

# System Software Crash Course

Samsung Research Russia  
Moscow 2019

Block C Compiler Construction  
6-7. YACC & Bootstrapping  
Eugene Zouev

# Automatic parser generation

- Top-down or bottom-up parsing?
- «Hand-made» or automated development?

|                   | By hand  | By tools                           |           |
|-------------------|--|------------------------------------|-----------|
| Top-down parsing  | Most often:<br><b>Recursive descent parser</b> | ANTLR<br>COCO<br>etc.              | ← LL(1)   |
| Bottom-up parsing | Rarely (too complicated)                       | Yacc, Bison,<br>Jay, GPPG,<br>etc. | ← LALR(1) |

# Yacc/Bison & clones

- **YACC** – Yet another compiler compiler 1970: based on C.
- **Bison** – Yacc version for для GNU: based on C.
- **GPPG** – Gardens Point Parser Generator: Yacc version for C# and .NET.
- **Jay** – Yacc version for Java.
- ...A lot of YACC clones for almost all popular languages.

All YACCs have **identical** parsing algorithm.

# Yacc/Bison: references (Russian)

**YACC - Yet Another Compiler Compiler**

[http://yacc.solotony.com/yacc\\_rus/index.html](http://yacc.solotony.com/yacc_rus/index.html)

Перевод оригинальной статьи (так себе, но понятно)

**Компилятор компиляторов Bison - первое знакомство**

[http://trpl.narod.ru/CC\\_Bison.htm](http://trpl.narod.ru/CC_Bison.htm)

**Bison - Генератор синтаксических анализаторов, совместимый с YACC**

[http://www.opennet.ru/docs/RUS/bison\\_yacc/bison\\_1.html](http://www.opennet.ru/docs/RUS/bison_yacc/bison_1.html)

Перевод официального руководства GNU

**Lex и YACC в примерах**

<http://rus-linux.net/lib.php?name=/MyLDP/algol/lex-yacc-howto.html>

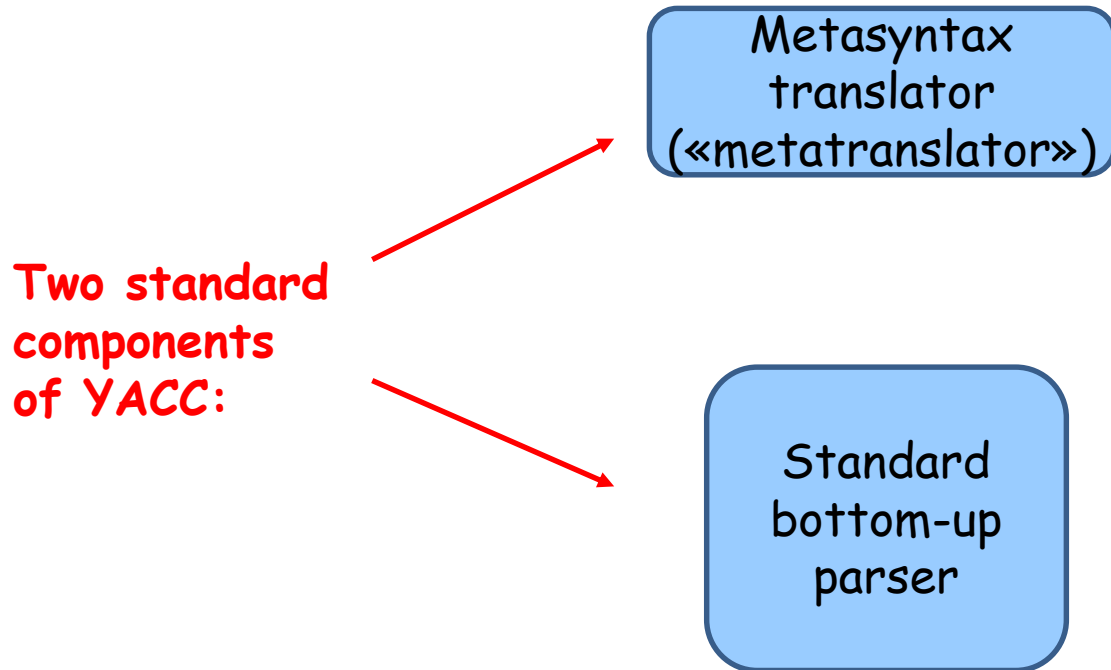
**Gardens Point Parser Generator**

YACC-совместимый генератор для C#; <http://gppg.codeplex.com/>

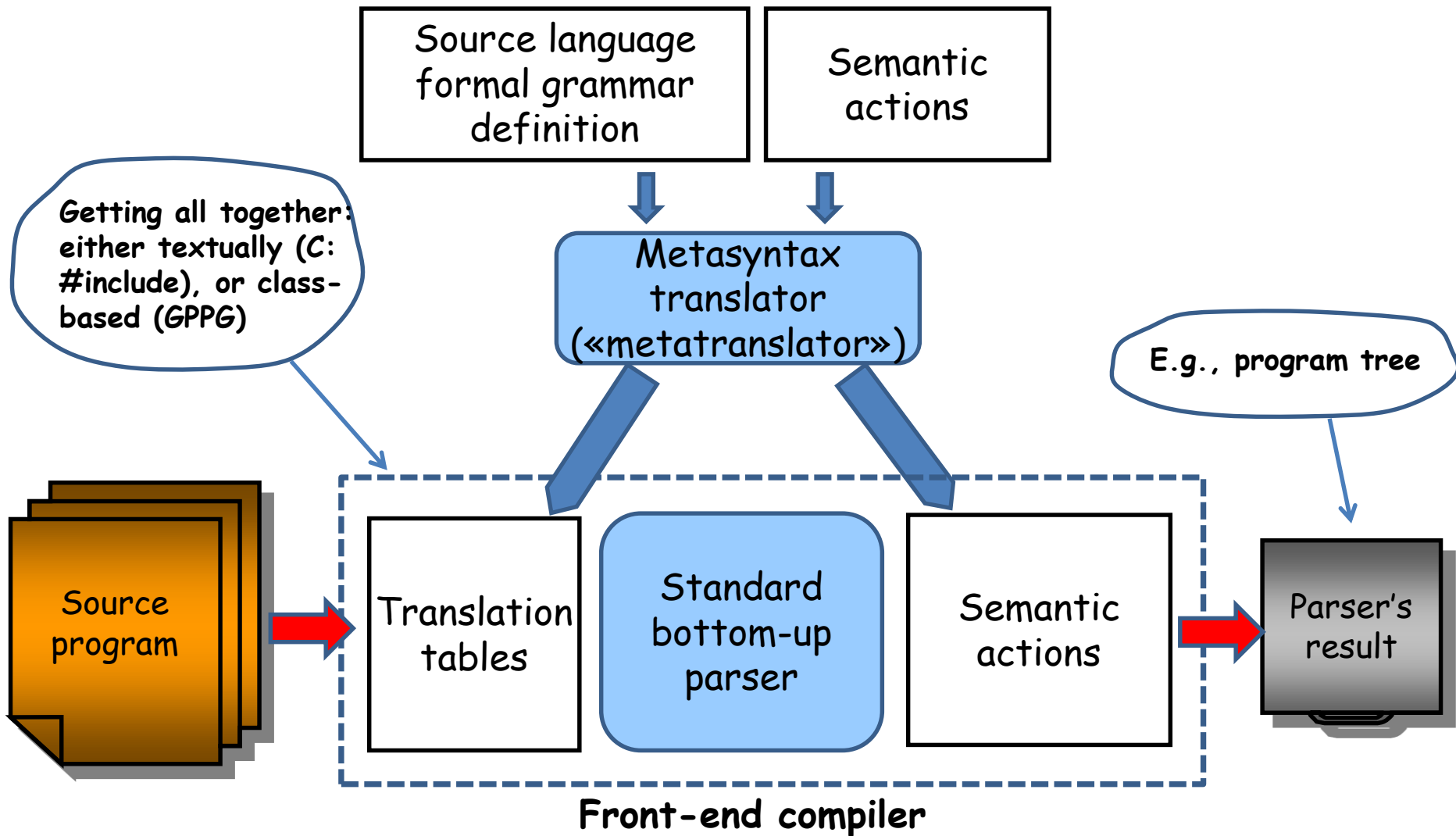
# Yacc/Bison & clones: features

- Generates **bottom-up** syntax parsers.
- Has its **own notation** (formalism) for grammar specification.
- Internally, the grammar is represented in a **table form**; the generated parser is table-driven.
- Source tokens should be generated by a separate lexical analyzer: either by a hand made analyzer or by Lex/Flex or compatible (Yacc uses integer token codes).
- Very good grammar readability.
- Separation the grammar from semantic actions.
- Rules with left recursion are allowed.
- Good standard support for error recovery.
- Hard to debug the grammar and to find ambiguities.

# Yacc based technology



# Yacc based technology



# Yacc: Grammar structure

Common declarations (implementation language)

%%

Declarations of token, types, associativity,...

Declaration of the **main rule**

%%

**Grammar rules** (together with semantic actions)

%%

Common declarations (implementation language)



# Grammar examples

- **ANSI C & Oberon grammars.**

**Are in DropBox:** Lecture 7–9 Addendum

- **C++ & C# grammars (a bit out of date).**

**In the Appendix of the Russian translation  
of the “Dragon Book” (1<sup>st</sup> edition).**

- **Following slides: Toy language grammar.**

**Are in DropBox:** Lecture 7–9 Addendum

# Example: Toy language grammar

```
// Identifiers & numbers
```

```
%token IDENTIFIER
```

```
%token NUMBER
```

```
// Keywords
```

```
%token IMPORT CLASS EXTENDS PRIVATE PUBLIC STATIC VOID IF ELSE
```

```
%token WHILE LOOP RETURN PRINT NULL NEW INT REAL
```

```
// Delimiters
```

```
%token LBRACE      // {
```

```
%token RBRACE      // }
```

```
%token LPAREN       // (
```

```
%token RPAREN       // )
```

```
%token LBRACKET     // [
```

```
%token RBRACKET     // ]
```

```
%token COMMA        // ,
```

```
%token DOT           // .
```

```
%token SEMICOLON    // ;
```

```
// Operator signs
```

```
%token ASSIGN        // =
```

```
%token LESS          // <
```

```
%token GREATER       // >
```

```
%token EQUAL         // ==
```

```
%token NOT_EQUAL     // !=
```

```
%token PLUS          // +
```

```
%token MINUS         // -
```

```
%token MULTIPLY      // *
```

```
%token DIVIDE        // /
```

```
%start CompilationUnit
```

Language  
alphabet

Gets converted to

```
enum Tokens
{
    IDENTIFIER,
    NUMBER,
    ...
};
```

Grammar  
main rule

# Example: Toy language grammar

```
CompilationUnit
: Imports ClassDeclarations
;

Imports
: /* empty */
| Import Imports
;

Import
: IMPORT IDENTIFIER SEMICOLON
;

ClassDeclarations
: /* empty */
| ClassDeclaration ClassDeclarations
;

ClassDeclaration
: CLASS IDENTIFIER SEMICOLON Extension ClassBody
| PUBLIC CLASS IDENTIFIER SEMICOLON Extension ClassBody
;

Extension
: /* empty */
| EXTENDS Identifier
;

ClassBody
: LBRACE RBRACE
| LBRACE ClassMembers RBRACE
;

ClassMembers
: ClassMember
| ClassMembers ClassMember
;
```

**Grammar:  
program &  
classes**

# Example: Toy language grammar

```
ClassMember
: FieldDeclaration
| MethodDeclaration
;

FieldDeclaration
: visibility Staticness Type IDENTIFIER SEMICOLON
;

visibility
: /* empty */
| PRIVATE
| PUBLIC
;

Staticness
: /* empty */
| STATIC
;

MethodDeclaration
: visibility Staticness MethodType IDENTIFIER Parameters Body
;

Parameters
: LPAREN RPAREN
| LPAREN ParameterList RPAREN
;

ParameterList
: Parameter
| ParameterList COMMA Parameter
;

Parameter
: Type IDENTIFIER ;
```

**Grammar:  
declarations**

# Example: Toy language grammar

```
MethodType
  : Type
  | VOID
  ;

Body
  : LBRACE LocalDeclarations Statements RBRACE
  ;

LocalDeclarations
  : LocalDeclaration
  | LocalDeclarations LocalDeclaration
  ;

LocalDeclaration
  : Type IDENTIFIER SEMICOLON
  ;
```

**Grammar:  
declarations**

# Example: Toy language grammar

Statements

```
: Statement
| Statements Statement
;
```

Statement

```
: Assignment | IfStatement | whileStatement | ReturnStatement
| CallStatement | PrintStatement | Block
;
```

Assignment

```
: LeftPart ASSIGN Expression SEMICOLON
;
```

LeftPart

```
: CompoundName
| CompoundName LBRACKET Expression RBRACKET
;
```

CompoundName

```
: IDENTIFIER
| CompoundName DOT IDENTIFIER
;
```

IfStatement

```
: IF LPAREN Relation RPAREN Statement
| IF LPAREN Relation RPAREN Statement ELSE Statement
;
```

whileStatement

```
: WHILE Relation LOOP Statement SEMICOLON
;
```

ReturnStatement

```
: RETURN SEMICOLON
| RETURN Expression SEMICOLON
;
```

**Grammar:  
statements**

# Example: Toy language grammar

callStatement

```
: CompoundName LPAREN RPAREN SEMICOLON  
| CompoundName LPAREN ArgumentList RPAREN SEMICOLON  
;
```

ArgumentList

```
: Expression  
| ArgumentList COMMA Expression  
;
```

PrintStatement

```
: PRINT Expression SEMICOLON  
;
```

Block

```
: LBRACE RBRACE  
| LBRACE Statements RBRACE  
;
```

**Grammar:  
statements**

# Example: Toy language grammar

```
Relation
: Expression
| Expression RelationalOperator Expression
;
RelationalOperator
: LESS | GREATER | EQUAL | NOT_EQUAL
;
Expression
: Term Terms
| AddSign Term Terms
;
AddSign
: PLUS | MINUS
;
Terms
: /* empty */
| AddSign Term Terms
;
Term
: Factor Factors
;
Factors
: /* empty */
| MultSign Factor Factors
;
MultSign
: MULTIPLY | DIVIDE
;
```

**Grammar:  
expressions**



# Example: Toy language grammar

Factor

```
: NUMBER
| LeftPart
| NULL
| NEW NewType
| NEW NewType LBRACKET Expression RBRACKET
;
```

NewType

```
: INT
| REAL
| IDENTIFIER
;
```

Type

```
: INT      ArrayTail
| REAL     ArrayTail
| IDENTIFIER ArrayTail
;
```

ArrayTail

```
: /* empty */
| LBRACKET RBRACKET
;
```

**Grammar:  
types**

# Toy grammar: comments

1. No means for expression repetitions (like in BNF format) in YACC notation; we have to use recursion instead.

```
ParameterList
    :
    | ParameterList COMMA Parameter
    ;

...
ArgumentList
    :
    | ArgumentList COMMA Expression
    ;

...
Statements
    :
    | Statements Statement
    ;
```

# Toy grammar: comments

2. Both right and left recursions are allowed and supported.

```
Expression
:      Term Terms
| AddSign Term Terms
;

AddSign
: PLUS | MINUS
;

Terms
: /* empty */
| AddSign Term Terms
;

Term
: Factor Factors
;

Factors
: /* empty */
| MultSign Factor Factors
;

MultSign
: MULTIPLY | DIVIDE
;
```

# Toy grammar: comments

## 3. Grouping is not supported; we have to add extra rules for grouping

```
AddSign
: PLUS | MINUS
;

Terms
: /* empty */
| AddSign Term Terms
;

Term
: Factor Factors
;

Factors
: /* empty */
| MultSign Factor Factors
;

MultSign
: MULTIPLY | DIVIDE
;
```



```
Terms
: /* empty */
| (PLUS|MINUS) Term Terms
;

Factors
: /* empty */
| (MULTIPLY | DIVIDE)
  Factor Factors
;
```

# Toy grammar translation

```
C:\Lectures\GPG 1.5.0\binaries>
```

```
gppg /conflicts "C:\Lectures\Lecture 8\Toy.yacc
```

Shift/Reduce conflict

Shift "IDENTIFIER": State-20 -> State-21

Reduce 30: MethodType -> Type

1

Shift/Reduce conflict

Shift "ELSE": State-87 -> State-88

Reduce 50: IfStatement -> IF, LPAREN, Relation, RPAREN, Statement

2

Shift/Reduce conflict

Shift "LBRACKET": State-120 -> State-122

Reduce 48: CompoundName -> IDENTIFIER

3

# Toy grammar translation

```
C:\Lectures\GPG 1.5.0\binaries>
```

```
gppg /conflicts "C:\Lectures\Lecture 8\Toy.yacc
```

Shift/Reduce conflict

```
Shift "IDENTIFIER":    State-20 -> State-21
```

```
Reduce 30:    MethodType -> Type
```

1

```
FieldDeclaration: Visibility Staticness Type . IDENTIFIER SEMICOLON
```

```
MethodType: Type .
```

```
FieldDeclaration
```

```
    : Visibility Staticness Type IDENTIFIER SEMICOLON  
    ;
```

```
...
```

```
MethodDeclaration
```

```
    : Visibility Staticness MethodType IDENTIFIER Parameters Body  
    ;
```

```
...
```

```
Type
```

```
    : ...  
    | IDENTIFIER ArrayTail  
    ;
```

```
MethodType
```

```
    : Type  
    | ...  
    ;
```

**Shift/Reduce conflicts are  
resolved in favor of Shift**

# Toy grammar translation

```
C:\Lectures\GPG 1.5.0\binaries>
```

```
gppg /conflicts "C:\Lectures\Lecture 8\Toy.yacc
```

Shift/Reduce conflict

Shift "ELSE": State-87 -> State-88

Reduce 50: IfStatement -> IF, LPAREN, Relation, RPAREN, Statement

2

IfStatement: IF LPAREN Relation RPAREN Statement .

IfStatement: IF LPAREN Relation RPAREN Statement . ELSE Statement

```
IfStatement
: IF LPAREN Relation RPAREN Statement
| IF LPAREN Relation RPAREN Statement ELSE Statement
;
```



```
IfStatement
: IF LPAREN Relation RPAREN Statement ElseTail
;
ElseTail
: /* empty */
| ELSE Statement
;
```

# Toy grammar translation

Shift/Reduce conflict

Shift "LBRACKET": State-120 -> State-122

Reduce 48: CompoundName -> IDENTIFIER

3

CompoundName: IDENTIFIER .

Type: IDENTIFIER . ArrayTail

```
C [ 10 ] = 7 ; // assignment
    ] a ;       // declaration
```

Assignment

```
: LeftPart ASSIGN Expression SEMICOLON
;
```

LeftPart

```
: CompoundName
| CompoundName LBRACKET Expression RBRACKET
;
```

Body

```
: LBRACE LocalDeclarations Statements RBRACE
;
```

LocalDeclaration

```
: Type IDENTIFIER SEMICOLON
;
```

Type

```
: ...
| IDENTIFIER ArrayTail
;
```

ArrayTail

```
: ...
| LBRACKET RBRACKET
;
```



# Toy grammar translation

Let's introduce an error to the grammar:

```
Body
    : LBRACE LocalDeclarations Statements RBRACE
    ;
...
Statement
    : LocalDeclaration
    | Assignment
    | IfStatement
    | whileStatement
    | ReturnStatement
    | CallStatement
    | PrintStatement
    | Block
    ;
```

Reduce/Reduce conflict in state 131 on symbol INT

Reduce 34: LocalDeclarations -> LocalDeclarations, LocalDeclaration

Reduce 38: Statement -> LocalDeclaration

Reduce/Reduce should be resolved by  
developer (by transforming the grammar)

# Toy grammar: semantic actions

```
Statements
:
| Statement { $$ = createStmtList($1); }
| Statements Statement { $$ = addStmtToList($1,$2); }
;

Statement
: Assignment | IfStatement | whileStatement | ReturnStatement
| ...
;

Assignment
: LeftPart ASSIGN Expression SEMICOLON { $$ = createAssign($1,$3); }
;

IfStatement
: IF LPAREN Relation RPAREN Statement
    { $$ = createIf($3,$5,NULL); }
| IF LPAREN Relation RPAREN Statement ELSE Statement
    { $$ = createIf($3,$5,$7); }
;

whileStatement
: WHILE Relation LOOP Statement SEMICOLON { $$ = createwhile($2,$4); }
;

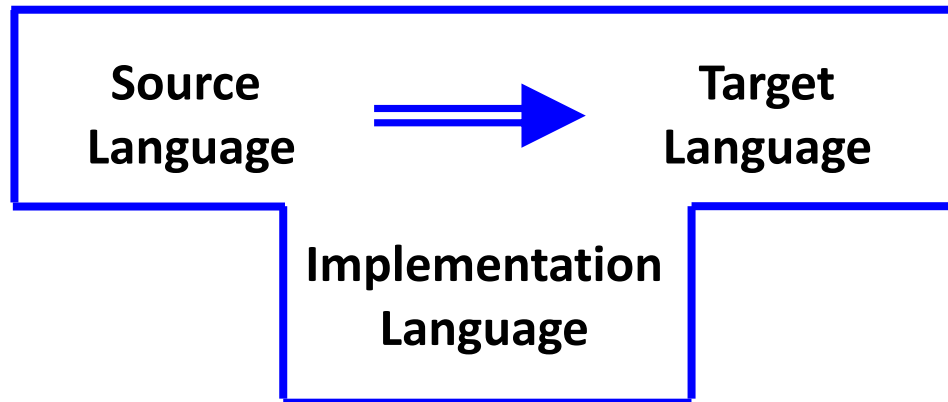
ReturnStatement
: RETURN SEMICOLON { $$ = createReturn(NULL); }
| RETURN Expression SEMICOLON { $$ = createReturn($2); }
;
```

# Bootstrapping Technology

- When implementation language & source language are **the same**. Examples: Ada; Eiffel; Scala.
- **Advantages:**  
More stable technology; supports graduate language & compiler improvement; no dependency on any third-party tools; the code of the compiler is an **excellent test** for both language and compiler itself.
- **Disadvantages:**  
A bit awkward technology; requires non-trivial management & powerful management tools (e.g., *ant*).

# Bootstrapping Technology

- Reference: Terence Pratt.
- Graphical notation ("T Notation"):



# Bootstrapping the Compiler

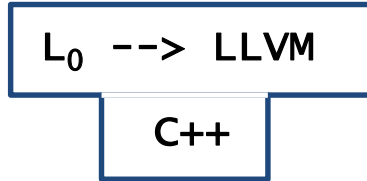
Initial development step

$L_0$

- Define a very simple subset of the target L language:  $L_0$

# Bootstrapping the Compiler

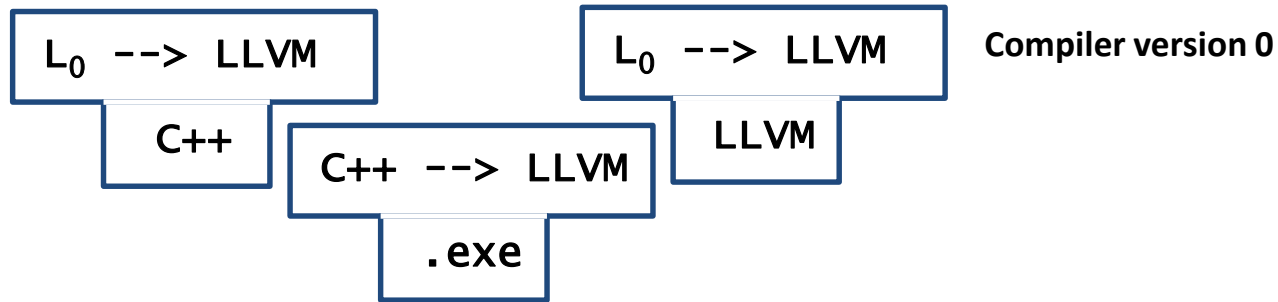
Initial development step



- Define a very simple subset of the target L language:  $L_0$
- Write the prototype compiler for  $L_0$  using a third party language with an existing compiler

# Bootstrapping the Compiler

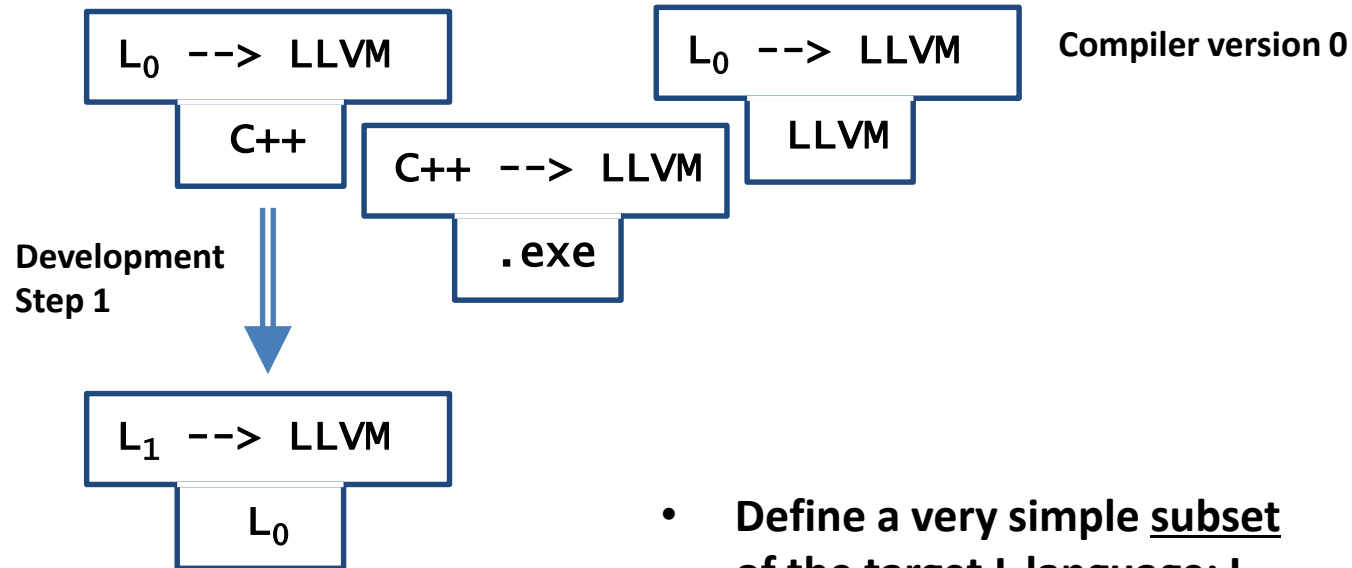
Initial development step



- Define a very simple subset of the target L language:  $L_0$
- Write the prototype compiler for  $L_0$  using a third party language with an existing compiler
- Compile the  $L_0$  compiler getting the 0th version of the target compiler.

# Bootstrapping the Compiler

Initial development step

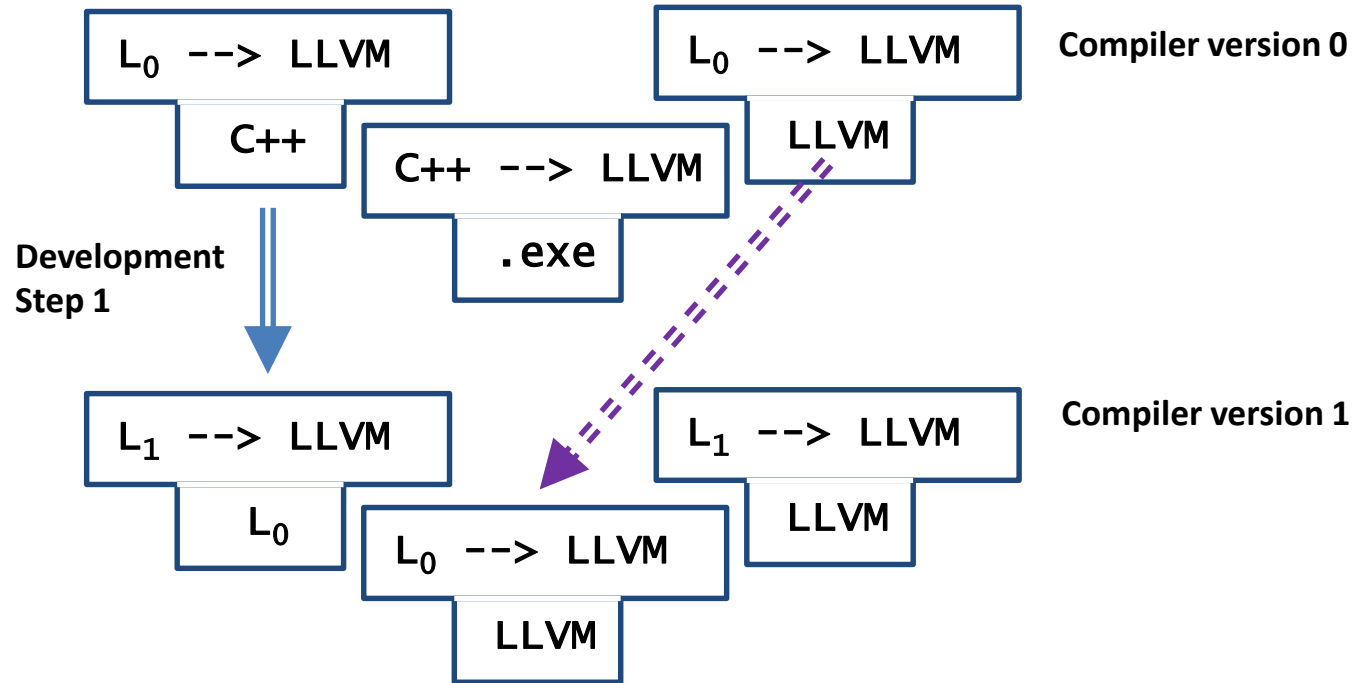


- Define a very simple subset of the target L language:  $L_0$
- Write the prototype compiler for  $L_0$  using a third party language with an existing compiler
- Compile the  $L_0$  compiler getting the 0th version of the target compiler.
- DevStep1: Create the next version of the language and (re)write the compiler in  $L_0$  language...



# Bootstrapping the Compiler

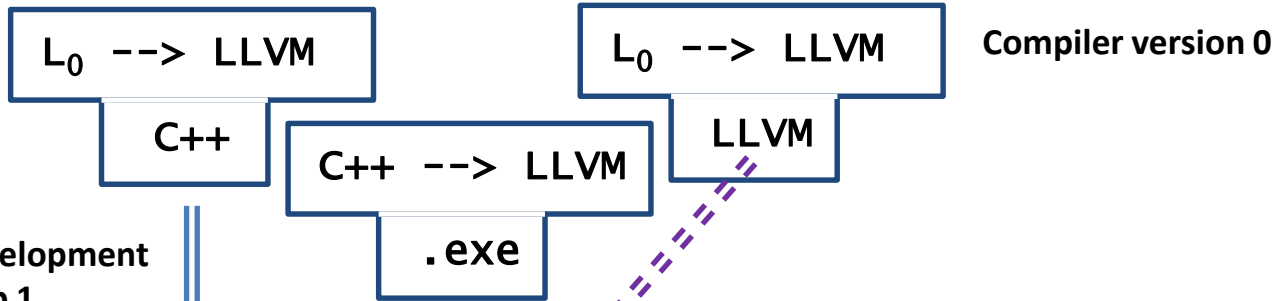
Initial development step



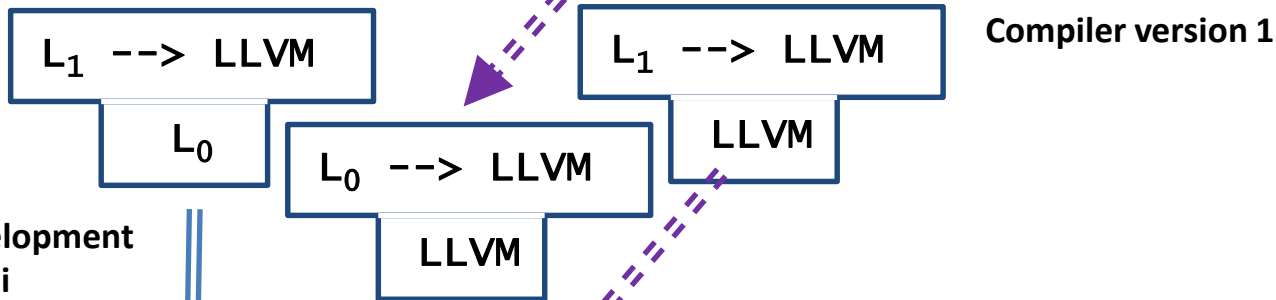
- **Dev Step1:**  
...and compile it using the previous compiler version!

# Bootstrapping the Compiler

Initial development step



Development Step 1



Development Step i

