

BI 前的悲哀

- 那些User所謂的**Big(?)** data

The Sadness before Business Intelligence:
When User says **Big** data...

Bingroom

About Me

1. Inspired by PyCon



2. Got **donated** in PyCon



3. **Contribute** to PyCon

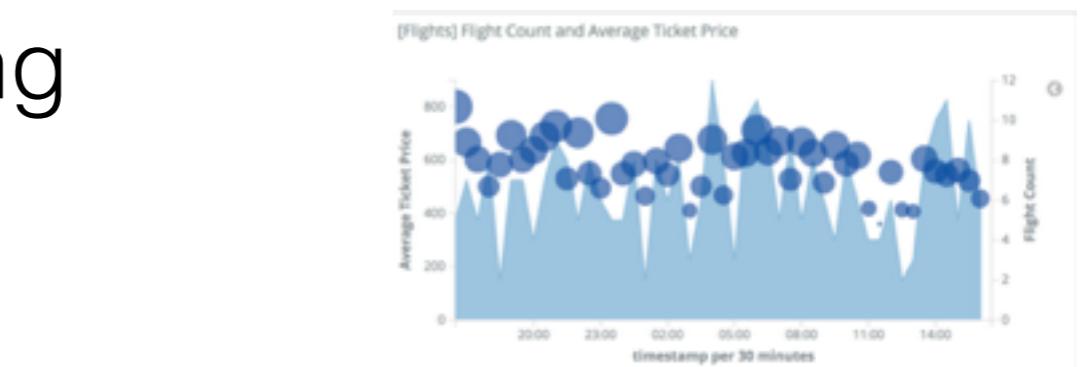
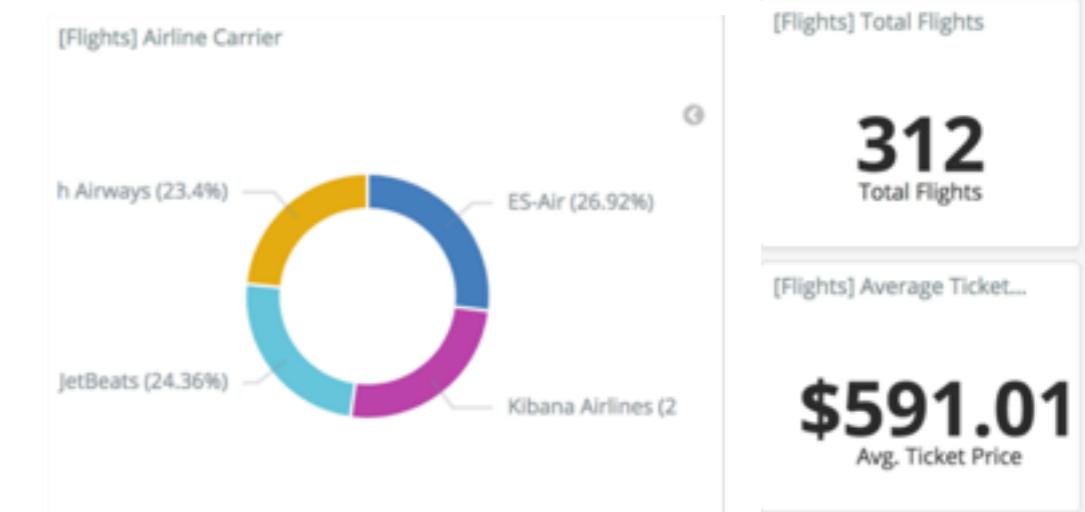
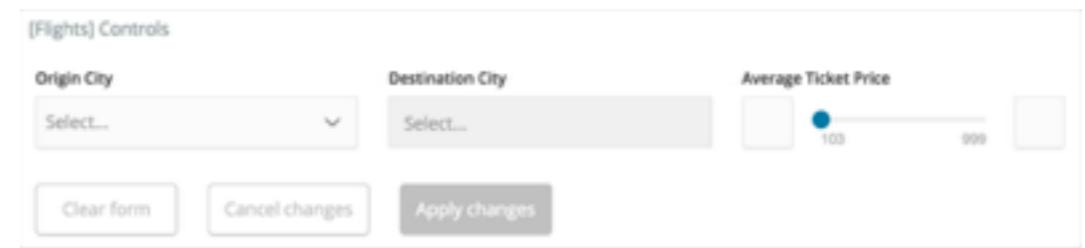


Outline

- BI
- ETL
- Data parser
 - Common document file extensions
 - Software for developers
- Case sharing

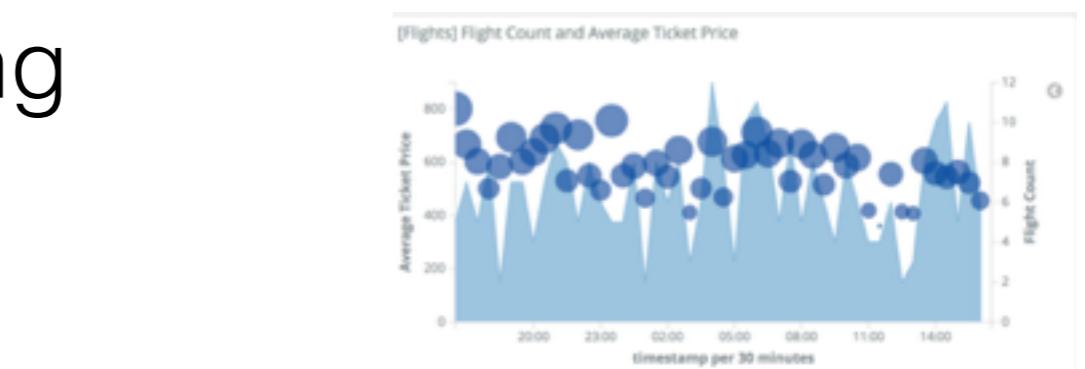
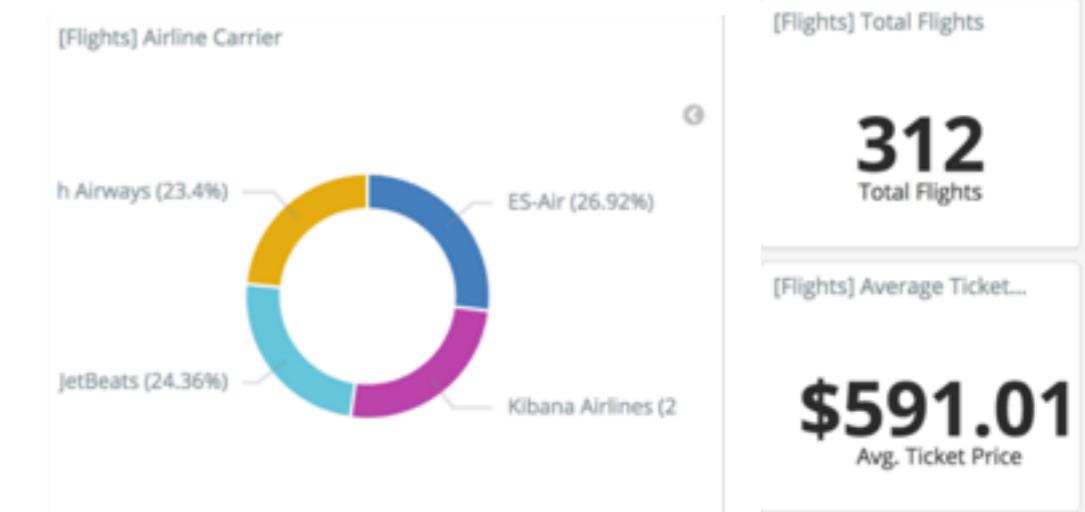
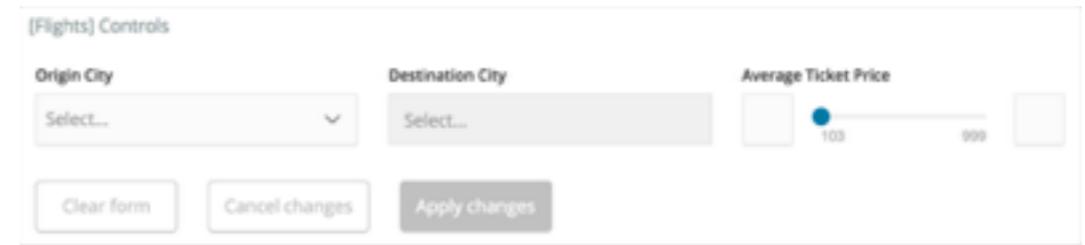
Business Intelligence

- Know business
- Decision making
- Real-time navigation & monitoring
- Make things easier & sharable



Business Intelligence

- Know business
- Decision making
- Real-time navigation & monitoring
- Make things easier & **sharable**

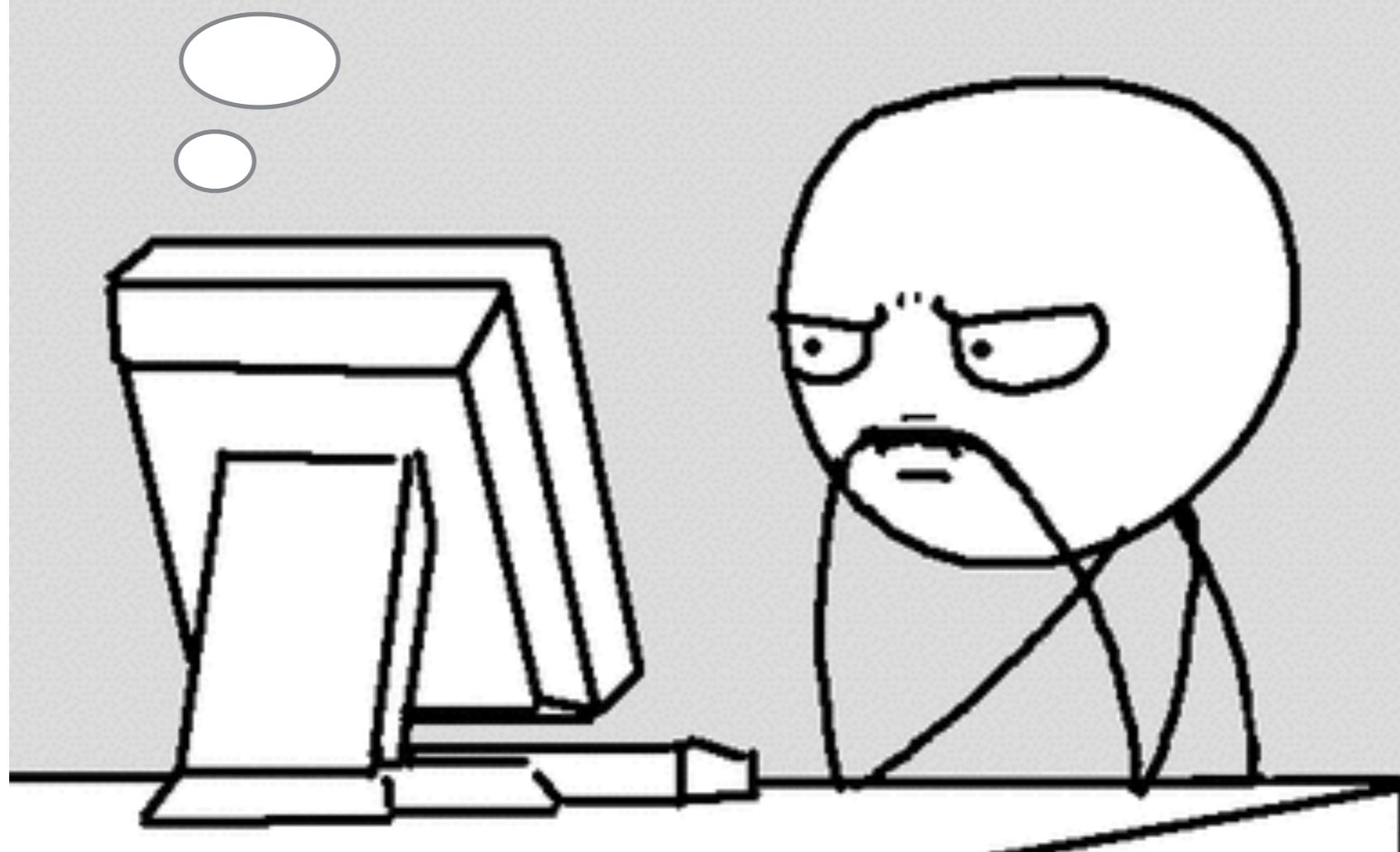


Before the Fall

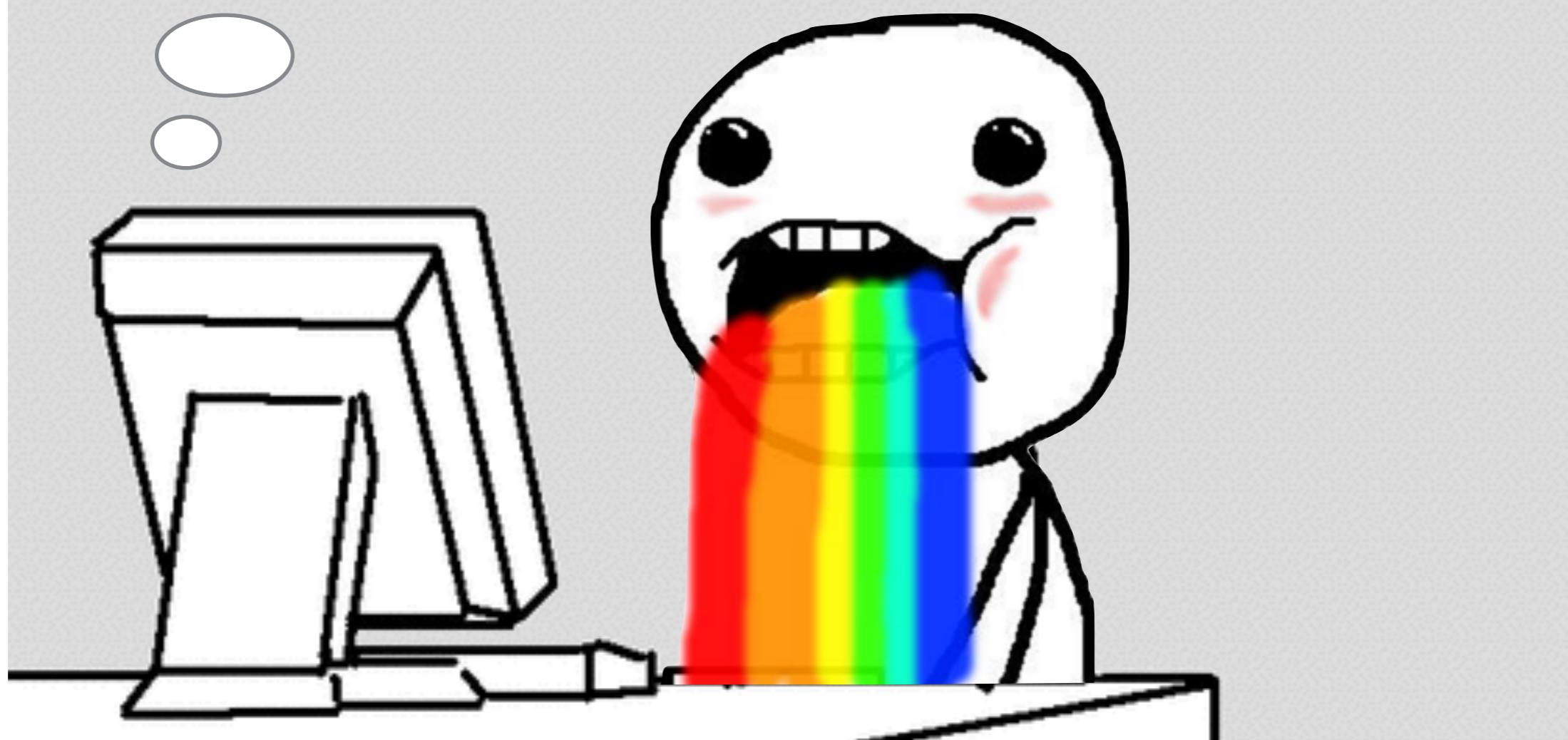
As a Data Engineer

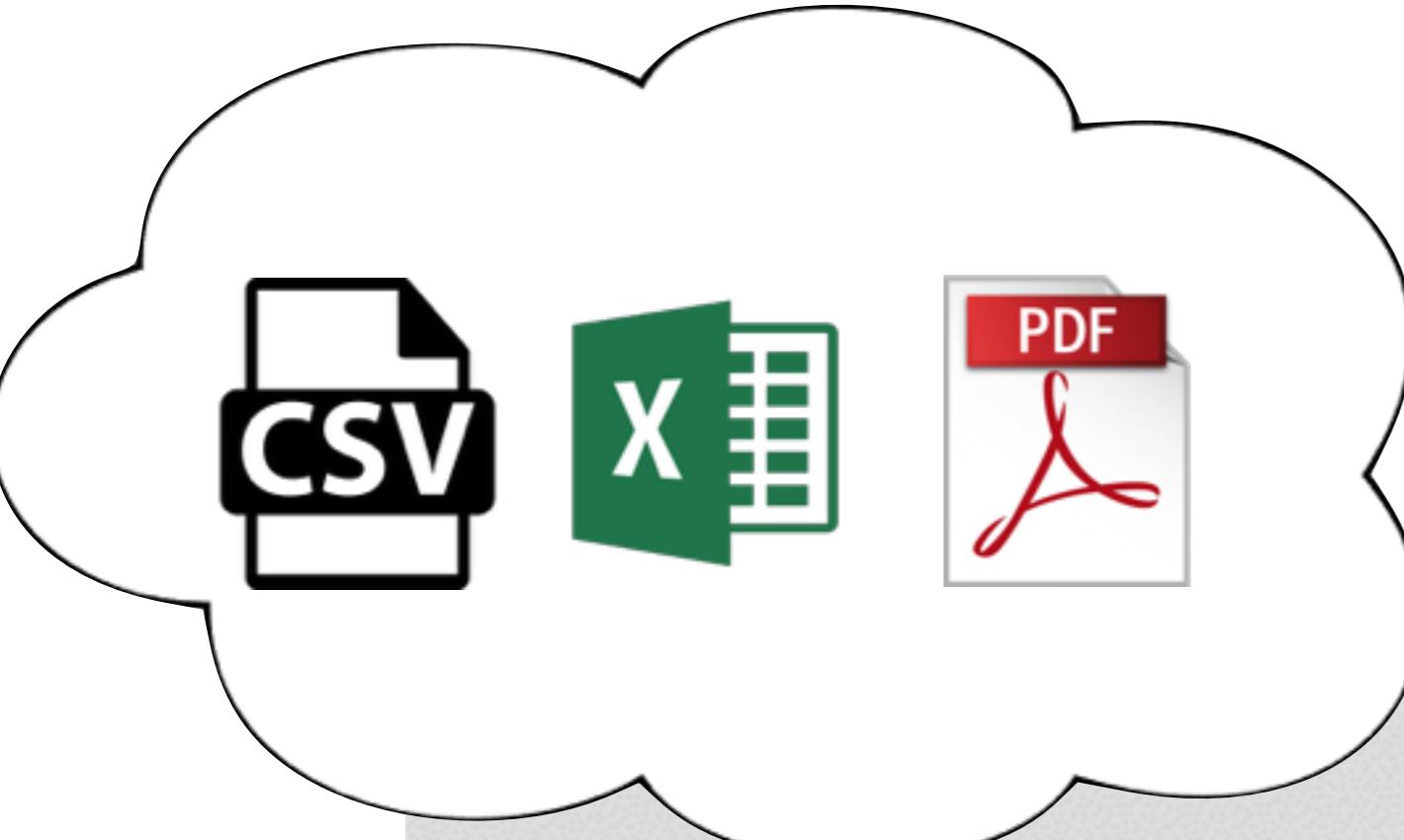


Google BigQuery



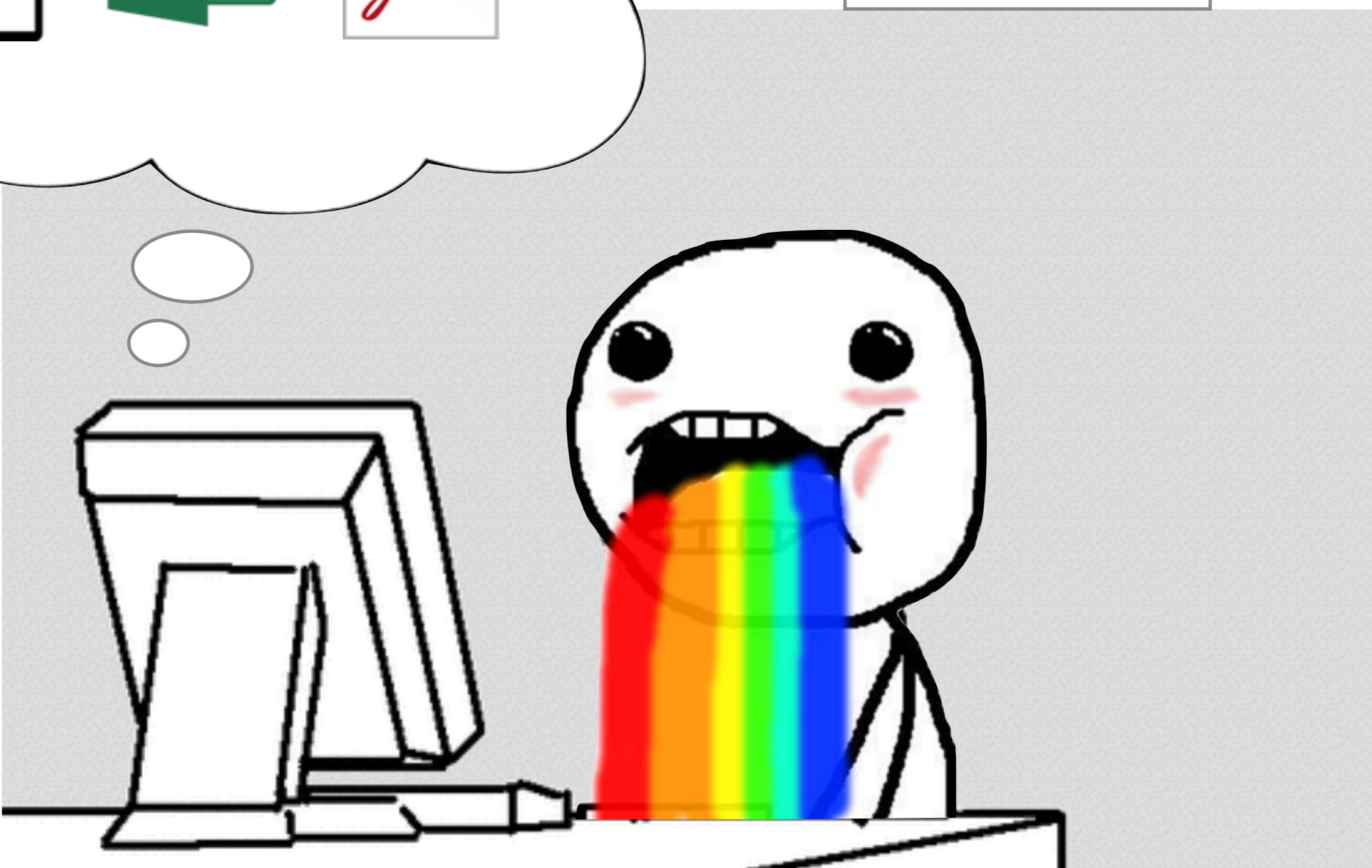
**THREE
WEEKS LATER**





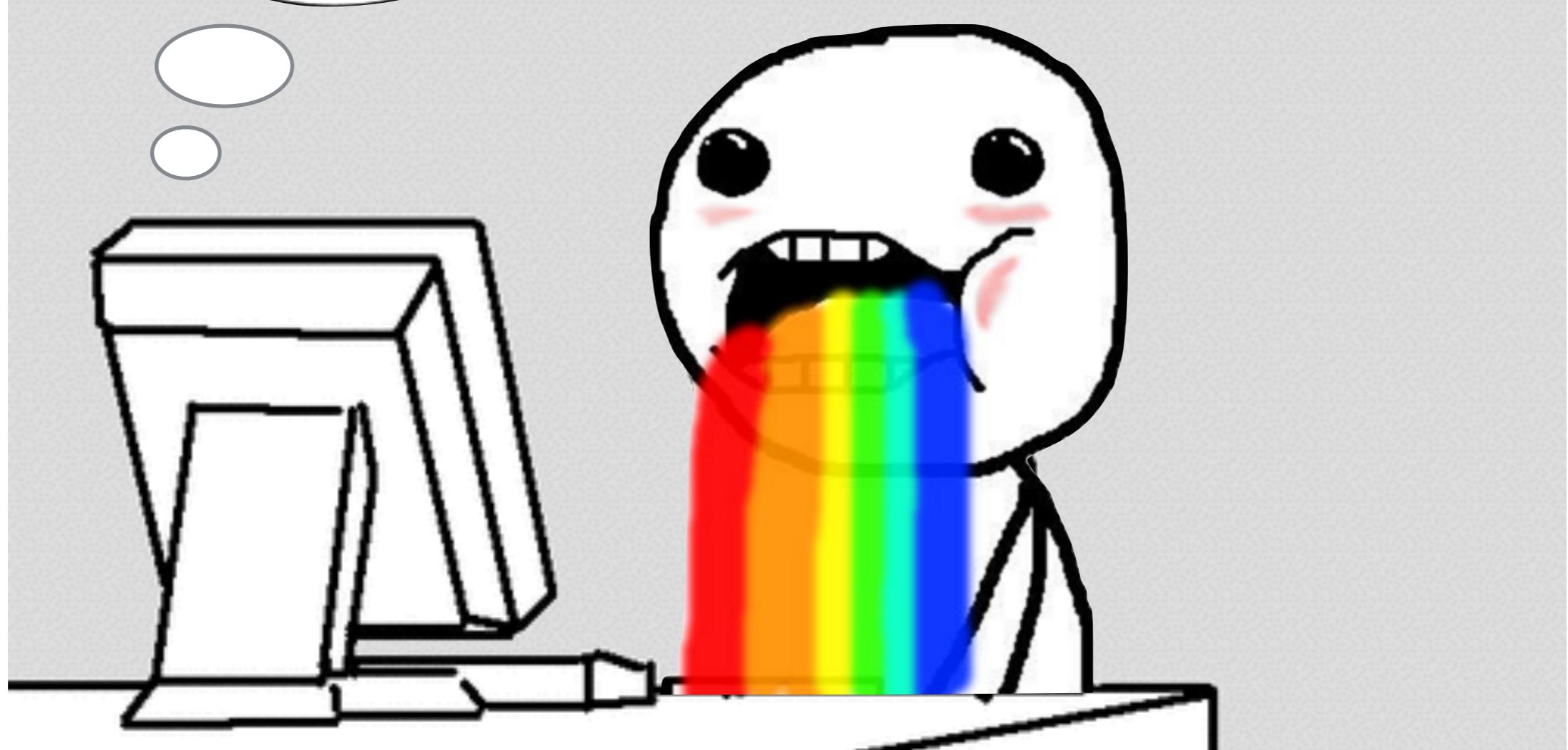
Big Data:

1. Volume
2. Variety
3. Velocity



Big(?) Data:

1. ~~Volume~~
2. Variety
3. Velocity



Extract-Transform-Load

Extract

- Load DB tables
- Read CSV files
- Web crawlers
- **import** requests

Extract-Transform-Load

Extract

- Load DB tables
- Read CSV files
- Web crawlers
- **import** requests

Transform

- Data clean
 - ✗ incomplete
 - ✗ inaccurate
 - ✗ irrelevant
- Data transform
 - ✓ re-index
 - ✓ aggregation
 - ✓ transpose

Extract-Transform-Load

Extract

- Load DB tables
- Read CSV files
- Web crawlers
- **import** requests

Transform

- Data clean
 - ✗ incomplete
 - ✗ inaccurate
 - ✗ irrelevant
- Data transform
 - ✓ re-index
 - ✓ aggregation
 - ✓ transpose

Load

- Data warehouse
- SQL / NoSQL
- **from sqlalchemy import ***

Extract-Transform-Load

Extract

- Load DB tables
- Read CSV files
- Web crawlers
- **import** requests

Transform

- Data clean
 - ✗ incomplete
 - ✗ inaccurate
 - ✗ irrelevant
- Data transform
 - ✓ re-index
 - ✓ aggregation
 - ✓ transpose

Load

- Data warehouse
- SQL / NoSQL
- **from sqlalchemy import ***



Data Transform (Preprocessing)



Don't reinvent the wheel!

A short recap...

Bingroom，迷途中的一個Python小書僮，由於之前Crawler寫太多，現在都在寫Parser還債。

Data Parser





Excel

- Similar to CSV - Lots of Python modules
- For Windows user: **win32com**
- Read only: **xlrd**
- Read & Write: **openpyxl** maybe **xlwings** for VBA lover?
- Scientific computing: **pandas**



Excel

Compatibility

	Win	Mac	Py2	Py3	xlsx	安装
win32com	✓	✗	✓	✓	✓	pip
xlwings	✓	✓	✓	✓	✓	pip
xlsxwriter	✓	✓	✓	✓	✓	pip
DataNitro	✓	✗	✓	✓	✓	安装包
pandas	✓	✓	✓	✓	✓	pip
openpyxl	✓	✓	✓	✓	✓	pip
xlutils	✓	✓	✓	✓	✗	pip

Task

	打开文档	新建文档	修改文档	保存文档
win32com	✓	✓	✓	✓
xlwings	✓	✓	✓	✓
xlsxwriter	✗	✓	✗	✓
pandas	✓	✗	✓	✓
openpyxl	✓	✓	✓	✓
xlutils	✓	✓	✓	✓

Efficiency

	写入用时(s)	读取用时(s)
win32com	3	2
xlwings	5	1.5
xlsxwriter	14	
DataNitro	21	
pandas	22	12
openpyxl	32	17
xlutils		11

CrossIn的编程教室

<https://zhuanlan.zhihu.com/p/23998083>



Excel

(In Practice)

- Background color of cell

Account	Region	2018	2018
		Jan	Feb
PyCon TW	TW	0	0
PyCon JP	JP	0	1

```
from openpyxl import load_workbook
wb = load_workbook('sample.xlsx')
ws = wb['工作表1']
cell_A5 = ws['A5'].fill.start_color.index[2:]
cell_C5 = ws['C5'].fill.start_color.index[2:]
cell_A5, cell_C5
('FFC000', '66FF66')
```

#FFC000

#66FF66



Excel

(In Practice)

- Transpose



Excel

(In Practice)

fixed_cols

excel_df

scan_cols

	Account	Region	2018-01	2018-02	2018-03	2018-04	2018-05
0	PyCon TW	TW	0	0	0	10	11
1	PyCon JP	JP	0	1	0	3	0

- Transpose

- Take apart df.columns
=> **fixed_cols** and **scan_cols**
- Fixed: keep value and repeat
- Scan: get value by key
then append to row

(Repeat) (key) (value)

Account Region Date Commit

0	PyCon TW	TW	2018-01	0
1	PyCon JP	JP	2018-01	0
2	PyCon TW	TW	2018-02	0
3	PyCon JP	JP	2018-02	1
4	PyCon TW	TW	2018-03	0
5	PyCon JP	JP	2018-03	0
6	PyCon TW	TW	2018-04	10
7	PyCon JP	JP	2018-04	3
8	PyCon TW	TW	2018-05	11



Excel

(In Practice)

- Transpose

- pandas.melt()



```
df = pd.melt(excel_df, id_vars=['Account', 'Region'], var_name='Date', value_name='Commit')
```

	Account	Region	Date	Commit
0	PyCon TW	TW	2018-01	0
1	PyCon JP	JP	2018-01	0
2	PyCon TW	TW	2018-02	0
3	PyCon JP	JP	2018-02	1



Excel

(In Practice)



CONFIRMED.

A photograph of a brown bear standing in a grassy, overgrown field. The bear is facing towards the right of the frame. Overlaid on the upper portion of the image is the word "CONFIRMED." in large, bold, black letters with a white outline.

I just reinvented a wheel!

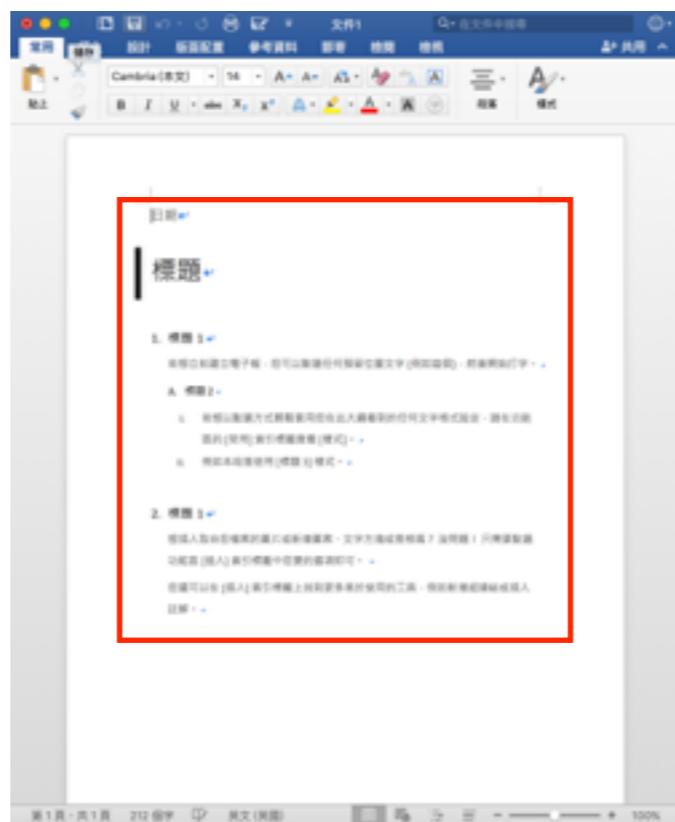
Word

- Before talking about



- Let's extract

first (or **text**, precisely)



Text Extraction

- On Linux: **libreoffice**

```
libreoffice --headless --convert-to "txt:Text (encoded):UTF8" mydocument.doc
```

<https://ask.libreoffice.org/en/question/2641/convert-to-command-line-parameter/>

- Cons:

- High memory usage in batch processing
- No root, no apt-get

Text Extraction

- **textract**: all-mighty text extractor
- **textract is not tesseract**
and
tesseract in textract

```
text = txtract.process(  
    'path/to/norwegian.pdf',  
    method='tesseract',  
    language='nor',  
)
```

COMPLETE
EDITION

DEFINITIVE
EDITION

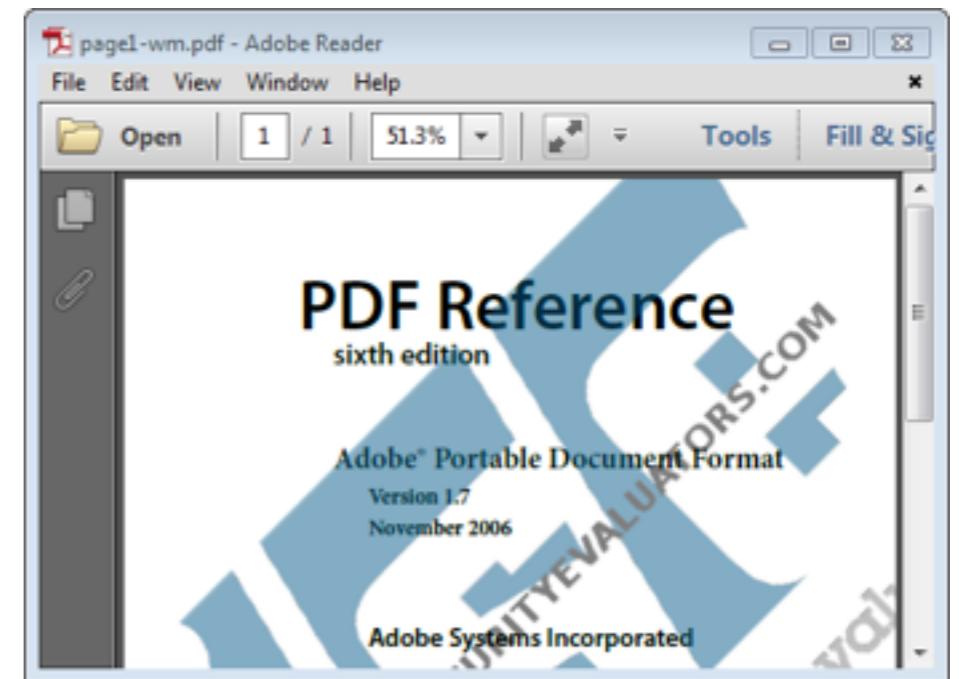
ROYAL EDITION

Currently supporting

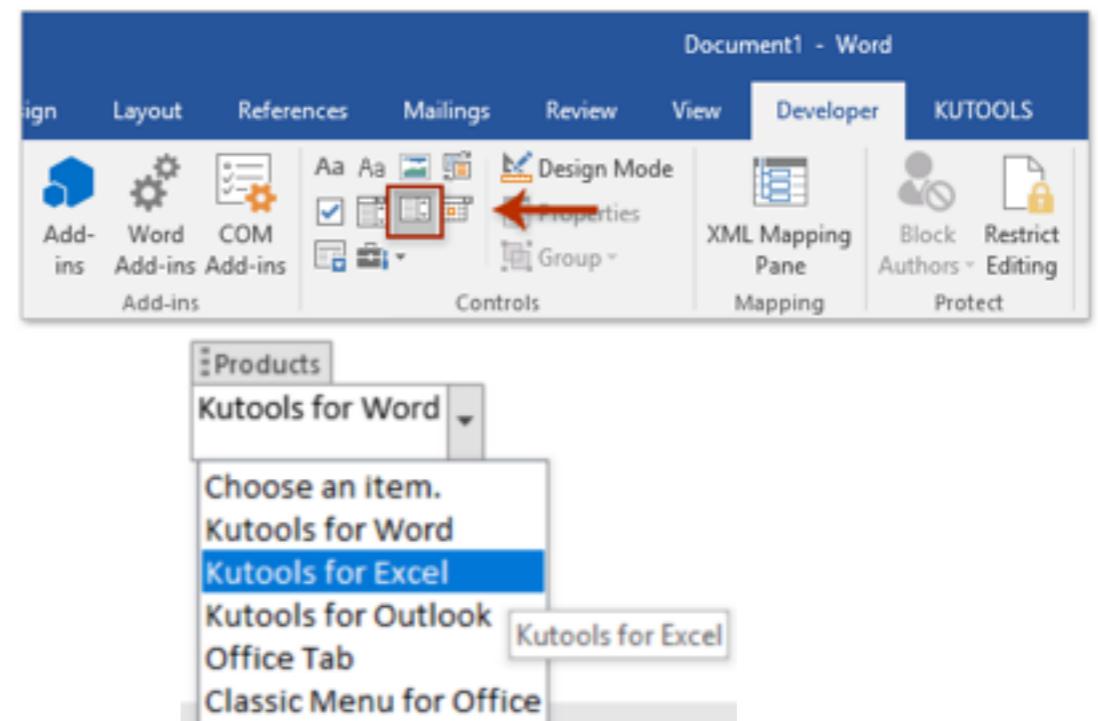
.csv	via python builtins
.doc	via antiword
.docx	via python-docx2txt
.eml	via python builtins
.epub	via ebooklib
.gif	via tesseract-ocr
.jpg	and .jpeg via tesseract-ocr
.json	via python builtins
.html	and .htm via beautifulsoup4
.mp3	via sox , SpeechRecognition , and pocketsphinx
.msg	via msg-extractor
.odt	via python builtins
.ogg	via sox , SpeechRecognition , and pocketsphinx
.pdf	via pdftotext (default) or pdfminer.six
.png	via tesseract-ocr
.pptx	via python-pptx
.ps	via ps2text
.rtf	via unrtf
.tiff	and .tif via tesseract-ocr
.txt	via python builtins
.wav	via SpeechRecognition and pocketsphinx
.xlsx	via xlrd
.xls	via xlrd

Text Extraction

- Challenge 1: Digital watermarking



- Challenge 2: Develop controls





Word

(In Practice)

- Text extraction
- Table extraction

Table Extraction

(Basic)

- **import python-docx**

▲ If you are using python 3x don't do `pip install docx` instead go for

102 `pip install python-docx`

▼ it is compatible with python 3x

official Document: <https://pypi.org/project/python-docx/>

▲ when I import `docx` I have this error:

54 >`File "/Library/Frameworks/Python.framework/Versions/3.3/lib/python3.3/site-packages/`
 `from exceptions import PendingDeprecationWarning`
 `ImportError: No module named 'exceptions'`

★ How to fix this error (`python3.3 , docx 0.2.4`)?

4

`python` `python-3.x` `python-docx`

Table Extraction

(Basic)

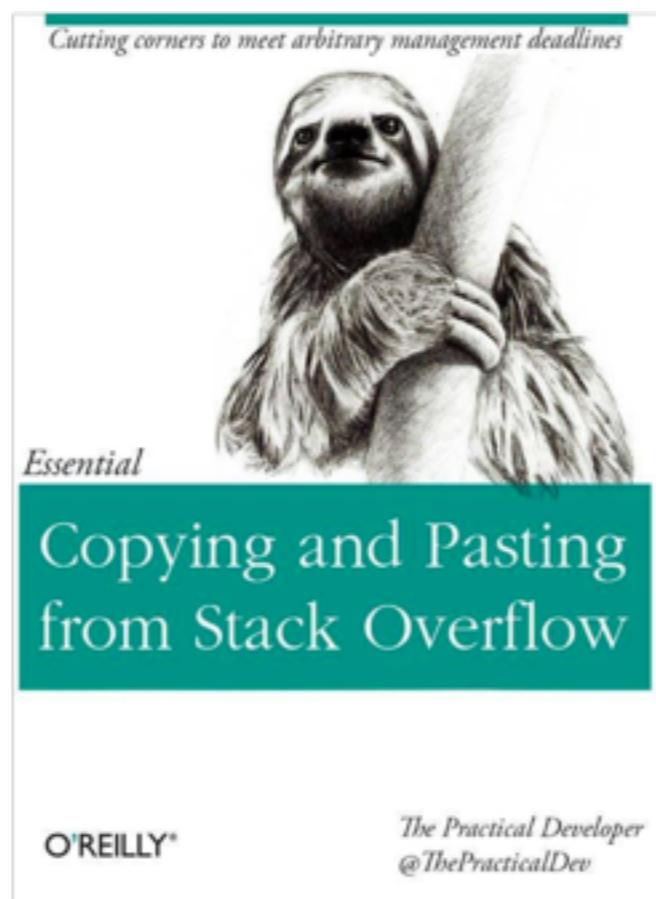
- import python-docx
- Parse XML

Result : Beautify XML

```
23      <w:widowControl w:val="1"/>
24      <w:tabs>
25          <w:tab w:val="left" w:pos="9214"/>
26          <w:tab w:val="left" w:pos="9639"/>
27      </w:tabs>
28      <w:spacing w:after="0" w:before="0" w:line="240"
29          w:lineRule="auto"/>
30      <w:ind w:right="107"/>
31          <w:contextualSpacing w:val="0"/>
32          <w:jc w:val="right"/>
33      </w:pPr>
34      <w:r w:rsidDel="00000000" w:rsidR="00000000" w
35          :rsidRPr="00000000">
36          <w:rPr>
37              <w:rFonts w:ascii="Arial Unicode MS" w:cs
38                  ="Arial Unicode MS" w:eastAsia="Arial
39                      Unicode MS" w:hAnsi="Arial Unicode MS"/>
40              <w:b w:val="1"/>
41              <w:sz w:val="24"/>
42              <w:szCs w:val="24"/>
43              <w:vertAlign w:val="baseline"/>
44              <w:rtl w:val="0"/>
45          </w:rPr>
46          <w:t xml:space="preserve">【機密】</w:t>
47      </w:r>
48      <w:r w:rsidDel="00000000" w:rsidR="00000000" w
49          :rsidRPr="00000000">
<w:rPr>
    <w:rtl w:val="0"/>
</w:rPr>
</w:r>
</w:p>
```

Table Extraction

(Medium)



```
$ INPUT | some scripts from Github | perfect.csv
```

<https://github.com/ivbeg/docx2csv>

Table Extraction

(Advanced)

- Add special text as **quotation** marks (" ") in target cell

Name ↕	Gender ↕	Age ↕	ID Number ↕													
			hi_i_am_head													
Address ↕				Phone ↕												
Mail ↕					Other ↕											hi_i_am_tail

Add text: "hi_i_am_head"

Add text: "hi_i_am_tail"



PDF

NATIONAL PARTNERSHIP FOR QUALITY AFTERSCHOOL LEARNING
www.sedl.org/afterschool/toolkits

AFTERSCHOOL TRAINING TOOLKIT

Tutoring to Enhance Science Skills

Tutoring Two: Learning to Make Data Tables

Sample Data for Data Tables

Use these data to create data tables following the Guidelines for Making a Data Table and Checklist for a Data Table.

Example 1: Pet Survey (GR 2-3)

Ms. Hubert's afterschool students took a survey of the 600 students at Morales Elementary School. Students were asked to select their favorite pet from a list of eight animals. Here are the results.

Lizard 25, Dog 250, Cat 115, Bird 50, Guinea pig 30, Hamster 45, Fish 75, Ferret 10

Example 2: Electromagnets—Increasing Coils (GR 3-5)

The following data were collected using an electromagnet with a 1.5 volt battery, a switch, a piece of #20 insulated wire, and a nail. Three trials were run. *Safety precautions in repeating this experiment include using safety goggles or safety spectacles and avoiding short circuits.*

Number of Coils	Number of Paperclips
5	3, 5, 4
10	7, 8, 6
15	11, 10, 12
20	15, 13, 14

Example 3: pH of Substances (GR 5-10)

The following are pH values of common household substances taken by three different teams using pH probes. *Safety precautions in repeating this experiment include hooded ventilation, chemical-splash safety goggles, gloves, and apron. Do not use bleach, ammonia, or strong acids with children.*

Lemon juice 2.4, 2.0, 2.2; Baking soda (1 Tbsp) in Water (1 cup) 8.4, 8.3, 8.7; Orange juice 3.5, 4.0, 3.4; Battery acid 1.0, 0.7, 0.5; Apples 3.0, 3.2, 3.5; Tomatoes 4.5, 4.2, 4.0; Bottled water 6.7, 7.0, 7.2; Milk of magnesia 10.5, 10.3, 10.6; Liquid hand soap 9.0, 10.0, 9.5; Vinegar 2.2, 2.9, 3.0; Household bleach 12.5, 12.5, 12.7; Milk 6.6, 6.5, 6.4; Household ammonia 11.5, 11.0, 11.5; Lye 13.0, 13.5, 13.4; and Sodium hydroxide 14.0, 14.0, 13.9; Anti-freeze 10.1, 10.9, 9.7; Windex 9.9, 10.2, 9.5; Liquid detergent 10.5, 10.0, 10.3; and Cola 3.0, 2.5, 3.2

Teaching tip: The pH scale is from 0 to 14. Have students make two data tables, one with the data as given and one with the pH scale 0 to 14 with the substances' average pH in rank order on the scale (Battery acid at the lower end and Sodium hydroxide at the upper end) or create a pH graphic organizer.

Example 4: Automobile Land Speed Records (GR 5-10)

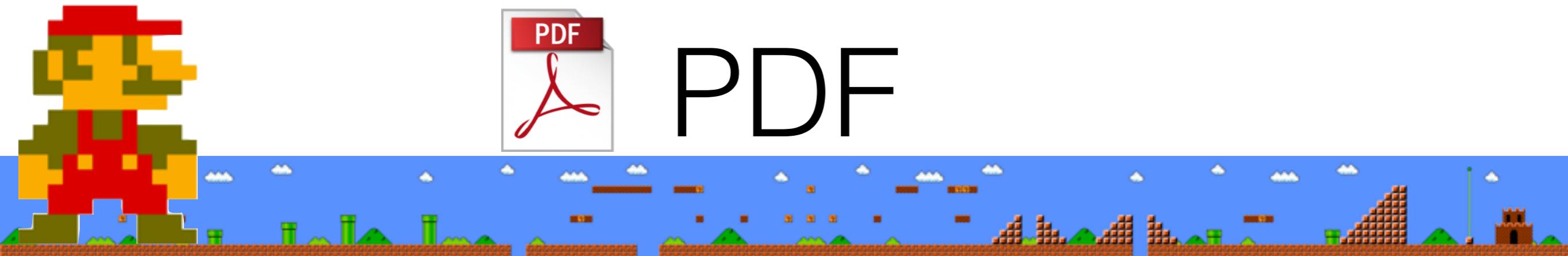
In the first recorded automobile race in 1898, Count Gaston de Chasseloup-Laubat of Paris, France, drove 1 kilometer in 57 seconds for an average speed of 39.2 miles per hour (mph) or 63.1 kilometers per hour (kph). In 1904, Henry Ford drove his Ford Arrow across frozen Lake St. Clair, MI, at an average speed of 91.4 mph. Now, the North American Eagle is trying to break a land speed record of 800 mph. The Federation Internationale de L'Automobile (FIA), the world's governing body for motor sport and land speed records, recorded the following land speed records. (Retrieved on February 5, 2006, from <http://www.landspeed.com/lsrcinfo.asp>.)

Speed (mph)	Driver	Car	Engine	Date
407.447	Craig Breedlove	Spirit of America	GE J47	8/5/63
413.199	Tom Green	Wingfoot Express	WE 346	10/2/64
434.22	Art Arfons	Green Monster	GE J79	10/5/64
468.719	Craig Breedlove	Spirit of America	GE J79	10/13/64
526.277	Craig Breedlove	Spirit of America	GE J79	10/15/65
536.712	Art Arfons	Green Monster	GE J79	10/27/65
555.127	Craig Breedlove	Spirit of America, Sonic 1	GE J79	11/2/65
576.553	Art Arfons	Green Monster	GE J79	11/7/65
600.601	Craig Breedlove	Spirit of America, Sonic 1	GE J79	11/15/65
622.407	Gary Gabelich	Blue Flame	Rocket	10/23/70
633.468	Richard Noble	Thrust 2	RR RG 146	10/4/83
763.035	Andy Green	Thrust SSC	RR Spey	10/15/97

Example 5: Distance and Time (GR 8-10)

The following data were collected using a car with a water clock set to release a drop in a unit of time and a meter stick. The car rolled down an inclined plane. Three trials were run. Create a data table with an average distance column and an average velocity column, create an average distance-time graph, and draw the best-fit line or curve. Estimate the car's distance traveled and velocity at six drops of water. Describe the motion of the car. Is it going at a constant speed, accelerating, or decelerating? How do you know?

Time (drops of water)	Distance (cm)
1	10, 11, 9
2	29, 31, 30
3	59, 58, 61
4	102, 100, 98
5	122, 125, 127



- Zone 1 - **tabula** (pip install tabula-py)

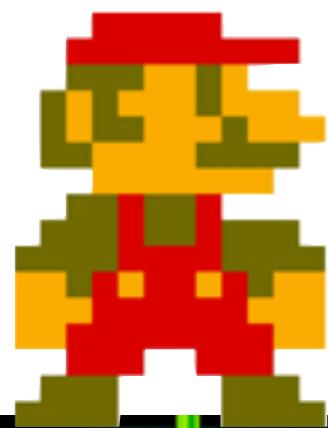
```
from tabula import read_pdf
df_list = read_pdf('sample.pdf', pages='all', multiple_tables=True)
for df in df_list:
    df
```

	0	1
0	Number of Coils	Number of Paperclips
1	5	3, 5, 4
2	10	7, 8, 6
3	15	11, 10, 12
4	20	15, 13, 14

	0	1	2	3	4
0	Speed (mph)	Driver	Car	Engine	Date
1	407.447	Craig Breedlove	Spirit of America	GE J47	8/5/63
2	413.199	Tom Green	Wingfoot Express	WE J46	10/2/64
3	434.22	Art Arfons	Green Monster	GE J79	10/5/64
4	468.719	Craig Breedlove	Spirit of America	GE J79	10/13/64
5	526.277	Craig Breedlove	Spirit of America	GE J79	10/15/65
6	536.712	Art Arfons	Green Monster	GE J79	10/27/65
7	555.127	Craig Breedlove	Spirit of America, Sonic 1	GE J79	11/2/65
8	576.553	Art Arfons	Green Monster	GE J79	11/7/65
9	600.601	Craig Breedlove	Spirit of America, Sonic 1	GE J79	11/15/65
10	622.407	Gary Gabelich	Blue Flame	Rocket	10/23/70
11	633.468	Richard Noble	Thrust 2	RR RG 146	10/4/83
12	763.035	Andy Green	Thrust SSC	RR Spey	10/15/97



PDF





PDF



WELCOME TO WRAP ZONE !



PDF



- Zone 2 - Decrypt password-protected pdf

- **pyPDF2**

```
from PyPDF2 import PdfFileReader  
  
fp = open(filename)  
pdfFile = PdfFileReader(fp)  
password = "mypassword"  
if pdfFile.isEncrypted:  
    try:  
        pdfFile.decrypt(password)  
        print('File Decrypted (PyPDF2)')
```

- **QPDF**: qpdf --decrypt input.pdf output.pdf

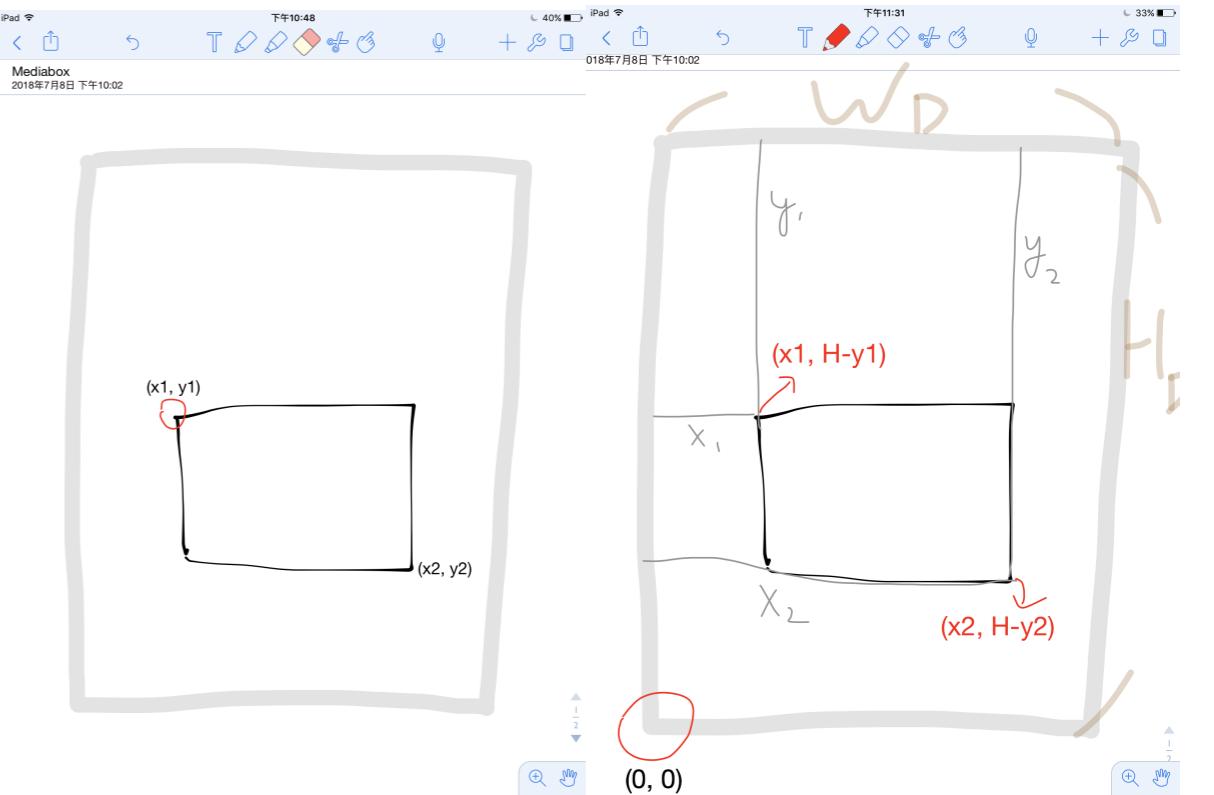
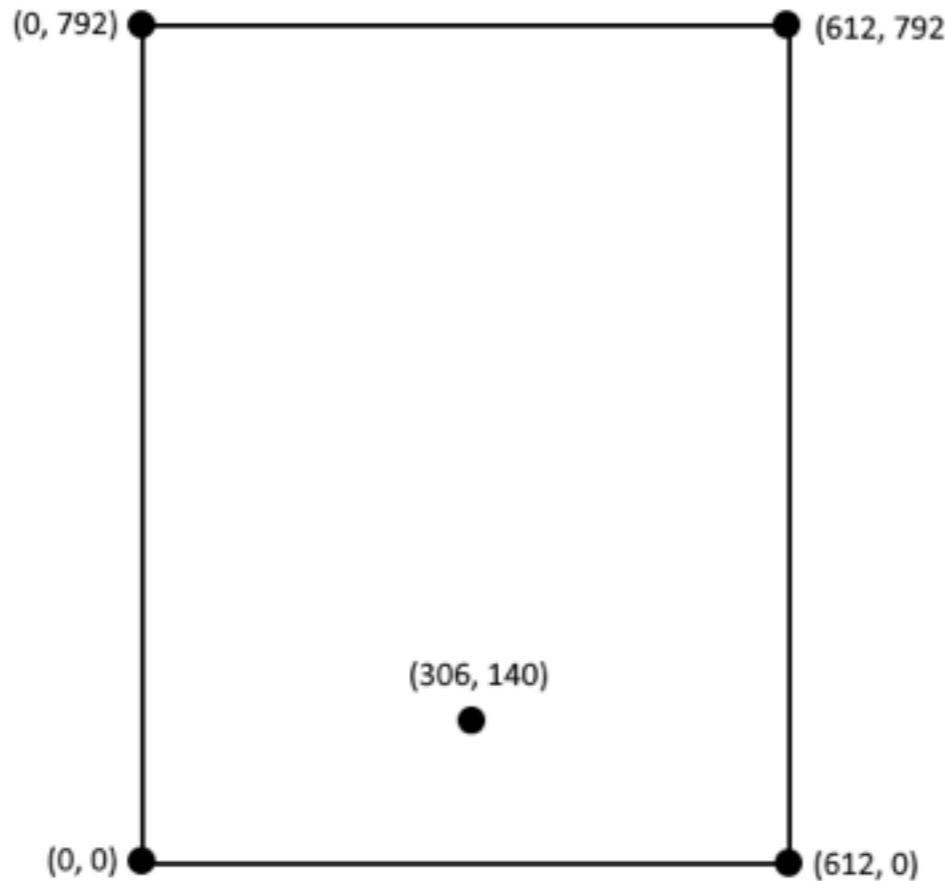
```
except:  
    command = ("cp "+ filename +  
              " temp.pdf; qpdf --password='' --decrypt temp.pdf " + filename  
              + "; rm temp.pdf")  
    os.system(command)  
    print('File Decrypted (qpdf)')  
    fp = open(filename)  
    pdfFile = PdfFileReader(fp)
```



PDF



- Zone 3 - Find table captions
 - 3-1 Find area coordinate

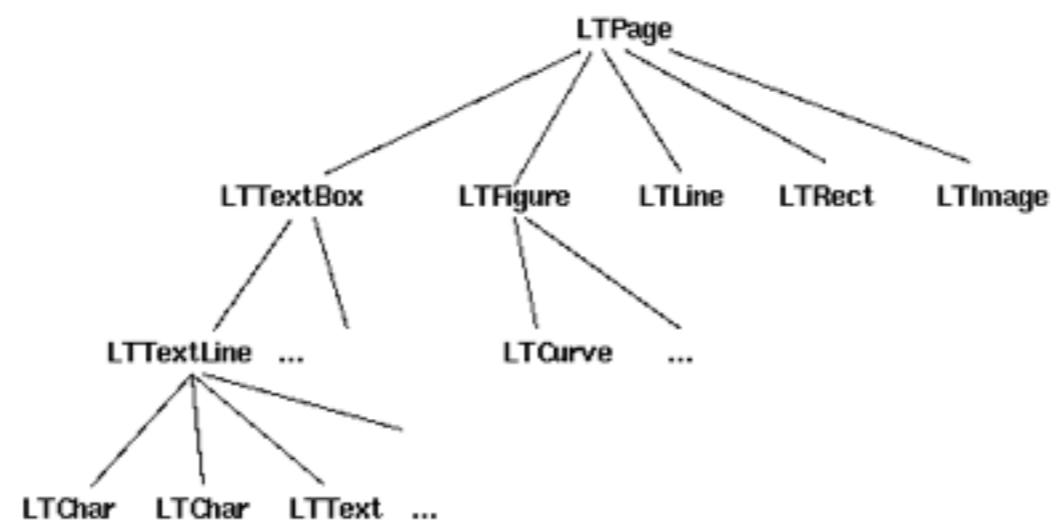
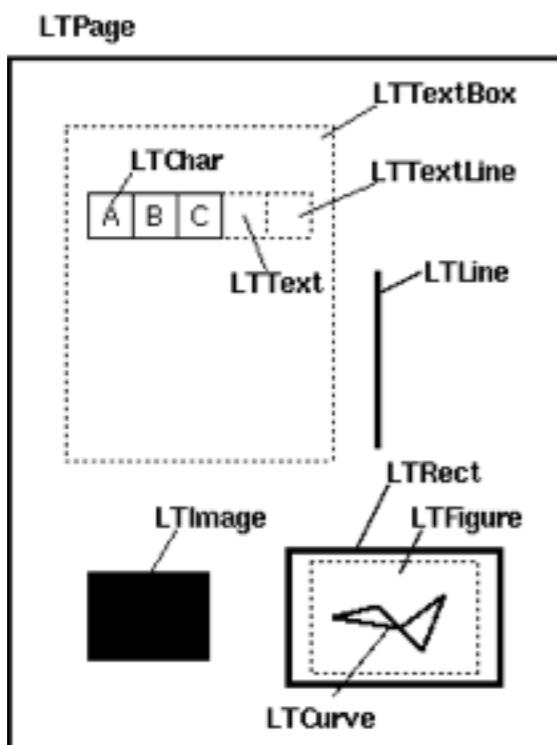




PDF



- Zone 3 - Find table captions
 - 3-1 Find area coordinate
 - 3-2 Tear down elements



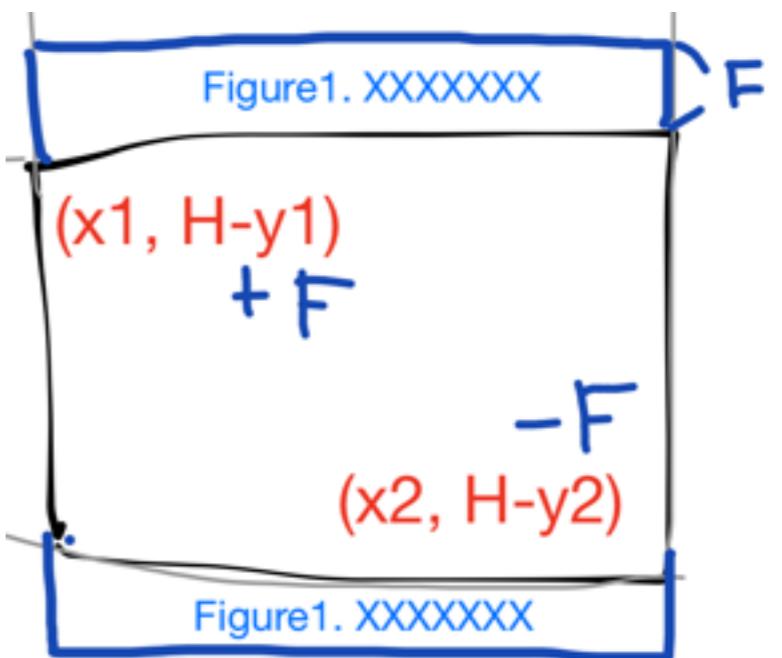
PDFMiner: <https://github.com/euske/pdfminer>



PDF



- Zone 3 - Find table captions
 - 3-1 Find area coordinate
 - 3-2 Tear down elements
 - 3-3 Get captions





PDF



- Zone 4 - Section tagging
 - Utilize OCR
https://github.com/WZBSocialScienceCenter/pdftabextract/blob/master/examples/catalogue_30s/catalog_30s_notebook.ipynb
 - ML approach
<https://github.com/HazyResearch/pdftotree>

ML Approach to tag sections

section_header

section header

TABLE I THERMAL CHARACTERISTICS ($T_A = 25^\circ\text{C}$ unless otherwise noted)

Characteristic	Symbol	Min	Typ	Max	Unit
OFF CHARACTERISTICS					
Collector = Emitter Breakdown Voltage ($I_C = 1.0 \text{ mA}$, $I_B = 0$)	$V_{(BR)CEO}$	BC546 BC547 BC549	63 45 30	- - -	- - -
Collector = Base Breakdown Voltage ($I_C = 100 \mu\text{A}$)	$V_{(BR)BBO}$	BC546 BC547 BC549	60 50 30	- - -	- - -
Emitter = Base Breakdown Voltage ($I_E = 10 \mu\text{A}$, $I_C = 0$)	$V_{(BR)BBO}$	BC546 BC547 BC549	6.0 6.0 6.0	- - -	- - -
Collector Cutoff Current ($V_{CE} = 70 \text{ V}$, $V_{BE} = 0$) ($V_{CE} = 50 \text{ V}$, $V_{BE} = 0$) ($V_{CE} = 35 \text{ V}$, $V_{BE} = 0$) ($V_{CE} = 30 \text{ V}$, $T_A = 125^\circ\text{C}$)	I_{CES}	BC546 BC547 BC549 BC546/547/549	- - - -	0.2 0.2 0.2 -	15 15 15 4.0
ON CHARACTERISTICS					
DC Current Gain ($I_C = 10 \mu\text{A}$, $V_{CE} = 5.0 \text{ V}$)	h_{FE}	BC547A BC546B/547B/548B BC548C	- - -	90 150 270	- - -
($I_C = 2.0 \text{ mA}$, $V_{CE} = 5.0 \text{ V}$)		BC546 BC547 BC548 BC547A BC546B/547B/548B BC547C/BC548C	110 110 110 110 200 420	- - - - 250 520	450 600 600 220 450 800
($I_C = 100 \text{ mA}$, $V_{CE} = 5.0 \text{ V}$)		BC547A/548A BC546B/547B/548B BC548C	- - -	120 180 300	- - -
Collector = Emitter Saturation Voltage ($I_C = 10 \text{ mA}$, $I_B = 0.5 \text{ mA}$) ($I_C = 100 \text{ mA}$, $I_B = 0.5 \text{ mA}$) ($I_C = 10 \text{ mA}$, $I_B = \text{See Note 1}$)	$V_{(CE)sat}$		- - -	0.09 0.2 0.3	0.25 0.6 0.6
Base = Emitter Saturation Voltage ($I_C = 10 \text{ mA}$, $I_B = 0.5 \text{ mA}$)	$V_{(BE)sat}$		-	0.7	-
Base = Emitter On Voltage ($I_C = 2.0 \text{ mA}$, $V_{CE} = 5.0 \text{ V}$) ($I_C = 10 \text{ mA}$, $V_{CE} = 5.0 \text{ V}$)	$V_{(BE)on}$		0.55 -	-	0.7 0.77
SMALL-SIGNAL CHARACTERISTICS					
Current = Gain - Bandwidth Product ($I_C = 10 \text{ mA}$, $V_{CE} = 5.0 \text{ V}$, $f = 100 \text{ MHz}$)	f_T	BC546 BC547 BC549	150 150 150	200 200 300	- - -
Output Capacitance ($V_{CE} = 10 \text{ V}$, $I_C = 0$, $f = 1.0 \text{ MHz}$)	C_{obs}		-	1.7	4.5
Input Capacitance ($V_{EB} = 0.5 \text{ V}$, $I_C = 0$, $f = 1.0 \text{ MHz}$)	C_{iss}		-	10	-
Small - Signal Current Gain ($I_C = 2.0 \text{ mA}$, $V_{CE} = 5.0 \text{ V}$, $f = 1.0 \text{ kHz}$)	h_{ie}	BC546 BC547/549 BC547A BC546B/547B/548B BC547C/BC548C	125 125 125 240 450	- - - 300 600	600 900 250 500 900
Noise Figure ($I_C = 0.2 \text{ mA}$, $V_{CE} = 5.0 \text{ V}$, $R_S = 2 \text{ k}\Omega$, $f = 1.0 \text{ kHz}$, $\Delta f = 200 \text{ Hz}$)	NF	BC546 BC547 BC549	- - -	2.0 2.0 2.0	10 10 10

1. In is value for which $i_C = 11$ mA at $V_{CE} = 1.0$ V

section header

Data Parser



SVN v.s. Git



SVN v.s. Git



GitPython

```
#!/usr/bin/env python

from git import *

repo = Repo("misctools")
o = repo.remotes.origin
o.pull()

master = repo.head.reference
print master.log()
```

SVN v.s. Git

PySVN



GitPython

```
#!/usr/bin/env python

from git import *

repo = Repo("misctools")
o = repo.remotes.origin
o.pull()

master = repo.head.reference
print master.log()
```

SVN v.s. Git

~~PySVN~~

glob



GitPython

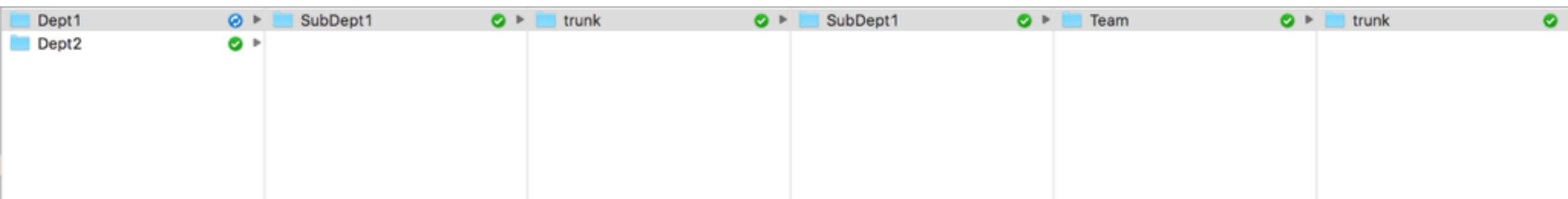
```
#!/usr/bin/env python

from git import *

repo = Repo("misctools")
o = repo.remotes.origin
o.pull()

master = repo.head.reference
print master.log()
```

Legacy Repo



Dept1 > SubDept1 > trunk > SubDept1 > Team > trunk > status > bug_list

http://hostaddress/svn/Dept1/SubDept1/trunk/SubDept1/Team/trunk/status/bug_list/bug.xls

path_exp = "**/Dept1/SubDept1/**/bug_list/**/*.*"

for fp in glob.glob(path_exp, recursive=True):

...

A screenshot of the JIRA web interface. The URL in the address bar is `bingroom.atlassian.net/projects/TEST/issues/?filter=allopenissues`.

The page title is "Projects / test". On the left sidebar, under "Issues and filters", the "Open issues" option is selected. The main content area shows a list of open issues for the "test summary" component:

Priority	Issue Key	Summary
TEST-2	TEST-2	test summary
TEST-1	TEST-1	test

On the right side, there are sections for "Advanced search", "test description" (with a red box around the "TEST-2" priority), "环境" (Environment) set to "None", and "Activity" (Activity) with a "B" badge.

- pip install jira

```
from jira import JIRA

HOST = 'your JIRA server'
jira = JIRA(HOST, basic_auth=('user', 'token'))
issue = jira.issue('issue key')
```



JIRA

(In Practice)

```
field_dict = {d['id']: d['name'] for d in jira.fields()}\nfield_dict
```

```
{'statuscategorychangedate': 'Status Category Changed',
'issuetype': 'Issue Type',
'timespent': 'Time Spent',
'project': 'Project',
'fixVersions': 'Fix versions',
'aggregatetimespent': 'Σ Time Spent',
'statusCategory': 'Status Category',
'resolution': 'Resolution',
'resolutiondate': 'Resolved',
'workratio': 'Work Ratio',
'lastViewed': 'Last Viewed',
'watches': 'Watchers',
'thumbnail': 'Images',
'created': 'Created',
'customfield_10020': 'Sprint',
'customfield_10021': 'Flagged',
'customfield_10022': '[CHART] Date of First Response',
'priority': 'Priority',
'customfield_10023': '[CHART] Time in Status',
'labels': 'Labels',
'customfield_10016': 'Story point estimate',
'customfield_10017': 'Issue color',
'customfield_10018': 'Parent Link',
'customfield_10019': 'Rank',
'timeestimate': 'Remaining Estimate',
'aggregatetimeoriginalestimate': 'Σ Original Estimate',
'versions': 'Affects versions'}
```

- Field name mapping



JIRA

(In Practice)

The screenshot shows a Jira issue page for an issue titled "test". The JSON representation of the issue is displayed on the right, with red arrows pointing from specific fields in the JSON to their corresponding values or descriptions on the page.

```
'assignee': None,  
'attachment': [],  
'comment': [(['高餅倫', 'gg'], ['高餅倫', '我也gg'])],  
'components': [],  
'created': '2019-09-17T23:54:37.734+0800',  
'creator': {'displayName': '高餅倫', 'emailAddress': 'joekaojoekao@gmail.com'},  
'description': 'test description',  
'duedate': None,  
'environment': None,  
'fixVersions': [],  
'issuelinks': [],  
'labels': [],  
'lastViewed': '2019-09-18T00:40:30.881+0800',  
'priority': 'Medium',  
'progress': {'progress': 0, 'total': 0},  
'project': 'test',  
'reporter': {'displayName': '高餅倫', 'emailAddress': 'joekaojoekao@gmail.com'},  
'resolution': None,  
'resolutiondate': None,  
'security': None,  
'status': 'To Do',  
'statuscategorychangedate': '2019-09-17T23:54:38.010+0800',  
'subtasks': [],  
'summary': 'test',  
'timeestimate': None,  
'timeoriginalestimate': None,  
'timespent': None,  
'timetracking': {},  
'type': 'Bug',  
'updated': '2019-09-18T00:40:36.632+0800',
```

Case 1: Low-cost Data Pipeline

Case 1: Low-cost Data Pipeline



Requirements

- Windows data, Linux backend
- Integration on heterogeneous files
- Error handling, early feedback
- User-friendly

Requirements

- Windows data, Linux backend
- Integration on heterogeneous files
- Error handling, early feedback
- User-friendly

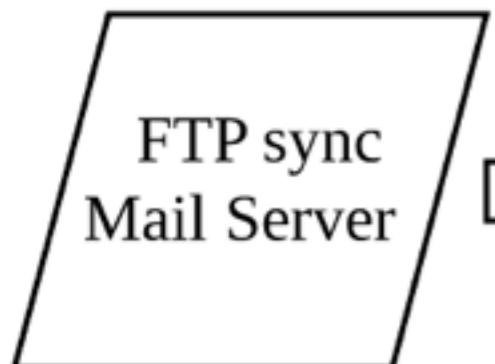


elasticsearch



Pipeline

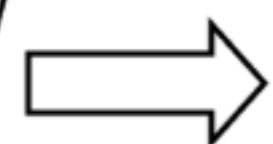
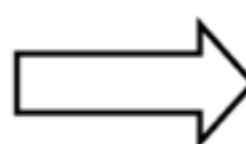
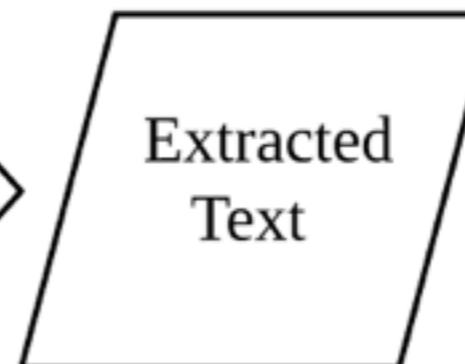
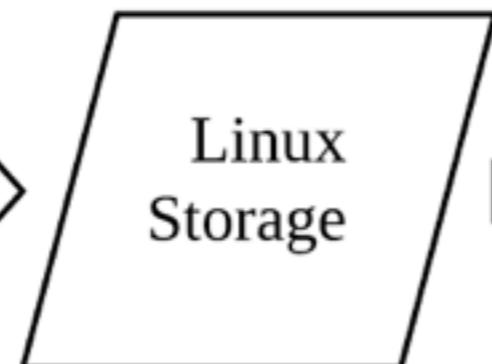
1. inotifywait



3. Elasticsearch



**2. Parser with
error handling**



Pipeline

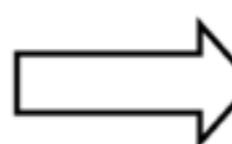
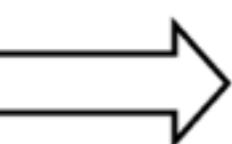
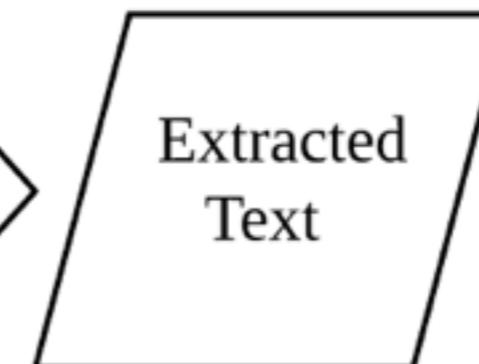
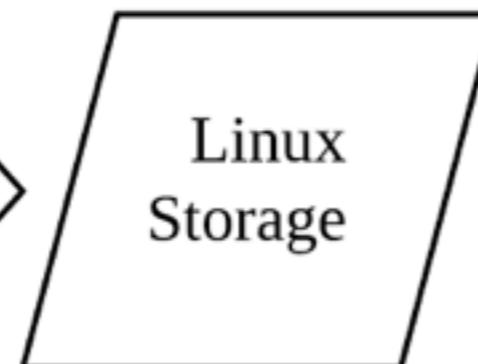
1. inotifywait



3. Elasticsearch



**2. Parser with
error handling**



(1) inotifywait

You should consider using `inotifywait`, as an example:

```
inotifywait -m /path -e create -e moved_to |  
while read path action file; do  
    echo "The file '$file' appeared in directory '$path' via '$action'"  
    # do something with the file  
done
```

(1) inotifywait

You should consider using `inotifywait`, as an example:

close_write

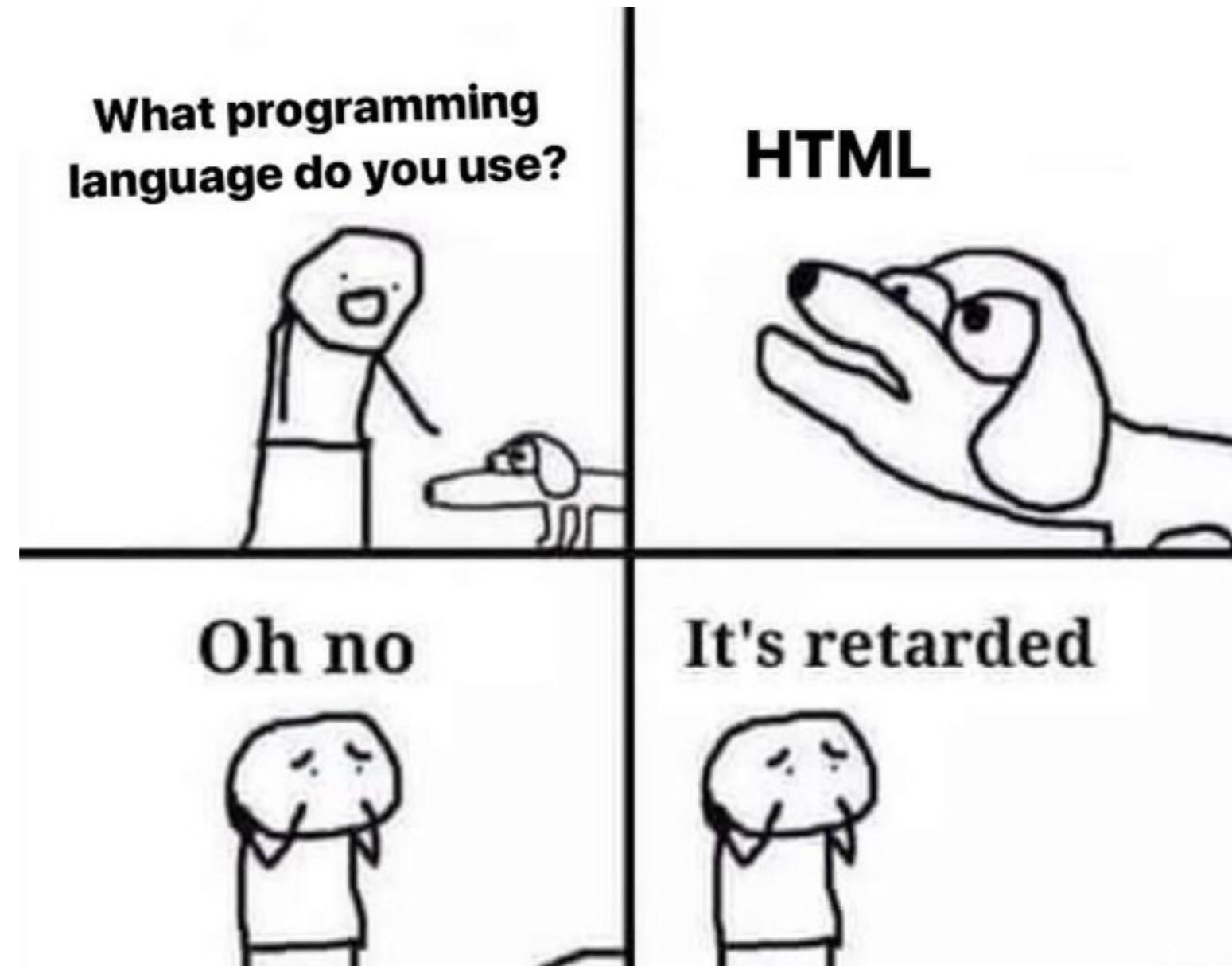
```
inotifywait -m /path -e create -e moved_to |  
while read path action file; do  
    echo "The file '$file' appeared in directory '$path' via '$action'"  
    # do something with the file  
done
```

(2) Error handling

- Massive **if-else** and **try-except**
- Report locations of bad cells when handling **Excel**

e.g. [Error] Wrong data type in column ‘**est. cost**’,
In Line **2, 5, 10**.
- Reply in email using **HTML**

(2) Error handling



(2) Error handling

```
import html
from htmltag import HTML, td, tr, th

table = ""
rows = [1, 2, 'x', 4]
for row in rows:
    s = ""
    if row == 'x':
        s += td(str(row), _class="bad-cell")
    else:
        s += td(str(row))

    table += tr(s)

html_str = """
<h3 class="error-msg">Invalid format: </h3>
<table align="center" cellpadding="2" cellspacing="0" width="5%">
    <tr><th>index</th></tr>
    <tr><td>1</td></tr><tr><td>2</td></tr><tr><td class="bad-cell">x</td></tr><tr><td>4</td></tr>
</table>""".replace("{", "{{").replace("}", "}}") \
.replace("[", "{").replace("]", "}").format(style, th('index'), table)

print(html.unescape(html_str))
```

```
<style>
table, th, td {
    border: 1px solid black;
}
td.bad-cell {
    width: 100%;
    background-color: #ff0000;
}
</style>
<h3 class="error-msg">Invalid format: </h3>
<table align="center" cellpadding="2" cellspacing="0" width="5%">
    <tr><th>index</th></tr>
    <tr><td>1</td></tr><tr><td>2</td></tr><tr><td class="bad-cell">x</td></tr><tr><td>4</td></tr>
</table>
```

Invalid format:	
index	
	1
	2
	x
	4

htmltag

<http://liftoff.github.io/htmltag/>

(2) Error handling



高餅倫

2018年9月3日 · 人 ·

...

catch exceptions: 5 min

**write HTML format error
message using python: 5 hr**

Case 2: Bug Search System

Case 2: Bug Search System

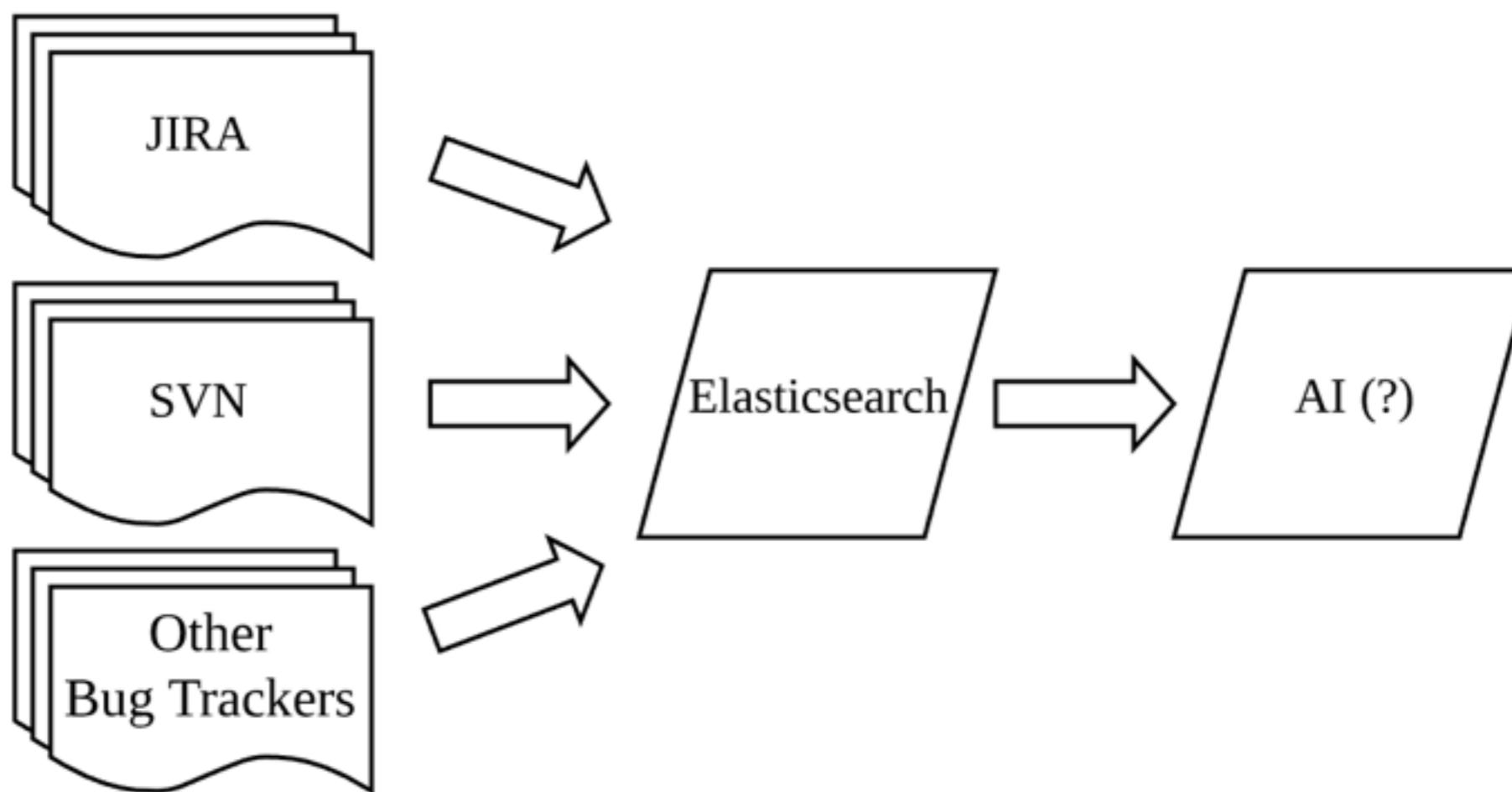
User explaining
Debug with Big data

me



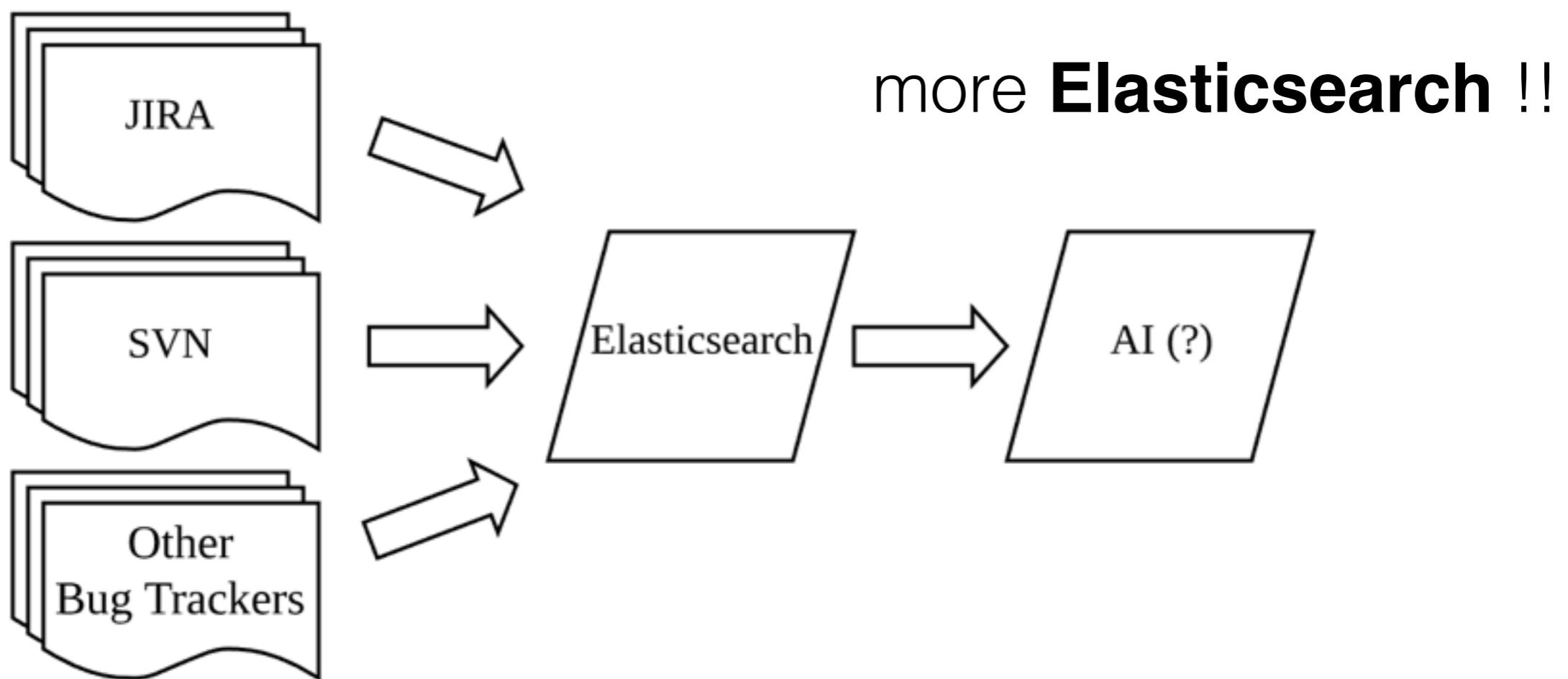
Pipeline

Recap: **textract**



Pipeline

Recap: **textract**





- NoSQL, schema-free or schema-*flexible*
- JSON everywhere
- Query DSL (Domain-specific Language)
- Search engine in nature

<https://www.elastic.co/guide/en/elasticsearch/reference/current/query-dsl.html>



elasticsearch

```
GET dota2/_search
{
  "query": {
    "filtered": {
      "filter": {
        "and": [
          {
            "term": {
              "heroname": "drow"
            }
          }
        ]
      }
    },
    "aggs": {}
  }
}

1  {
2    "took": 2,
3    "timed_out": false,
4    "_shards": {
5      "total": 5,
6      "successful": 5,
7      "failed": 0
8    },
9    "hits": {
10      "total": 1,
11      "max_score": 1,
12      "hits": [
13        {
14          "_index": "dota2",
15          "_type": "agust",
16          "_id": "AU-N0y9iNzuaTKNspu3X",
17          "_score": 1,
18          "_source": {
19            "heroname": "DROW RANGER",
20            "class": "Ranged - Carry",
21            "skills": {
22              "main": [
23                "Frost Arrows",
24                "Gust",
25                "Precision Aura"
26              ],
27              "ulti": "Marksmanship"
28            },
29            "ability": "agility",
30            "items": [
31              "mitem1",
32              "mitem2",
33              "mitem3",
34              "mitem4",
35              "mitem5"
36            ],
37            "dateofplay": "01/01/2012 00:00:00",
38            "experienced": "false"
39          }
40        }
41      ]
42    }
43  }
```

Approaches to AI 0.5

- UI side
 - Google-like searching experience
 - ✓ Space separator
 - ✓ Ranked search result
 - ✓ Near-real-time indexing

Approaches to AI 0.5

- Backend side:
 - Modularizing query strategy

```
174     query_all = {  
175         "query": {  
176             "bool": {  
177                 "should": [  
178                     {  
179                         "bool": {  
180                             "must_not": {  
181                                 "exists": {  
182                                     "field": "created"  
183                                 }  
184                             }  
185                         }  
186                     },  
187                     {  
188                         "bool": {  
189                             "must": [  
190                                 {  
191                                     "range": {  
192                                         "created": {  
193                                             "gte": date_from,  
194                                             "lte": date_to  
195                                         }  
196                                     }  
197                                 },  
198                                 {  
199                                     "terms": {  
200                                         "status.keyword": status  
201                                     }  
202                                 }  
203                             ]  
204                         }  
205                     }  
206                 ]  
207             },  
208             "_source": {  
209                 "includes": search_field,  
210                 "excludes": []  
211             },  
212             "sort": sort_list  
213         }  
214     }
```



Approaches to AI 0.5

- Backend side:
 - Modularizing query strategy

```
project_boost = {
    "function_score": {
        "query": {
            "match": {
                "project": term
            }
        },
        "boost": "200",
        "random_score": {},
        "boost_mode": "multiply"
    }
}
```

```
match_phrase_prefix = {
    "multi_match": {
        "type": "phrase_prefix",
        "query": term,
        "fields": [
            "_id",
            "summary",
            "description",
            "project"
        ]
    }
}
```

```
query_string_wildcard = {
    "query_string": {
        "query": "*{}*".format(re.sub(r'\W+', ' ', term)),
        "fields": [
            "issue_key",
            "summary",
            "description",
            "project"
        ]
    }
}
```

Approaches to AI 0.5

- Backend side:
 - Modularizing query strategy

Domain-based labeling

Utilizing existent features
ex: JQL

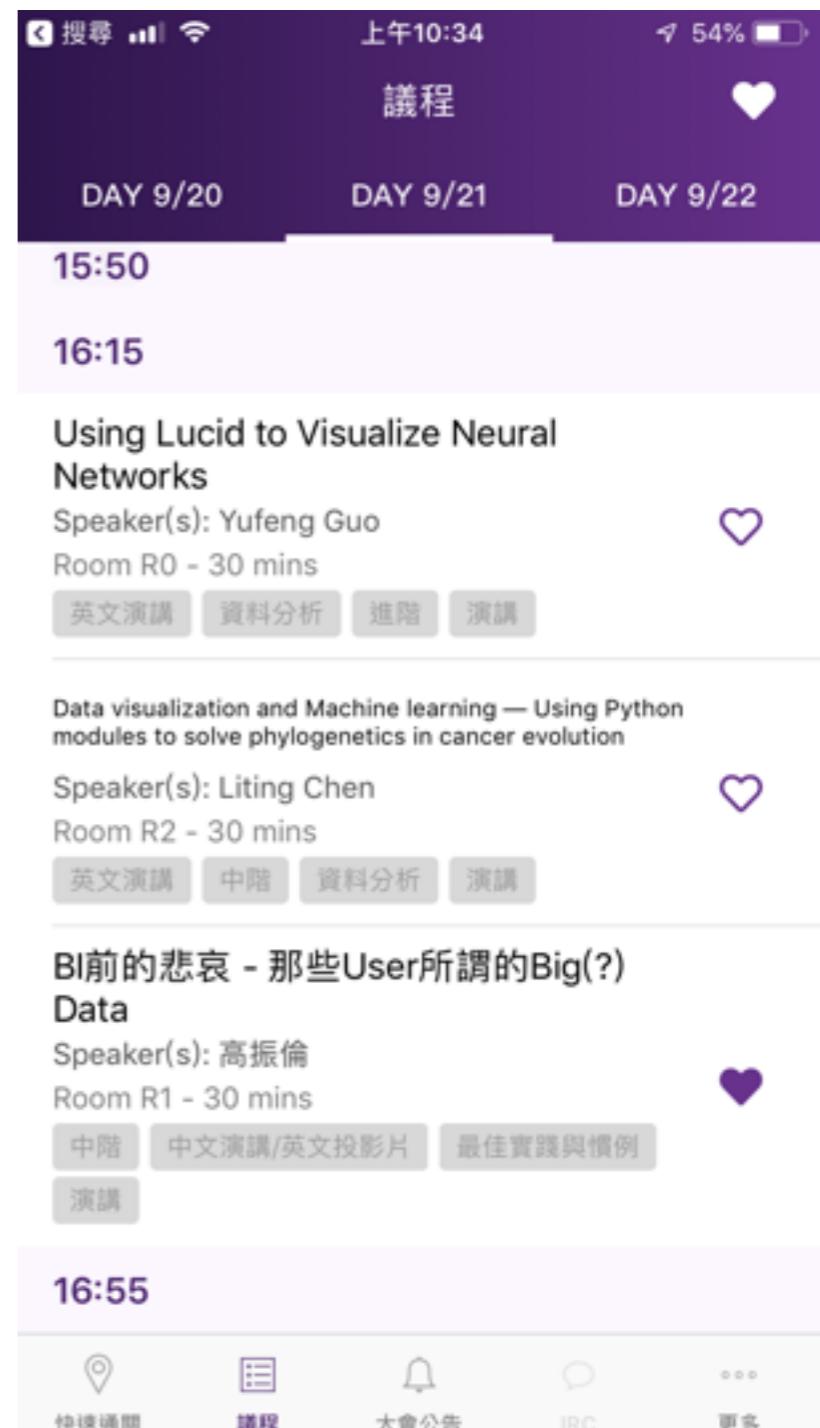
MOTHER OF THE HORIZONTAL SCROLLBAR

```
238 ":{  
239     "must": [  
240         {  
241             "range": {  
242                 "created": {  
243                     "gte": date_from,  
244                     "lte": date_to  
245                 }  
246             }  
247         },  
248         {  
249             "terms": {  
250                 "status.keyword": status  
251             }  
252         },  
253         {  
254             "bool": {  
255                 "should": [match_issue, match_phrase, match_word, match_phrase_prefix, query_string_wildcard, match_phrase_prefix_post, query_string_wildcard_post, match_phrase_pref  
256             ]  
257         }  
258     }  
259 }
```

Approaches to AI 0.5

- Backend side:

- Modularizing query strategy
- Domain-based labeling



Approaches to AI 0.5

- Backend side:
 - Modularizing query strategy
 - Domain-based labeling
 - Utilizing existent features
ex: JQL



JQL Cheat Sheet

A simple query in JQL (also known as a "clause") consists of a field, followed by an operator, followed by one or more values or functions. For example:

project = Test
field operator value

To perform a more complex query, you can link clauses together with keywords.

project = TEST AND assignee in (currentUser())
field operator value keyword field operator function



```
jql = 'summary ~ test'  
issue_list = jira.search_issues(jql)  
issue_list
```

```
[<JIRA Issue: key='TEST-2', id='10001'>,  
<JIRA Issue: key='TEST-1', id='10000'>]
```

**Friends: How did you write this
code so beautifully ?**

Me(Proudly):

**Friends: How did you write this
code so beautifully ?
Me(Proudly):**



AI 0.5: Must have

- Meet user expectations
- UAT (User Acceptance Testing)
- Hourly build

AI 0.5: Nice to have

- User rating 
- Auto-complete (Completion suggester)
- Semantic analysis

autocompl
autocomplete
autocomplete c++
autocomplete off
autocompletetextview
autocomplete google

Grab-and-Go

- Not Best practice but Test practice
- Not local / global optimum but User optimum
- Show your respect for frontend colleagues