

# Hands-on Workshop : Python Programming

SSN ACM Student Chapter

# AGENDA FOR THE DAY

- Intro to Python
- Basics of Python (3 Coding Sessions)
- Lunch Break
- Python Specific Constructs and Features (2 Coding Sessions)
- After Session:
  - ◆ Google Image Search API with Python
  - ◆ Scraping “www.netlingo.com” for Web abbreviations (eg. LOL, ROFL, etc.)
  - ◆ Parser for Social Network Text data

WHAT MOTIVATED ME ?

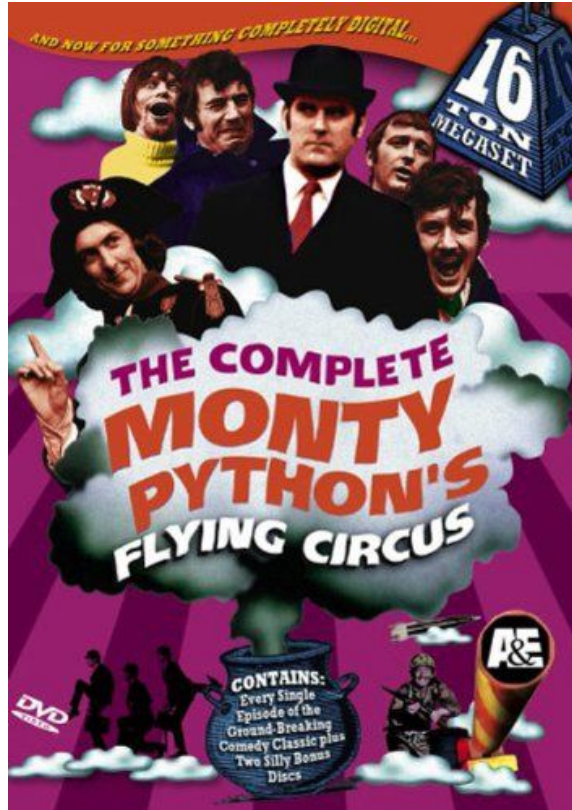
# FORMAL INTRO

- Created by **Guido Van Rossum**.
- Extremely Simple + Intuitive + Minimalistic + Readable + Expressive => **Very Pythonic !!**
- Interpreted Language
- Multi-Paradigm programming language
  - Object Oriented
  - Structural
  - And Functional
- Dynamic Typing, Dynamic Name Resolution, Cycle Detecting Garbage Collector



Guido Van Rossum  
**Worked in,**  
Google - 2005 - 2012  
DropBox - 2013 - Now

# WHY THE NAME "PYTHON"?

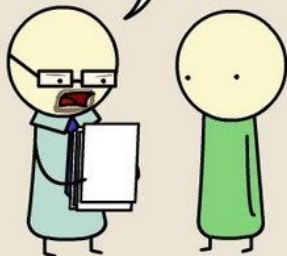


# WHY LEARN PYTHON?

# PYTHON IS EASY!

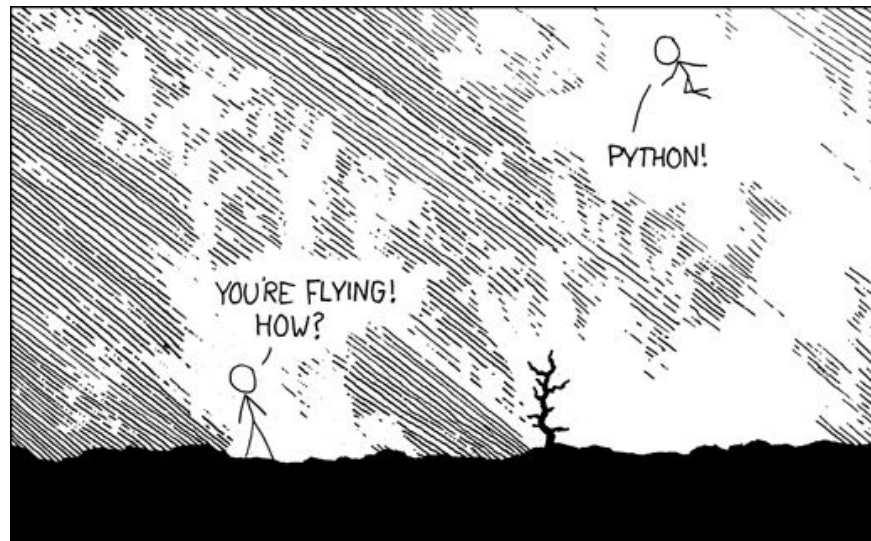
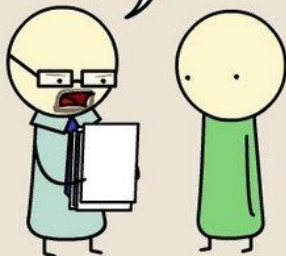
## PYTHON

THIS IS PLAGIARISM.  
YOU CAN'T JUST "IMPORT ESSAY."



## JAVA

I'M TWO PAGES IN AND I STILL  
HAVE NO IDEA WHAT YOU'RE SAYING.



I LEARNED IT LAST  
NIGHT! EVERYTHING  
IS SO SIMPLE!  
HELLO WORLD IS JUST  
print "Hello, world!"

I DUNNO...  
DYNAMIC TYPING?  
WHITESPACE?

COME JOIN US!  
PROGRAMMING  
IS FUN AGAIN!  
IT'S A WHOLE  
NEW WORLD  
UP HERE!

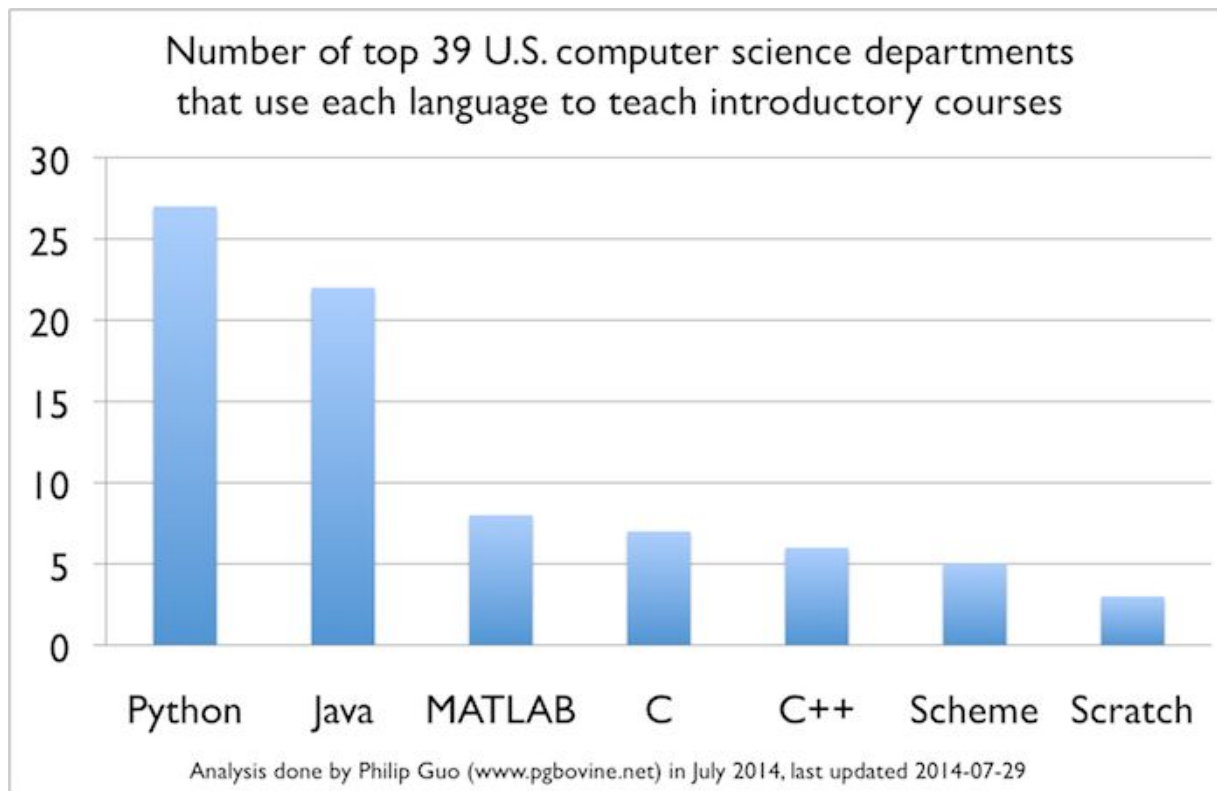
BUT HOW ARE  
YOU FLYING?

I JUST TYPED  
import ontigravity  
THAT'S IT?

... I ALSO SAMPLED  
EVERYTHING IN THE  
MEDICINE CABINET  
FOR COMPARISON.

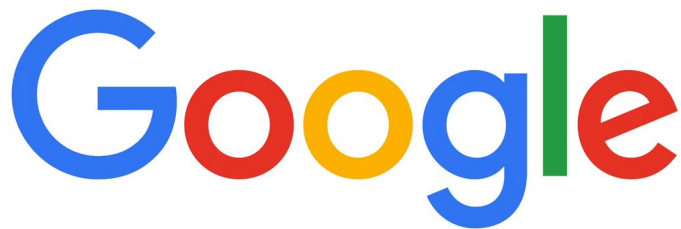
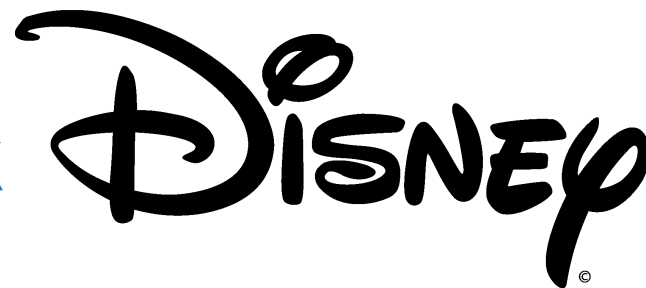
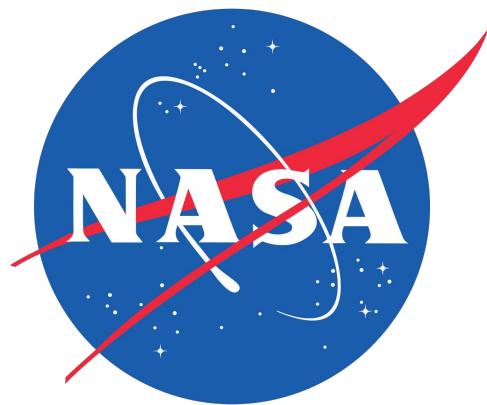
BUT I THINK THIS  
IS THE PYTHON.

# PYTHON IS THE “FIRST” LANGUAGE!





WHO USES PYTHON?



# PYTHON IS EFFICIENT!

```
1  def read_in_chunks(file_object, chunk_size=1024):  
2      """Lazy function (generator) to read a file piece by piece.  
3      Default chunk size: 1k."""  
4      while True:  
5          data = file_object.read(chunk_size)  
6          if not data:  
7              break  
8          yield data  
9  
10  
11  f = open('really_big_file.dat')  
12  for piece in read_in_chunks(f):  
13      process_data(piece)  
14
```

# GENERATORS IN C!

```
1  def pow2s():  
2      ...yield 1  
3      ...for i in map((lambda x: 2 * x), pow2s()):  
4          ...yield i  
5        
6        
7  def mymap(func, iteration):  
8      ...for i in iteration:  
9          ...yield func(i)  
10     
```

```
1  #include <stdio.h>  
2  #include <stdlib.h>  
3    
4  struct gen { // generic structure, the base of all generators  
5      int (*next)();  
6      int continue_from;  
7  };  
8    
9  typedef int (*fptr)();  
10   
11 struct mapgen { // structure for map  
12     int (*next)();  
13     int continue_from; // not really required, provided for compatibility  
14     fptr f;  
15     struct gen *g;  
16 };  
17   
18 int map_next(struct mapgen *p) { // next function for map  
19     return p->f(p->g->next(p->g));  
20 }  
21   
22 struct gen *map(fptr f, struct gen *g) { // constructor for map iterator  
23     struct mapgen *p = (struct mapgen *)malloc(sizeof(struct mapgen));  
24     p->next = map_next;  
25     p->continue_from = 0;  
26     p->f = f;  
27     p->g = g;  
28     return (struct gen *)p;  
29 }  
30   
31 struct pow2s { // structure  
32     int (*next)();  
33     int continue_from;  
34     struct gen *g;  
35 };  
36   
37 int times2(int x) { // anonymous lambda is translated into this  
38     return 2*x;  
39 }  
40   
41 struct gen *pow2() { // forward declaration of constructor  
42   
43     int pow2next(struct pow2s *p) { // next function for iterator  
44         switch(p->continue_from) {  
45             case 0:  
46                 p->continue_from = 1;  
47                 return 1;  
48             case 1:  
49                 p->g = map(times2, pow2());  
50                 p->continue_from = 2;  
51                 return p->g->next(p->g);  
52             case 2:  
53                 p->continue_from = 2;  
54                 return p->g->next(p->g);  
55             }  
56         }  
57     }  
58     struct gen *pow2() { // constructor for pow2  
59         struct pow2s *p = (struct pow2s *)malloc(sizeof(struct pow2s));  
60         p->next = pow2next;  
61         p->continue_from = 0;  
62         return (struct gen *)p;  
63     }  
64   
65     int main() {  
66         int i;  
67         struct gen *p = pow2();  
68         for(i=0; i<10; i++)  
69             printf("%d ", p->next(p));  
70         printf("\n");  
71     }  
72 }
```

# PYTHON IS “FAST” !

## NATIVE PYTHON

```
import numpy as np

def pairwise_python(X, D):
    M = X.shape[0]
    N = X.shape[1]
    for i in range(M):
        for j in range(M):
            d = 0.0
            for k in range(N):
                tmp = X[i, k] - X[j, k]
                d += tmp * tmp
            D[i, j] = np.sqrt(d)
```

```
In [2]: import numpy as np
In [3]: X = np.random.random((1000, 3))
In [4]: D = np.empty((1000, 1000))
In [5]: %timeit pairwise_python(X, D)
1 loops, best of 3: 12.1 s per loop
```

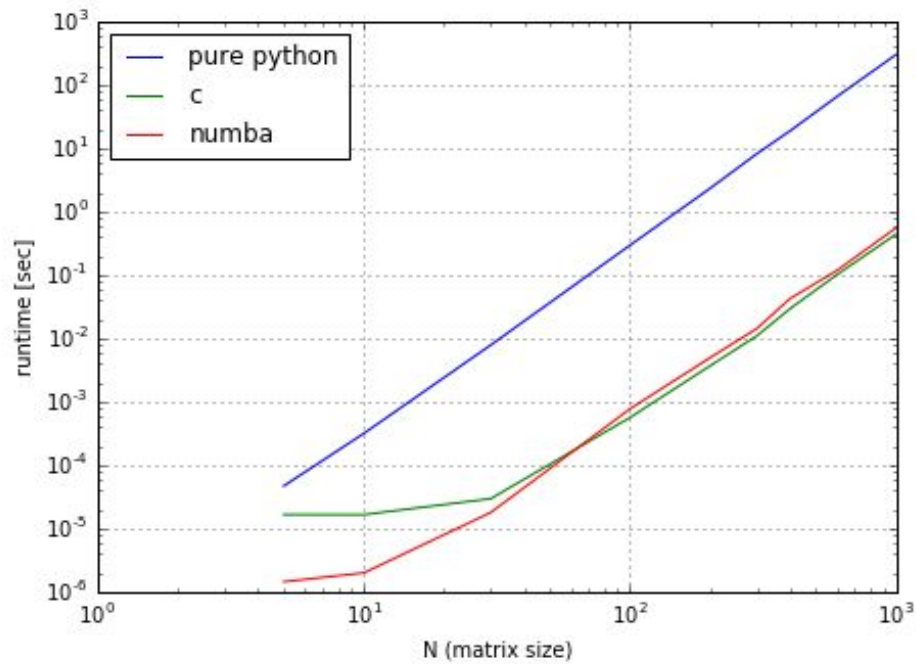
## “NUMBA” PYTHON

```
import numpy as np
from numba import double
from numba.decorators import jit

@jit(arg_types=[double[:, :], double[:, :]])
def pairwise_numba(X, D):
    M = X.shape[0]
    N = X.shape[1]
    for i in range(M):
        for j in range(M):
            d = 0.0
            for k in range(N):
                tmp = X[i, k] - X[j, k]
                d += tmp * tmp
            D[i, j] = np.sqrt(d)
```

```
In [2]: import numpy as np
In [3]: X = np.random.random((1000, 3))
In [4]: D = np.empty((1000, 1000))
In [5]: %timeit pairwise_numba(X, D)
100 loops, best of 3: 15.5 ms per loop
```

# BENCHMARKING

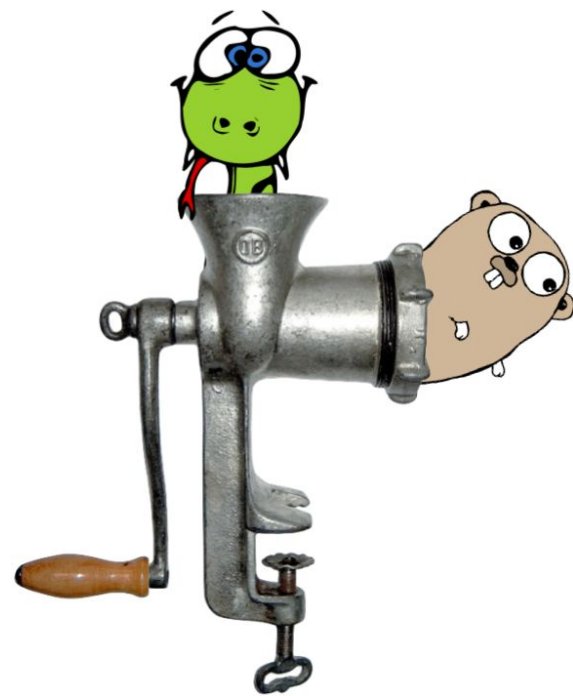


Source: [IBM developerWorks Blog](#)

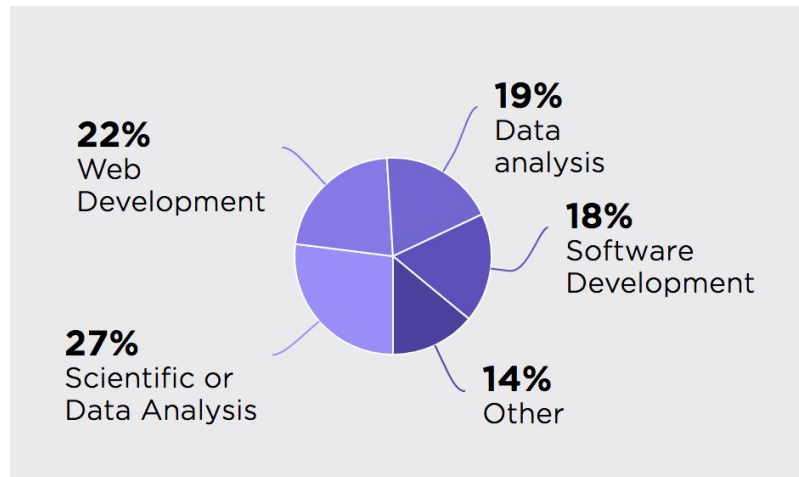
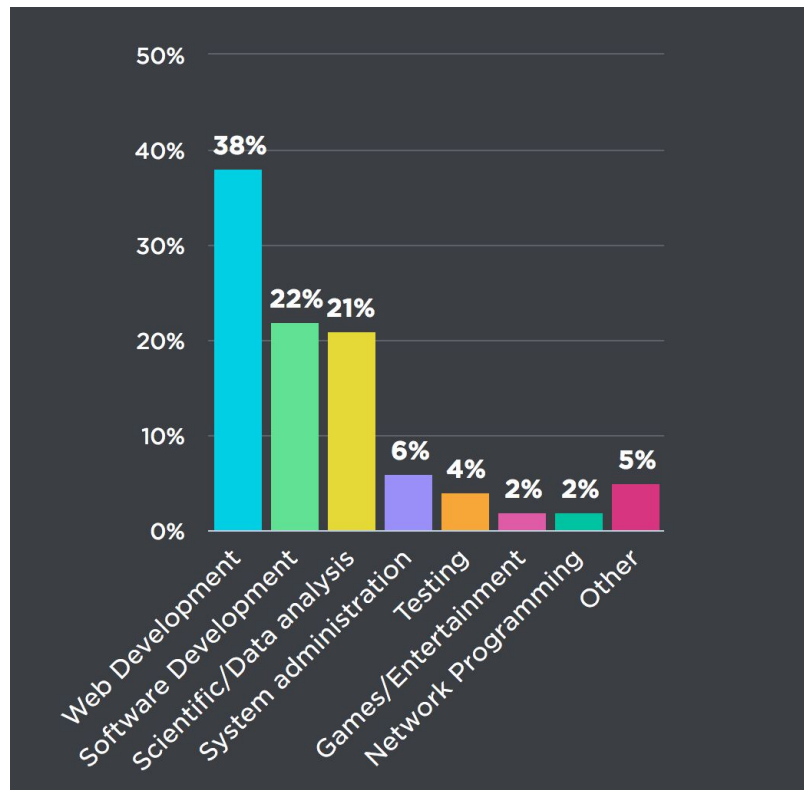
PYTHON IS NOT PYTHON ANYMORE!



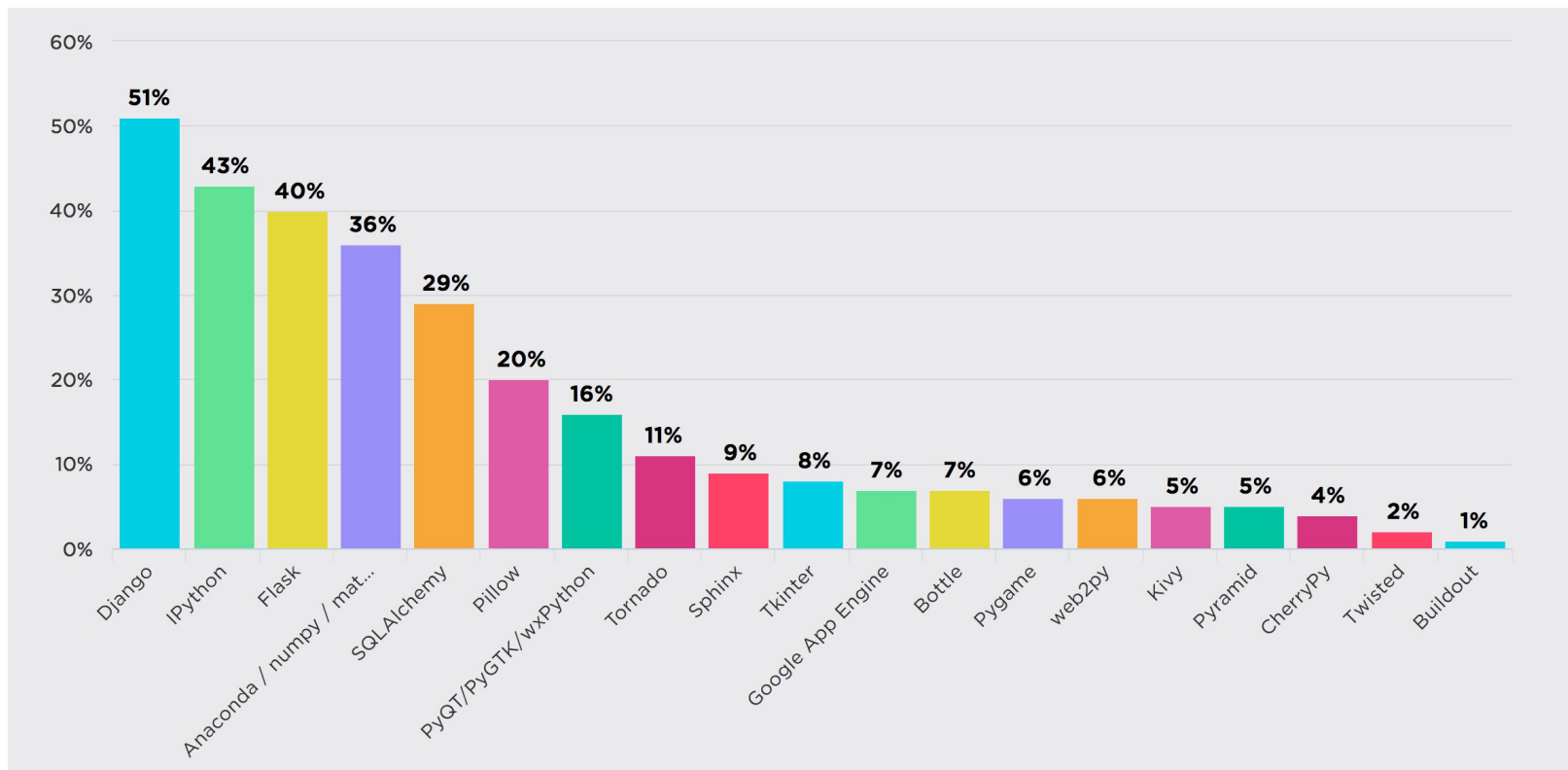
IronPython



# PYTHON IS EVERYWHERE!



# PYTHON IS OPEN SOURCE!





# ARE YOU ALL SET?

## → Python Environment Setup

- Refer to the **“Setup.txt”** file in the **“Py Environment Setup”** directory.
- **Anaconda:** It is a completely free Python distribution that has 195+ of the most popular Python packages for science, mathematics, engineering, data analysis.

→ If, in case you have problems with setting up, try the **“Alternate”** in **“Py Environment Setup”** folder.

→ All the **‘Time to Code ...!’** slides have supporting files in the **“./Py Script Files/”** folder.

# HELLO WORLD!

Java

```
public class Main {  
    public static void main(String[] args) {  
        System.out.println("hello world");  
    }  
}
```

Python

```
print "hello world";
```

# INDENTATION ...

- Python uses tabs/spaces for representing blocks rather than curly braces or keywords - off-side rule (Standard is - four spaces).
- Codestyling Paradigms: PEP-8, Flake-8, etc.

```
1  # Normal C Syntax
2  int x = 1;
3  if(x == 1){
4      printf("x is 1.");
5  }
6
7  # Python Syntax
8  x = 1
9  if x == 1:
10     # indented four spaces
11     print("x is 1.")
```

# BASICS OF PYTHON

- Variables and Data types
- Data Structures
- Operators
- Conditions
- Loops
- Functions
- Classes and Objects
- Modules and Packages

# VARIABLES AND DATA TYPES

→ Variables are,

- Objects
- Dynamically Typed - No need to declare variables or their types prior to use.

→ Basic Types

- Numbers
  - Integers
  - Floating Point
  - Complex
- Strings

# STRINGS

```
1 s = 'hi'
2 print s[1]      ## i
3 print len(s)    ## 2
4 print s + ' there'  ## hi there
5
6 pi = 3.14
7 ##text = 'The value of pi is ' + pi      ## NO, does not work
8 text = 'The value of pi is ' + str(pi)    ## yes
9
10 raw = r'this\t\n and that'
11 print raw      ## this\t\n and that
12
13 multi = """It was the best of times.
14 It was the worst of times."""
15
```

# STRING OPERATIONS

- `s.lower()`, `s.upper()`
- `s.strip()`
- `s.isalpha()/s.isdigit()/s.isspace()`
- `s.startswith('other')/s.endswith('other')`
- `s.find('other')`
- `s.replace('old', 'new')`
- `s.split('delim')`
- `s.join(list)`

```
>>> "\n Harry Potter ".strip()
'Harry Potter'
>>>
>>> "Harry Potter".lower()
'harry potter'
>>> "Harry Potter".upper()
'HARRY POTTER'
>>>
>>>
>>> "abc".isalpha()
True
>>> "abc ".isalpha()
False
>>>
>>>
>>> "123".isdigit()
True
>>> "123a".isdigit()
False
>>>
>>>
>>> " \t\n".isspace()
True
>>> " 1\t\n".isspace()
False
```

```
>>>
>>> "Star Wars".startswith("S")
True
>>> "Star Wars".startswith("s")
False
>>>
>>>
>>> "Star Wars".find("Wars")
5
>>> "Star Wars".find("wars")
-1
>>>
>>>
>>> "Star Wars".replace("Wars", "Battle")
'Star Battle'
>>>
>>>
>>> "Star Wars".split(" ")
['Star', 'Wars']
>>>
>>>
[>>> "||".join(["Monty", "Python's", "Circus"])
'Monty||Python's||Circus'
```



# STRING SLICES

Hello

0 1 2 3 4  
-5 -4 -3 -2 -1

```
Python 3.4.3 [Anaconda 2.2.0 (64-bit)] (default, Jun  4 2015, 15:29:08)
[GCC 4.4.7 20120313 (Red Hat 4.4.7-1)] on linux
Type "help", "copyright", "credits" or "license" for more information.
>>> str = "Hello SSN!"
>>> str[0]
'H'
>>> str[1:3]
'el'
>>> str[-1]
'!'
>>> str[::-1]
'!NSS olleH'
>>> str[-5:]
' SSN!'
```

# STRING FORMATTING

```
1 # % Operator
2 # This prints out "Hello, John!"
3 name = "John"
4 print("Hello, %s!" % name)
5
6 # This prints out "3 little pigs come out or I'll huff and puff and blow down"
7 text = "%d little pigs come out or I'll %s and %s and %s" % (3, 'huff', 'puff', 'blow down')
8 print(text)
```

```
1 # This prints out: 'First Value-Second Value'~
2 '{0}-{1}'.format('First Value', 'Second Value')~
3 ~
4 # This prints out: 'My name is Chris. I like Machine Learning'~
5 'My name is {name}. I like {interest}'.format(name='Chris', interest="Machine Learning")~
6 ~
```

# STRING ENCODING

```
10  # Encoding
11  >> ustring = u'A unicode \u018e string \xf1'
12  ## (ustring from above contains a unicode string)
13  >> s = ustring.encode('utf-8')
14  >> 'A unicode \xc6\x8e string \xc3\xb1'
15  >> print(s)                                ## bytes of utf-8 encoding
16  >> t = unicode(s, 'utf-8')                 ## Convert bytes back to a unicode string
17  >> t == ustring
18  >> True
```

ISO-8859-1 is graphic character set that is a superset of “UTF-8”

# DATA STRUCTURES

<b>List</b>	Mutable, Mixed Types
<b>Dictionary</b>	Mutable, Mixed Types (key and value)
<b>Set</b>	Mutable, Hashable, Unordered, Mixed Types
<b>Tuple</b>	Mutable, Mixed Types

# PYTHON LISTS

- Lists are similar to arrays but with mixed-type support.
- They have the `len()` function and `[]` to access data, similar to strings.

```
1  >> colors = ['red', 123, 12.34]
2  >> print colors[0]
3  >> 'red'
4  >> print colors[2]
5  >> 12.34
6  >> print len(colors)
7  >> 3
8  >> A = [2, 3]
9  >> B = [4, 5]
10 >> A + B # append
11 >> [2, 3, 4, 5]
```

# "FOR" AND "IN" IN LISTS

```
>>> squares = [1, 4, 9, 16]
>>> sum = 0
>>> for num in squares:
...     sum += num
...
>>> print (sum)
30
>>> list = ['larry', 'curly', 'moe']
>>> if 'curly' in list:
...     print("Found !")
...
Found !
```

# LIST METHODS

- `list.append(elem)`
- `list.extend(list2)`
- `list.index(elem)`
- `list.remove(elem)`
- `list.sort()`
- `list.reverse()`
- `list.pop(index)`
- `list.insert(index, elem)`

```
>>> list = ['larry', 'curly', 'moe']
>>> list.append('shemp')           ## append elem at end
>>> list.insert(0, 'xxx')          ## insert elem at index 0
>>> list.extend(['yyy', 'zzz'])    ## add list of elems at end
>>> print (list)
['xxx', 'larry', 'curly', 'moe', 'shemp', 'yyy', 'zzz']
>>> print (list.index('curly'))
2
>>> list.remove('curly')           ## search and remove that element
>>> list.pop(1)                   ## removes and returns 'larry'
'larry'
>>> print (list)
['xxx', 'moe', 'shemp', 'yyy', 'zzz']
>>>
```

**List Slicing works the same as Strings.**

```
>>>
>>> a = ["one", "two", "three"]
>>> a.append("four")
>>> a
['one', 'two', 'three', 'four']
>>>
>>>
>>> a = ["one", "two", "three"]
>>> a.append(["four", "five"])
>>> a
['one', 'two', 'three', ['four', 'five']]
>>>
>>>
>>> a = ["one", "two", "three"]
>>> a.extend(["four", "five"])
>>> a
['one', 'two', 'three', 'four', 'five']
>>>
>>>
>>> a = ["one", "two", "three"]
>>> print(a.index("two"))
1
>>>
>>>
```

```
>>> a = ["one", "two", "three"]
>>> a.remove("one")
>>> a
['two', 'three']
>>>
>>>
>>> a = [5, -1, 10]
>>> a.sort()
>>> a
[-1, 5, 10]
>>>
>>>
>>> a = ["one", "two", "three"]
>>> a.reverse()
>>> a
['three', 'two', 'one']
>>>
>>>
>>> a = ["one", "two", "three"]
>>> a.pop(2)
'three'
>>>
>>>
>>> a = ["one", "two", "three"]
>>> a.insert(1, "four")
[>>> a
['one', 'four', 'two', 'three']
```



# TIME TO CODE ...!

- **Edit Words:** Given a word  $W$ , perform the following edit operations and output a single unique list of edited words.
- Eg. “pokemon”
  - **Splits** – ‘p’, ‘okemon’; ‘po’, ‘kemon’; ‘pok’, ‘emon’; etc.
  - **Deletes** – ‘pkemon’, ‘poemon’, ‘pokmon’, etc.
  - **Transposes** – ‘pkoemon’, ‘poekmon’, ‘pokmeon’, etc.
  - **Replaces** – ‘pakemon’, ‘pobemon’, ‘pokcmon’, etc.
  - **Inserts** – ‘paokemon’, ‘pobkemon’, ‘pokcemon’, etc.

```
[>>>
>>> input_word = "pokemon"
>>> output_word_list = list()
>>> for i in range(1, len(input_word)):
...     output_word_list.extend([input_word[:i], input_word[i:]])
...
[>>> print(output_word_list)
['p', 'okemon', 'po', 'kemon', 'pok', 'emon', 'poke', 'mon', 'pokem', 'on', 'pokemo', 'n']
[>>>
[>>>
```

# PYTHON DICTIONARIES

- Efficient key/value hash tables
- {} are used to denote dict

```
1  ## Can build up a dict by starting with the the empty dict {}
2  ## dict[key] = value-for-that-key
3  dict = {}
4  dict['a'] = 'alpha'
5  dict['g'] = 'gamma'
6  dict['o'] = 'omega'
7  # (or)
8  dict = {'a': 'alpha', 'g': 'gamma', 'o': 'omega'}
9  print dict    ## {'a': 'alpha', 'o': 'omega', 'g': 'gamma'}
10
11 print dict['a']    ## Simple lookup, returns 'alpha'
12 dict['a'] = 6      ## Put new key/value into dict
13 'a' in dict        ## True
14 ## print dict['z']    ## Throws KeyError
15 if 'z' in dict: print dict['z']    ## Avoid KeyError
16 print dict.get('z')  ## None (instead of KeyError)
```

# CONTINUED...

```
1  ## By default, iterating over a dict iterates over its keys.
2  ## Note that the keys are in a random order.
3  for key in dict: print (key) # prints a g o
4  ## Exactly the same as above
5  for key in dict.keys(): print (key)
6
7  ## Get the .keys() list:
8  print (dict.keys()) ## ['a', 'o', 'g']
9
10 ## Likewise, there's a .values() list of values
11 print (dict.values()) ## ['alpha', 'omega', 'gamma']
12
13 ## .items() is the dict expressed as (key, value) tuples
14 print (dict.items()) ## [('a', 'alpha'), ('o', 'omega'), ('g', 'gamma')]
15
16 for k, v in dict.items(): print (k + '>' + v)
17 ## a > alpha      o > omega      g > gamma
```

# MUTABLE VS IMMUTABLE

Class	Description	Immutable?
<b>bool</b>	Boolean value	✓
<b>int</b>	integer (arbitrary magnitude)	✓
<b>float</b>	floating-point number	✓
<b>list</b>	mutable sequence of objects	
<b>tuple</b>	immutable sequence of objects	✓
<b>str</b>	character string	✓
<b>set</b>	unordered set of distinct objects	
<b>frozenset</b>	immutable form of set class	✓
<b>dict</b>	associative mapping (aka dictionary)	

# MUTABLE VS IMMUTABLE

```
>>> def func1(val):  
...     val += 'bar'  
...  
>>> x = 'foo'  
>>> x  
'foo'  
>>> func1(x)  
>>> x  
'foo'  
>>>  
>>>  
>>> def func2(val):  
...     val += [3, 2, 1]  
...  
>>> x = [1, 2, 3]  
>>> x  
[1, 2, 3]  
>>> func2(x)  
[>>> x  
[1, 2, 3, 3, 2, 1]
```

# OPERATORS

```
1  # Just as any other programming languages, the addition, subtraction, multiplication,  
2  # and division operators can be used with numbers.  
3  >> 1 + 2 * 3 / 4.0  
4  2.5  
5  >> 11 % 3 # remainder  
6  >> 2  
7  # Using two multiplication symbols makes a power relationship.  
8  >> 7 ** 2 # squared  
9  >> 49  
10 >> 2 ** 3 # cubed  
11 >> 8  
12 # Python supports concatenating strings using the addition operator:  
13 >> "hello" + " " + "world"  
14 >> 'hello world'  
15 # Python also supports multiplying strings to form a string with a repeating sequence:  
16 >> "hello" * 10  
17 >> 'hellohellohellohellohellohellohellohellohello'  
18 >> print ([1,2,3] * 3)  
19 >> [1, 2, 3, 1, 2, 3, 1, 2, 3]
```

# DEL OPERATOR

```
1  var = 6
2  del var  # var no more!
3
4  list = ['a', 'b', 'c', 'd']
5  del list[0]    ## Delete first element
6  del list[-2:]  ## Delete last two elements
7  print (list)   ## ['b']
8
9  dict = {'a':1, 'b':2, 'c':3}
10 del dict['b']   ## Delete 'b' entry
11 print (dict)   ## {'a':1, 'c':3}
```



# CONDITIONS

```
>>> x = 2
>>> x == 2
True
>>> x == 3
False
>>> x < 3
True
>>> x != 2
False
```

The "and" and "or" boolean operators allow building complex boolean expressions.

```
1 name = "John"
2 age = 23
3 if name == "John" and age == 23:
4     print ("Your name is John, and you are also 23 years old.")
5
6 if name == "John" or name == "Rick":
7     print ("Your name is either John or Rick.")
```

'Is' vs '=='

```
>>> # 'is' compares object level identity
... # '==' compares value level identity
...
>>> a = [1, 2, 3]
>>> b = a
>>> a
[1, 2, 3]
>>> b
[1, 2, 3]
>>>
>>>
>>> a is b
True
>>> a == b
True
>>>
>>>
>>> a = [1, 2, 3]
>>> b = list(a)
>>> a
[1, 2, 3]
>>> b
[1, 2, 3]
>>>
>>> a is b
False
[>>> a == b
True
```

# IF ... ELSE ...

```
if <statement is true>:  
    <do something>  
    ....  
    ....  
elif <another statement is true>: # else if  
    <do something else>  
    ....  
    ....  
else:  
    <do another thing>  
    ....  
    ....
```

# LOOPS - FOR, WHILE

```
1  primes = [2, 3, 5, 7]
2  for prime in primes:
3      print prime
4  # 2 3 5 7
5
6  for x in range(5):
7      print x
8  # 0 1 2 3 4
9
10 count = 0
11 while count < 5:
12     print count
13     count += 1  # This is the same as count = count + 1
14 # 0 1 2 3 4
```

# "BREAK" AND "CONTINUE"

```
1  # Prints out 0,1,2,3,4
2
3  count = 0
4  while True:
5      print (count)
6      count += 1
7      if count >= 5:
8          break
9
10 # Prints out only odd numbers - 1,3,5,7,9
11 for x in xrange(10):
12     # Check if x is even
13     if x % 2 == 0:
14         continue
15     print (x)
```

# TIME TO CODE ...!

- **Download Progress:** Open “**download-progress**”, you can find a ‘**download\_1.log**’ and ‘**download\_2.log**’ files in it.
- You have to read this file and sort all of the log entries based on time in ascending order. Each log entry occupies one line as is in Dict format.
- Then, you will have to look-up the ‘**progress\_percentage**’ key in each of the Dict entries and verify whether the download has been smooth.
- Smooth download (1,2,3,3,10,90,100) => Output ‘True’
- Non-smooth download (1,2,3,10,8,90,100) => Output ‘False’

```
{"timestamp": "2017-02-24 01:12:03.234892", "progress_percentage": 10}  
{"timestamp": "2017-02-24 01:12:04.234892", "progress_percentage": 11}
```

# HINTS...

```
# read from file~
logs = open("download_1.log").readlines()~
~
# sorting logs based on time~
sorted_logs = sorted(logs, key=lambda x: x["timestamp"])~
~
# conveting string to variables~
import ast~
log_str = '''{"timestamp": "2017-02-24 01:12:03.234892", "progress_percentage": 10}'''~
log_dict = ast.literal_eval(log_str)~
print(log_dict["progress_percentage"]) # 10~
~
# converting string date time datetime object~
from datetime import datetime~
time1 = datetime.strptime("2017-02-24 01:12:03.234892", "%Y-%m-%d %H:%M:%S.%f")~
time2 = datetime.strptime("2017-02-25 01:12:03.234892", "%Y-%m-%d %H:%M:%S.%f")~
print(time1 < time2) # True~
```

# FUNCTIONS

```
1  def my_function():
2      print ("Hello From My Function!")
3
4  def my_function_with_args(username, greeting):
5      print ("Hello, %s , From My Function!, I wish you %s" % (username, greeting))
6
7  def sum_two_numbers(a, b):
8      return a + b
9
10 # print a simple greeting
11 my_function()
12
13 #prints - "Hello, John, From My Function!, I wish you a great year!"
14 my_function_with_args("John", "a great year!")
15
16 # after this line x will hold the value 3!
17 x = sum_two_numbers(1,2)
```



# CLASSES AND OBJECTS

```
1  class Person:
2      def __init__(self, n):
3          self.name = n
4
5      def printMessage(self):
6          print ("Hello Mr./Ms. " + self.name)
7
8  # Create an object
9  p1 = Person("John")
10 p2 = Person("Angel")
11
12 # Function call
13 p1.printMessage() # Hello Mr./Ms. John
14 p2.printMessage() # Hello Mr./Ms. Angel
```

# MODULES AND PACKAGES

- **Modules** are Python files with a `.py` extension that implements a set of functions; Collection of modules is a **Package**.
- Installation of Python Packages
  - Use **`pip install <Module-or-Package-Name>`**
  - Download Source from PyPi and install using, **`python setup.py install`**
- ***“import” is the keyword used.***
- `# import the library`  
`import urllib.request`
- `# use it`  
`urllib.request.urlopen(...)`

# EXPLORING BUILT-IN MODULES

- “*dir*” and “*help*” functions are used to explore.

```
>>> from urllib import request
>>> dir(request)
['AbstractBasicAuthHandler', 'AbstractDigestAuthHandler', 'AbstractHTTPHandler',
 'BaseHandler', 'CacheFTPHandler', 'ContentTooShortError', 'DataHandler', 'FTPHa
ndler', 'FancyURLopener', 'FileHandler', 'HTTPBasicAuthHandler', 'HTTPCookieProc
essor', 'HTTPDefaultErrorHandler', 'HTTPDigestAuthHandler', 'HTTPError', 'HTTPEr
rorProcessor', 'HTTPHandler', 'HTTPPasswordMgr', 'HTTPPasswordMgrWithDefaultReal
m', 'HTTPRedirectHandler', 'HTTPSHandler', 'MAXFTPCACHE', 'OpenerDirector', 'Pro
xyBasicAuthHandler', 'ProxyDigestAuthHandler', 'ProxyHandler', 'Request', 'URLER
ror', 'URLopener', 'UnknownHandler', '_all_', '_builtins_', '_cached_', '_
_doc_', '_file_', '_loader_', '_name_', '_package_', '_spec_', '_ver
sion_', '_cut_port_re', '_ftperrors', '_have_ssl', '_localhost', '_noheaders',
 '_opener', '_parse_proxy', '_proxy_bypass_macosx_sysconf', '_randombytes', '_saf
e_gethostbyname', '_thishost', '_url_tempfiles', '_addclosehook', '_addinfourl',
'_base64', '_bisect', '_build_opener', '_collections', '_contextlib', '_email', '_ftpcac
he', '_ftperrors', '_ftpwrapper', '_getproxies', '_getproxies_environment', '_hashlib
', '_http', '_install_opener', '_io', '_localhost', '_noheaders', '_os', '_parse_http_l
ist', '_parse_keqv_list', '_pathname2url', '_posixpath', '_proxy_bypass', '_proxy_byp
ass_environment', '_quote', '_re', '_request_host', '_socket', '_splitattr', '_splitho
st', '_splitpasswd', '_splitport', '_splitquery', '_splittag', '_splitttype', '_splitus
er', '_splitvalue', '_ssl', '_sys', '_tempfile', '_thishost', '_time', '_to bytes', '_un
quote', '_unquote_to bytes', '_unwrap', '_url2pathname', '_urlcleanup', '_urljoin', '_
urlopen', '_urlparse', '_urlretrieve', '_urlsplit', '_urlunparse', '_warnings']
>>> help(request)
```

# WRITING MODULES AND PACKAGES

- Modules - Create a new “.py” with the Module name and all necessary functions inside and import it using the python file name.
- Packages
  - Are namespaces that contain multiple packages and Modules themselves.
  - They are simply directories with a special file named “**`__init__.py`**”, this can be empty for now.

## Folder Structure:

```
-> Foo
    --> bar1.py
    --> bar2.py
    --> __init__.py
-> spam.py
```

```
# inside spam.py import Foo package
>> from Foo import bar1
>> import Foo.bar2
```

# SOME PACKAGES I HAVE USED:

- PyCaffe, Keras, etc - Deep Learning Tools
- Scikit-Learn - Machine Learning Tools
- Nltk - Natural Language Processing Tools
- Gensim - Topic Modelling Library
- Scrappy - Web Scraping
- Numpy, Scipy - Mathematical and Scientific Library
- Flask, Cherrypy - lightweight web server
- Tweepy - Twitter API
- Twokenize - Tweets Tokenizer Tool.
- Etc...

# TIME TO CODE ... !

- Open the folder “**top-10**”.
- Try to parse the contents of the “**content.txt**” file and print the top 10 highly frequent words in the file.
- Try using looping, dicts (clue: defaultdict, counter), list slicing, files io, sorting, etc.
- Hint (File IO and Sorting):

```
# file open
f = open("filename","r") # r for Read mode
lines = f.readlines() # to get the lines in file as a list
f.close()

# sorting dict based on "value"
import operator
sorted_dict = dict(sorted(unsorted_dict.items(), key = operator.itemgetter(1)))
```

# PYTHON SPECIFIC CONSTRUCTS AND FEATURES

- Generators
- List Comprehensions
- Multiple Function Arguments
- Regular Expressions
- Exception Handling
- Sets
- Serialization - JSON, Pickle
- Lambda Operator
- Decorators

# GENERATORS

Simple functions which return an iterable set of items, one at a time, in a special way.

```
import random

def lottery():

    # returns 6 numbers between 1 and 40
    for i in range(6):
        yield random.randint(1, 40)

    # returns a 7th number between 1 and 15
    yield random.randint(1, 15)

for random_number in lottery():
    print("And the next number is... %d !" % random_number)
```



# LIST COMPREHENSIONS

Very powerful tool, which creates a new list based on another list, in a single, readable line.

```
1 sentence = "the quick brown fox jumps over the lazy dog"
2 words = sentence.split()
3 word_lengths = []
4 ✓ for word in words:
5     ✓ if word != "the":
6         word_lengths.append(len(word))
7 print(word_lengths)
8 # >> [5, 5, 3, 5, 4, 4, 3]
```

# MUCH CLEAN AND SIMPLE VERSION

```
10 sentence = "the quick brown fox jumps over the lazy dog"
11 words = sentence.split()
12 word_lengths = [len(word) for word in words if word != "the"]
13 print(word_lengths)
14 # >> [5, 5, 3, 5, 4, 4, 3]
```

# FUN ! - FIND THE N-GRAMS FROM A LIST OF WORDS

```
def find_ngrams(listI, n):  
    grams = []  
    for i in zip(*[listI[j:] for j in range(n)]):  
        grams.append(" ".join(i))  
    return grams
```

# MULTIPLE FUNCTION ARGUMENTS

```
# Normal Function Call
```

```
def myfunction(first, second, third):  
    # do something with the 3 variables  
    print(first + second + third)
```

```
myfunction(1, 2, 3)
```

```
# >> 6
```

```
def foo(first, second, third, *theRest):  
    print "First: %s" % first  
    print "Second: %s" % second  
    print "Third: %s" % third  
    print "And all the rest... %s" % list(theRest)
```

```
foo(1, 2, 3, 4, 5, 6, 7)
```

```
foo(1, 2, 3, 4, 5)
```

```
First: 1
```

```
Second: 2
```

```
Third: 3
```

```
And all the rest... [4, 5, 6, 7]
```

```
First: 1
```

```
Second: 2
```

```
Third: 3
```

```
And all the rest... [4, 5]
```

# MORE ...

```
def bar(first, second, third, **options):  
    if options.get("action") == "sum":  
        print "The sum is: %d" % (first + second + third)    # The sum is: 6  
  
    if options.get("number") == "first":  
        return first  
  
result = bar(1, 2, 3, action="sum", number="first")  
print "Result: %d" % result    # Result: 1
```

\*ARGS, \*\*KWARGS

```
>>> def foo(a, *args):
...     print(a, args)
...
>>> foo(1, 2, 3, 4)
(1, (2, 3, 4))
>>>
>>> foo(1, *[2, 3, 4])
(1, (2, 3, 4))
>>>
>>>
>>> def bar(a, **kwargs):
...     print(a, kwargs)
...
>>> bar(1, b=2, c=3, d=4)
(1, {'c': 3, 'b': 2, 'd': 4})
>>>
>>>
>>> def foo_bar(a, *args, **kwargs):
...     print(a, args, kwargs)
...
[>>> foo_bar(1, 2, 3, 4, b=5, c=6, d=7)
(1, (2, 3, 4), {'c': 6, 'b': 5, 'd': 7})
```

# REGULAR EXPRESSIONS

```
1  import re
2
3  line = "Cats are smarter than dogs"
4
5  matchObj = re.match(r'(.*) are (.*?) .*', line, re.M | re.I)
6
7  if matchObj:
8      print "matchObj.group() : %s" % matchObj.group()
9      print "matchObj.group(1) : %s" % matchObj.group(1)
10     print "matchObj.group(2) : %s" % matchObj.group(2)
11 else:
12     print("No match!!")
13
14 # RESULT
15 # matchObj.group() : Cats are smarter than dogs
16 # matchObj.group(1) : Cats
17 # matchObj.group(2) : smarter
```

**Further Reading:** [http://www.tutorialspoint.com/python/python\\_reg\\_expressions.htm](http://www.tutorialspoint.com/python/python_reg_expressions.htm)

# EXCEPTION HANDLING

```
1 def do_stuff_with_number(n):  
2     print(n)  
3  
4 the_list = (1, 2, 3, 4, 5)  
5  
6 for i in range(7):  
7     try:  
8         do_stuff_with_number(the_list[i])  
9     except IndexError: # Raised when accessing a non-existing index of a list  
10        do_stuff_with_number(0)
```

1  
2  
3  
4  
5  
0  
0  
[Finished in 0.039s]



# SETS

Similar to lists with no duplicate entries provided all entries must be hashable.

```
>>> print (set("my name is Eric and Eric is my name".split()))  
{'my', 'and', 'name', 'is', 'Eric'}
```

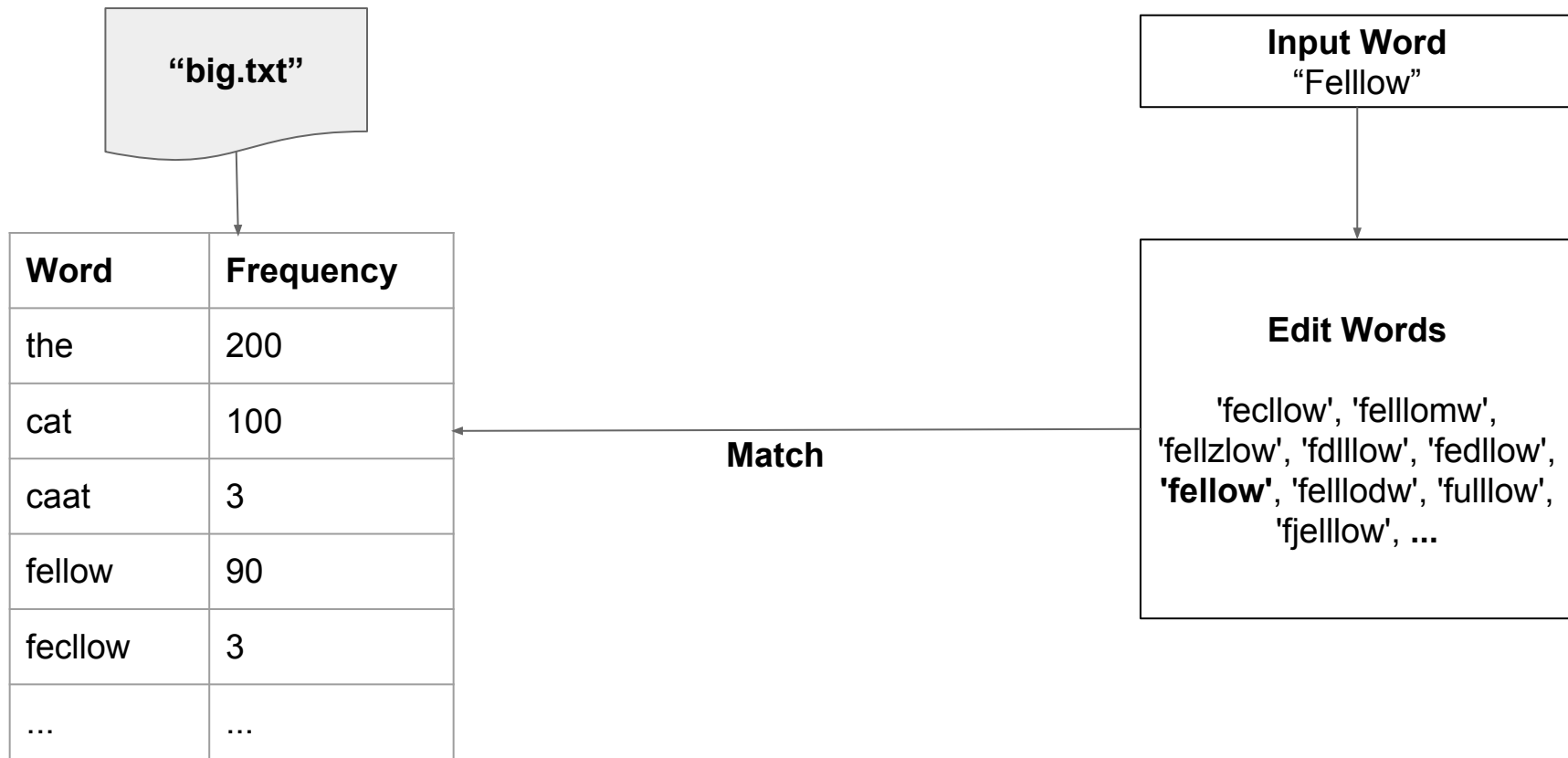
# EXAMPLE:

```
1  # list of participants in events A and B
2  a = set(["Jake", "John", "Eric"])
3  b = set(["John", "Jill"])
4
5  # To find out which members attended both events
6  a.intersection(b)  # set(['John'])
7  b.intersection(a)  # set(['John'])
8
9  # To find out which members attended only one of the events
10 a.symmetric_difference(b)  # set(['Jill', 'Jake', 'Eric'])
11 b.symmetric_difference(a)  # set(['Jill', 'Jake', 'Eric'])
12
13 # To find out which members attended only one event and not the other
14 a.difference(b)  # set(['Jake', 'Eric'])
15 b.difference(a)  # set(['Jill'])
16
17 # To receive a list of all participants
18 a.union(b)  # set(['Jill', 'Jake', 'John', 'Eric'])
```

# TIME TO RECODE ...!

- **Part 1:** Try to convert all the **Edit Words** for loop constructs to use **List Comprehension** and produce a unique list with **Set**.
- **Part 2:** Go into “**spell-corrector**”,
  - Try the command ``python spell.py fellow``, it will output ``fellow``
  - We need to extend **Edit Words** to build this spell corrector!

# HOW TO BUILD A SPELL CORRECTOR ?



# SERIALIZATION

```
1  # JSON
2  import json
3
4  dict = {'a': 'alpha', 'o': 'omega', 'g': 'gamma'}
5  json_string = json.dumps(dict)
6  print(json_string)
7
8  print(json.loads(json_string))
9
10 # PICKLE
11 import pickle
12 pickled_string = pickle.dumps([1, 2, 3, "a", "b", "c"])
13 print(pickle.loads(pickled_string))
14
```

# SAVING OBJECTS AND DATA STRUCTURES - PERSISTENCE

```
15 # Saving DataStructures as .pkl binary files
16 def save_obj(obj, name):
17     with open(name + '.pkl', 'wb') as f:
18         pickle.dump(obj, f, protocol=2)
19
20 def load_obj(name):
21     with open(name + '.pkl', 'rb') as f:
22         return pickle.load(f)
23
24 dict = {'a': 'alpha', 'o': 'omega', 'g': 'gamma'}
25 save_obj(dict, "dictSave")
```

# LAMBDA OPERATOR

- Anonymous Functions
- Functional Programming - Pass functions to other functions to do stuff.

```
1  # Using lambda
2  addTwo = lambda x: x+2
3  addTwo(2)
4  # 4
5
6  # the above is similar to
7  def addTwo(x):
8      return x+2
9  addTwo(2)
10 # 4
```

# USE

```
# Use:
mult3 = filter(lambda x: x % 3 == 0, [1, 2, 3, 4, 5, 6, 7, 8, 9])
print(mult3)
# [3, 6, 9]

# Equivalent to
new = []
for i in [1, 2, 3, 4, 5, 6, 7, 8, 9]:
    if i % 3 == 0:
        new.append(i)
print(new)
# [3, 6, 9]
```



# DECORATORS

```
1  from functools import wraps
2
3
4  def authorize(func):
5      @wraps(func)
6      def check_credentials(*args, **kwargs):
7          # check the login credentials
8          return func(*args, **kwargs)
9      return check_credentials
10
11
12 @authorize
13 def apiCall(username, password):
14     # do operation
15     pass
```

# TIME TO CODE ... !

- Open the folder **“imdb”**.
- Try to read from all the text files in **“MovieReviews”** folder.
- Try generating a list of all possible unique bi-grams and their frequencies in numbers and save the frequency map in .pkl format for persistence.
- Code should be resumable. i.e,(If we add new files from **“subMovieReviews”**) into **“MovieReviews”** and compile again, the code should only parse the new files and update the frequencies accordingly and update the .pkl binaries rather than overwriting.

# HINTS

- To get list of all files in a folder
  - `import os`
  - `os.listdir("/path/to/folder/")`
- Bi-Grams: They are phrases with exactly 2 consecutive words, from any piece of text.
  - eg. "India is my country!"
  - **All Bi-Grams** - "India is", "is my", "my country!".
- Use "sets" to find unique files list (use it with `os.listdir`)
- Use pickle to store and load frequency map
  - use `DefaultDict` from `Collections`

PYTHON WORK AT QUBE CINEMA TECHNOLOGIES!

# WEB SCRAPING WITH PYTHON

- Libraries we will be using
  - **Urllib** - open and read information from urls
  - **Beautifulsoup** - Parse HTML documents (similar to DOM in JS)
- Before building a scraper (“**web-scraping**”), let’s see how to download images from Google Image Search using Google Images API(“**google-image-api**”).

Parser for Social Network Text data

SUGGESTIONS?

# Thank You!

Satish Palaniappan

Contact:

Email: [tpsathish95@gmail.com](mailto:tpsathish95@gmail.com)

Phone: +91 9488515784

GitHub: <https://github.com/tpsathish95>

LinkedIn: <https://in.linkedin.com/in/satishpalaniappan>

*Except as otherwise noted, the contents of this workshop are licensed under the [Creative Commons Attribution 3.0 License](#), and code samples are licensed under the [Apache 2.0 License](#).*