# ANTLR 3

Mark Volkmann
mark@ociweb.com
Object Computing, Inc.
2008

---

# ANTLR Overview

▸ ANother Tool for Language Recognition
  ▸ written by Terence Parr in Java

▸ Easier to use than most/all similar tools

▸ Supported by ANTLRWorks
  ▸ graphical grammar editor and debugger
  ▸ written by Jean Bovet using Swing

▸ Used to implement
  ▸ "real" programming languages
  ▸ domain-specific languages (DSLs)

▸ http://www.antlr.org
  ▸ download ANTLR and ANTLRWorks here
  ▸ both are free and open source
  ▸ docs, articles, wiki, mailing list, examples

I'm a professor at the University of San Francisco.

**Ter**

I worked with Ter as a masters student there.

**Jean**

ANTLR 3

# ANTLR Documentation

# ANTLR Overview …

▸ **Uses EBNF grammars**
  - ▸ Extended Backus-Naur Form
  - ▸ can directly express optional and repeated elements
  - ▸ supports subrules (parenthesized groups of elements)

> BNF grammars require more verbose syntax to express these.

▸ **Supports many target languages for generated code**
  - ▸ Java, Ruby, Python, Objective-C, C, C++ and C#

▸ **Provides infinite lookahead**
  - ▸ most parser generators don't
  - ▸ used to choose between rule alternatives

▸ **Plug-ins available for IDEA and Eclipse**

# ANTLR Overview ...

▸ Three main use cases

> We'll explain actions and rewrite rules later.

▸ 1) Implementing "**validators**" | no actions or rewrite rules

> ▸ generate code that validates that input obeys grammar rules

▸ 2) Implementing "**processors**" | actions but no rewrite rules

> ▸ generate code that validates and processes input
>
> ▸ could include performing calculations, updating databases, reading configuration files into runtime data structures, ...
>
> ▸ our Math example coming up does this

▸ 3) Implementing "**translators**" | actions containing println and/or rewrite rules

> ▸ generate code that validates and translates input into another format such as a programming language or bytecode

---

# Projects Using ANTLR

▸ Programming languages

> ▸ Boo
>> ▸ http://boo.codehaus.org
>
> ▸ Groovy
>> ▸ http://groovy.codehaus.org
>
> ▸ Mantra
>> ▸ http://www.linguamantra.org
>
> ▸ Nemerle
>> ▸ http://nemerle.org
>
> ▸ XRuby
>> ▸ http://xruby.com

▸ Other tools

> ▸ Hibernate
>> ▸ for its HQL to SQL query translator
>
> ▸ Intellij IDEA
>
> ▸ Jazillian
>> ▸ translates COBOL, C and C++ to Java
>
> ▸ JBoss Rules (was Drools)
>
> ▸ Keynote (Apple)
>
> ▸ WebLogic (Oracle)
>
> ▸ too many more list!

> See showcase and testimonials at http://antlr.org/showcase/list and http://www.antlr.org/testimonial/.

# Books



▸ "ANTLR Recipes"? in the works

  ▸ another Pragmatic Programmers book from Terence Parr

---

# Other DSL Approaches

▸ Languages like Ruby and Groovy
  are good at implementing DSLs, but ...

▸ The DSLs have to live within
  the syntax rules of the language

▸ For example

  ▸ dots between object references and method names

  ▸ parameters separated by commas

  ▸ blocks of code surrounded by { ... } or do ... end

▸ What if you don't want these
  in your language?

# Conventions

▸ ANTLR grammar syntax makes frequent use of the characters [ ] and { }

▸ In these slides

  ▸ when describing a placeholder, I'll use italics

  ▸ when describing something that's optional, I'll use `item?`

---

# Some Definitions

```
character
stream
   │
   ▼
 Lexer
   │
   ▼
 token
 stream
   │
   ▼
 Parser
   │
   ▼
  AST
   │
   ▼
 Tree
 Parser
   │
   ▼
template
calls
   │
   ▼
 text
output
```

▸ Lexer

  ▸ converts a stream of characters to a stream of tokens

▸ Parser

  | Token objects know their start/stop character stream index, line number, index within the line, and more. |

  ▸ processes a stream of tokens, possibly creating an AST

▸ Abstract Syntax Tree (AST)

  ▸ an intermediate tree representation of the parsed input that

    ▸ is simpler to process than the stream of tokens

    ▸ can be efficiently processed multiple times

▸ Tree Parser

  ▸ processes an AST

▸ StringTemplate

  ▸ a library that supports using templates with placeholders for outputting text (for example, Java source code)

# General Steps

▸ Write grammar
  ▸ can be in one or more files

▸ Optionally write StringTemplate templates

▸ Debug grammar with ANTLRWorks

▸ Generate classes from grammar
  ▸ these validate that text input conforms to the grammar and execute target language "actions" specified in the grammar

▸ Write application that uses generated classes

▸ Feed the application
  text that conforms to the grammar

---

# Let's Create A Language!

▸ Features
  ▸ run on a file or interactively
  ▸ get help – **? or help**
  ▸ one data type, double
  ▸ assign values to variables – **a = 3.14**
  ▸ define polynomial functions – **f(x) = 3x^2 – 4x + 2**
  ▸ print strings, numbers, variables and function evaluations –
    **print "The value of f for " a " is " f(a)**
  ▸ print the definition of a function and its derivative –
    **print "The derivative of " f() " is " f'()**
  ▸ list variables and functions –
    **list variables** and **list functions**
  ▸ add/subtract functions – **h = f – g** ←

    | Input: |
    | --- |
    | f(x) = 3x^2 – 4 |
    | g(y) = y^2 – 2y + 1 |
    | h = f – g |
    | print h() |
    | |
    | **Output:** |
    | h(x) = 2x^2 + 2x – 5 |

  ▸ the function variables don't have to match
  ▸ exit – **exit** or **quit**

# Example Input/Output

```
a = 3.14

f(x) = 3x^2 - 4x + 2

print "The value of f for " a " is " f(a)


print "The derivative of " f() " is " f'()


list variables
list functions


g(y) = 2y^3 + 6y - 5
h = f + g
print h()
```

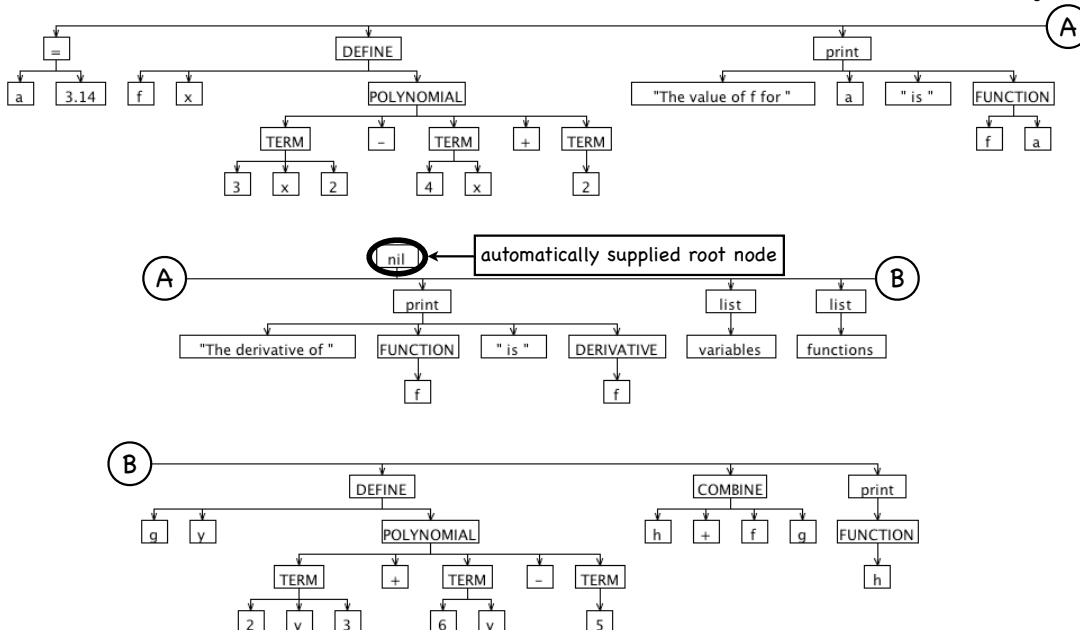```
The value of f for 3.14 is 19.0188
The derivative of f(x) = 3x^2 - 4x + 2
is f'(x) = 6x - 4
# of variables defined: 1
a = 3.14
# of functions defined: 1
f(x) = 3x^2 - 4x + 2
h(x) = 2x^3 + 3x^2 + 2x - 3
```

ANTLR 3

# Example AST

drawn by
ANTLRWorks

ANTLR 3

# Important Classes

Key:
provided
generated
written

```
                    ┌──────────────────┐
                    │  BaseRecognizer  │
                    └──────────────────┘
                       △    △    △
          ┌─────────┐  ┌──────────┐  ┌──────────────┐
          │  Lexer  │  │  Parser  │  │  TreeParser  │
          └─────────┘  └──────────┘  └──────────────┘
              △             △              △
       ┌──────────────┐ ┌──────────────┐ ┌──────────────┐     ┌──────────────┐
       │  MathLexer   │ │  MathParser  │ │   MathTree   │────→│   Function   │
       └──────────────┘ └──────────────┘ └──────────────┘     └──────────────┘
                │              │       ╲       │  ╲            │              │
                │              │         ╲     │    ╲────→┌──────────────┐    │
                │              │           ╲   │         │  Polynomial  │←───┘
                │              │             ╲ │         └──────────────┘
                │              │       ┌─────────────┐          │
                └──────────────┴───────│  Processor  │     ┌─────────┐
                                       └─────────────┘     │  Term   │
                                                           └─────────┘
```

---

# Grammar Actions

▸ Add to the generated code

▸ @*grammar-type*::header { ... }

> ▸ inserts contained code before the class definition
>
>> ▸ commonly used to specify a package name
>> and import classes in other packages

> *grammar-type*
> must be lexer,
> parser (the default)
> or treeparser

▸ @*grammar-type*::members { ... }

> ▸ inserts field declarations and methods inside the class definition
>
> ▸ commonly used to
>
>> ▸ define constants and attributes accessible to
>> all rule methods in the generated class
>>
>> ▸ define methods used by multiple rule actions
>>
>> ▸ override methods in the superclasses of the generated classes
>>
>>> ▸ useful for customizing error reporting and handling

# Lexer Rules

- Need one for every kind of token to be processed in parser grammar
- Name must start uppercase
  - typically all uppercase
- Assign a token name to
  - a single literal string found in input
  - a selection of literal strings found in input
  - one or more characters and ranges of characters
    - can use cardinality indicators ?, * and +
- Can refer to other lexer rules
- "fragment" lexer rules
  - do not result in tokens
  - are only referenced by other lexer rules

Regular expressions aren't supported.

The next lexer rule used is the one that matches the most characters. If there is a tie, the one listed first is used, so order matters!

See LETTER and DIGIT rules in the upcoming example.

---

# Whitespace & Comments

- Handled in lexer rules
- Two common options
  - throw away - `skip();`
  - write to a different "channel" - `$channel = HIDDEN;`

The ANTLRWorks debugger input panel doesn't display skipped characters, but does display hidden ones.

constant defined in BaseRecognizer

- Examples

Don't skip or hide NEWLINEs if they are used as statement terminators.

```
WHITESPACE: (' ' | '\t')+ { $channel = HIDDEN; };
NEWLINE: ('\r'? '\n')+;
SINGLE_COMMENT: '//' ~('\r' | '\n')* NEWLINE { skip(); };
MULTI_COMMENT
options { greedy = false; }
  : '/*' .* '*/' NEWLINE? { skip(); };
```

The greedy option defaults to true, except for the patterns .* and .+, so it doesn't need to be specified here. When true, the lexer matches as much input as possible. When false, it stops when input matches the next element.

# Our Lexer Grammar

```
(lexer) grammar MathLexer;

@header { package com.ociweb.math; }
```

We want the generated lexer class to be in this package.

```
APOSTROPHE: '\''; // for derivative
ASSIGN: '=';
CARET: '^'; // for exponentiation
FUNCTIONS: 'functions'; // for list command
HELP: '?' | 'help';
LEFT_PAREN: '(';
LIST: 'list';
PRINT: 'print';
RIGHT_PAREN: ')';
SIGN: '+' | '-';
VARIABLES: 'variables'; // for list command


NUMBER: INTEGER | FLOAT;
fragment FLOAT: INTEGER '.' '0'..'9'+;
fragment INTEGER: '0' | SIGN? '1'..'9' '0'..'9'*;
```

ANTLR 3

# Our Lexer Grammar …

```
NAME: LETTER (LETTER | DIGIT | '_')*;
STRING_LITERAL: '"' NONCONTROL_CHAR* '"';

fragment NONCONTROL_CHAR: LETTER | DIGIT | SYMBOL | SPACE;
fragment LETTER: LOWER | UPPER;
fragment LOWER: 'a'..'z';
fragment UPPER: 'A'..'Z';
fragment DIGIT: '0'..'9';
fragment SPACE: ' ' | '\t';

// Note that SYMBOL omits the double-quote character,
// digits, uppercase letters and lowercase letters.
fragment SYMBOL: '!' | '#'..'/' | ':'..'@' | '['..'`' | '{'..'~';

// Windows uses \r\n. UNIX and Mac OS X use \n.
// To use newlines as a terminator,
// they can't be written to the hidden channel!
NEWLINE: ('\r'? '\n')+;
WHITESPACE: SPACE+ { $channel = HIDDEN; };
```

ANTLR 3

# Token Specification

▸ The lexer creates tokens
  for all input character sequences
  that match lexer rules

▸ It can be useful to create other tokens that

  ▸ don't exist in the input (imaginary)

    ▸ often serve to group other tokens

  ▸ have a better name than is found in the input

| See all the uppercase token names in the AST diagram on slide 14. |
|---|

▸ Do this with a token specification
  in the parser grammar

| We need this for the imaginary tokens DEFINE, POLYNOMIAL, TERM, FUNCTION, DERIVATIVE and COMBINE. |
|---|

  ▸ **tokens {**
      ***imaginary-name;***
      ***better-name = 'input-name';***
      **. . .**
    **}**

---

# Rule Syntax

| only for lexer rules |
|---|

**fragment?** *rule-name arguments*?

**(returns** *return-values***)?**

*throws-spec*?

*rule-options*?

*rule-attribute-scopes*?

| add code before and/or after code in the generated method for this rule |
|---|

→ *rule-actions*?

  **: *token-sequence-1***

  **| *token-sequence-2***

  **...**

  **;**

| Each element in these alternative sequences can be followed by an action which is target language code in curly braces. The code is executed immediately after a preceding element is matched by input. |
|---|

***exceptions-spec*?**

| to customize exception handling for this rule |
|---|

# Creating ASTs

▸ Requires grammar option `output = AST;`

▸ Approach #1 - Rewrite rules

  ▸ appear <u>after a rule alternative</u>
  ▸ the recommended approach in most cases
  ▸ `-> ^(parent child-1 child-2 ... child-n)`

  | can't use both approaches in the same rule alternative! |
  | --- |

▸ Approach #2 - AST operators

  ▸ appear <u>in a rule alternative</u>, immediately after tokens
  ▸ works best for sequences like mathematical expressions
  ▸ operators
    ▸ ^ - make new root node for all child nodes at the same level
    ▸ none - make a child node of current root node
    ▸ ! - don't create a node ⟵──────── often used for bits of syntax that aren't needed in the AST such as parentheses, commas and semicolons
  ▸ parent^ '('! child-1 child-2 ... child-n ')'!

---

# Declaring Rule Arguments
# and Return Values

```
rule-name[type1 name1, type2 name2, ...]
returns [type1 name1, type2 name2, ...] :
  ...
;
```

| arguments |
| --- |

| return values; can have more than one |
| --- |

ANTLR generates a class to use as the return type of the generated method for the rule.

Instances of this class hold all the return values.

The generated method name matches the rule name.

The name of the generated return type class is the rule name with "_return" appended.

# Our Parser Grammar

```
parser grammar MathParser;

options {
    output = AST;
    tokenVocab = MathLexer;
}

tokens {
    COMBINE;
    DEFINE;
    DERIVATIVE;
    FUNCTION;
    POLYNOMIAL;
    TERM;
}

@header { package com.ociweb.math; }
```

We're going to output an AST.

We're going to use the tokens defined in our MathLexer grammar.

These are imaginary tokens that will serve as parent nodes for grouping other tokens in our AST.

We want the generated parser class to be in this package.

---

# Our Parser Grammar ...

```
// This is the "start rule".
script: statement* EOF;

statement: assign | define | interactiveStatement | combine | print;

interactiveStatement: help | list;

assign: NAME ASSIGN value terminator -> ^(ASSIGN NAME value);

value: NUMBER | NAME | functionEval;

functionEval
  : fn=NAME LEFT_PAREN (v=NUMBER | v=NAME) RIGHT_PAREN -> ^(FUNCTION $fn $v);

// EOF cannot be used in lexer rules, so we made this a parser rule.
// EOF is needed here for interactive mode where each line entered ends in EOF
// and for file mode where the last line ends in EOF.
terminator: NEWLINE | EOF;
```

AST operator

EOF is a predefined token that represents the end of input. The start rule should end with this.

An expression starting with "->" is called a "**rewrite rule**".

Examples:
a = 19
a = b
a = f(2)
a = f(b)

Parts of rule alternatives can be assigned to variables (ex. **fn** & **v**) that are used to refer to them in rule actions. Alternatively rule names (ex. **NAME**) can be used.

Examples:
f(2)
f(b)

When parser rule alternatives contain literal strings, they are converted to references to automatically generated lexer rules.
For example, we could eliminate the ASSIGN lexer rule and change **ASSIGN** to '=' in this grammar.
The rules in this grammar don't use literal strings.

# Our Parser Grammar ...

```
define
  : fn=NAME LEFT_PAREN fv=NAME RIGHT_PAREN ASSIGN
    polynomial[$fn.text, $fv.text] terminator
    -> ^(DEFINE $fn $fv polynomial);


// fnt = function name text; fvt = function variable text
polynomial[String fnt, String fvt]
  : term[$fnt, $fvt] (SIGN term[$fnt, $fvt])*
    -> ^(POLYNOMIAL term (SIGN term)*);
```

Examples:
**f(x) = 3x^2 - 4**
**g(y) = y^2 - 2y + 1**

To get the text value from a variable that refers to a Token object, use "**$var.text**".

Examples:
**3x^2 - 4**
**y^2 - 2y + 1**

ANTLR 3

---

# Our Parser Grammar ...

```
// fnt = function name text; fvt = function variable text
term[String fnt, String fvt]
  // tv = term variable
  : c=coefficient? (tv=NAME e=exponent?)?
    // What follows is a validating semantic predicate.
    // If it evaluates to false, a FailedPredicateException will be thrown.
    { tv == null ? true : ($tv.text).equals($fvt) }?
    -> ^(TERM $c? $tv? $e?)
  ;
  catch [FailedPredicateException fpe] {
    String tvt = $tv.text;
    String msg = "In function \"" + fnt +
      "\" the term variable \"" + tvt +
      "\" doesn't match function variable \"" + fvt + "\".";
    throw new RuntimeException(msg);
  }


coefficient: NUMBER;


exponent: CARET NUMBER -> NUMBER;
```
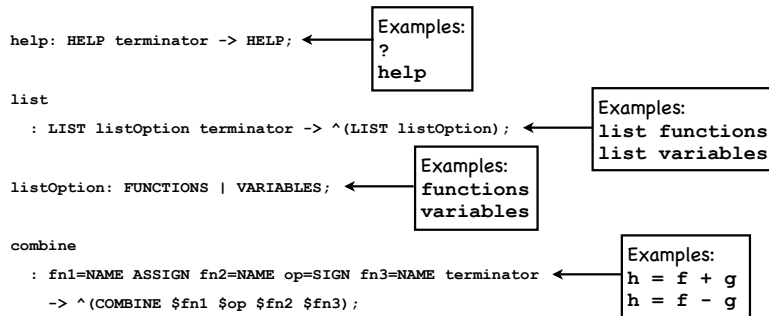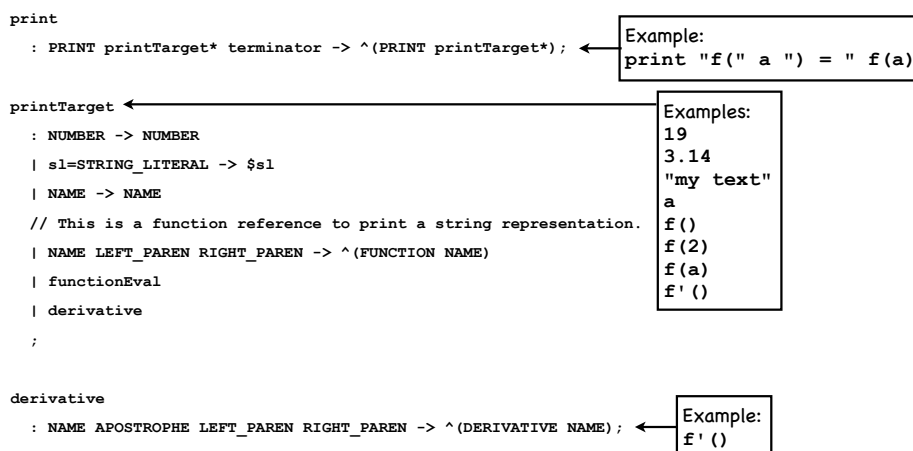
Examples:
**4**
**4x**
**x^2**
**4x^2**

Term variables must match their function variable.
This catches bad function definitions such as f(x) = 2y.

Example:
**^2**

ANTLR 3

# Our Parser Grammar ...

```
help: HELP terminator -> HELP;
```
Examples:
```
?
help
```

```
list
    : LIST listOption terminator -> ^(LIST listOption);
```
Examples:
```
list functions
list variables
```

```
listOption: FUNCTIONS | VARIABLES;
```
Examples:
```
functions
variables
```

```
combine
    : fn1=NAME ASSIGN fn2=NAME op=SIGN fn3=NAME terminator
      -> ^(COMBINE $fn1 $op $fn2 $fn3);
```
Examples:
```
h = f + g
h = f - g
```

---

# Our Parser Grammar ...

```
print
    : PRINT printTarget* terminator -> ^(PRINT printTarget*);
```
Example:
```
print "f(" a ") = " f(a)
```

```
printTarget
    : NUMBER -> NUMBER
    | sl=STRING_LITERAL -> $sl
    | NAME -> NAME
    // This is a function reference to print a string representation.
    | NAME LEFT_PAREN RIGHT_PAREN -> ^(FUNCTION NAME)
    | functionEval
    | derivative
    ;
```
Examples:
```
19
3.14
"my text"
a
f()
f(2)
f(a)
f'()
```

```
derivative
    : NAME APOSTROPHE LEFT_PAREN RIGHT_PAREN -> ^(DERIVATIVE NAME);
```
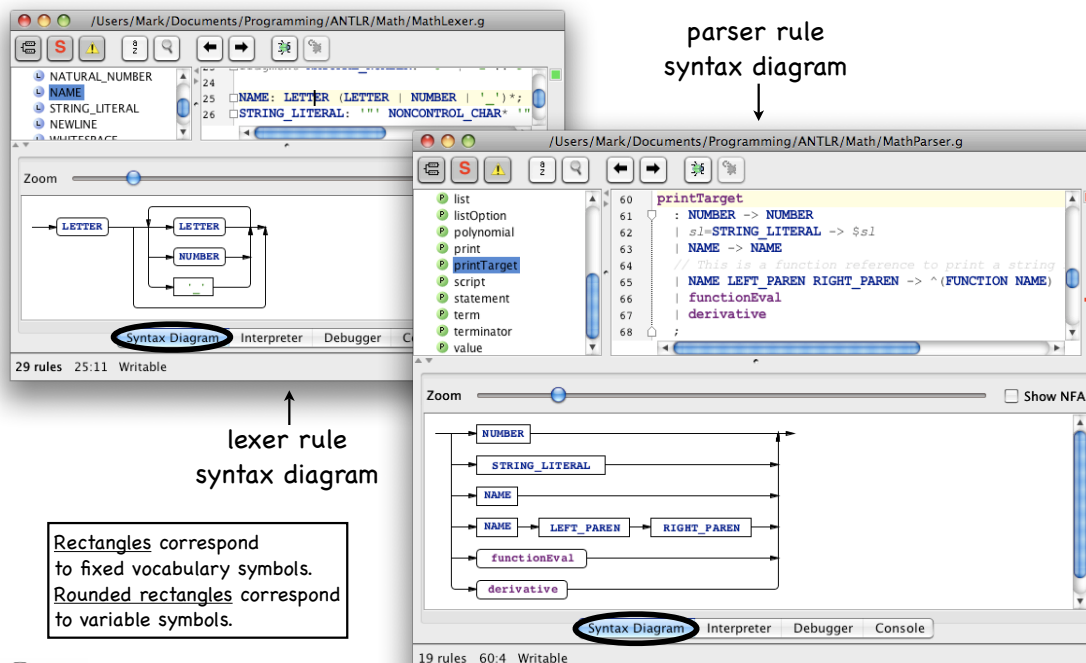Example:
```
f'()
```

# ANTLRWorks

▸ A graphical grammar editor and debugger

▸ Features

    ▸ highlights grammar syntax errors

    ▸ checks for grammar errors beyond the syntax variety

        ▸ such as conflicting rule alternatives

    ▸ displays a syntax diagram for the selected rule

    ▸ debugger can step through creation of parse trees and ASTs

ANTLR 3

---

# ANTLRWorks ...



parser rule
syntax diagram

lexer rule
syntax diagram

Rectangles correspond
to fixed vocabulary symbols.
Rounded rectangles correspond
to variable symbols.

ANTLR 3

# ANTLRWorks ...

grammar check result

requesting a grammar check

ANTLR 3

# ANTLRWorks Debugger

▸ **Simple when lexer and parser rules are combined in a single grammar file**

  ▸ press Debug toolbar button

  ▸ enter input text or select an input file

  ▸ select start rule

    ▸ allows debugging a subset of grammar

  ▸ press OK button

ANTLR 3

# ANTLRWorks Debugger ...

▸ At the bottom of the ANTLRWorks window



ASTExample
Start Rule: script

---

# ANTLRWorks Debugger ...

▸ A bit more complicated when
lexer and parser rules are in separate files

See the ANTLR Wiki page
"When do I need to use remote debugging?" at
http://www.antlr.org/wiki/pages/viewpage.action?pageId=5832732

▸ We'll demonstrate this
after we see the Java code that
ties all the generated classes together

  ▸ see slides 53-56

# Using Rule Return Values

These are examples from our tree grammar.

The code in curly braces is a rule "action" written in the target language, in this case Java.

```
printTarget
  : NUMBER { out($NUMBER); }
  | STRING_LITERAL {
      String s = unescape($STRING_LITERAL.text);
      out(s.substring(1, s.length() - 1)); // remove quotes
    }
  | NAME { out(getVariable($NAME)); }
  | ^(FUNCTION NAME) { out(getFunction(NAME)); }
  | functionEval { out($functionEval.result); }
  | derivative // handles own output
  ;
```

"unescape", "out", "getFunction", "getVariable", "evalFunction" and "toDouble" are methods we wrote in the tree grammar coming up.

```
functionEval returns [double result]
  : ^(FUNCTION fn=NAME v=NUMBER) {
      $result = evalFunction($fn, toDouble($v));
    }
  | ^(FUNCTION fn=NAME v=NAME) {
      $result = evalFunction($fn, getVariable($v));
    }
  ;
```

---

# Rule Actions

▸ Add code before and/or after
  the generated code
  in the method generated for a rule

  ▸ can be used for AOP-like wrapping of methods

▸ **@init { ... }**

  ▸ inserts contained code before generated code

    ▸ can be used to declare local variables used in actions of rule alternatives

  ▸ used in our tree parser **polynomial** and **term** rules ahead

▸ **@after { ... }**

  ▸ inserts contained code after generated code

# Attribute Scopes

▸ **Data is shared between rules in two ways**
  - ▸ passing parameters and/or returning values
  - ▸ using attributes

  > same as options to share data between Java methods in the same class

▸ **Attributes can be accessible to**
  - ▸ a single rule using `@init` to declare them
  - ▸ a rule and all rules invoked by it - <u>rule scope</u>
  - ▸ all rules that request the named <u>global scope</u> of the attributes

▸ **Attribute scopes**
  - ▸ define collections of attributes that can be accessed by multiple rules
  - ▸ two kinds, <u>global</u> and <u>rule</u> scopes

---

# Attribute Scopes ...

▸ **Global scopes**
  - ▸ named scopes defined outside any rule
  - ▸ define with
    ```
    scope name {
      type variable;
      . . .
    }
    ```
    > Use an @init rule action to initialize attributes.
  - ▸ request access to the scope in a rule with
    ```
    scope name;
    ```
    > To access multiple scopes, list them separated by spaces.
  - ▸ rule actions access variables in the scope with
    ```
    $name::variable
    ```

▸ **Rule scopes**
  - ▸ unnamed scopes defined inside a rule
  - ▸ define with
    ```
    scope {
      type variable;
      ...
    }
    ```
  - ▸ rule actions in the defining rule and rules invoked by it access attributes in the scope with
    ```
    $rule-name::variable
    ```

# Our Tree Grammar

```
(tree) grammar MathTree;

(options) {
    ASTLabelType = CommonTree;
    tokenVocab = MathParser;
}

(@header) {
    package com.ociweb.math;

    import java.util.Map;
    import java.util.TreeMap;
}

(@members) {
    private Map<String, Function> functionMap = new TreeMap<String, Function>();
    private Map<String, Double> variableMap = new TreeMap<String, Double>();
```

We're going to process an AST whose nodes are of type CommonTree.

We're going to use the tokens defined in both our MathLexer and MathParser grammars. The MathParser grammar already includes the tokens defined in the MathLexer grammar.

We want the generated parser class to be in this package.

We're using TreeMaps so the entries are sorted on their keys which is desired when listing them.

# Our Tree Grammar …

```
private void define(Function function) {
    functionMap.put(function.getName(), function);
}

private Function getFunction(CommonTree nameNode) {
    String name = nameNode.getText();
    Function function = functionMap.get(name);
    if (function == null) {
        String msg = "The function \"" + name + "\" is not defined.";
        throw new RuntimeException(msg);
    }
    return function;
}

private double evalFunction(CommonTree nameNode, double value) {
    return getFunction(nameNode).getValue(value);
}
```

This adds a Function to our function Map.

This retrieves a Function from our function Map whose name matches the text of a given AST tree node.

This evaluates a function whose name matches the text of a given AST tree node for a given value.

# Our Tree Grammar ...

```
private double getVariable(CommonTree nameNode) {
    String name = nameNode.getText();
    Double value = variableMap.get(name);
    if (value == null) {
        String msg = "The variable \"" + name + "\" is not set.";
        throw new RuntimeException(msg);
    }
    return value;
}


private static void out(Object obj) {
    System.out.print(obj);
}

private static void outln(Object obj) {
    System.out.println(obj);
}
```

This retrieves the value of a variable from our variable Map whose name matches the text of a given AST tree node.

These just shorten the code for print and println calls.

# Our Tree Grammar ...

```
private double toDouble(CommonTree node) {
    double value = 0.0;
    String text = node.getText();
    try {
        value = Double.parseDouble(text);
    } catch (NumberFormatException e) {
        throw new RuntimeException("Cannot convert \"" + text + "\" to a double.");
    }
    return value;
}


private static String unescape(String text) {
    return text.replaceAll("\\\\n", "\n");
}

} // @members
```

This converts the text of a given AST node to a double.

This replaces all escaped newline characters in a String with unescaped newline characters. It is used to allow newline characters to be placed in literal Strings that are passed to the print command.

# Our Tree Grammar ...

```
script: statement*;


statement: assign | combine | define | interactiveStatement | print;


interactiveStatement: help | list;


assign: ^(ASSIGN NAME v=value) { variableMap.put($NAME.text, $v.result); };
```

could also use $value here

This adds a variable to the variable map.

```
value returns [double result]
  : NUMBER { $result = toDouble($NUMBER); }
  | NAME { $result = getVariable($NAME); }
  | functionEval { $result = $functionEval.result; }
  ;
```

This returns a value as a double. The value can be a number, a variable name or a function evaluation.

```
functionEval returns [double result]
  : ^(FUNCTION fn=NAME v=NUMBER) {
      $result = evalFunction($fn, toDouble($v));
    }
  | ^(FUNCTION fn=NAME v=NAME) {
      $result = evalFunction($fn, getVariable($v));
    }
  ;
```

This returns the result of a function evaluation as a double.

ANTLR 3

---

# Our Tree Grammar ...

```
define
  : ^(DEFINE name=NAME variable=NAME polynomial) {
      define(new Function($name.text, $variable.text, $polynomial.result));
    }
  ;


polynomial returns [Polynomial result]
scope { Polynomial current; }
@init { $polynomial::current = new Polynomial(); }
  : ^(POLYNOMIAL term[""] (s=SIGN term[$s.text])*) {
      $result = $polynomial::current;
    }
  ;
```

This builds a Function object and adds it to the function map.

This builds a Polynomial object and returns it.

The "current" attribute in this rule scope is visible to rules invoked by this one, such as term.

There can be no sign in front of the first term, so "" is passed to the term rule.
The coefficient of the first term can be negative.
The sign between terms is passed to subsequent invocations of the term rule.

ANTLR 3

# Our Tree Grammar …

```
term[String sign]
@init { boolean negate = "-".equals(sign); }
  : ^(TERM coefficient=NUMBER) {
      double c = toDouble($coefficient);
      if (negate) c = -c; // applies sign to coefficient
      $polynomial::current.addTerm(new Term(c));
    }
  | ^(TERM coefficient=NUMBER? variable=NAME exponent=NUMBER?) {
      double c = coefficient == null ? 1.0 : toDouble($coefficient);
      if (negate) c = -c; // applies sign to coefficient
      double exp = exponent == null ? 1.0 : toDouble($exponent);
      $polynomial::current.addTerm(new Term(c, $variable.text, exp));
    }
  ;
```

This builds a Term object and adds it to the current Polynomial.

---

# Our Tree Grammar …

```
help
  : HELP {
      outln("In the help below");
      outln("* fn stands for function name");
      outln("* n stands for a number");
      outln("* v stands for variable");
      outln("");
      outln("To define");
      outln("* a variable: v = n");
      outln("* a function from a polynomial: fn(v) = polynomial-terms");
      outln("  (for example, f(x) = 3x^2 - 4x + 1)");
      outln("* a function from adding or subtracting two others: " +
        "fn3 = fn1 +|- fn2");
      outln("  (for example, h = f + g)");
      outln("");
      outln("To print");
      // some lines omitted for space
      outln("To exit: exit or quit");
    }
  ;
```

This outputs help on our language which is useful in interactive mode.

# Our Tree Grammar …

```
list
  : ^(LIST FUNCTIONS) {
      outln("# of functions defined: " + functionMap.size());
      for (Function function : functionMap.values()) {
        outln(function);
      }
    }
  | ^(LIST VARIABLES) {
      outln("# of variables defined: " + variableMap.size());
      for (String name : variableMap.keySet()) {
        double value = variableMap.get(name);
        outln(name + " = " + value);
      }
    }
  ;
```

> This lists all the functions or variables that are currently defined.

---

# Our Tree Grammar …

```
combine
  : ^(COMBINE fn1=NAME op=SIGN fn2=NAME fn3=NAME) {
      Function f2 = getFunction(fn2);
      Function f3 = getFunction(fn3);
      if ("+".equals($op.text)) {
        define(f2.add($fn1.text, f3));
      } else if ("-".equals($op.text)) {
        define(f2.subtract($fn1.text, f3));
      } else {
        // This should never happen since SIGN is defined to be either "+" or "-".
        throw new RuntimeException(
          "The operator \"" + $op +
          " cannot be used for combining functions.");
      }
    }
  ;
```

> This adds or subtracts two functions to create a new one.

> "$fn1.text" is the name of the new function to create.

# Our Tree Grammar ...

```
print
  : ^(PRINT printTarget*)
      { System.out.println(); };
```

This prints a list of printTargets then prints a newline.

```
printTarget
  : NUMBER { out($NUMBER); }
```

This prints a single printTarget without a newline.

```
  | STRING_LITERAL {
      String s = unescape($STRING_LITERAL.text);
      out(s.substring(1, s.length() - 1)); // removes quotes
    }
  | NAME { out(getVariable($NAME)); }
  | ^(FUNCTION NAME) { out(getFunction($NAME)); }
  | functionEval { out($functionEval.result); }
  | derivative
  ;
```

on slide 46

This prints the derivative of a function. This also could have been done in place in the printTarget rule.

```
derivative
  : ^(DERIVATIVE NAME) {
      out(getFunction($NAME).getDerivative());
    }
  ;
```

---

# Using Generated Classes

▸ Our manually written Processor class

- ▸ uses the generated classes
  - ▸ MathLexer extends Lexer
  - ▸ MathParser extends Parser
  - ▸ MathTree extends TreeParser

  Lexer, Parser and TreeParser extend BaseRecognizer

- ▸ uses other manually written classes
  - ▸ Function
  - ▸ Polynomial
  - ▸ Term
- ▸ supports two modes
  - ▸ batch - see processFile method
  - ▸ interactive - see processInteractive method

# Processor.java

```java
package com.ociweb.math;

import java.io.*;
import java.util.Scanner;
import org.antlr.runtime.*;
import org.antlr.runtime.tree.*;

public class Processor {

    public static void main(String[] args) throws IOException, RecognitionException {
        if (args.length == 0) {
            new Processor().processInteractive();
        } else if (args.length == 1) { // name of file to process was passed in
            new Processor().processFile(args[0]);
        } else { // more than one command-line argument
            System.err.println("usage: java com.ociweb.math.Processor [file-name]");
        }
    }
```

ANTLR 3

# Processor.java …

```java
    private void processFile(String filePath) throws IOException, RecognitionException {
        CommonTree ast = getAST(new FileReader(filePath));
        //System.out.println(ast.toStringTree()); // for debugging
        processAST(ast);
    }

    private CommonTree getAST(Reader reader) throws IOException, RecognitionException {
        MathParser tokenParser = new MathParser(getTokenStream(reader));
        MathParser.script_return parserResult = tokenParser.script(); // start rule method
        reader.close();
        return (CommonTree) parserResult.getTree();
    }

    private CommonTokenStream getTokenStream(Reader reader) throws IOException {
        MathLexer lexer = new MathLexer(new ANTLRReaderStream(reader));
        return new CommonTokenStream(lexer);
    }

    private void processAST(CommonTree ast) throws RecognitionException {
        MathTree treeParser = new MathTree(new CommonTreeNodeStream(ast));
        treeParser.script(); // start rule method
    }
```

ANTLR 3

# Processor.java ...

```java
private void processInteractive() throws IOException, RecognitionException {
    MathTree treeParser = new MathTree(null); // a TreeNodeStream will be assigned later
    Scanner scanner = new Scanner(System.in);

    while (true) {
        System.out.print("math> ");
        String line = scanner.nextLine().trim();
        if ("quit".equals(line) || "exit".equals(line)) break;
        processLine(treeParser, line);
    }
}
```

# Processor.java ...

```java
private void processLine(MathTree treeParser, String line) throws RecognitionException {
    // Run the lexer and token parser on the line.
    MathLexer lexer = new MathLexer(new ANTLRStringStream(line));
    MathParser tokenParser = new MathParser(new CommonTokenStream(lexer));
    MathParser.statement_return parserResult = tokenParser.statement(); // start rule method

    // Use the token parser to retrieve the AST.
    CommonTree ast = (CommonTree) parserResult.getTree();
    if (ast == null) return; // line is empty

    // Use the tree parser to process the AST.
    treeParser.setTreeNodeStream(new CommonTreeNodeStream(ast));
    treeParser.statement(); // start rule method
}

} // end of Processor class
```

We can't create a new instance of MathTree
for each line processed because
it maintains the variable and function Maps.

# ANTLRWorks Debugger

▸ Let's demonstrate using remote debugging
  which is necessary when lexer and parser
  rules are in separate grammar files

  ▸ edit build.properties to include **-debug** in `tool.options`

  ▸ **ant clean run**

    ▸ the run target in build.xml tells it to parse the file "simple.math"

  ▸ start ANTLRWorks

  ▸ open the parser grammar file

  ▸ select Debugger ... Debug Remote...

  ▸ press "Connect" button

  ▸ debug as usual

ANTLR 3

---

# ANTLRWorks Debugger ...

ANTLR 3

# References

- ANTLR
  - http://www.antlr.org
- ANTLRWorks
  - http://www.antlr.org/works
- StringTemplate
  - http://www.stringtemplate.org
  - http://www.codegeneration.net/
    tiki-read_article.php?articleId=65 and 77
- My slides and code examples
  - http://www.ociweb.com/mark - look for "ANTLR 3"

OBJECT COMPUTING, INC.

ANTLR 3