



XSEDE[16]

Tutorial – An Introduction to Brown Dog

An Elastic Data Cyberinfrastructure for Autocuration
and Digital Preservation

Monday, July 18, 2016

1:00 pm – 1:30 pm | Introduction to Brown Dog | Kenton McHenry

How to Use Brown Dog Services

1:30 pm – 1:45 pm	Hands-on Outline Environment Set-up	Smruti Padhy
1:45 pm – 2:05 pm	Demonstration of BD Fiddle	Jong Sung Lee
2:05 pm – 2:55 pm	Create Applications Using BD Services	
	- Conversion Example	Jong Sung Lee
	- Extraction/Indexing/Retrieval Example	Marcus Slavenas
	- Combined Conversion and Extraction Example	Marcus Slavenas
2:55 pm – 3:10 pm	Break	

How to Add Your Tool to Brown Dog Services

3:10 pm – 3:50 pm	How to Write a Converter	Kenton McHenry
3:50 pm – 4:30 pm	How to Write an Extractor	Smruti Padhy
4:30 pm – 5:00 pm	Wrap Up and Feedback	Team



Linux Cheat sheet



Basic Commands

Basic Terminal Shortcuts

CTRL L = Clear the terminal
 CTRL D = Logout
 SHIFT Page Up/Down = Go up/down the terminal
 CTRL A = Cursor to start of line
 CTRL E = Cursor the end of line
 CTRL U = Delete left of the cursor
 CTRL K = Delete right of the cursor
 CTRL W = Delete word on the left
 CTRL Y = Paste (after CTRL U,K or W)
 TAB = auto completion of file or command
 CTRL R = reverse search history
 !! = repeat last command
 CTRL Z = stops the current command (resume with fg in foreground or bg in background)

Basic Terminal Navigation

ls -a = list all files and folders
 ls <folderName> = list files in folder
 ls -lh = Detailed list, Human readable
 ls -l *.jpg = list jpeg files only
 ls -lh <fileName> = Result for file only

cd <folderName> = change directory
 if folder name has spaces use " "
 cd / = go to root
 cd .. = go up one folder, tip: ../../../

du -h: Disk usage of folders, human readable
 du -ah: " " " files & folders, Human readable
 du -sh: only show disc usage of folders

pwd = print working directory

man <command> = shows manual (RTFM)

File Permissions

chown = change the owner of a file
 ex --> chown bob hello.txt
 chown user:bob report.txt = changes the user owning report.txt to 'user' and the group owning it to 'bob'
 -R = recursively affect all the sub folders
 ex --> chown -R bob:bob /home/Daniel

chmod = modify user access/permission - simple way
 u = user
 g = group
 o = other

d = directory (if element is a directory)
 l = link (if element is a file link)
 r = read (read permissions)
 w = write (write permissions)
 x = eXecute (only useful for scripts and programs)

Basic file manipulation

cat <fileName> = show content of file
 (less, more)

head = from the top
 -n <#oflines> <fileName>

tail = from the bottom
 -n <#oflines> <fileName>

mkdir = create new folder
 mkdir myStuff ..
 mkdir myStuff/pictures/ ..

cp image.jpg newimage.jpg = copy and rename a file
 cp image.jpg <folderName>/ = copy to folder
 cp image.jpg folder/sameImageNewName.jpg
 cp -R stuff otherStuff = copy and rename a folder
 cp *.txt stuff/ = copy all of *<file type> to folder

mv file.txt Documents/ = move file to a folder
 mv <folderName> <folderName2> = move folder in folder
 mv filename.txt filename2.txt = rename file
 mv <fileName> stuff/newfileName
 mv <folderName>/ .. = move folder up in hierarchy

rm <fileName> .. = delete file (s)
 rm -i <fileName> .. = ask for confirmation each file
 rm -f <fileName> = force deletion of a file
 rm -r <foldername>/ = delete folder

touch <fileName> = create or update a file

ln file1 file2 = physical link
 ln -s file1 file2 = symbolic link

File Permissions (continued)

'+' means add a right
 '-' means delete a right
 '=' means affect a right

ex --> chmod g+w someFile.txt
 (add to current group the right to modify someFile.txt)

more info: man chmod

Accessing BD Services locally installed in the VM















- To generate tokens in locally installed bdfiddle

username: **bdtut.user@xsede.org**

password: **bdtutXy71816de**

- While using `key_token.py` to obtain tokens to access local bdfiddle, you need to type in the following values
 - Server [bd-api.ncsa.illinois.edu]: localhost:8080
 - Protocol [https]: http
 - Username: bdtut.user@xsede.org
 - Password: bdtutXy71816de

Index of /tutorials

<u>Name</u>	<u>Last modified</u>	<u>Size</u>	<u>Description</u>
 Parent Directory		-	
 01-BrownDog.pptx	2016-07-08 12:27	78M	
 02-Setup.pptx	2016-07-08 12:27	6.6M	
 03-BDFiddle.pptx	2016-07-08 12:27	19M	
 04-Conversions.pptx	2016-07-08 12:27	3.3M	
 05-Extractions.pptx	2016-07-08 11:35	5.0M	
 06-Combined.pptx	2016-07-08 12:27	3.7M	
 07-AddingConverters.pptx	2016-07-13 14:43	5.3M	
 08-AddingExtractors.pptx	2016-07-08 12:27	6.1M	
 09-ToolsCatalog-updated.pptx	2016-07-08 12:27	4.1M	
 Template.pptx	2016-07-08 12:27	3.2M	
 handouts.zip	2016-07-14 16:21	1.5M	
 old/	2016-07-14 16:22	-	
 samplefiles.zip	2016-07-14 16:21	4.8M	

Apache/2.4.7 (Ubuntu) Server at browndog.ncsa.illinois.edu Port 80

<http://browndog.ncsa.illinois.edu/tutorials/>

File: URL:

Token (Generate Token)


Select Options:

Convert:

Extract:

Output:


Code Snippets:

 Terminal

 Python
bd.py

 Python
Native

 R
bd.r

 Matlab
bd.m

 JS
JavaScript
bd.js

```
1  #!/usr/bin/python -u
2  import bd
3  import time
4  import os
5  from os.path import isfile, join, splitext, basename
6
7  #parameters to access BD services
8  bds = 'https://bd-api.ncsa.illinois.edu'
9
10 # key and token from the Brown Dog API Gateway service
11 api_key = ""
12 token = ""
13
14 # path to the folder to process
15 input_path = ""
16
17 # path to the folder to store output
18 # add time.time() to make the output folder unique for each run
19 output_path = ""+str(int(time.time()))
20
21 # output format
22 output_format = "pdf"
23
24 # Create a output folder
25 os.makedirs(output_path)
26
27 # create a list of files that is not a folder
28 onlyfiles = [join(input_path, f) for f in os.listdir(input_path) if
29 isfile(join(input_path, f))]
30
31 for (input_file) in onlyfiles:
32     filename, file_extension = splitext(basename(input_file))
33     input_format = file_extension[1:]
34
35     print 'File: ' + basename(input_file) + ', ' + input_format + ' to ' + output_format
36
37     # getting possible output format by the input format
38     outputs = bd.outputs(bds, input_format, token)
39
40     # if the output format is not supported, skip the conversion
41     if not (output_format in outputs):
42         print output_format, 'is not supported for', input_format
43         continue
44
45     # craete a output file name with path
46     output_file = output_path + '/' + filename + "." + output_format
47
48     # do the conversion
49     bd.convert(bds, input_file, output_format, output_file, token, 60, True)
50
```

```
1  #!/usr/bin/python -u
2  import bd
3  import os
4  from os.path import isfile, join, splitext, basename
5
6  # this is the method to search dictionary by key
7  def findItem(obj, key):
8      if key in obj: return obj[key]
9      for k, v in obj.items():
10         if isinstance(v, dict):
11             item = findItem(v, key)
12             if item is not None:
13                 return item
14
15  # this is the method to search key with given string
16  def search(values, searchFor):
17      found = []
18      for k in values:
19         for v in values[k]:
20             if searchFor in v.lower():
21                 found.append(k)
22      return found
23
24  #parameters to access BD services
25  bds = 'https://bd-api.ncsa.illinois.edu'
26
27  # key and token from the Brown Dog API Gateway service
28  token = ""
29
30  # path to the folder to process
31  input_path = ""
32
33  # create a list of files that is not a folder
34  onlyfiles = [join(input_path, f) for f in os.listdir(input_path) if
35               isfile(join(input_path, f))]
36
37  text_store = {}
38
39  # processing the files in the given folder
40  for (input_file) in onlyfiles:
41      filename, file_extension = splitext(basename(input_file))
42      input_format = file_extension[1:]
43
44      print 'Processing file: ' + basename(input_file)
45
46      # do the extraction
47      metadata = bd.extract(bds, input_file, token, 120)
48
49      # process the metadata
50      # find the "OCR" results and store
51      for m in metadata['metadata.jsonld']:
52         txt = findItem(m, 'ocr_simple')
53         if not (txt is None):
54             text_store[basename(input_file)] = txt
55
56  # find the keyword from the ocr results
57  print search(text_store, 'information')
```

```
1  #!/usr/bin/python -u
2  import bd
3  import time
4  import os
5  from os.path import isfile, join, splitext, basename
6
7  #parameters to access BD services
8  bds = 'https://bd-api.ncsa.illinois.edu'
9
10 # token from the Brown Dog API Gateway service
11 token = ""
12
13 # path to the folder to process
14 input_file = ""
15
16 # path to the folder to store output
17 output_path = ""+str(int(time.time()))
18
19 # output format
20 output_format = "png"
21
22 def main():
23     # create a output folder
24     os.makedirs(output_path)
25
26     # convert the image file to png
27     output = convert(input_file, output_format, output_path, bds, token)
28
29     # extract metadata from the converted file
30     metadata = extract(output, bds, token)
31
32     # print the "tag" section for to see the face detected
33     print metadata['tags']
34
35 def convert(input_file, output_format, output_path, bds, token):
36     filename, file_extension = splitext(basename(input_file))
37     input_format = file_extension[1:]
38
39     print 'File: ' + basename(input_file) + ', ' + input_format + ' to ' + output_format
40
41     # getting possible output format by the input format
42     outputs = bd.outputs(bds, input_format, token)
43
44     # craete a output file name with path
45     output_file = output_path + '/' + filename + "." + output_format
46
47     # do the conversion
48     return bd.convert(bds, input_file, output_format, output_file, token, 60, True)
49
50 def extract(extract_file, bds, token):
51     print 'Processing File: ' + extract_file
52     extract_output = bd.extract(bds, extract_file, token, 120)
53     return extract_output
54
55 if __name__ == "__main__":
56     main()
57
```


Template - ImageMagick_convert.sh

```
#!/bin/sh
#ImageMagick (v6.5.2)
#image
#bmp, dib, eps, fig, gif, ico, jpg, jpeg
#bmp, dib, eps, gif, jpg, jpeg, jp2, pcd, pdf, pgm

# NOTE: First four lines of comments above MUST follow pattern:
#Full software name (Version)
#Data types dealt with
#Comma separated list of inputs accepted
#Comma separated list of outputs accepted

output_filename=$(basename "$2")
output_format="${output_filename##*.}"

#Output PGM files as ASCII
if [ "$output_format" = "pgm" ]; then
    convert "$1" -compress none "$2"
else
    convert "$1" "$2"
fi
```

Exercise (/home/ubuntu/git/bd-tutorial-problems/07-AddingConverters/)

1. Write a wrapper script

```
nano sciencify_convert.sh
```

2. Modify the SoftwareServer configuration (i.e. copy the script and edit aliases)

```
cp sciencify_convert.sh /home/ubuntu/git/polyglot/scripts/sh
```

```
nano /home/ubuntu/git/polyglot/scripts/sh/.aliases.txt
```

3. Start the SoftwareServer (you can do this part in the bd-tmux pane)

```
cd /home/ubuntu/git/polyglot/
```

```
bin/SoftwareServerRestlet.sh
```

4. Test it with the BD CLI

```
cd /home/ubuntu/git/bdcli
```

```
bd -o science.txt examples/sonnet80.txt
```

```
1  #!/usr/bin/env python
2  import logging
3  import numpy
4  import csv
5  from config import *
6  import pycowder.extractors as extractors
7
8  def main():
9      global extractorName, messageType, rabbitmqExchange, rabbitmqURL, logger,
10         registrationEndpoints
11
12         # Set logging
13         logging.basicConfig(format='%(asctime)-15s %(levelname)-7s : %(name)s -
14         %(message)s', level=logging.WARN)
15         logging.getLogger('pycowder.extractors').setLevel(logging.INFO)
16         logger = logging.getLogger('extractor')
17         logger.setLevel(logging.DEBUG)
18
19         extractors.setup(extractorName=extractorName, messageType=messageType,
20         rabbitmqURL=rabbitmqURL,
21         rabbitmqExchange=rabbitmqExchange)
22
23         # Register extractor
24         extractors.register_extractor(registrationEndpoints)
25
26         # Connect to extractor and wait for files
27         extractors.connect_message_bus(extractorName=extractorName, messageType=messageType,
28         rabbitmqURL=rabbitmqURL,
29         rabbitmqExchange=rabbitmqExchange,
30         processFileFunction=process_file,
31         checkMessageFunction=None)
32
33 def process_file(parameters):
34     global field_name
35
36     inputfile = parameters['inputfile']
37
38     mean_desc = 'Mean total monthly precipitation'
39     median_desc = 'Median total monthly precipitation'
40     standard_dev_desc = 'Standard deviation of total monthly precipitation'
41
42     # Initialize a records list and field name
43     records = list()
44     field_name = 'TPCP'
45
46     # Read csv file contents into a variable
47     csv_file = open(inputfile, 'rb')
48     data_table = csv.DictReader(csv_file)
49
50     # Read each row of the table and append the data into records array
51     for data_row in data_table:
52         records.append(float(data_row[field_name]))
53
54     # Calculate mean, median and standard deviation and upload it back as metadata to
55     Clowder
56     mean = numpy.mean(records)
57     median = numpy.median(records)
58     std = numpy.std(records)
59
60     # Context url
61     context_url = "https://clowder.ncsa.illinois.edu/contexts/metadata.jsonld"
62
63     # Store results as metadata
64     metadata = {
65         "@context": [context_url,
```

```
60         {
61             "mean": "http://clowder.ncsa.illinois.edu/" + extractorName +
62                 "#mean",
63             "median": "http://clowder.ncsa.illinois.edu/" + extractorName
64                 + "#median",
65             "std": "http://clowder.ncsa.illinois.edu/" + extractorName +
66                 "#std",
67             "value": "http://clowder.ncsa.illinois.edu/" + extractorName +
68                 "#value",
69             "description": "http://clowder.ncsa.illinois.edu/" +
70                 extractorName + "#description"
71         }],
72     "attachedTo": {"resourceType": "file", "id": parameters["fileid"]},
73     "agent": {
74         "@type": "cat:extractor",
75         "extractor_id":
76             "https://clowder.ncsa.illinois.edu/clowder/api/extractors/" + extractorName
77     },
78     "content": {
79         "mean": {
80             "value": round(mean, 3),
81             "description": mean_desc
82         },
83         "median": {
84             "value": round(median, 3),
85             "description": median_desc
86         },
87         "std": {
88             "value": round(std, 3),
89             "description": standard_dev_desc
90         }
91     }
92 }
93
94 # upload metadata
95 extractors.upload_file_metadata_jsonld(mdata=metadata, parameters=parameters)
96 logger.info("Uploaded metadata %s", metadata)
97
98 if __name__ == '__main__':
99     main()
```

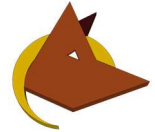
```
1  # =====
2  #
3  # In order for this extractor to run according to your preferences,
4  # the following parameters need to be set.
5  #
6  # Some parameters can be left with the default values provided here - in that
7  # case it is important to verify that the default value is appropriate to
8  # your system. It is especially important to verify that paths to files and
9  # software applications are valid in your system.
10 #
11 # =====
12
13 import os
14
15 # Name to show in rabbitmq queue list
16 extractorName = os.getenv('RABBITMQ_QUEUE', "nca.csv_stats")
17
18 # URL to be used for connecting to rabbitmq
19 rabbitmqURL = os.getenv('RABBITMQ_URI', "amqp://guest:guest@localhost:%2f")
20
21 # Name of rabbitmq exchange
22 rabbitmqExchange = os.getenv('RABBITMQ_EXCHANGE', "clowder")
23
24 # Type of files to process
25 messageType = "*.file.text.csv", "extractors." + extractorName
26
27 # Trust certificates, set this to false for self signed certificates
28 sslVerify = os.getenv('RABBITMQ_SSLVERIFY', False)
29
30 # Endpoints and keys for registering extractor information in CSV format.
31 registrationEndpoints = os.getenv('REGISTRATION_ENDPOINTS',
32
33                                     "http://localhost:9000/clowder/api/extractors?key=key1,"
34                                     "
35                                     "http://host2:9000/api/extractors?key=key2")
```

```
1  {
2    "@context": "http://clowder.ncsa.illinois.edu/contexts/extractors.jsonld",
3    "name": "ncsa.csv_stats",
4    "version": "1.0",
5    "description": "A simple JSON-LD extractor that calculates mean, median, and
6    standard deviation of a column in a CSV file.",
7    "author": "Sandeep Satheesan <sandeeps@illinois.edu>",
8    "contributors": [],
9    "contexts": [],
10   "repository": {
11     "repType": "git",
12     "repUrl":
13       "https://opensource.ncsa.illinois.edu/stash/scm/cats/extractors-tutorial.git"
14   },
15   "external_services": [],
16   "dependencies": [],
17   "bibtext": [],
18   "dataset": ["gov.noaa.ncdc:C00841"]
19 }
```


XSEDE Tutorial Session: July 18, 2016

Introduction to Brown Dog

Tutorial Feedback



Thank you for attending our tutorial session. As a part of our on-going effort to improve the utility and usability of the tools we are developing we are asking attendees to answer the following questions. Note that all answers are anonymous. Comments are welcome anywhere!

1. What is your current job role/title?

2. What level would you rate yourself in terms of writing code (in any language)?

1	2	3	4	5
expert		proficient		none

What programming languages are you familiar with?

3. How accessible was the content of presentations?

1	2	3	4	5
very		neutral		difficult to follow

Any suggestions as to what could be improved?

4. Please rate the perceived ease or difficulty with regards to calling/using brown dog via the API.

1	2	3	4	5
easy		neutral		difficult

5. Please rate the perceived ease or difficulty with regards to calling/using brown dog via any one of the sample clients.

1	2	3	4	5
easy		neutral		difficult

6. Please rate the perceived ease or difficulty with regards to adding a new converter or extractor.

1	2	3	4	5
easy		neutral		difficult

7. Was there anything in particular that you found confusing or counterintuitive? If so, what?



XSEDE Tutorial Session: July 18, 2016

Introduction to Brown Dog

Tutorial Feedback

8. Do you have any suggestions for improving any of the interfaces or additional functionality you would like to see? Please describe:

9. Were there other material/topics you would have liked to have had presented? If so, what?

10. Is there material you think should be deleted? If so, what?

11. Are there any other relevant comments you want to share?

(Optional)

Name: _____ Email: _____

Thank you for your feedback!
Please feel free to write additional comments below.