CSE110A: Compilers

April 12, 2023



• Topics:

- Finishing regular expressions
- Using regular expression's in scanners
 - Exact match scanner
 - Start-of-string Scanner
 - Named group matcher

Announcements

- Not feeling well 😣
 - Remote class today and possibly Friday
- Homework has not been released. Trying to get it out today or tomorrow
- We have office hours this week; come see us!
 - Mine will be remote tomorrow

Announcements

• Other office hours:

Neal is available on Tuesdays and Thursdays, 6:30PM - 7:30PM, virtual.

Announcements

- Please enroll in Piazza!
 - Only ~68 students are enrolled. There are ~100 students in the class. People are asking good clarification questions that will help you
- Threads on docker questions, etc.

Quiz

Integer RE

The following RE is a good candidate for non-negative integers: "[0-9]*"

 \bigcirc True

 \bigcirc False

Integer RE

The following RE is a good candidate for non-negative integers: "[0-9]*"

⊖ True

 \bigcirc False

Does the "" match the RE?

Fundamental RE operators

All regular expressions can be expressed in terms of concatenation or choice operators

⊖ True

⊖ False

Fundamental RE operators

- Fundamental RE operators are:
 - Concatenate: put the regexes next to each other
 - "|" : Choice: one or the other
 - "*": Repeat: 0 or more copies
- Practically:
 - a* roughly is the same as "" | "a" | "aa" | "aaa" ...
 - in theory, REs can accept strings of arbitrary length (not infinite strings though).
 - in practice, strings have a reasonable bound. Repeat (*) is a good abstraction though!

RE examples

which of the following strings	s do NOT match ac* b	*
--------------------------------	----------------------	---

"" (empty string)		
🗌 ab		
🗌 acac		
🗌 bbb		

RE examples

ac*|b*

• ""

Let's work through them

- "ab"
- "acac"
- "acccc"
- "bbb"

RE experiences

Have you used regular expressions before? If so, in what language or tool did you use them, and for what application?

Review

• Some syntactic sugar and useful interfaces

- strict repeat operator: +
- one or more repeats (the * operator is 0 or more repeats)
- derivation: "r+" = "rr*"
- Let's revisit binary numbers and decimal numbers

"(0|1)+"

- Ranges:
 - digits [0-9]
 - alpha [a-z], [A-Z]
- Derivation: [0-9] = "1|2|3|4|5|6|7|8|9"
- Lets try C style IDs: "[a-zA-Z][0-9a-zA-Z]*"
- Hexadecimal numbers: "0x[0-9a-fA-F]+"

- optional operator ?
 - optional characters
- "r?" = "|r"
- Example: "ab?"
- Let's do simple floating point numbers: "[0-9]+(\.[0-9]+)?"

- any character '.'
- example using email (this is probably too general!)
- ".*@.*\.com"

Using REs

- What if we want either the domain or user name from the email?
- We can use groups!
 - use ()s to deliminate groups
- "(.*)@(.*\.com)"
- Index the resulting object with [1] and [2] to get to the user name and domain respectively

Using REs

- you can give groups id names rather than using indices
- "<mark>(?P<name></mark>.+)@<mark>(?P<domain></mark>.+\.com)"

Review

• Why do we want REs?

Naïve Scanner

simple string stream, peek/eat model

class NaiveScanner:

ID	=	[characters]
NUM	=	[numbers]
ASSIGN	=	<i>"="</i>
PLUS	=	"+"
MULT	=	<i>''</i> * <i>''</i>
IGNORE	=	[""]

Shortcomings of Naïve scanner

- IDs with numbers in them?
 - x1, y1, etc.
 - how would you solve?
- Numbers with a decimal point in them?
 - 4.5, 9999.99998
 - how would you solve this?
- Two character operators:
 - ++, +=
 - how would you solve this?

We need a new token definition language

Can we express these tokens using REs?

- ARTICLE
- NOUN
- VERB
- ADJECTIVE

- = {The, A, My, Your}
- = {Dog, Car, Computer}
 - = {Ran, Crashed, Accelerated}
- = {Purple, Spotted, Old}

Can we express these tokens using REs? Yes!



- NOUN
- VERB
- ADJECTIVE

- = "The|A|Mine|Your"
- = "Dog|Car|Computer"
- = "Ran|Crashed|Accelerated"
- = "Purple|Spotted|Old"

Let's write our tokens as regular expressions

• For our simple programming language

ID	=	[characters]
NUM	=	[numbers]
ASSIGN	=	"="
PLUS	=	"+"
MULT	=	<i>"</i> * <i>"</i>
IGNORE	=	["","\n"]

Let's write our tokens as regular expressions

• For our simple programming language

ID	=	[characters]
NUM	=	[numbers]
ASSIGN	=	<i>"="</i>
PLUS	=	"+"
MULT	=	<i>''</i> * <i>''</i>
IGNORE	=	["","\n"]

ID =
$$"[a-z]+"$$

NUM = $"[0-9]+"$
ASSIGN = $"="$
PLUS = $"+"$
MULT = $"*"$
IGNORE = $" \mid \backslash n"$

Some benefits of REs? Let's try adding some extensions:

Let's write our tokens as regular expressions

• For our simple programming language

ID	=	[characters]		
NUM	=	[numbers]		
ASSIGN	=	<i>"="</i>		
PLUS	=	"+"		
MULT	=	<i>"</i> * <i>"</i>		
IGNORE	=	["","\n"]		

ID = "
$$[a-z]+"$$

NUM = " $[0-9]+"$
ASSIGN = "="
PLUS = "+"
MULT = "*"
IGNORE = " $|n"$

Some benefits of REs? Let's try adding some extensions: * increment operator?

* digits in IDs?

Finishing up last lecture

• A few final thoughts:

RE examples

- What can REs not do?
- Nested structures, such as parathesis matching:
 - Try doing arithmetic expressions
 - You will not be able to match ()s
- Classical example: REs cannot capture same number of repeats:
 - A{N}B{N}
- REs cannot parse HTML!!!
 - One of the most upvoted answers on stackoverflow!
 - <u>https://stackoverflow.com/questions/1732348/regex-match-open-tags-except-xhtml-self-contained-tags/1732454#1732454</u>

How to implement an RE matcher?

- Overview: first you have to parse the RE...
 - Chicken and egg problem
 - The language of REs is not a regular language. It is context free (because it has ()s)
 - But once you can parse the RE, there are several options

How to implement an RE matcher?

- parsing with derivatives
 - We discuss this in CSE211
 - Elegant solution, but difficult to make fast
- Convert to an automata
 - Learn more about this CSE103
 - A cool website
 - <u>https://ivanzuzak.info/noam/webapps/fsm_simulator/</u>

New material for today

- Using RE matchers to build scanners
 - Exact match (EM) scanners
 - Start-of-string (SOS) scanners
 - named group (NG) scanners

New material for today

• Using RE matchers to build scanners

- Exact match (EM) scanners
- Start-of-string (SOS) scanners
- named group (NG) scanners

The problem

 How do we move from an RE match to performing lexical analysis on a string

ID	=	"[a-z]+"
NUM	=	"[0-9]+"
ASSIGN	=	<i>"="</i>
PLUS	=	"+"
MULT	=	<i>''</i> * <i>''</i>
IGNORE	=	" \\n"
SEMI	=	// • // /

"variable = 50 + 30 * 20;"

The problem

 How do we move from an RE match to performing lexical analysis on a string

ID	=	"[a-z]+"
NUM	=	"[0-9]+"
ASSIGN	=	<i>"="</i>
PLUS	=	"+"
MULT	=	"*"
IGNORE	=	" \n"
SEMI	=	// • //

"variable = 50 + 30 * 20;"

[(ID, "variable"), (ASSIGN, "="), (NUM, "50"), (PLUS, "+"), (NUM, "30"), (MULT, "*"), (NUM, "20"), (SEMI, ";")]
• How do we move from an RE match to performing lexical analysis on a string

ID	=	"[a-z]+"
NUM	=	"[0-9]+"
ASSIGN	=	"="
PLUS	=	"+"
MULT	=	<i>"</i> * <i>"</i>
IGNORE	=	" \n"
SEMI	=	// • // /

Do these match?

• How do we move from an RE match to performing lexical analysis on a string

ID	=	"[a-z]+"
NUM	=	"[0-9]+"
ASSIGN	=	<u>"="</u>
PLUS	=	<mark>"+"</mark>
MULT	=	<u>" * "</u>
IGNORE	=	" \n"
SEMI	=	"•" 7

Do any of the tokens match?

• How do we move from an RE match to performing lexical analysis on a string

ID	=	"[a-z]+"
NUM	=	<mark>"[0-9]+"</mark>
ASSIGN	=	<u>"="</u>
PLUS	=	<mark>"+"</mark>
MULT	=	<mark>" * "</mark>
IGNORE	=	" \n"
SEMI	=	"•"

What if we start "peeking" characters

• How do we move from an RE match to performing lexical analysis on a string

ID	=	"[a-z]+"
NUM	=	"[0-9]+"
ASSIGN	=	<i>"="</i>
PLUS	=	"+"
MULT	=	<i>''</i> * <i>''</i>
IGNORE	=	" \n"
SEMI	=	" • " •

Match!

• How do we move from an RE match to performing lexical analysis on a string

ID	=	"[a-z]+"
NUM	=	"[0-9]+"
ASSIGN	=	<i>"="</i>
PLUS	=	"+"
MULT	=	<i>''</i> * <i>''</i>
IGNORE	=	" \n"
SEMI	=	// • // /

Match! <mark>(ID, "v")</mark>

"variable = 50 + 30 * 20;"

• How do we move from an RE match to performing lexical analysis on a string

ID	=	"[a-z]+"
NUM	=	"[0-9]+"
ASSIGN	=	<i>"="</i>
PLUS	=	"+"
MULT	=	"*"
IGNORE	=	" \n"
SEMI	=	// • // /

Match! (ID, "v") but what is the issue?

 How do we move from an RE match to performing lexical analysis on a string

ID	=	"[a-z]+"
NUM	=	"[0-9]+"
ASSIGN	=	<i>"="</i>
PLUS	=	"+"
MULT	=	<i>''</i> * <i>''</i>
IGNORE	=	" \n"
SEMI	=	// • // /

Match! (ID, "v") but what is the issue? Not the longest match

• How do we move from an RE match to performing lexical analysis on a string

So what's our strategy?

ID	=	"[a-z]+"
NUM	=	"[0-9]+"
ASSIGN	=	<i>"="</i>
PLUS	=	"+"
MULT	=	"*"
IGNORE	=	" \\n"
SEMI	=	// • // /

"variable = 50 + 30 * 20;"

New material for today

- Using RE matchers to build scanners
 - Exact match (EM) scanners
 - Start-of-string (SOS) scanners
 - named group (NG) scanners

• Start with the whole string, remove one character at the end until a match is found. Then return the lexeme

ID	=	"[a-z]+"
NUM	=	"[0-9]+"
ASSIGN	=	<i>"="</i>
PLUS	=	"+"
MULT	=	<i>''</i> * <i>''</i>
IGNORE	=	" \\n"
SEMI	=	// • // /

"variable = 50 + 30 * 20;"

• Start with the whole string, remove one character at the end until a match is found. Then return the lexeme

ID	=	"[a-z]+"
NUM	=	"[0-9]+"
ASSIGN	=	<i>"="</i>
PLUS	=	"+"
MULT	=	<i>''</i> ★ <i>''</i>
IGNORE	=	" \n"
SEMI	=	// • // /

start with the whole string

• Start with the whole string, remove one character at the end until a match is found. Then return the lexeme

ID = "[a-z]+"
NUM = "[0-9]+"
ASSIGN = "="
PLUS = "+"
MULT = "*"
IGNORE = " \\n"
SEMI = ";"

Try to match with all the tokens

start with the whole string

• Start with the whole string, remove one character at the end until a match is found. Then return the lexeme

Try to match with all the tokens. No match.

ID	=	"[a-z]+"
NUM	=	" [0-9]+"
ASSIGN	=	<u>"="</u>
PLUS	=	<mark>"+"</mark>
MULT	=	" * "
IGNORE	=	" \n"
SEMI	=	" • " 7

start with the whole string

• Start with the whole string, remove one character at the end until a match is found. Then return the lexeme

Try to match with all the tokens. No match.

ID	=	"[a-z]+"
NUM	=	" [0-9]+"
ASSIGN	=	<u>"="</u>
PLUS	=	<mark>"+"</mark>
MULT	=	" * "
IGNORE	=	" \n"
SEMI	=	" <mark>''</mark>

Try with one character chopped from back

• Start with the whole string, remove one character at the end until a match is found. Then return the lexeme

Try to match with all the tokens. No match.

So on

ID	=	"[a-z]+"
NUM	=	"[0-9]+"
ASSIGN	=	<u>"="</u>
PLUS	=	<mark>"+"</mark>
MULT	=	" * "
IGNORE	=	" \n"
SEMI	=	" <mark>,</mark> "

"variable = 50 + 30 * 20;"

• Start with the whole string, remove one character at the end until a match is found. Then return the lexeme

Try to match with all the tokens. No match.

So on

ID	=	"[a-z]+"
NUM	=	" [0-9]+ "
ASSIGN	=	<u>"="</u>
PLUS	=	<mark>"+"</mark>
MULT	=	" * "
IGNORE	=	" \n"
SEMI	=	" • "

"variable = 50 + 30 * 20;"

• Start with the whole string, remove one character at the end until a match is found. Then return the lexeme

Try to match with all the tokens. No match.

ID	=	"[a-z]+"
NUM	=	<u>"[0-9]+"</u>
ASSIGN	=	<u>"="</u>
PLUS	=	<mark>"+"</mark>
MULT	=	" * "
IGNORE	=	" \n"
SEMI	=	" • " 7

Where do find a match?



• Start with the whole string, remove one character at the end until a match is found. Then return the lexeme

ID	=	"[a-z]+"
NUM	=	"[0-9]+"
ASSIGN	=	<u>"="</u>
PLUS	=	<mark>"+"</mark>
MULT	=	<u>" * "</u>
IGNORE	=	" \n"
SEMI	=	"•" 7

we can match id

at this point



• Start with the whole string, remove one character at the end until a match is found. Then return the lexeme

we can match id

ID	=	"[a-z]+"
NUM	=	" [0-9]+ "
ASSIGN	=	<u>"="</u>
PLUS	=	" + "
MULT	=	" * "
IGNORE	=	" \n"
SEMI	=	"•"

at this point

Return the lexeme

• Start with the whole string, remove one character at the end until a match is found. Then return the lexeme

ID	=	"[a-z]+"
NUM	=	<mark>"[0-9]+"</mark>
ASSIGN	=	<u>"="</u>
PLUS	=	<mark>"+"</mark>
MULT	=	" * "
IGNORE	=	" \n"
SEMI	=	"•"

Chop the string

• Start with the whole string, remove one character at the end until a match is found. Then return the lexeme

ID	=	"[a-z]+"
NUM	=	" [0-9]+"
ASSIGN	=	<u>"="</u>
PLUS	=	<mark>"+"</mark>
MULT	=	" * "
IGNORE	=	" \n"
SEMI	=	"•" 7

Start the process over

• Start with the whole string, remove one character at the end until a match is found. Then return the lexeme

ID	=	"[a-z]+"
NUM	=	" [0-9]+ "
ASSIGN	=	<u>"="</u>
PLUS	=	<mark>"+"</mark>
MULT	=	" * "
IGNORE	=	" \n"
SEMI	=	"•" 7

Start the process over Where is our next match?

code for exact match scanner

• Provided in your homework

- Pros
- Cons

- Pros
 - Uses an exact RE matcher. Many RE match algorithms are exact!
- Cons
 - SLOW! Each lexeme requires many many many calls to each RE match!

New material for today

- Using RE matchers to build scanners
 - Exact match (EM) scanners
 - Start-of-string (SOS) scanners
 - named group (NG) scanners

• We will use a new RE match function

re. fullmatch(pattern, string, flags=0) ¶

If the whole *string* matches the regular expression *pattern*, return a corresponding match object. Return None if the string does not match the pattern; note that this is different from a zero-length match.

re.match(pattern, string, flags=0)

If zero or more characters at the beginning of *string* match the regular expression *pattern*, return a corresponding match object. Return None if the string does not match the pattern; note that this is different from a zero-length match.

• The match API gives us a match starting at the beginning of the string

ID	=	"[a-z]+"
NUM	=	"[0-9]+"
ASSIGN	=	<i>"="</i>
PLUS	=	"+"
MULT	=	" * "
IGNORE	=	" \n"
SEMI	=	// • // /

"variable = 50 + 30 * 20;"

• The match API gives us a match starting at the beginning of the string

ID	=	"[a-z]+"
NUM	=	"[0-9]+"
ASSIGN	=	<u>"="</u>
PLUS	=	<mark>"+"</mark>
MULT	=	" * "
IGNORE	=	" \n"
SEMI	=	" • "

Feed full string into each token definition

• The match API gives us a match starting at the beginning of the string

ID	=	"[a-z]+"
NUM	=	" [0-9]+"
ASSIGN	=	<u>"="</u>
PLUS	=	<mark>"+"</mark>
MULT	=	" * "
IGNORE	=	" \n"
SEMI	=	" <mark>'''</mark>

Feed full string into each token definition

We get 1 match. We can return the lexeme

• The match API gives us a match starting at the beginning of the string

ID	=	"[a-z]+"
NUM	=	"[0-9]+"
ASSIGN	=	<u>"="</u>
PLUS	=	<mark>"+"</mark>
MULT	=	" * "
IGNORE	=	" \n"
SEMI	=	// • // /

Chop the string

We get 1 match. We can return the lexeme

• The match API gives us a match starting at the beginning of the string

ID	=	"[a-z]+"
NUM	=	<mark>"[0-9]+"</mark>
ASSIGN	=	<u>"="</u>
PLUS	=	<mark>"+"</mark>
MULT	=	" * "
IGNORE	=	" \n"
SEMI	=	"•"

Chop the string

We get 1 match. We can return the lexeme

• The match API gives us a match starting at the beginning of the string

ID	=	"[a-z]+"
NUM	=	" [0-9]+"
ASSIGN	=	<u>"="</u>
PLUS	=	<mark>"+"</mark>
MULT	=	" * "
IGNORE	=	" \n"
SEMI	=	"•"

What about the next one

• The match API gives us a match starting at the beginning of the string

ID	=	"[a-z]+"
NUM	=	"[0-9]+"
ASSIGN	=	<u>"="</u>
PLUS	=	<mark>"+"</mark>
MULT	=	" * "
IGNORE	=	" \n"
SEMI	=	" • "

What about the next one

1 match: IGNORE

• The match API gives us a match starting at the beginning of the string

ID	=	"[a-z]+"
NUM	=	" [0-9]+"
ASSIGN	=	<u>"="</u>
PLUS	=	<mark>"+"</mark>
MULT	=	" * "
IGNORE	=	" \n"
SEMI	=	"•"

Chop the string

1 match: IGNORE

• The match API gives us a match starting at the beginning of the string

ID	=	"[a-z]+"
NUM	=	<mark>"[0-9]+"</mark>
ASSIGN	=	<u>"="</u>
PLUS	=	<mark>"+"</mark>
MULT	=	" * "
IGNORE	=	" \n"
SEMI	=	"•"

Chop the string

1 match: IGNORE
Consideration

How to scan this string?

"CSE110A"

LETTERS	=	"[A-Z]+"
NUM	=	"[0-9]+"
CLASS	=	"CSE110A"

Consideration

How to scan this string?

Try to match on each token

"CSE110A"

LETTERS	=	"[A-Z]+"
NUM	=	" [0-9]+ "
CLASS	=	"CSE110A"

Consideration

How to scan this string?

Try to match on each token

"CSE110A"

LETTERS	=	"[A-Z]+"
NUM	=	"[0-9]+"
CLASS	=	"CSE110A"

Two matches: LETTERS: "CSE" CLASS: "CSE110A"

Which one do we choose?

Consideration

How to scan this string?

Try to match on each token

"CSE110A"

LETTERS	=	"[A-Z]+"
NUM	=	"[0-9]+"
CLASS	=	"CSE110A"

Two matches: LETTERS: "CSE" CLASS: "CSE110A"

> Which one do we choose? The longest one!

After each pass through token REs we have to measure match length

Consideration

How to scan this string?

Try to match on each token

"CSE110A"

LETTERS	=	"[A-Z]+"
NUM	=	"[0-9]+"
CLASS	=	"CSE110A"

Two matches: LETTERS: "CSE" CLASS: "CSE110A"

> Which one do we choose? The longest one!

After each pass through token REs we have to measure match length

Why didn't we have to do this for the exact match Scanner?

• One more consideration

Within 1 RE, how does this match?

"CSE110A"

CLASS = "CSE | 110A | CSE110A"

• One more consideration

Within 1 RE, how does this match?

"CSE110A"

CLASS = "CSE | 110A | CSE110A"

Returns "CSE", but this isn't what we want!!!

• One more consideration

Within 1 RE, how does this match?

"CSE110A"

CLASS = "CSE | 110A | CSE110A"

Returns "CSE", but this isn't what we want !!!

When using the SOS Scanner: A token definition either should not:

- contain choices where one choice is a prefix of another
- order choices such that the longest choice is the first one

• One more consideration

Within 1 RE, how does this match?

"CSE110A"

Returns "CSE", but this isn't what we want !!!

When using the SOS Scanner: A token definition either should not:

- contain choices where one choice is a prefix of another
- order choices such that the longest choice is the first one

CLASS = "CSE110A | 110A | CSE"

- Pros
- Cons

- Pros
 - Much faster than EM scanner. Only 1 call to each RE per token() call
- Cons
 - Depends on an efficient implementation of match()
 - Typically provided in most RE libraries (for this exact reason)
 - Requires some care in token definitions and prefixes

New material for today

- Using RE matchers to build scanners
 - Exact match (EM) scanners
 - Start-of-string (SOS) scanners
 - named group (NG) scanners

We're going to optimize this to 1 RE call! It can really help if you have many tokens

- Pros
 - Much faster than EM scanner. Only 1 call to each RE per token() call
- Cons
 - Depends on an efficient implementation of match()
 - Typically provided in most RE libraries (for this exact reason)
 - Requires some care in token definitions and prefixes

• We will still use the match API call

re. fullmatch(pattern, string, flags=0) ¶

If the whole *string* matches the regular expression *pattern*, return a corresponding match object. Return None if the string does not match the pattern; note that this is different from a zero-length match.

re.match(pattern, string, flags=0)

If zero or more characters at the beginning of *string* match the regular expression *pattern*, return a corresponding match object. Return None if the string does not match the pattern; note that this is different from a zero-length match.

- Start out with token definitions
- Merge them into one RE definition

ID	=	"[a-z]+"	
NUM	=	"[0-9]+"	SINGLE_RE =
ASSIGN	=	<i>"="</i>	
PLUS	=	"+"	
MULT	=	<i>"</i> * <i>"</i>	
IGNORE	=	" \n"	
SEMI	=	// • // /	

- Start out with token definitions
- Merge them into one RE definition

ID	=	"[a-z]+"
NUM	=	"[0-9]+"
ASSIGN	=	<i>"="</i>
PLUS	=	"+"
MULT	=	<i>''</i> * <i>''</i>
IGNORE	=	" \\n"
SEMI	=	// • //

SINGLE_RE =
$$"[a-z]+"$$

- Start out with token definitions
- Merge them into one RE definition

ID	=	"[a-z]+"
NUM	=	"[0-9]+"
ASSIGN	=	<i>"="</i>
PLUS	=	"+"
MULT	=	<i>''</i> * <i>''</i>
IGNORE	=	" \\n"
SEMI	=	// • //

SINGLE_RE = "
$$([a-z]+)$$
"

- Start out with token definitions
- Merge them into one RE definition

SINGLE_RE = "(
$$[a-z]+$$
) ($[0-9]+$)"

- Start out with token definitions
- Merge them into one RE definition

- Start out with token definitions
- Merge them into one RE definition

Give each group a name corresponding to its token

ID	=	"[a-z]+"
NUM	=	"[0-9]+"
ASSIGN	=	<i>"="</i>
PLUS	=	"+"
MULT	=	"*"
IGNORE	=	" \\n"
SEMI	=	// • // /

- Start out with token definitions
- Merge them into one RE definition

It's a giant RE, but you can construct it automatically

to implement token()

Try to match the whole string to the single RE

"variable = 50 + 30 * 20;"

to implement token()

Try to match the whole string to the single RE

"variable = 50 + 30 * 20;"

Check the group dictionary in the result

{"ID" : "variable" "NUM" : None "ASSIGN" : None "PLUS" : None "MULT" : None "IGNORE" : None "SEMI" : None}

"variable = 50 + 30 * 20;"

Try to match the whole string to the single RE

NG Scanner

to implement token()

 $\frac{\text{SINGLE RE}}{\text{SINGLE RE}} = "(?P < ID > [a - z] +)$ (?P<NUM>[0-9]+) (?P<ASSIGN>=) (?P<PLUS>+) (?P<MULT>*) (?P < IGNORE > | n)(?P<SEMI>;)"

'"ID" : "variable" "NUM" : None "ASSIGN" : None "PLUS" : None "MULT" : None "IGNORE" : None "SEMI" : None}

Try to match the whole string to the single RE

to implement token()

Try to match the whole string to the single RE

"variable = 50 + 30 * 20;"

{ "ID"	:	"variable"
"NUM"	:	None
"ASSIGN"	:	None
"PLUS"	:	None
"MULT"	:	None
"IGNORE"	:	None
"SEMI"	:	None}

Return the lexeme (ID, "variable")

Return the lexeme (ID, "variable")

{ "ID"	:	"variable"
"NUM"	:	None
"ASSIGN"	:	None
"PLUS"	:	None
"MULT"	:	None
"IGNORE"	:	None
"SEMI"	:	None}

chop!

(?P<MULT>*)

(?P<SEMI>;)"

(?P < IGNORE > | n)

to implement token()

to implement token()

chop!

• Recall from SOS scanner:

LETTERS

NUM

CLASS

How to scan this string?

=	"[A-Z]+"
=	"[0-9]+"

= "CSE110A"

"CSE110A"

• Convert to a single RE

How to scan this string?

SINGLE_RE = "
 (?P<LETTERS>([A-Z]+)|
 (?P<NUM>([0-9]+)|
 (?P<CLASS>CSE110A)"

"CSE110A"

• Convert to a single RE

How to scan this string?

"CSE110A"

SINGLE_RE = "
(?P<LETTERS>([A-Z]+)|
(?P<NUM>([0-9]+)|
(?P<CLASS>CSE110A)"

What do we think the dictionary will look like?

• Convert to a single RE

How to scan this string?

SINGLE_RE = "
(?P<LETTERS>([A-Z]+)|
(?P<NUM>([0-9]+)|
(?P<CLASS>CSE110A)"

"CSE110A"

{"LETTERS"	:	"CSE"
"NUM"	:	None
"CLASS"	:	None
}		

• Convert to a single RE

SINGLE_RE = "
(?P<LETTERS>([A-Z]+)|
(?P<NUM>([0-9]+)|
(?P<CLASS>CSE110A)"

What does this mean?

Tokens should not contain prefixes of each other

OR

• Tokens that share a common prefix should be ordered such that the longer token comes first

"CSE1	L10A"
-------	-------

{"LETTERS"	:	"CSE"
"NUM"	:	None
"CLASS"	:	None
}		

• Careful with these tokens

Ensure that you provide them in the right order so that the longer one is first!

- Pros
- Cons

- Pros
 - FAST! Only 1 RE call per token()
- Cons
 - Requires a named group RE library
 - inter-token interactions need to be considered
Scanners we have discussed

- Naïve Scanner
- RE based scanners
 - Exact match (EM) scanners
 - Start-of-string (SOS) scanners
 - named group (NG) scanners

Which one to use? Complex decision with performance, expressivity, and token requirements

On Friday

- We will discuss token actions and how to use them to implement keywords and line numbers
- We will discuss a classic scanner generator: lex
- See you on Friday!