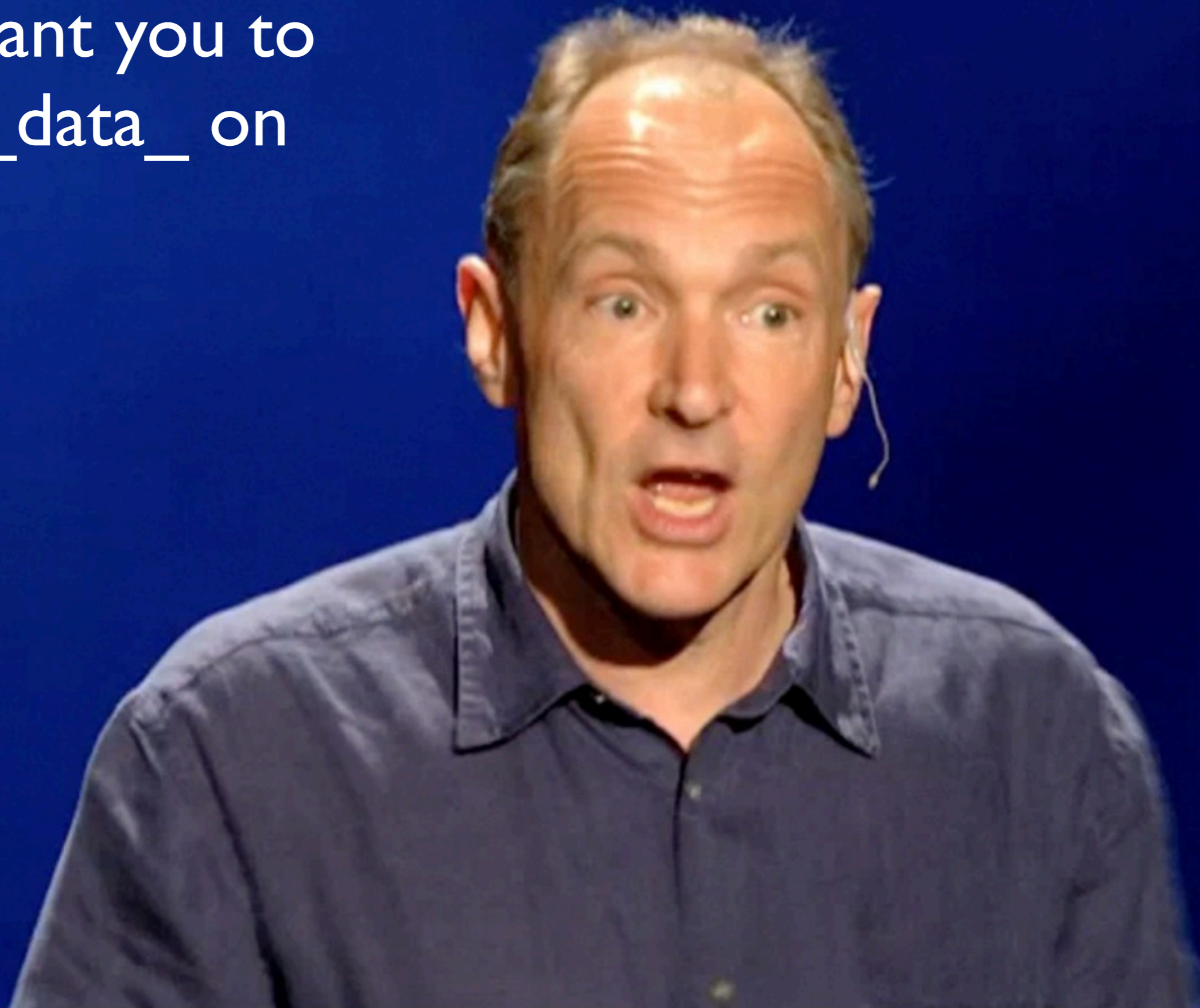


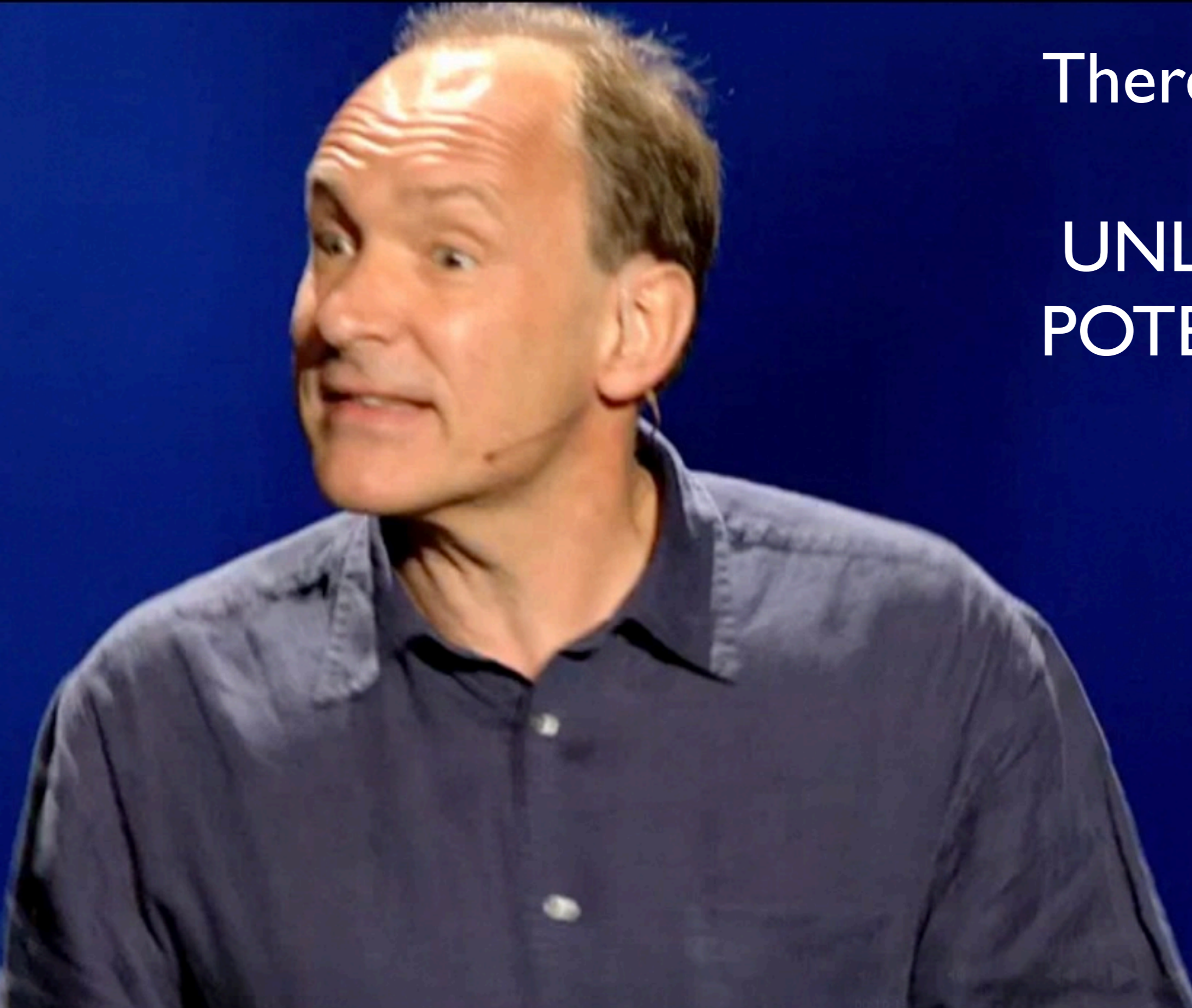
# Scraping with Python for Fun & Profit

@PyCon India 2010  
<http://in.pycon.org>

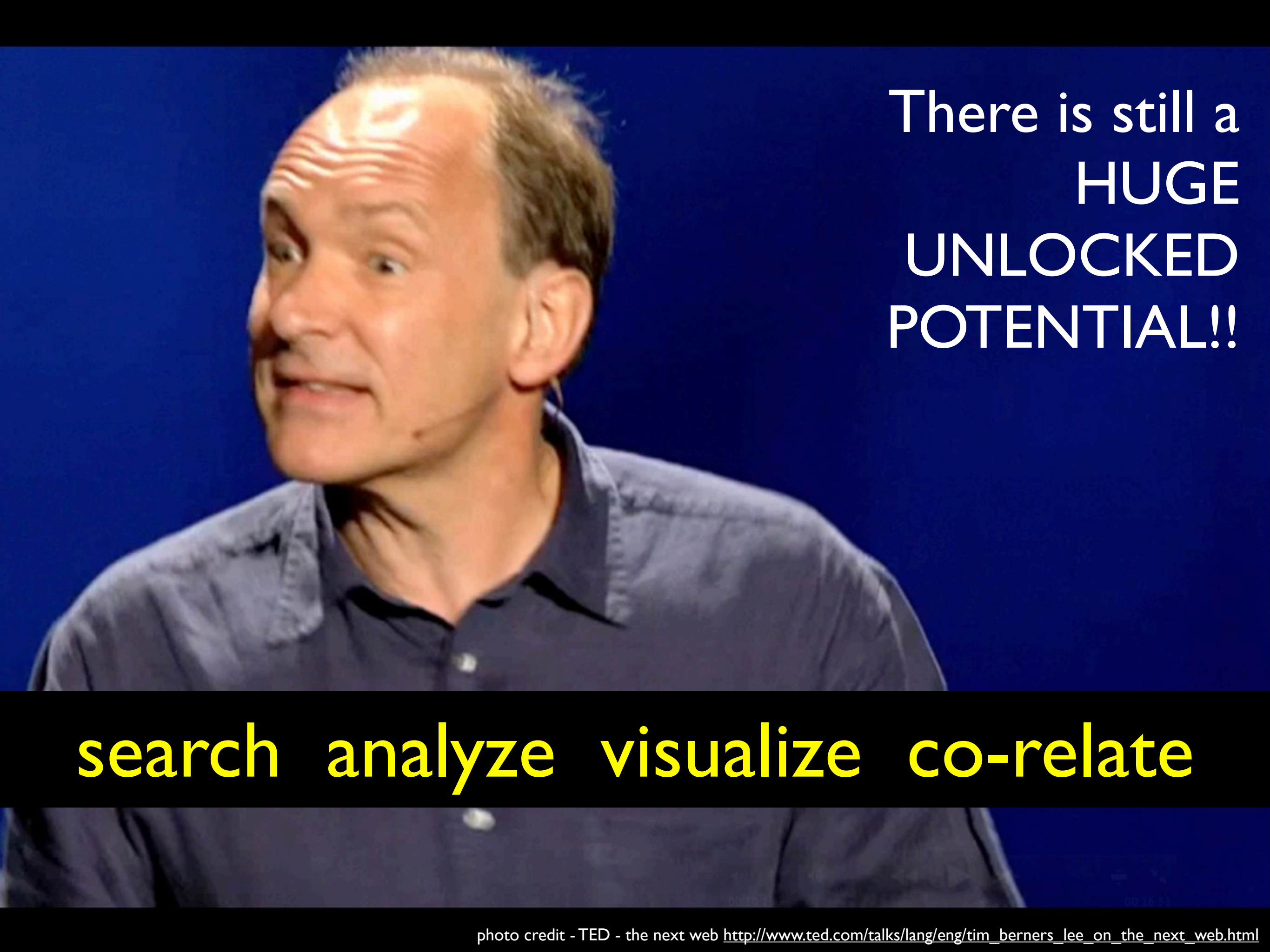
Abhishek Mishra  
[hello@ideamonk.com](mailto:hello@ideamonk.com)

Now... I want you to  
put your `_data_` on  
the web!





There is still a  
**HUGE**  
**UNLOCKED**  
**POTENTIAL!!**

A photograph of Tim Berners-Lee, a middle-aged man with thinning hair, wearing a dark blue button-down shirt. He is looking slightly to his left with a thoughtful expression. The background is a solid dark blue.

There is still a  
**HUGE**  
**UNLOCKED**  
**POTENTIAL!!**

**search analyze visualize co-relate**

Haven't got DATA  
on the web as  
DATA!



# Haven't got DATA on the web as DATA!



A screenshot of the GOI web directory website. The page is titled "GOI web directory" and shows a navigation menu on the left with categories like "Union Government", "States &amp; UTs", "Districts", "Legislature", "Judiciary", and "International". The main content area displays a list of districts in Karnataka, including Bagalkot, Bangalore Rural, Bangalore Urban, Belgaum, Bellary, Bidar, Bijapur, Chamarajanagar, Chickmagalur, Chikballapur, Chitradurga, Dakshina Kannada, Davangere, Dharwad, Gadag, Gulbarga, Hassan, Haveri, Kodagu, Kolar, Koppal, Mandya, Mysore, Raichur, Ramnagara, Shimoga, Tumkur, Udupi, Uttara Kannada (Karwar), and Yadgir. The page also features a search bar, a "Participate" section, and a "Last Updated" date of Sep 8, 2010.

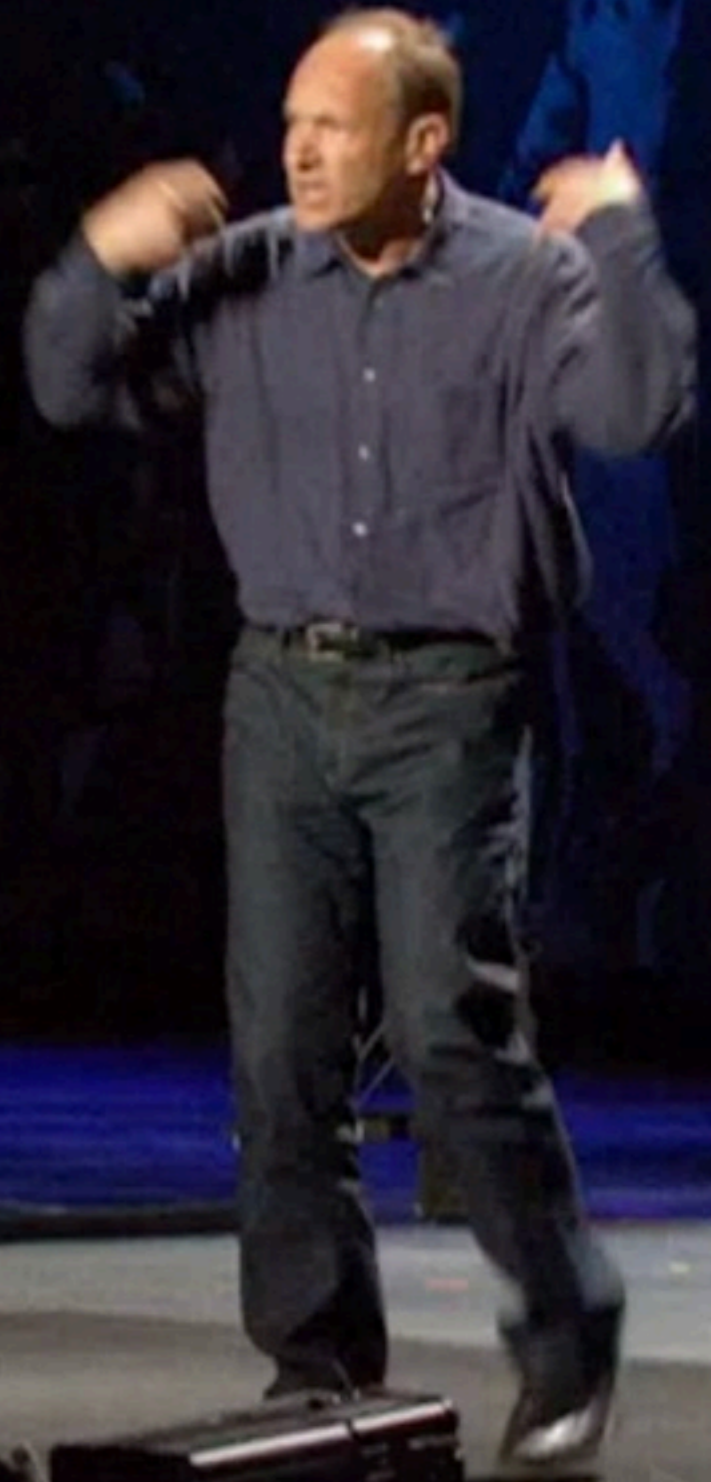
# Haven't got DATA on the web as DATA!



```
648
649
650 <ul style="display:none" id="pKR009" class="bg_link">
651
652 <li>
653
654 <a
655 href="javascript:openChild('http://goidirectory.nic.in
/sitecounter.php?id=1439','win2');"
class="heading_light_url"
title="http://chickmagalur.nic.in">Chickmagalur
District, Karnataka</a>
</li>
</ul>
```

# FRUSTRATION !!!

1. will take a lot of time
2. the data isn't consumable
3. you want it right now!
4. missing source





# Solution ?



**Solution ?**

**just scrape it out!**

A hand is shown using a white plastic scraper to remove a layer of white, fibrous material from a wooden cutting board. The material being scraped is thick and has a web-like, fibrous texture. The background is a light-colored wooden surface.

# Solution ?

just **scrape** it out!

1. analyze request
2. find patterns
3. extract data

# Solution ?

just **scrape** it out!

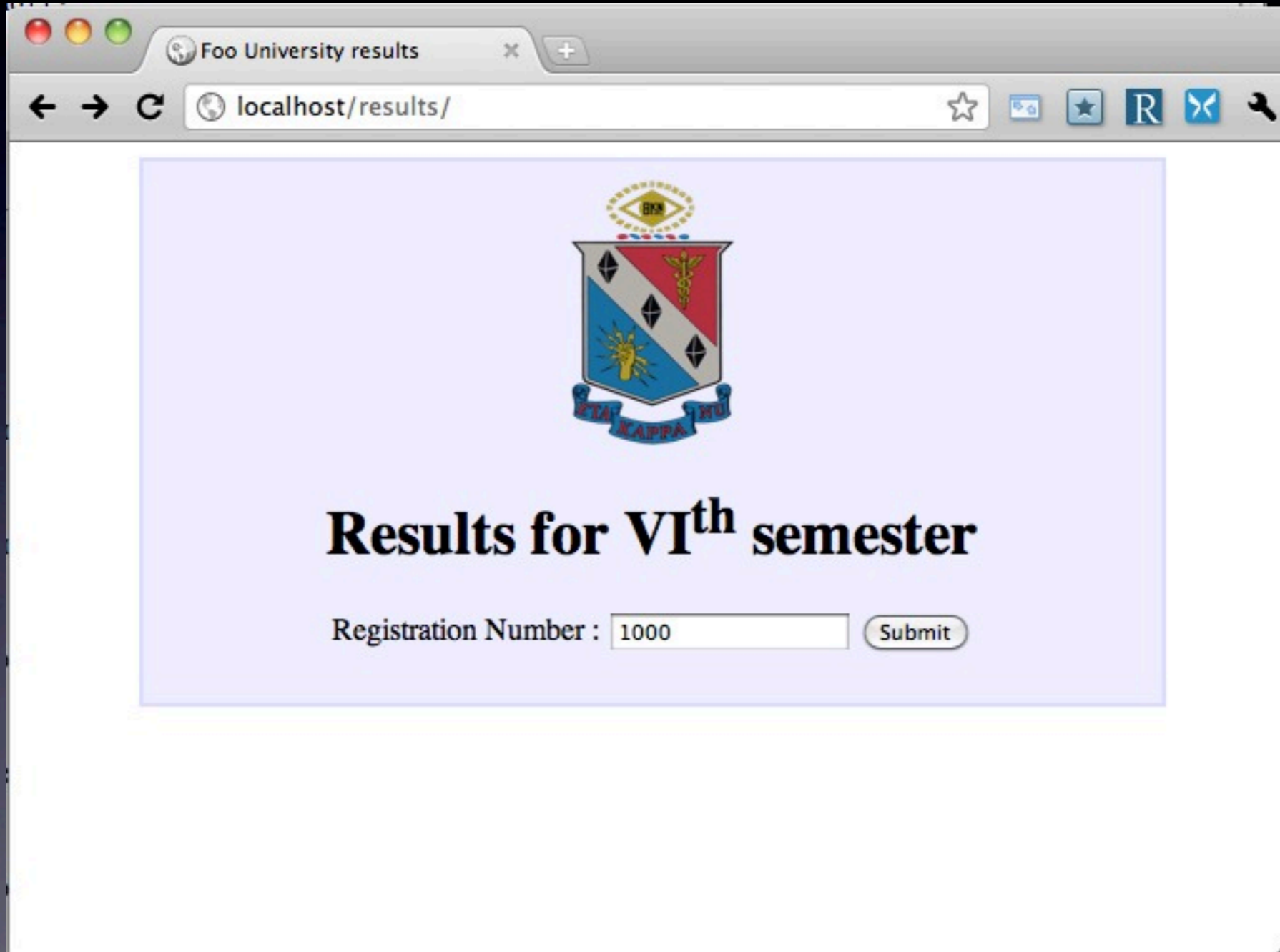
1. analyze request
2. find patterns
3. extract data

Use cases - **good** **bad** **ugly**

# Tools of Trade


- **urllib2**  
passing your data, getting back results
- **BeautifulSoup**  
parsing it out of the entangled web
- **Mechanize**  
programmatic web browsing
- **Scrapy**  
a web scraping framework

# Our objective



Foo University results

localhost/results/



## Results for VI<sup>th</sup> semester

Registration Number :

# Our objective

```
<body>
<div id="results">

<h1>Lohit Narayan
</h1>
<div>
Registration Number : 1234<br />
<table>
<tr>
<td>Object Oriented Programming</td>
<td> B+</td>
</tr>
<tr>
<td>Probability and Statistics</td>
<td> B+</td>
</tr>
<tr>
<td><b>CGPA</b></td>
<td><b>9.3</b></td>
</tr>
</table>
```

Name	Reg	CGPA
====	===	====
Lohit Narayan,	1234,	9.3
...	...	...

# urllib2

- Where to send?
- What to send?
- Answered by -  
firebug,  
chrome developer tools,  
a local debugging proxy



# urllib2

```
import urllib
import urllib2

url = "http://localhost/results/results.php"

data = { 'regno' : '1000' }
data_encoded = urllib.urlencode(data)

request = urllib2.urlopen( url, data_encoded )
html = request.read()
```

# urllib2

```
import urllib  
import urllib2
```

```
url = "http://localhost/results/results.php"
```

```
data = { 'regno' : '1000' }  
data_encoded = urllib.urlencode(data)
```

```
request = urllib2.urlopen( url, data_encoded )  
html = request.read()
```

```
'<html>\n<head>\n<title>Foo University results</title>\n<style type="text/css">\n<tbody>\n<tr>\n<td> B- </td>\n<td> Humanities Elective II </td>\n<td> A </td>\n<td> Discrete Mathematics </td>\n<td> C </td>\n<td> Environmental Science and Engineer  
ing </td>\n<td> B </td>\n<td> Data Structures </t  
d>\n<td> D </td>\n<td> Science Elective I </td>\n<td> D </td>\n<td> Object Oriented Programmin  
g </td>\n<td> B </td>\n<td> Probability and Stat  
istics </td>\n<td> A </td>\n<td> <b>CGPA</b> </td>\n<td> <b>9.1</b> </td>\n</tbody>\n</html>\n'
```

# urllib2

```
import urllib    # for urlencode
import urllib2   # main urllib2 module

url = "http://localhost/results/results.php"

data = { 'regno' : '1000' }
# ^ data as an object

data_encoded = urllib.urlencode(data)
# a string, key=val pairs separated by '&'

request = urllib2.urlopen( url, data_encoded )
# ^ returns a network object

html = request.read()
# ^ reads out its contents to a string
```

# urllib2

```
import urllib
import urllib2

url = "http://localhost/results/results.php"
data = {}
scraped_results = []

for i in xrange(1000, 3000):
    data['regno'] = str(i)
    data_encoded = urllib.urlencode(data)
    request = urllib2.urlopen(url, data_encoded)
    html = request.read()
    scraped_results.append(html)
```

```
'<html>\n<head>\n<title>Foo University results</title>\n<style type="text/css">\n<#results {\n  margin:0px 60px;\n  background: #efefff;\n  padding:10px;\n  border:2px solid #dddff;\n  text-align:center;\n}</style>\n</head>\n<body>\n<div id="results">\n  \n  <h1>Rahim Kumar</h1>\n  <div>\n    <tRegistration Number : 1000<br/>\n    <table>\n      <tr>\n        <td>Humanities Elective II</td>\n        <td> B-</td>\n      </tr>\n      <tr>\n        <td>Humanities Elective III</td>\n        <td> A</td>\n      </tr>\n      <tr>\n        <td>Discrete Mathematics</td>\n        <td> C-</td>\n      </tr>\n      <tr>\n        <td>Environmental Science and Engineering</td>\n        <td> B</td>\n      </tr>\n      <tr>\n        <td>Data Structures</td>\n        <td> D</td>\n      </tr>\n      <tr>\n        <td>Science Elective I</td>\n        <td> D-</td>\n      </tr>\n      <tr>\n        <td>Object Oriented Programming</td>\n        <td> B</td>\n      </tr>\n      <tr>\n        <td>Probability and Statistics</td>\n        <td> A-</td>\n      </tr>\n      <tr>\n        <td><b>CGPA</b></td>\n        <td><b>9.1</b></td>\n      </tr>\n    </table>\n  </div>\n</div>\n</body>\n</html>\n'
```

Thats alright,  
But this is still meaningless.



# BeautifulSoup



# BeautifulSoup

- Python HTML/XML Parser
- Wont choke you on **bad markup**
- auto-detects encoding
- easy tree traversal
- find all links whose url contains 'foo.com'



# BeautifulSoup

```
>>> import re
>>> from BeautifulSoup import BeautifulSoup

>>> html = """<a href="http://foo.com/foobar">link1</a> <a
href="http://foo.com/beep">link2</a> <a href="http://
foo.co.in/beep">link3</a>"""

>>> soup = BeautifulSoup(html)
>>> links = soup.findAll('a', {'href': re.compile
(".*foo.com.*")})
>>> links
[<a href="http://foo.com/foobar">link1</a>, <a href="http://
foo.com/beep">link2</a>]

>>> links[0]['href']
u'http://foo.com/foobar'

>>> links[0].contents
[u'link1']
```

# BeautifulSoup

```
>>> html = """<html>\n <head>\n  <title>\n    Foo University results\n  </title>\n  <style type="text/css">\n    #results {\n      \t\t\tmargin:0px 60px;\n      \t\t\tbackground: #efefff;\n      \t\t\tpadding:10px;\n      \t\t\tborder:2px solid #ddddff;\n      \t\t\ttext-align:center;\n      \t\t\t}\n    </style>\n  </head>\n  <body>\n    <div id="results">\n      \n      <h1>\n        Bipin Narayan\n      </h1>\n      <div>\n        Registration Number : 1000\n        <br />\n        <table>\n          <tr>\n            <td>\n              Humanities Elective II\n            </td>\n            <td>\n              C-\n            </td>\n          </tr>\n          <tr>\n            <td>\n              Humanities Elective III\n              Discrete Mathematics\n              Environmental Science and Engineering\n              Data Structures\n              I\n            </td>\n            <td>\n              B-\n            </td>\n          </tr>\n          <tr>\n            <td>\n              B+\n            </td>\n            <td>\n              A\n            </td>\n          </tr>\n          <tr>\n            <td>\n              <b>\n                8\n              </b>\n            </td>\n            <td>\n              <b>\n                CGPA\n              </b>\n            </td>\n          </tr>\n        </table>\n      </div>\n    </div>\n  </body>\n</html>"""
```

# BeautifulSoup

```
>>> html = """<html>\n <head>\n <title>\n  Foo University results\n </title>\n <style type="text/css">\n  #results {\n    \t\t\tmargin:0px 60px;\n    \t\t\tbackground: #efefff;\n    \t\t\tpadding:10px;\n    \t\t\tborder:2px solid #ddddff;\n    \t\t\ttext-align:center;\n    \t\t\t}\n </style>\n </head>\n <body>\n <div id="results">\n \n <h1>\n  Bipin Narayan\n </h1>\n <div>\n  Registration Number : 1000\n <br />\n <table>\n <tr>\n <td>\n  Humanities Elective II\n </td>\n <td>\n  C-\n </td>\n </tr>\n <tr>\n <td>\n  Humanities Elective III\n </td>\n <td>\n  B-\n </td>\n </tr>\n <tr>\n <td>\n  Discrete Mathematics\n </td>\n <td>\n  C\n </td>\n </tr>\n <tr>\n <td>\n  Environmental Science and Engineering\n </td>\n <td>\n  C\n </td>\n </tr>\n <tr>\n <td>\n  Data Structures\n </td>\n <td>\n  A\n </td>\n </tr>\n <tr>\n <td>\n  I\n </td>\n <td>\n  B-\n </td>\n </tr>\n <tr>\n <td>\n <td>\n  B+\n </td>\n </tr>\n <tr>\n <td>\n <td>\n  A\n </td>\n </tr>\n <tr>\n <td>\n <td>\n  <b>\n  8\n </b>\n </td>\n </tr>\n </table>\n </div>\n </div>\n </body>\n</html>"""
```

```
>>> soup.h1
```

```
<h1>Bipin Narang</h1>
```

```
>>> soup.h1.contents
```

```
[u'Bipin Narang']
```

```
>>> soup.h1.contents[0]
```

```
u'Bipin Narang'
```

# BeautifulSoup

```
>>> html = """<html>\n <head>\n  <title>\n    Foo University results\n  </title>\n  <style type="text/css">\n    #results {\n      \t\t\tmargin:0px 60px;\n      \t\t\tbackground: #efefff;\n      \t\t\tpadding:10px;\n      \t\t\tborder:2px solid #ddddff;\n      \t\t\ttext-align:center;\n      \t\t\t}\n    </style>\n  </head>\n  <body>\n    <div id="results">\n      \n      <h1>\n        Bipin Narayan\n      </h1>\n      <div>\n        Registration Number : 1000\n        <br />\n        <table>\n          <tr>\n            <td>\n              Humanities Elective II\n            </td>\n            <td>\n              C-\n            </td>\n          </tr>\n          <tr>\n            <td>\n              Humanities Elective III\n              Discrete Mathematics\n              Environmental Science and Engineering\n              Data Structures\n              I\n            </td>\n            <td>\n              B-\n            </td>\n          </tr>\n          <tr>\n            <td>\n              B+\n            </td>\n            <td>\n              A\n            </td>\n          </tr>\n          <tr>\n            <td>\n              <b>\n                8\n              </b>\n            </td>\n            <td>\n              <b>\n                CGPA\n              </b>\n            </td>\n          </tr>\n        </table>\n      </div>\n    </div>\n  </body>\n</html>"""
```

```
>>> soup.div.div.contents[0]
u'\n\t\tRegistration Number : 1000\t\t'
>>> soup.div.div.contents[0].split(':')
[u'\n\t\tRegistration Number ', u' 1000\t\t']
>>> soup.div.div.contents[0].split(':')[1].strip()
u'1000'
```

# BeautifulSoup

```
>>> html = """<html>\n <head>\n  <title>\n    Foo University results\n  </title>\n  <style type="text/css">\n    #results {\n      \t\t\tmargin:0px 60px;\n      \t\t\tbackground: #efefff;\n      \t\t\tpadding:10px;\n      \t\t\tborder:2px solid #ddddff;\n      \t\t\ttext-align:center;\n      \t\t\t}\n    </style>\n  </head>\n  <body>\n    <div id="results">\n      \n      <h1>\n        Bipin Narayan\n      </h1>\n      <div>\n        Registration Number : 1000\n        <br />\n        <table>\n          <tr>\n            <td>\n              Humanities Elective II\n            </td>\n            <td>\n              C-\n            </td>\n          </tr>\n          <tr>\n            <td>\n              Humanities Elective III\n              Discrete Mathematics\n              Environmental Science and Engineering\n              Data Structures\n              I\n            </td>\n            <td>\n              B-\n            </td>\n          </tr>\n          <tr>\n            <td>\n              B+\n            </td>\n            <td>\n              A\n            </td>\n          </tr>\n          <tr>\n            <td>\n              <b>\n                8\n              </b>\n            </td>\n            <td>\n              <b>\n                CGPA\n              </b>\n            </td>\n          </tr>\n        </table>\n      </div>\n    </div>\n  </body>\n</html>"""
```

```
>>> soup.find('b')
```

```
<b>CGPA</b>
```

```
>>> soup.find('b').findNext('b')
```

```
<b>8</b>
```

```
>>> soup.find('b').findNext('b').contents[0]
```

```
u'8'
```

# urllib2 + BeautifulSoup

```
import urllib
import urllib2
from BeautifulSoup import BeautifulSoup

url = "http://localhost/results/results.php"
data={}
for i in xrange(1000,1100):
    data['regno'] = str(i)
    html = urllib2.urlopen( url, urllib.urlencode(data) ).read()
    soup = BeautifulSoup(html)
    name = soup.h1.contents[0]
    regno = soup.div.div.contents[0].split(':')[1].strip()
    # ^ though we already know current regno ;) == str(i)
    cgpa = soup.find('b').findNext('b').contents[0]
    print "%s,%s,%s" % (regno, name, 'cgpa')
```

# Mechanize



# Mechanize



Stateful web browsing in Python

after Andy Lester's WWW:Mechanize

fill forms  
follow links  
handle cookies  
browser history



# Mechanize

```
>>> import mechanize
>>> br = mechanize.Browser()
>>> response = br.open("http://google.com")
>>> br.title()
'Google'
>>> br.geturl()
'http://www.google.co.in/'
>>> print response.info()
Date: Fri, 10 Sep 2010 04:06:57 GMT
Expires: -1
Cache-Control: private, max-age=0
Content-Type: text/html; charset=ISO-8859-1
Server: gws
X-XSS-Protection: 1; mode=block
Connection: close
content-type: text/html; charset=ISO-8859-1
>>> br.open("http://yahoo.com")
>>> br.title()
'Yahoo! India'
>>> br.back()
>>> br.title()
'Google'
```

# Mechanize

## filling forms

```
>>> import mechanize
>>> br = mechanize.Browser()
>>> response = br.open("http://in.pycon.org/2010/account/
login")
>>> html = response.read()
>>> html.find('href="/2010/user/ideamonk">Abhishek I')
-1

>>> br.select_form(name="login")
>>> br.form["username"] = "ideamonk"
>>> br.form["password"] = "*****"
>>> response2 = br.submit()
>>> html = response2.read()
>>> html.find('href="/2010/user/ideamonk">Abhishek I')
10186
```

# Mechanize

## filling forms

```
>>> import mechanize
>>> br = mechanize.Browser()
>>> response = br.open("http://in.pycon.org/2010/account/
login")
>>> html = response.read()
>>> html.find('href="/2010/user/ideamonk">Abhishek I')
-1

>>> br.select_form(name="login")
>>> br.form["username"] = "ideamonk"
>>> br.form["password"] = "*****"
>>> response2 = br.submit()
>>> html = response2.read()
>>> html.find('href="/2010/user/ideamonk">Abhishek I')
10186
```

# Mechanize

logging into gmail

```
<!-- login box -->
▼ <div id="login">
  ▶ <script>...</script>
  ▶ <style type="text/css">...</style>
  ▶ <style type="text/css">...</style>
  ▼ <form id="gaia_loginform" action="https://
    "return(gaia_onLoginSubmit());">
    ▼ <div id="gaia_loginbox">
      ▼ <table class="form-noindent" cellpadding="0" cellspacing="0">
        ▼ <tbody>
          ▼ <tr>
            ▼ <td valign="top" style="text-align: left; padding-right: 10px;">
```

# Mechanize

logging into gmail

```
▼ <form id="gaia_loginform" action="https:  
  "return(gaia_onLoginSubmit());">
```

the form has no name!  
cant **select\_form** it

# Mechanize

## logging into gmail

```
>>> br.open("http://mail.google.com")
>>> login_form = br.forms().next()
>>> login_form["Email"] = "ideamonk@gmail.com"
>>> login_form["Passwd"] = "*****"
>>> response = br.open( login_form.click() )
>>> br.geturl()
'https://mail.google.com/mail/?shva=1'
>>> response2 = br.open("http://mail.google.com/mail/h/")
>>> br.title()
'Gmail - Inbox'
>>> html = response2.read()
```

# Mechanize

## reading my inbox

```
>>> soup = BeautifulSoup(html)
>>> trs = soup.findAll('tr', {'bgcolor': '#ffffff'})
>>> for td in trs:
...     print td.b
...
<b>madhukargunjan</b>
<b>NIHI</b>
<b>Chamindra</b>
<b>KR2 blr</b>
<b>Imran</b>
<b>Ratnadeep Debnath</b>
<b>yasir</b>
<b>Glenn</b>
<b>Joyent, Carly Guarcello</b>
<b>Zack</b>
<b>SahanaEden</b>
...
>>>
```

# Mechanize



Can scrape my gmail chats !?



# Scrapy



a **Framework** for web scraping

# Scrapy



uses XPath  
navigate through elements

# Scrapy

## XPath and Scrapy

```
<h1> My important chunk of text </h1>  
//h1/text()
```

```
<div id="description"> Foo bar beep &gt; 42 </div>  
//div[@id='description']
```

```
<div id='report'> <p>Text1</p> <p>Text2</p> </div>  
//div[@id='report']/p[2]/text()
```

```
// select from the document, wherever they are
```

```
/ select from root node
```

```
@ select attributes
```

scrapy.selector.XmlXPathSelector

scrapy.selector.HtmlXPathSelector

<http://www.w3schools.com/xpath/default.asp>

# Scrapy



## Interactive Scraping Shell!

```
$ scrapy-ctl.py shell http://my/target/website
```

# Scrapy



## Step 1

define a model to store **Items**

# Scrapy



## Step II

create your **Spider** to extract **Items**

# Scrapy



## Step III

write a **Pipeline** to store them

# Scrapy

Lets write a Scrapy spider



# SpojBackup

- Spoj? Programming contest website
- 4021039 code submissions
- 81386 users
- web scraping == backup tool
- uses Mechanize
- <http://github.com/ideamonk/spojbackup>

# SpojBackup

- <https://www.spoj.pl/status/ideamonk/signedlist/>

```
*****
*      Signed document issued by Sphere Online Judge  (http://www.spoj.pl)
*****

----- BEGIN OF DATA -----
Submissions by [ideamonk] in contest [SPOJ].  DoI: 2010-09-10 12:18:26.
/-----\
```

ID	DATE	PROBLEM	RESULT	TIME	MEM	LNG
3198919	2010-01-22 01:00:45	BOXES	AC	0.07	2544	C++
3023085	2009-11-29 14:54:36	ACODE	AC	0.46	5832	PYT
3023081	2009-11-29 14:54:15	ACODE	WA	0.68	5844	PYT
3023078	2009-11-29 14:53:38	ACODE	AC	0.47	5836	PYT
3023066	2009-11-29 14:50:22	ACODE	AC	0.37	5560	PYT
3023061	2009-11-29 14:50:00	ACODE	CE	0.00	0	PYT
3023030	2009-11-29 14:46:01	ACODE	AC	0.54	5556	PYT
3022963	2009-11-29 14:36:16	ACODE	RE	0.14	4124	PYT

- <https://www.spoj.pl/files/src/save/3198919>
- Code walkthrough

# Conclusions



# Conclusions

Python - large set of tools for scraping

Choose them wisely

Do not steal



# Conclusions

Python - large set of tools for scraping

Choose them wisely

Do not steal

## License and Site Access

IMDb grants you a limited license to access and make personal use of this site and not to download (other than page caching) or modify it, or any portion of it, except with express written consent of IMDb. This site or any portion of this site may not be reproduced, duplicated, copied, sold, resold, visited, or otherwise exploited for any commercial purpose without express written consent of IMDb. This license does not include any resale or commercial use of this site or its contents or any derivative use of this site or its contents.

**Robots and Screen Scraping:** You may not use data mining, robots, screen scraping, or similar data gathering and extraction tools on this site, except with our express written consent as noted below.

**Framing:** You may not frame or utilize framing techniques to enclose any trademark, logo, or other proprietary information (including images, text, page layout, or form) of IMDb without express written consent.

# Conclusions

Python - large set of tools for scraping

Choose them wisely

Do not steal

Use the cloud - [http://docs.picloud.com/advanced\\_examples.html](http://docs.picloud.com/advanced_examples.html)

Share your scrapers - <http://scraperwiki.com/>

Thanks to @amatix , flickr, you, etc

@PyCon India 2010  
<http://in.pycon.org>

Abhishek Mishra  
[hello@ideamonk.com](mailto:hello@ideamonk.com)