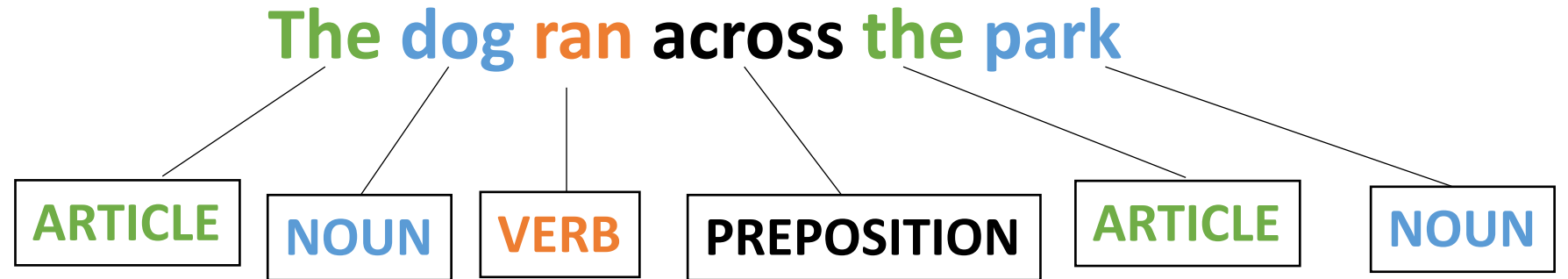


CSE110A: Compilers

April 8, 2023



- **Topics:**

- *Lexical Analysis:*
 - Shortcomings of naïve scanner
- *Regular expressions:*
 - Recursive definition
 - Syntactic sugar
 - groups

Announcements

- HW 1 will be released by midnight tonight
 - You have what you need to start working on part 1
 - You will have what you need for part 2 after Wednesday
 - You will have what you need for part 3 after Friday
- The TAs are trying a new gradescope approach
 - Should make life easier for everyone
 - Let us know if there are issues
- Due one week from today (by midnight)
- We will have office hours this week, come see us!

Announcements

- TA Tutor hours

<https://sorensenucsc.github.io/CSE110A-sp2024/overview.html#office-hours>

TA Office Hours:

Day and Time	TA	location/Zoom Link
Tuesday 3 - 4 PM	Rithik	BE-151
Wednesday 4 - 5 PM	Sakshi	E2-216
Thursday (TBD)	Sakshi	E2-216

Mentoring Hours:

Day and Time	Mentor	location/Zoom Link
Monday 11 AM - 12 PM	Kaushal	Zoom
Monday 12:30 - 1:30 PM	Ryan	in-person (location TBD)
Monday 4 - 5 PM	Ananth	in-person (location TBD)
Tuesday 5 - 6 PM	Kaushal	in-person (location TBD)
Wednesday 1 - 2 PM	Jack	Zoom
Thursday 11:30 AM - 12:30 PM	Ryan	Zoom
Friday 11 AM - 12 PM	Jack	in-person (location TBD)
Friday 4 - 5 PM	Ananth	Zoom

Quiz

Scanner API

The scanner member function "token" returns a list of the tokens that can recognize

☐ True

☐ False

Programs for Lexical Analysis

Scanner (sometimes called lexer)

Defined by a list of tokens and definitions:

- ARTICLE

- NOUN

- VERB

- ADJECTIVE

Tokens

= {The, A, My, Your}

= {Dog, Car, Computer}

= {Ran, Crashed, Accelerated}

= {Purple, Spotted, Old}

Tokens Definitions

Original program:

Lex

[https://en.wikipedia.org/wiki/Lex_\(software\)](https://en.wikipedia.org/wiki/Lex_(software))

Popular implementations

Flex

Scanner API

```
# Constructor, generates a Scanner  
s = ScannerGenerator(tokens)  
  
# The string we want to do  
# lexical analysis on  
s.input("My Old Computer Crashed")  
  
# Returns the next lexeme  
s.token()
```

```
> s = ScannerGenerator(tokens)
> s.input("My Old Computer Crashed")
> s.token()
(ARTICLE, "My")
> s.token()
(ADJECTIVE, "Old")
> s.token()
(NOUN, "Computer")
> s.token()
(VERB, "Crashed")
> s.token()
None
```


Scanning vs. Parsing

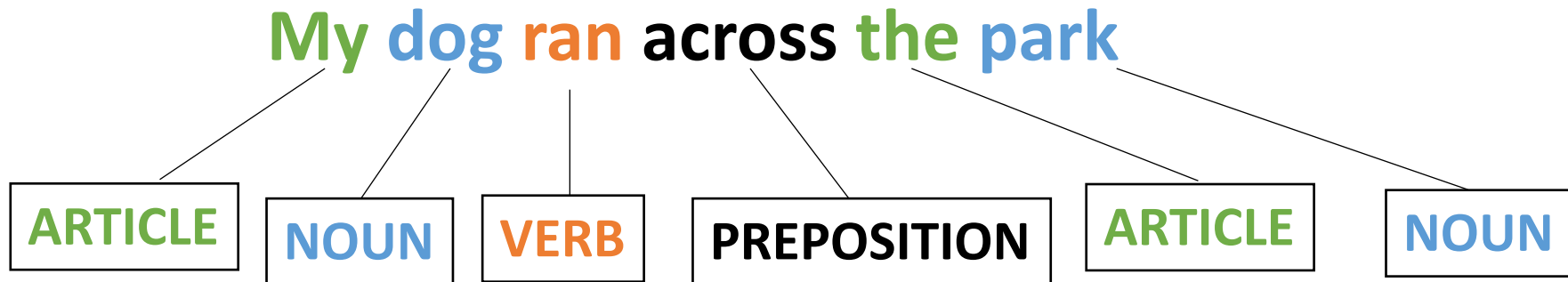
A scanner should make sure that the sequence of lexemes is valid, e.g. the scanner should make sure two numbers are separated by a valid operator.

☐ True

☐ False

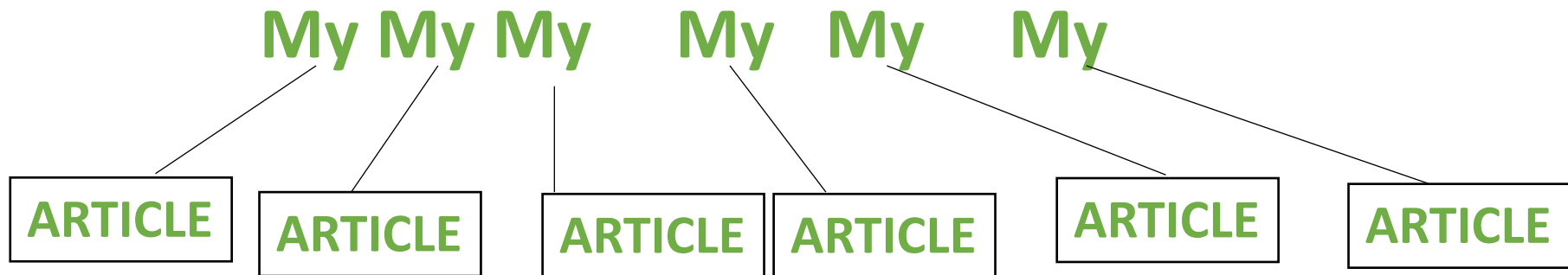
Parsing is the first step in a compiler

- How do we parse a sentence in English?



Parsing is the first step in a compiler

- How do we parse a sentence in English?



Lexical analysis doesn't care about the order of tokens. Just so long as there are valid tokens.

Programs for Lexical Analysis

Scanner (sometimes called lexer)

Defined by a list of tokens and definitions:

- ARTICLE

- NOUN

- VERB

- ADJECTIVE

= {The, A, My, Your}

= {Dog, Car, Computer}

= {Ran, Crashed, Accelerated}

= {Purple, Spotted, Old}

Tokens

Tokens Definitions

Original program:

Lex

[https://en.wikipedia.org/wiki/Lex_\(software\)](https://en.wikipedia.org/wiki/Lex_(software))

Popular implementations

Flex

Parsing is the first step in a compiler

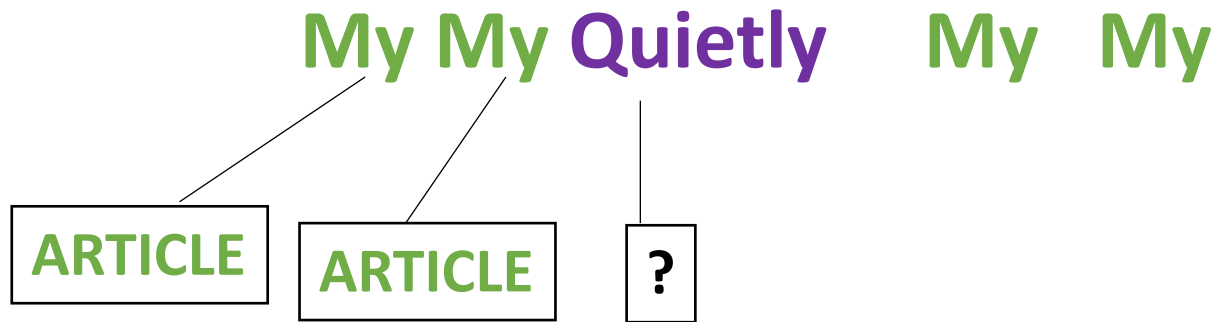
- How do we parse a sentence in English?

My My Quietly My My My

What happens here?

Parsing is the first step in a compiler

- How do we parse a sentence in English?



What happens here?

Scanner error here. Many scanners stop and report the error location

Parsing is the first step in a compiler

- How do we parse a sentence in English?



What happens here?

Scanner error here. Some scanners try to recover and keep going (difficult, and requires ad hoc rules)

Scanning vs. Parsing

A scanner should make sure that the sequence of lexemes is valid, e.g. the scanner should make sure two numbers are separated by a valid operator.

☐ True

☐ False

False! The order of tokens will be checked by the parser later on!

Scanning a simple PL statement

How many lexemes do you think the following statement should have?

```
for (int i = 0; i <=5; i++)
```

What lexemes do you think they should be?

Scanning a simple PL statement

```
for (int i = 0; i <= 5; i++)
```

Scanning a simple PL statement

```
for (int i = 0; i <= 5; i++)
```

```
[ (ID, "for"),      (PAR, "("), (ID, "int"), (ID, "i"),  
  (ASSIGN, "="),    (NUM, "0"),  (SEMI, ";"), (ID, "i"),  
  (LE, "<="),        (NUM, "5"),  (SEMI, ";"), (ID, "i"),  
  (INCR, "++"),     (PAR, ")") ]
```

Scanning a simple PL statement

```
for (int i = 0; i <= 5; i++)
```

```
[ (ID, "for"), (PAR, "("), (ID, "int"), (ID, "i"),  
  (ASSIGN, "="), (NUM, "0"), (SEMI, ";"), (ID, "i"),  
  (LE, "<="), (NUM, "5"), (SEMI, ";"), (ID, "i"),  
  (INCR, "++"), (PAR, ")") ]
```

Why not: "<" and "=" separately?

Scanning a simple PL statement

```
for (int i = 0; i <= 5; i++)
```

```
[ (ID, "for"), (PAR, "("), (ID, "int"), (ID, "i"),  
  (ASSIGN, "="), (NUM, "0"), (SEMI, ";"), (ID, "i"),  
  (LE, "<="), (NUM, "5"), (SEMI, ";"), (ID, "i"),  
  (INCR, "++"), (PAR, ")") ]
```

Should these be the same token?

Scanning a simple PL statement

```
for (int i = 0; i <= 5; i++)
```

```
[ (ID, "for"), (LPAR, "("), (ID, "int"), (ID, "i"),  
  (ASSIGN, "="), (NUM, "0"), (SEMI, ";"), (ID, "i"),  
  (LE, "<="), (NUM, "5"), (SEMI, ";"), (ID, "i"),  
  (INCR, "++"), (RPAR, ")") ]
```

Should these be the same token? Probably not

Review

Naïve implementation

- A scanner that implements

ID	=	[characters]
NUM	=	[numbers]
ASSIGN	=	"="
PLUS	=	"+"
MULT	=	"*"
IGNORE	=	[" "]

Naïve implementation

Building block:

```
class StringStream:
    def __init__(self, input_string):
        self.string = input_string

    def is_empty(self):
        return len(self.string) == 0

    def peek_char(self):
        if not self.is_empty():
            return self.string[0]
        return None

    def eat_char(self):
        self.string = self.string[1:]
```

Naïve implementation

First step in implementing the scanner

```
class NaiveScanner:

    def __init__(self, input_string):
        self.ss = StringStream(input_string)

    def token(self):

        while self.ss.peek_char() in IGNORE:
            self.ss.eat_char()

        if self.ss.is_empty():
            return None
```

ID	=	[characters]
NUM	=	[numbers]
ASSIGN	=	"="
PLUS	=	"+"
MULT	=	"*"
IGNORE	=	[" "]

Naïve implementation

First step in implementing the scanner

```
class NaiveScanner:

    def token(self):
        ...
        if self.ss.peek_char() == "+":
            value = self.ss.peek_char()
            self.ss.eat_char()
            return ("ADD", value)

        if self.ss.peek_char() == "*":
            value = self.ss.peek_char()
            self.ss.eat_char()
            return ("MULT", value)
```

ID	=	[characters]
NUM	=	[numbers]
ASSIGN	=	"="
PLUS	=	"+"
MULT	=	"*"
IGNORE	=	[" "]

Naïve implementation

First step in implementing the scanner

```
class NaiveScanner:

    def token(self):
        ...
        if self.ss.peek_char() in NUMS:
            value = ""
            while self.ss.peek_char() in NUMS:
                value += self.ss.peek_char()
                self.ss.eat_char()
            return ("NUM", value)
```

ID	=	[characters]
NUM	=	[numbers]
ASSIGN	=	"="
PLUS	=	"+"
MULT	=	"*"
IGNORE	=	[" "]

Schedule

- Naïve Parser:
 - Code demo and discussion
- Regular expressions

Code Demo

Shortcomings of Naïve scanner

- Any thoughts?

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?

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?

*Things get really hacky
really quickly!*

*Creates
a bad design that is
not easily extended
or maintained*

How do we solve this?

A new token definition language:

- **Regular expressions**
- Tokens will be defined using regular expressions
- Scanners can then utilize regular expression matchers

Benefits:

- Extensible design
 - easy to add new tokens, modify existing definitions
- Modular
 - Scanner can utilize common regex libraries

Cons:

How do we solve this?

A new token definition language:

- **Regular expressions**
- Tokens will be defined using regular expressions
- Scanners can then utilize regular expression matchers

Benefits:

- Extensible design
 - easy to add new tokens, modify existing definitions
- Modular
 - Scanner can utilize common regex libraries

Cons:

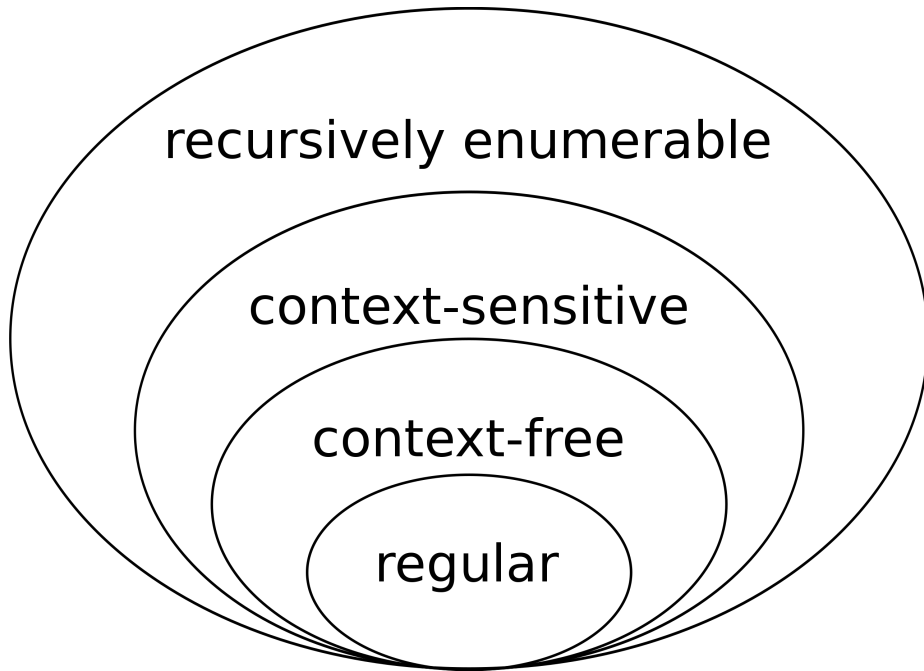
- Token definitions are restricted to regular languages
- Potentially slower
- Regular expression matchers are complicated

Regular expressions

Some theory:

- Given a language L , a string s is either part of that language or not
 - Integers are a language: “5”, “6”, “-7” is in the language. “abc” is not.
- Languages are grouped into families depending on how “hard” it is to determine if a string is part of that language.

Regular expressions

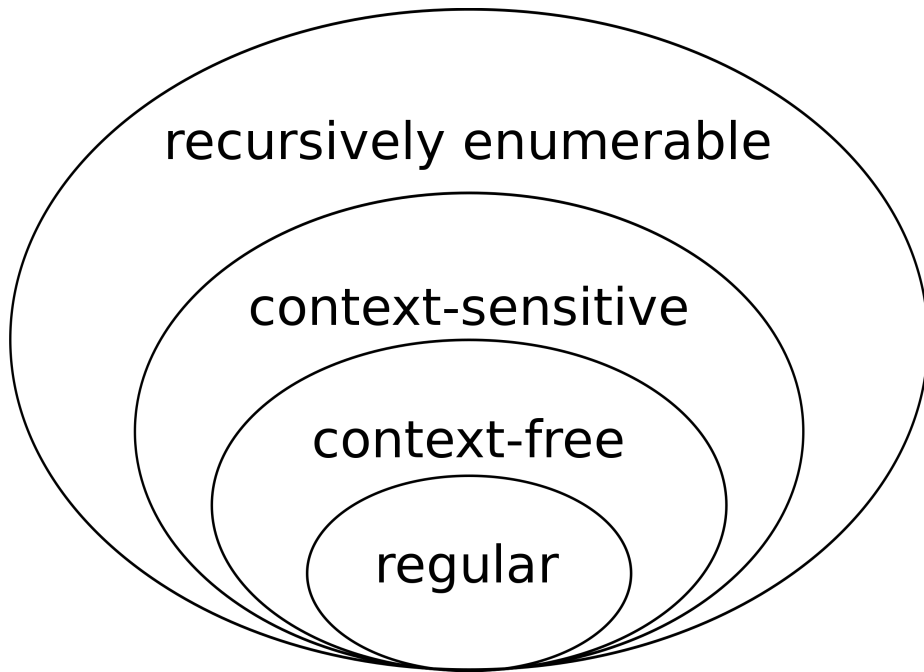


The simplest languages are regular. We will use regular languages as our token language.

We will use the next level: context-free, as the language for our parser.

Higher levels are interesting, but not as useful in compilers. Why?

Regular expressions



The simplest languages are regular. We will use regular languages as our token language.

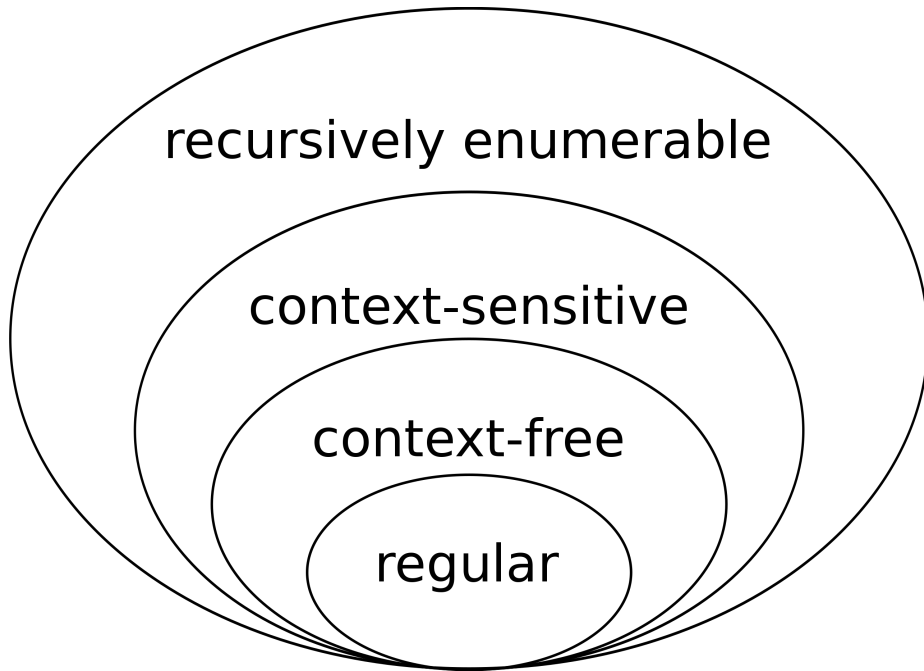
We will use the next level: context-free, as the language for our parser.

Higher levels are interesting, but not as useful in compilers. Why?

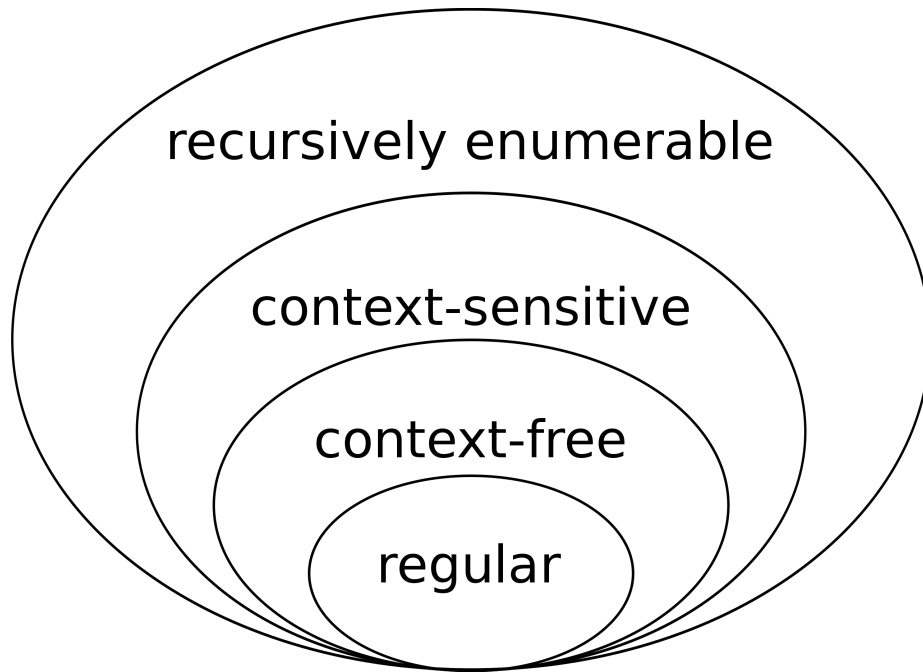
Because deciding if a string is in a recursively enumerable language is undecidable.

Regular expressions

What is a regular language?



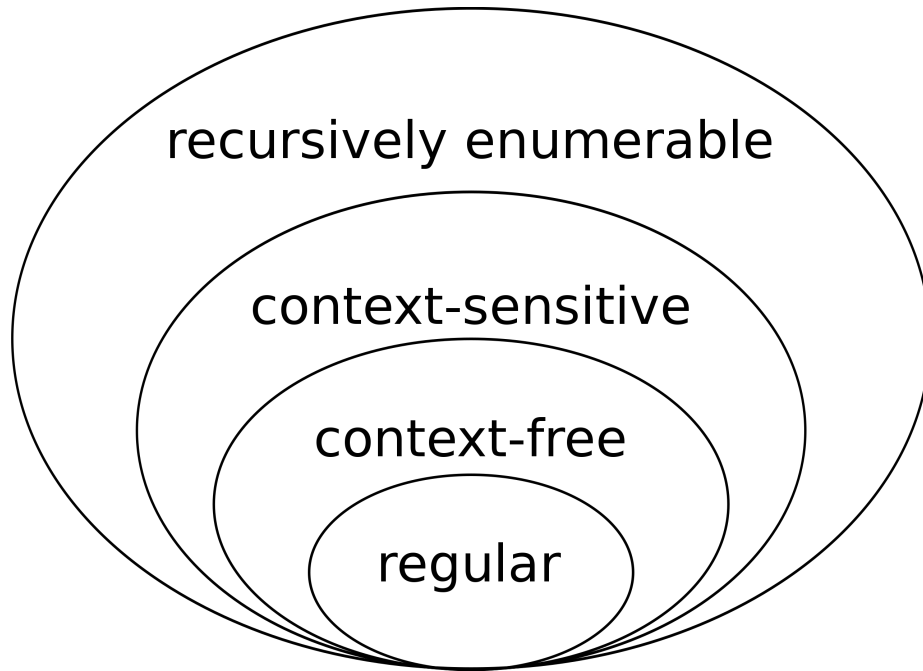
Regular expressions



What is a regular language?

For this class: *A regular language is a language that can be expressed as a regular expression.*

Regular expressions



What is a regular language?

For this class: *A regular language is a language that can be expressed as a regular expression.*

What is a regular expression?

Regular expressions

- We will define regular expressions (RE) recursively
- We will show examples at each step.
- And show to match them in Python
 - *A string matches an RE if it belongs to the regular language defined by the RE*
 - Python has a great RE matching library

Regular expressions

```
# import the library
```

```
import re
```

```
# pattern is a string representing the RE
```

```
# the function reports whether string matches RE
```

```
re.fullmatch(pattern, string)
```

Regular expressions

- **We will define regular expressions (RE) recursively**
- Like any recursive function, we can start with the base case:

a regular expression can be a single character or the empty string

Regular expressions

- **We will define regular expressions (RE) recursively**
- Like any recursive function, we can start with the base case:

a regular expression can be a single character or the empty string

Example:

```
ASSIGN = "="  
PLUS   = "+"
```

Python:

```
import re  
re.fullmatch("=", "=")  
  
re.fullmatch("+", "+") # what happens here?
```

Regular expressions

- When we define regular expressions, some characters are special.
 - They are operators in the regular expression language
 - If we want to use them as a character, then we need to "escape them" with a \
 - "+" happens to be one of those characters

<https://riptutorial.com/regex/example/15848/what-characters-need-to-be-escaped->

Python:

```
import re
re.fullmatch("=", "=")

re.fullmatch("\+", "+") # what happens here?
```

Regular expressions

- **We will define regular expressions (RE) recursively**
- Like any recursive function, we can start with the base case:

*a regular expression can be a single character or the **empty string***

Python:

```
import re  
re.fullmatch("", "")
```

*Not super useful for us,
but useful for the theory*

Regular expressions

- First recursive case: **concatenation**
- Two REs can be concatenated by simply writing them in sequence:
 - RE1 = "a", RE2 = "b"
 - concatenated it is: RE12 = "ab"
- This allows us to build words

Example:

```
FOR    = "for"  
WHILE  = "while"
```

Python:

```
import re  
re.fullmatch("for", "for")  
re.fullmatch("a+b", "a+b") # what happens here?
```


Can we define these tokens yet?

- ARTICLE
- NOUN
- VERB
- ADJECTIVE

Tokens

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

Tokens Definitions

Can we define these tokens yet? No, we need one more operator

- ARTICLE
- NOUN
- VERB
- ADJECTIVE

Tokens

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

Tokens Definitions

Regular expressions

- Second recursive operator: **choice** (sometimes called "union", or "or")
- Two REs can be choiced together using the "|" operator
 - RE1 = "a", RE2 = "b"
 - The choice is: RE1|2 = "a|b"
 - Matches either

Example:

```
OP      = "*" | "+"
CMP     = "==" | "<=" | ">="
```

Python:

```
import re
re.fullmatch("*|+", "+")
re.fullmatch("==|<=|>=", "==")
```

Can we define these tokens yet?

- ARTICLE
- NOUN
- VERB
- ADJECTIVE

Tokens

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

Tokens Definitions

Can we define these tokens yet? Yes!

- ARTICLE
- NOUN
- VERB
- ADJECTIVE

Tokens

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

Tokens Definitions

Can we define these tokens yet?

```
ID      = [characters]
NUM      = [numbers]
ASSIGN   = "="
PLUS     = "+"
MULT     = "*"
IGNORE   = [" "]
```

Can we define these tokens yet? No!

```
ID      = [characters]
NUM      = [numbers]
ASSIGN   = "="
PLUS     = "+"
MULT     = "*"
IGNORE   = [" "]
```

Regular expressions

- Last recursive operator: **Repeat**
- Unary operator: *****
 - RE1 = "a"
 - Repeat RE1 zero or more times: "a*"

Example:

```
RE1    = "a*"
RE2    = "a*|b*"
RE3    = "a|b"
```

Python:

```
import re
re.fullmatch("a*|b*", "aaa")
re.fullmatch("a*|b*", "")
```


Regular expressions

- Last recursive operator: **Repeat**
- Unary operator: *****
 - RE1 = "a"
 - Repeat RE1 zero or more times: "a*"

Example:

```
RE1    = "a*"
RE2    = "a*|b*"
RE3    = "a|b"
```

Precedence?

Python:

```
import re
re.fullmatch("a*|b*", "aaa")
re.fullmatch("a*|b*", "")
```

Regular expressions

- These are the theoretical foundational operators.
- Most languages give syntactic sugar to make common cases easier
- Most languages also break the theory
 - Perl regexes are extremely complicated
 - https://www.perlmonks.org/?node_id=809842
 - Python regexes (with recursion) are can capture context free languages
 - <https://www.npopov.com/2012/06/15/The-true-power-of-regular-expressions.html#matching-context-free-languages>

Regular expressions

- strict repeat operator: +
- one or more repeats (the * operator is 0 or more repeats)
- derivation: "r+" = "rr*"

Regular expressions

- 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:
- Hexadecimal numbers:

Regular expressions

- 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]"

Regular expressions

- optional operator ?
 - optional characters
- “r?” = “|r”
- Example: “ab?”

Regular expressions

- optional operator ?
 - optional characters
- “r?” = “|r”
- Example: “ab?”
- Let’s do simple floating point numbers:

Regular expressions

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

Regular expressions

- any character ‘.’
- example using email (this is probably too general!)

Regular expressions

- 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 delimitate 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
- “(?P<name>.+)(?P<domain>+\\.com)”

REs are good for?

- Scanning large amounts of documents quickly, looking for:
 - Websites
 - Email
 - Profiling numbers
 - Variable usages
 - **What else?**

RE examples

- **What can REs not do?**
- Nested structures, such as parenthesis 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!
 - <https://stackoverflow.com/questions/1732348/regex-match-open-tags-except-xhtml-self-contained-tags/1732454#1732454>

For your homework

- You'll need to write tokens for a simple programming language, including:

ID	=	[characters]
NUM	=	[numbers]
ASSIGN	=	"="
PLUS	=	"+"
MULT	=	"*"
IGNORE	=	[" "]

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 sensitive (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
 - https://ivanzuzak.info/noam/webapps/fsm_simulator/

How to use REs in a scanner implementation?

- We will discuss next class
- See you on Wednesday!