the table of main scala ICodes and their corresponding push/pop operations.

ICode	Operations	Stack Operations
THIS	loads "this" on top of the stack	push 1 <sup>1</sup>
CONSTANT	loads a constant on the stack	push 1
LOAD_ARRAY_ITEM	loads an element of an array	pop 2, push 1
LOAD_LOCAL	load a local variable on the stack	push 1
LOAD_FIELD	load a field on the stack	pop 1, push 1
LOAD_MODULE	load a module on the stack	pop 1
STORE_ARRAY_ITEM	store a value into an array at a specified index	pop 3
STORE_LOCAL	store a value into a local variable	pop 1
STORE_FIELD	store a value into a field	pop 2
STORE_THIS	store a value into the 'this' pointer	pop 1
CALL_PRIMITIVE	call a primitive function	pop n <sup>2</sup> , push 1
CALL_METHOD	call a method	pop n, push 1
NEW	create a new instance of a class	pop n, push 1
CREATE_ARRAY	create an array	pop n, push 1
SWITCH	a switch instruction	pop 1
CJUMP	jump according to the result of comparing two values	pop 2
CZJUMP	jump according to the result of comparing with zero	pop 1
THROW	throws an exception	pop 1
DROP	drop one value from the stack	pop 1
DUP	duplicate the top of the stack	pop 1, push 2
LOAD_EXCEPTION	load an exception	pop all, push 1

With these operations on the stack, we can know exactly each value on the stack is produced and consumed by which instruction. Thus data dependences can be represented using a sequence of push/pop operations. Again we demonstrate with the running example to see how these push/pop operations are extracted from the ICode.

First let's see how the simple `add` function is translated into ICode.

```

5  def add(x: Int (INT)): Int {
6    locals: value x
7    startBlock: 1
8    blocks: [1]
9
10   1:
11     3 LOAD_LOCAL(value x)
12     3 CONSTANT(1)
13     3 CALL_PRIMITIVE(Arithmetic(ADD,INT))
14     3 RETURN(INT)
15
16   }
```

Our simple `add` function is translated into 4 ICodes. First, it loads the parameter `x` onto stack, then a constant 1. After that, the `ADD` primitive operation is called. Finally, the function returns. The number 3 before these ICodes is the line number of corresponding source code.

Here's the ICode for the assignment `val addResult = add(8)`

```

25   7 THIS(Add)
26   7 CONSTANT(8)
27   7 CALL_METHOD Add.add (dynamic)
28   7 STORE_LOCAL(value addResult)
```

Again it can be represented by 4 lines of ICode. First it loads “this” onto the stack, this is needed to call `add` method. Note that `add` is a method which consumes an object of `Add` class. Since here it is called inside the `Add` class, so “this” object will be used. Then it loads the constant 8, as the argument for `add` method. The third ICode calls `add` method. After this ICode, the ICodes of `add` method will be executed. When the `add` call finishes, store its result into `addResult`.

Putting it all together, here's the sequence of ICodes being called when the assignment `val a = add(8)` is executed, with push/pop operations for each ICode shown on the right.

In our push/pop operations above, `push <something>` means putting `<something>` on the stack. `pop <something>` means storing the top of stack into `<something>` then pop the stack. Sometimes we need `>>>` to represent data dependences directly. For example `a b >>> c` means the value of `c` is determined by `a` and `b`.

## 4.2 Obtain Data Dependences at Runtime from ICodes

Now we know how push/pop operations that exactly describes data dependences can be extracted from ICodes. In this section let's focus on how to insert “log” function calls into ICodes.

<sup>1</sup>push 1 means pushing one value/object onto stack, similarly hereinafter.

<sup>2</sup>n means the number of arguments, similarly hereinafter.

Step	ICode	Operation
1	THIS(ADD)	push @this
2	CONSTANT(8)	push 8
3	CALL_METHOD Add.add (dynamic)	pop x, pop @this
4	LOAD_LOCAL(value x)	push x
5	CONSTANT(1)	push 1
6	CALL_PRIMITIVE(Arithmetic(ADD, INT))	pop @operand1 pop @operand2 @operand1 @operand2 >>> @result push @result
7	RETURN	
8	CALL_METHOD Add.add (dynamic)	push @result
9	STORE_LOCAL(value addResult)	pop addResult

We modify the icode phase of scala compiler 2.11 to insert additional ICodes for each original ICode. These additional ICodes calls the “log” function, passing push/pop operations to the analyzing unit.

Following shows the ICodes to call “log” function.

```
LOAD_FIELD ScalaTrace.logger
CONSTANT("#1 Add.x @")
CONSTANT("Hello.scala,line-2")
CALL_METHOD Logger.log (dynamic)
```

Our “log” function is a method of Logger class. The program and our analyzing tool will interfere through ScalaTrace module, which is a final java class. logger is a static field of ScalaTrace.

To call ScalaTrace.logger.log, ScalaTrace.logger should be loaded first. This is done by the first LOAD\_FIELD ICode. Then two String arguments, line number and push/pop operation, needed by log should be put on stack. We can do this using two CONSTANT ICodes. Finally, CALL\_METHOD ICode calls Logger.log, using ScalaTrace.logger as the object to call this method, with the two String constants as arguments.

For each ICode generated from the original program, these additional ICodes will be added before them, to passing the position and push/pop operations of original ICodes to our analyzing module. Following figure shows ICodes of our add function, with additional ICodes.

### 4.3 Analyze Data Dependences Extracted from ICodes

From section 4.2 we know how push/pop operations and code positions are passed to our analyzing module. This section focus on how to analyze these information are used to analyze data dependences, and find out dead codes.

#### 4.3.1 Build Data Dependence Graph

Here’s the basic idea of our analyzing algorithm. Our analyzing modules maintains a stack to imitate every push/pop operation of the program’s original ICodes. In this way, whenever an ICode pop something from the stack as an



operator, our analyzing module knows exactly which object it gets. For example, at some time, an object `ObjA` is push on the stack in source code position `file1,line-21`. Later on, `ObjA` is popped from the stack in source code position `file2,line-34`. This means `file2,line-34` has a data dependence on `file1,line-21`.

In this way, the data dependences between all lines of source codes in the original programs can be detected. Thus we can build a dependence graph in memory to maintain such relationships. Vertices in the graph corresponds to lines in the source code. Vertices are connected by directed edges. If there's an edge from Vertex A to Vertex B, that means source code line of Vertex B depends on that of Vertex A.

TODO Sample Dependence Graph

### 4.3.2 Search for Dead Codes

To search for dead codes within the dependence graph, a target must first be assigned. A target is an variable name in the program that the user considers as the final result of the program. For example, suppose the program is our Scala compiler, then the target can be `jclassBytes`, the variable storing the final java byte code, which is the final product of Scala compiler.

The target should be assigned before the program is executed. Then after the program lanches, our analyzing module will build the dependence graph gradually. Then, when the program is writing the target variable, a push operation `push <target>` will be passed to the analyzing module, along with its source code line. Then vertex corresponding to this source code line will be marked as useful, as well as all descendants of this vertex.

### 4.3.3 Get the Result

The last question is how to obtain the useless source code lines. Since we only mark source code lines as useful or useless, a straight forward idea is we can simply dump all the source code lines and all the useful source code lines into two separate files when the compilation finishes, then take the difference of these two files to get the useless lines.

However, to do this the program must be able to inform the analyzing module before it terminates so that the analyzing module can dump the result to files. As the program terminates, all the result in our analyzing module will disappear. This requires special treatment to inserting function calls other than “log” to code segments that the program will execute before terminating. It is extremely hard to identify such code segments at compile time.

Thus whenever a new source code line is executed or being marked as useful, the analyzing module has to write it to the file. Because all the actions of our analyzing module is activated function calls from the program, and it doesn't know when the program will terminate, it has to write the result incrementally.

## 5 Experiments and Results