

Geaux

Project 1

Project 1 - Lexer

- A lexer is a program that takes in text as input and spits out “tokens”
- Tokens are just java objects that represents the original text

int => INT

string => STRING

const => CONST

Why is this useful?

- A lexer is supposed to make sense of unstructured text
 - We want to *collect* strings of chars into a singular token
- In the previous cases “int” => new Token(INT)
 - “int” is 3 chars
 - INT is a singular java object
 - ... kind of useful, I guess...
- For String literals however:
 - “Hello World” => STRING_LITERAL(“Hello World”)
 - Way more useful

ANTLR4

- We are using ANTLR4
- ANTLR4 is a Lexer and Parser Generator
- Today, we will use it as a Lexer Generator
- In project 2, we will generate a Parser

```
2 lexer grammar gLexer;
3
4 @header {
5     package Parse.antlr_build;
6 }
7
8
9 @members {
10     StringBuilder sb;
11     private int stringToInt(String target) {
12         // TODO Implement me!
13         return 0;
14     }
15 }
16
17 fragment ALPHA
18     : [A-Za-z]
19 ;
20 fragment DIGIT
21     : [0-9]
22 ;
23
24 ADD
25     : '+'
26 ;
```

```
2 // Generated from Parse/gLexer.g4 by ANTLR 4.13.2
3
4     package Parse antlr_build;
5
6 import org.antlr.v4.runtime.Lexer;
7 import org.antlr.v4.runtime.CharStream;
8 import org.antlr.v4.runtime.Token;
9 import org.antlr.v4.runtime.TokenStream;
10 import org.antlr.v4.runtime.*;
11 import org.antlr.v4.runtime.atn.*;
12 import org.antlr.v4.runtime.dfa.DFA;
13 import org.antlr.v4.runtime.misc.*;
14
15 @SuppressWarnings({"all", "warnings", "unchecked", "unused"})
16 public class gLexer extends Lexer {
17     static { RuntimeMetaData.checkVersion("4.13.2", RuntimeMe
18
19     protected static final DFA[] _decisionToDFA;
20     protected static final PredictionContextCache _sharedCont
21     new PredictionContextCache();
22     public static final int
23     ADD=1;
24     public static String[] channelNames = {
25     "DEFAULT_TOKEN_CHANNEL", "HIDDEN"
26     };
27
28     public static String[] modeNames = {
29     "DEFAULT_MODE"
30     };
31
32     private static String[] makeRuleNames() {
33     return new String[] {
34     "ALPHA", "DIGIT", "ADD"
35     };
36     }
37     public static final String[] ruleNames = makeRuleNames();
38
39     private static String[] makeLiteralNames() {
40     return new String[] {
41     null, "'+'"
42     };
43     }
44     private static final String[] _LITERAL_NAMES = makeLitera
45     private static String[] makeSymbolicNames() {
```

Using ANTLR

- ANTLR is 2 things:
 - A generator:
 - We provide the grammar, and it generates the java file
 - A runtime library
 - The generated file needs certain imports provided by ANTLR's library

```
6 import org.antlr.v4.runtime.Lexer;
7 import org.antlr.v4.runtime.CharStream;
8 import org.antlr.v4.runtime.Token;
9 import org.antlr.v4.runtime.TokenStream;
10 import org.antlr.v4.runtime.*;
11 import org.antlr.v4.runtime.atn.*;
12 import org.antlr.v4.runtime.dfa.DFA;
13 import org.antlr.v4.runtime.misc.*;
```

How do we use ANTLR?

- I have provided two things:
 - A shell script for operating ANTLR
 - An “ANTLR Driver”
 - This is a program we write that “Drives” the lexer
 - Once the lexer java file is generated, we still need to call it!
- I have also included the ANTLR tool, and runtime library
 - It’s in the lib folder!

The project - Geaux

- We are making a lexer for our C-like language: Geaux.
- Our lexer will be similar to a C lexer, with some deviations
- C lexer documentation can be found here:
 - <https://port70.net/~nsz/c/c89/c89-draft.html#3.1>

constant:

~~floating constant~~

~~integer-constant~~

~~enumeration-constant~~

~~character-constant~~

integer-constant:

decimal-constant integer-suffix_{opt}

octal-constant integer-suffix_{opt}

hexadecimal-constant integer-suffix_{opt}

string-literal:

" s-char-sequence_{opt} "

L "s-char-sequence_{opt} "

s-char-sequence:

s-char

s-char-sequence s-char

s-char:

any member of the source chara

escape-sequence

Keywords, Operators, Punctuators

- examples/keywords.g contains all of the ones we will be implementing.
- It includes a couple keywords not in the C standard, such as fun and var.

Preprocessing

- We are not implementing anything from the preprocessing section

Grading

- Your project will be expected to have the same output as my project
 - A binary of my version is provided
- I will be using a test case generator to generate 10 random test cases. You will lose a point for each test case where our outputs don't line up
 - This system is new, so I will be going back over the results manually to verify that they are reasonable
- You have been provided 100 generated tests

Some notes

- Opt means optional
- If you have any question, better to ask than not.

Grading

- Some sticking points:
 - String escape sequences must not show up as “hello \n world”
 - I expect to see:

hello

world

- Just make sure it matches my output

Video

1. How well does your project work?
 - a. What's wrong with it, if anything?
 - b. Show it working on the 10 random test cases
2. What was the experience like doing the project?
 - a. Did you have trouble with ANTLR?
 - b. How long did it take you?
 - c. What strategies for group coding, AI use, time management, etc led to success, and which ones didn't work as well?
3. Not long, 5 minutes

AI

- ANTLR documentation is... not good (imo)
- I used AI to explain to me how ANTLR works
 - It worked incredibly well
- I don't want to make you do the same if you don't want to
- I provided an example grammar that demonstrates what you will need to know to do this project

DEMO