

# Big Data with ADAMS

Tweet, tweet, tweet

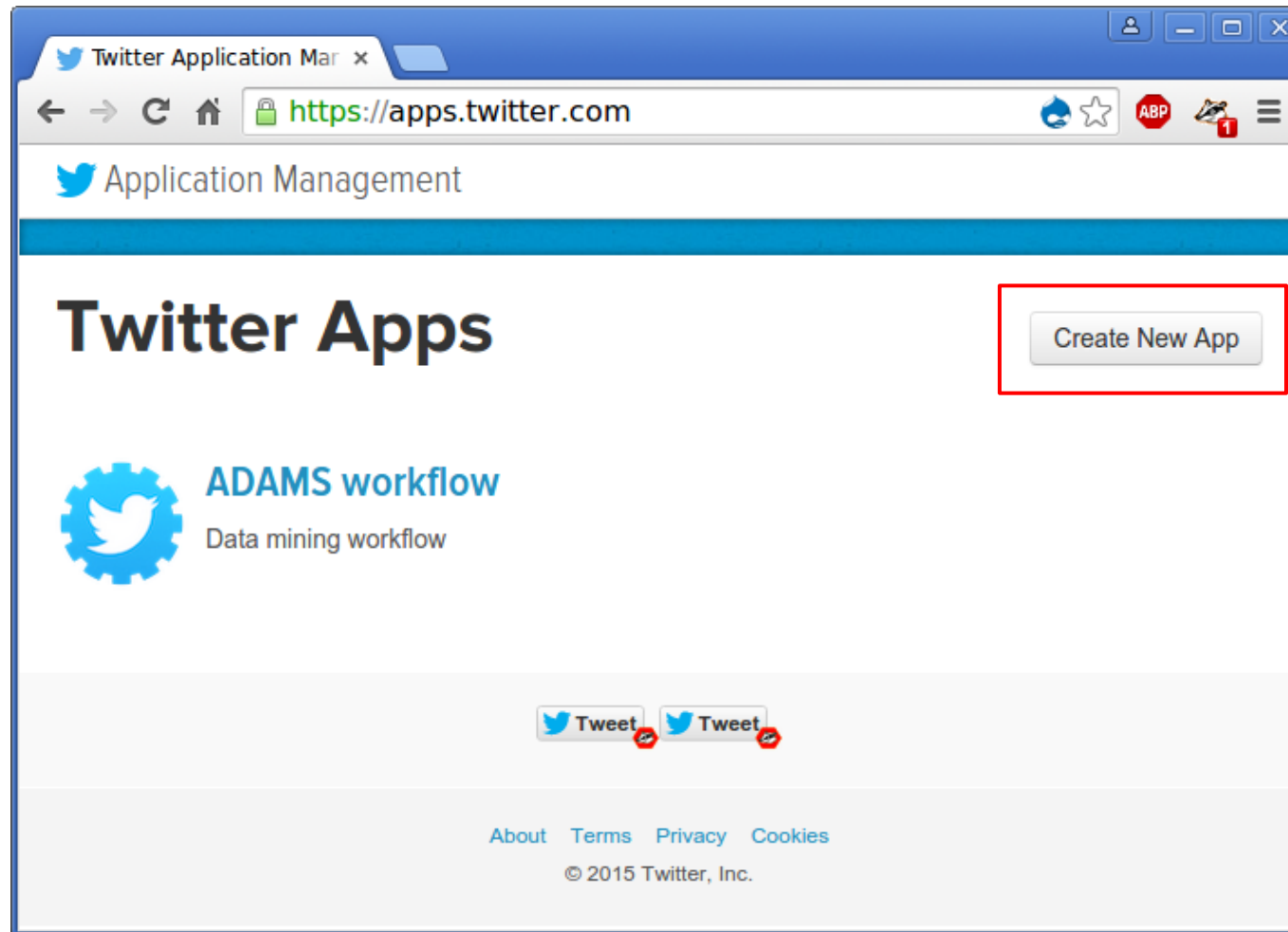
# Collecting tweets

---

- Twitter allows searches of public tweets
- Twitter offers access to tweets in real time
- ADAMS uses [twitter4j.org](https://twitter4j.org) to access Twitter
- Requires setting up an App  
<https://apps.twitter.com>
- Don't worry, it's not that hard...

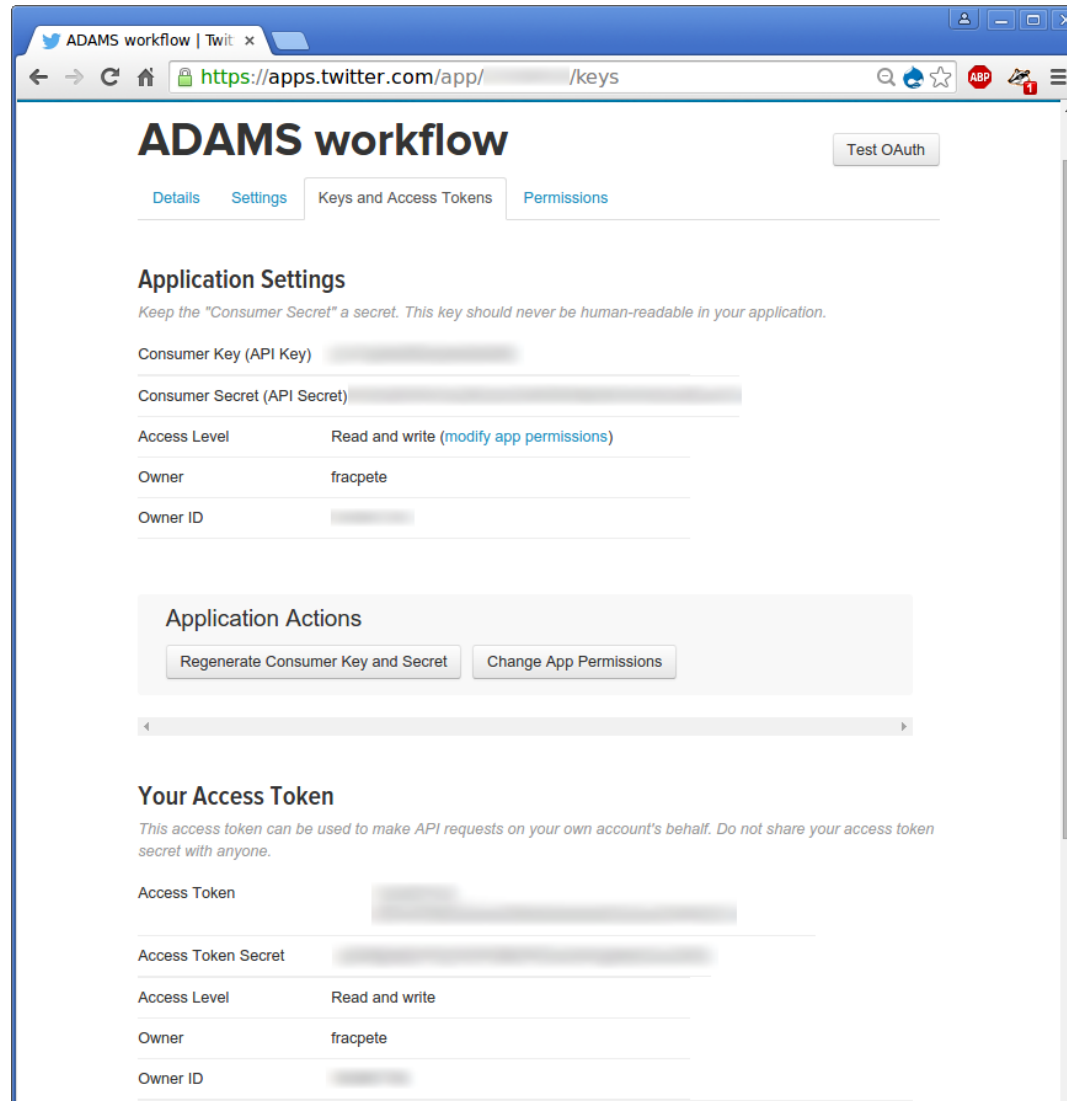
# Settings things up

- Create an app



# Settings things up (2)

- Set up tokens (consumer and access)



The screenshot shows the 'Keys and Access Tokens' page for an application named 'ADAMS workflow'. The page is divided into two main sections: 'Application Settings' and 'Your Access Token'.

**Application Settings**

Keep the "Consumer Secret" a secret. This key should never be human-readable in your application.

Consumer Key (API Key) [redacted]

Consumer Secret (API Secret) [redacted]

Access Level: Read and write ([modify app permissions](#))

Owner: fracpete

Owner ID: [redacted]

**Application Actions**

[Regenerate Consumer Key and Secret](#) [Change App Permissions](#)

**Your Access Token**

This access token can be used to make API requests on your own account's behalf. Do not share your access token secret with anyone.

Access Token [redacted]

Access Token Secret [redacted]

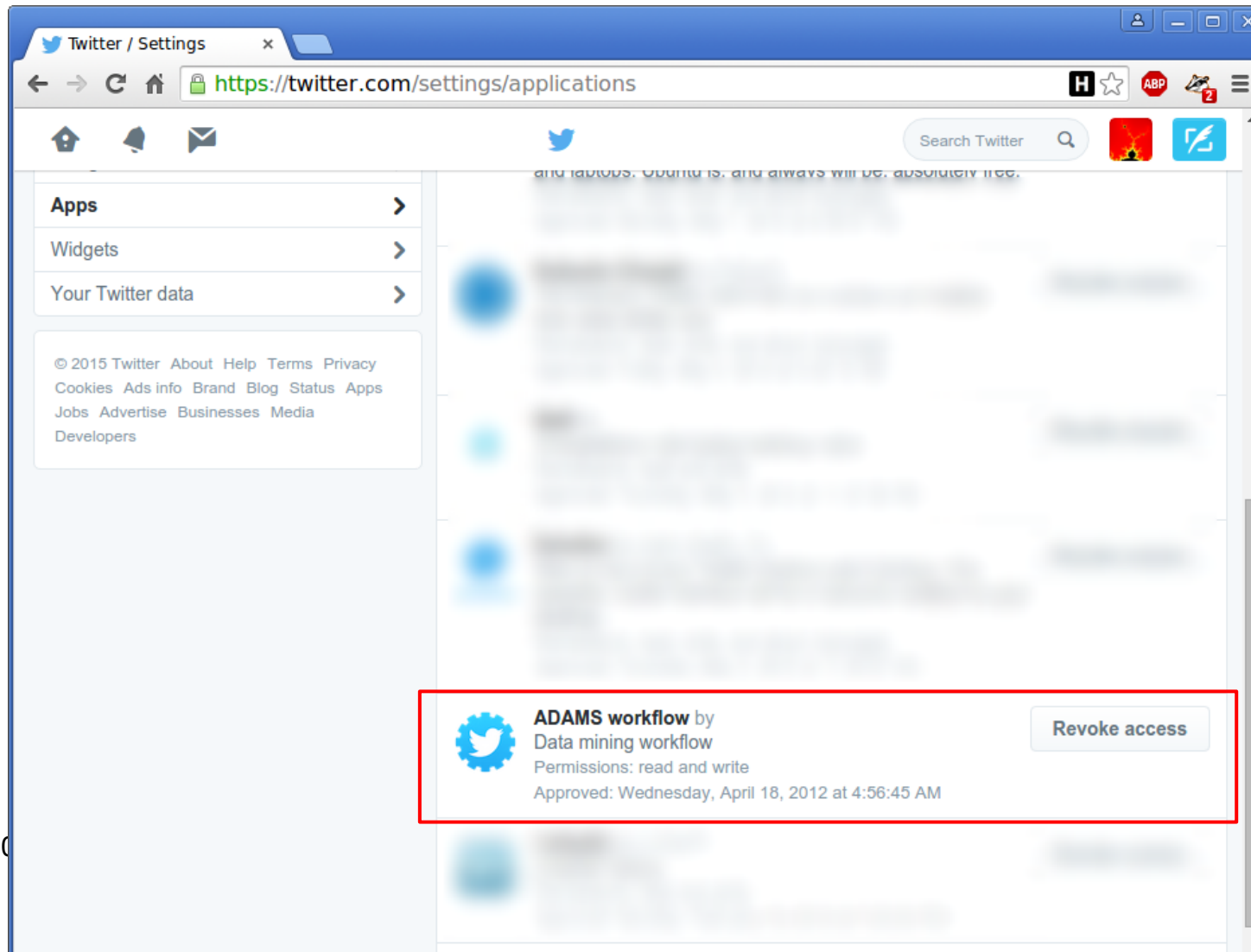
Access Level: Read and write

Owner: fracpete

Owner ID: [redacted]

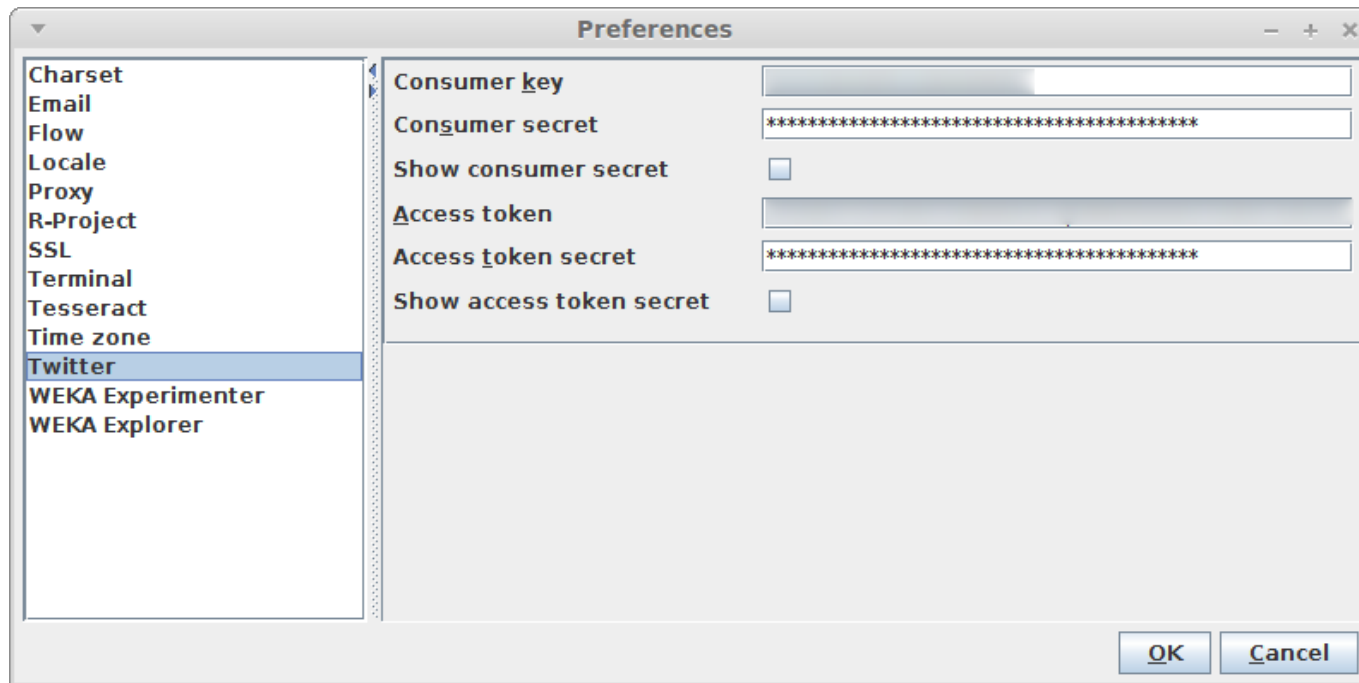
# Settings things up (3)

- App should show up in your profile settings





# Settings things up (4)

- Finally, fill in Twitter preferences in ADAMS



# User queries

---

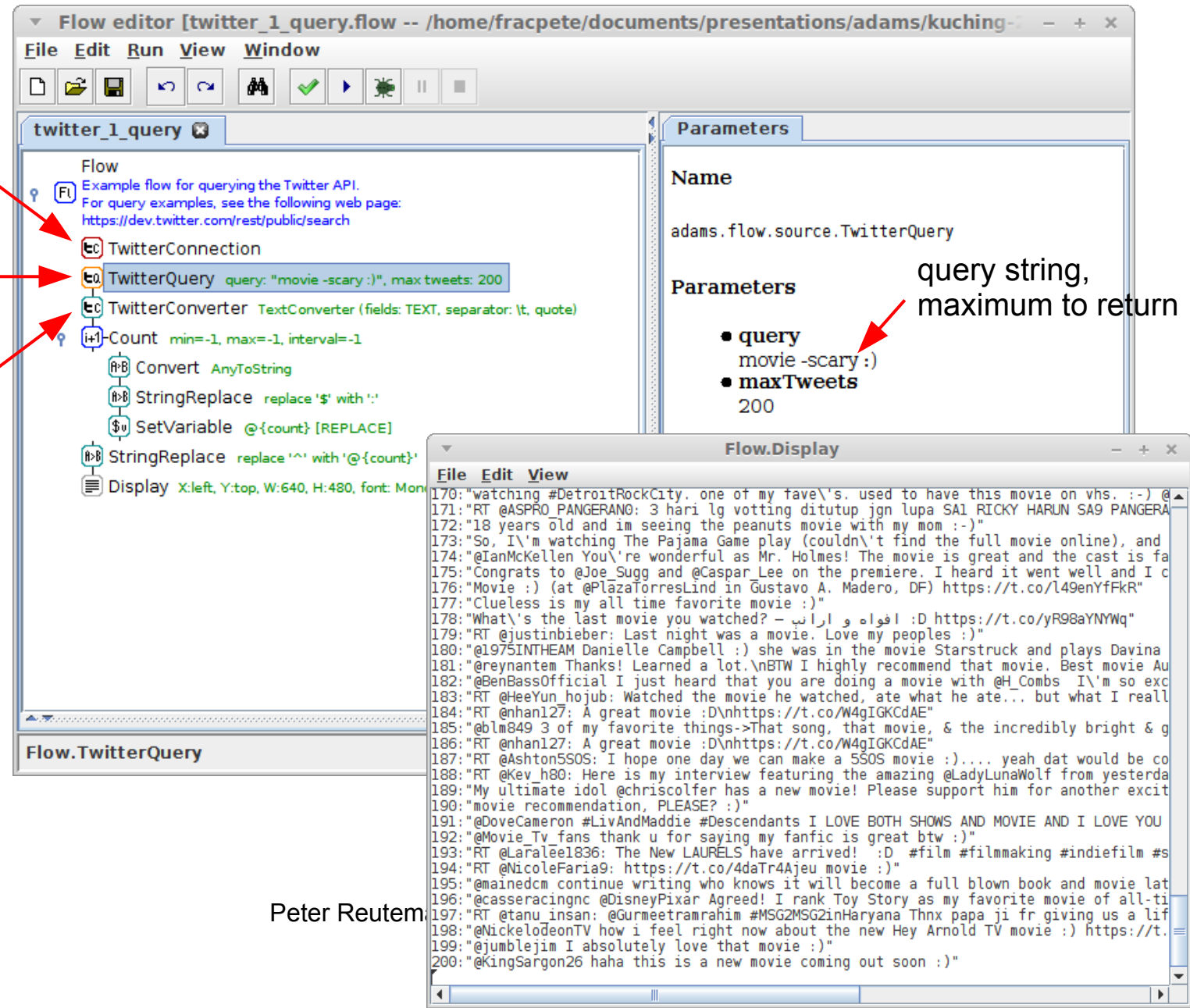
- You can query twitter using something like:  
movie -scary :)  
[tweets containing “movie” and “:)” but not “scary”]
- Actors to use
  -  TwitterQuery
  -  TwitterConverter
- Use case?
  - when looking for specific keywords

# User queries

connection settings;  
can override global  
preference settings

send query

convert tweet  
into textual format



The screenshot shows a flow editor window titled "Flow editor [twitter\_1\_query.flow -- /home/fracpete/documents/presentations/adams/kuching-]". The flow is named "twitter\_1\_query" and contains the following steps:

- Flow: Example flow for querying the Twitter API. For query examples, see the following web page: <https://dev.twitter.com/rest/public/search>
- TwitterConnection
- TwitterQuery query: "movie -scary :)", max tweets: 200
- TwitterConverter TextConverter (fields: TEXT, separator: {t, quote})
- Count min=-1, max=-1, interval=-1
- Convert AnyToString
- StringReplace replace '\$' with ':'
- SetVariable @-{count} [REPLACE]
- StringReplace replace '^' with '@-{count}'
- Display X:left, Y:top, W:640, H:480, font: Monospace

The Parameters panel on the right shows the following settings:

- Name: adams.flow.source.TwitterQuery
- Parameters:
  - query: movie -scary :)
  - maxTweets: 200


The Flow.Display window shows the results of the query, displaying a list of tweets. The first few tweets are:

```
170: "watching #DetroitRockCity. one of my fave\'s. used to have this movie on vhs. :-) @
171: "RT @ASPRO PANGERANO: 3 hari lg votting ditutup jgn lupa SA1 RICKY HARUN SA9 PANGERA
172: "18 years old and im seeing the peanuts movie with my mom :-)"
173: "So, I\'m watching The Pajama Game play (couldn\'t find the full movie online), and
174: "@IanMcKellen You\'re wonderful as Mr. Holmes! The movie is great and the cast is fa
175: "Congrats to @Joe_Sugg and @Caspar Lee on the premiere. I heard it went well and I c
176: "Movie :) (at @PlazaTorresLind in Gustavo A. Madero, DF) https://t.co/l49enYfKkR"
177: "Clueless is my all time favorite movie :)"
178: "What\'s the last movie you watched? - افواه و اراب :D https://t.co/yR98aYNYWq"
179: "RT @justinbieber: Last night was a movie. Love my peoples :)"
180: "@1975INTHEAM Danielle Campbell :) she was in the movie Starstruck and plays Davina
181: "@reynantem Thanks! Learned a lot.\nBTW I highly recommend that movie. Best movie Au
182: "@BenBassOfficial I just heard that you are doing a movie with @HCombs I\'m so exc
183: "RT @HeeYun hojub: Watched the movie he watched, ate what he ate... but what I reall
184: "RT @nhan127: A great movie :D\nhttps://t.co/W4gIGKcDAE"
185: "@blm849 3 of my favorite things->That song, that movie, & the incredibly bright & g
186: "RT @nhan127: A great movie :D\nhttps://t.co/W4gIGKcDAE"
187: "RT @Ashton5SOS: I hope one day we can make a 5SOS movie :)... yeah dat would be co
188: "RT @Kev_h80: Here is my interview featuring the amazing @LadyLunaWolf from yesterda
189: "My ultimate idol @chriscolfer has a new movie! Please support him for another excit
190: "movie recommendation, PLEASE? :)"
191: "@DoveCameron #LivAndMaddie #Descendants I LOVE BOTH SHOWS AND MOVIE AND I LOVE YOU
192: "@Movie_Tv fans thank u for saying my fanfic is great btw :)"
193: "RT @LaFaleel1836: The New LAURELS have arrived! :D #film #filmmaking #indiefilm #s
194: "RT @NicoleFaria9: https://t.co/4daTr4Ajeu movie :)"
195: "@mainedcm continue writing who knows it will become a full blown book and movie lat
196: "@casseracingnc @DisneyPixar Agreed! I rank Toy Story as my favorite movie of all-ti
197: "RT @tanu insan: @Gurmeetramrahim #MSG2MSG2inHaryana Thnx papa ji fr giving us a lif
198: "@NickelodeonTV how i feel right now about the new Hey Arnold TV movie :) https://t.
199: "@jumblejim I absolutely love that movie :)"
200: "@KingSargon26 haha this is a new movie coming out soon :)"
```



# Listening

---

- Rather than posting queries, listen to tweets in real time
- Public access to 1% sample of tweets
  - “garden hose” vs “fire hose”
- Sample bias?
  - <http://arxiv.org/abs/1212.1684>
- Use case?
  - create tweet archives for repeatable experiments
  - capture “the moment”
- Actor
  -  TwitterListener

convert tweets into  
spreadsheet format  
and save to CSV  
**“archive tweets”**



# Continuous archiving

---

- “Fleeting” nature of tweets makes it hard to repeat experiments
- Archival and replay of tweets solves this
- Next flow shows
  - continuous archival
  - compressed CSV spreadsheets
  - one archive per day

# Continuous archival

create filename based  
on current day

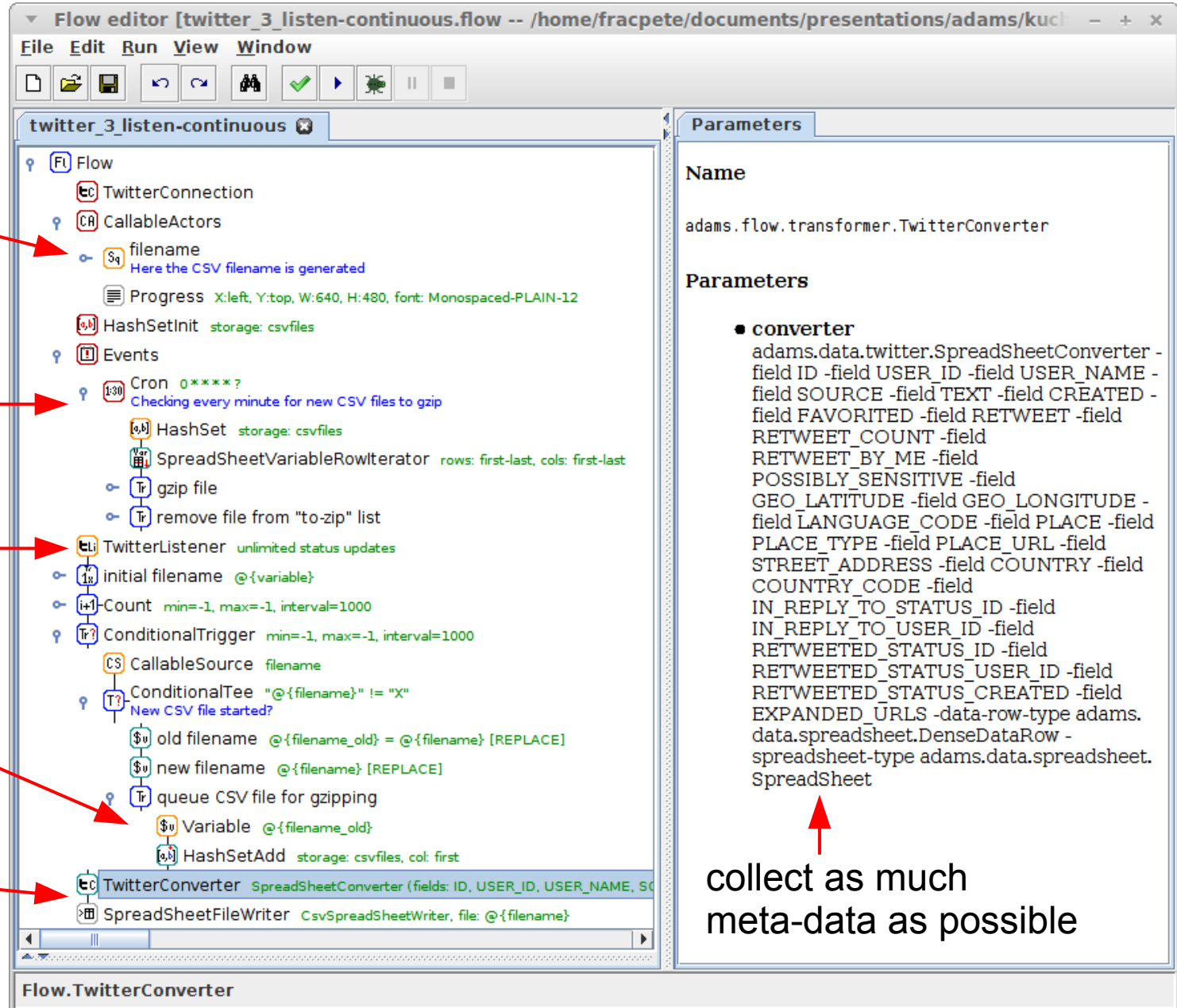
check "queue"  
whether new  
spreadsheet to  
compress

listen till flow  
gets stopped

"queue" new  
spreadsheet  
to be gzip'ed

create spreadsheet  
row and append file

26 & 27 Nov 2015



Flow editor [twitter\_3\_listen-continuous.flow -- /home/fracpete/documents/presentations/adams/kuc] - + x

File Edit Run View Window

twitter\_3\_listen-continuous

Flow

- TwitterConnection
- CallableActors
- filename  
Here the CSV filename is generated
- Progress X:left, Y:top, W:640, H:480, font: Monospaced-PLAIN-12
- HashSetInit storage: csvfiles
- Events
- Cron 0\* \* \* \* ?  
Checking every minute for new CSV files to gzip
- HashSet storage: csvfiles
- SpreadSheetVariableRowIterator rows: first-last, cols: first-last
- gzip file
- remove file from "to-zip" list
- TwitterListener unlimited status updates
- initial filename @ {variable}
- Count min=-1, max=-1, interval=1000
- ConditionalTrigger min=-1, max=-1, interval=1000
- CallableSource filename
- ConditionalTee " @ {filename} " != "X"  
New CSV file started?
- old filename @ {filename\_old} = @ {filename} [REPLACE]
- new filename @ {filename} [REPLACE]
- queue CSV file for gzipping
- Variable @ {filename\_old}
- HashSetAdd storage: csvfiles, col: first
- TwitterConverter SpreadsheetConverter (fields: ID, USER\_ID, USER\_NAME, SOURCE, TEXT, CREATED, FAVORITED, RETWEET, RETWEET\_COUNT, RETWEET\_BY\_ME, POSSIBLY\_SENSITIVE, GEO\_LATITUDE, GEO\_LONGITUDE, LANGUAGE\_CODE, PLACE, PLACE\_TYPE, PLACE\_URL, STREET\_ADDRESS, COUNTRY, COUNTRY\_CODE, IN\_REPLY\_TO\_STATUS\_ID, IN\_REPLY\_TO\_USER\_ID, RETWEETED\_STATUS\_ID, RETWEETED\_STATUS\_USER\_ID, RETWEETED\_STATUS\_CREATED, EXPANDED\_URLS, data-row-type adams.data.spreadsheet.DenseDataRow - spreadsheet-type adams.data.spreadsheet.SpreadSheet)
- SpreadSheetFileWriter CsvSpreadSheetWriter, file: @ {filename}

Parameters

Name

adams.flow.transformer.TwitterConverter

Parameters




- converter**
  - adams.data.twitter.SpreadSheetConverter - field ID -field USER\_ID -field USER\_NAME -field SOURCE -field TEXT -field CREATED -field FAVORITED -field RETWEET -field RETWEET\_COUNT -field RETWEET\_BY\_ME -field POSSIBLY\_SENSITIVE -field GEO\_LATITUDE -field GEO\_LONGITUDE -field LANGUAGE\_CODE -field PLACE -field PLACE\_TYPE -field PLACE\_URL -field STREET\_ADDRESS -field COUNTRY -field COUNTRY\_CODE -field IN\_REPLY\_TO\_STATUS\_ID -field IN\_REPLY\_TO\_USER\_ID -field RETWEETED\_STATUS\_ID -field RETWEETED\_STATUS\_USER\_ID -field RETWEETED\_STATUS\_CREATED -field EXPANDED\_URLS -data-row-type adams.data.spreadsheet.DenseDataRow - spreadsheet-type adams.data.spreadsheet.SpreadSheet

Flow.TwitterConverter

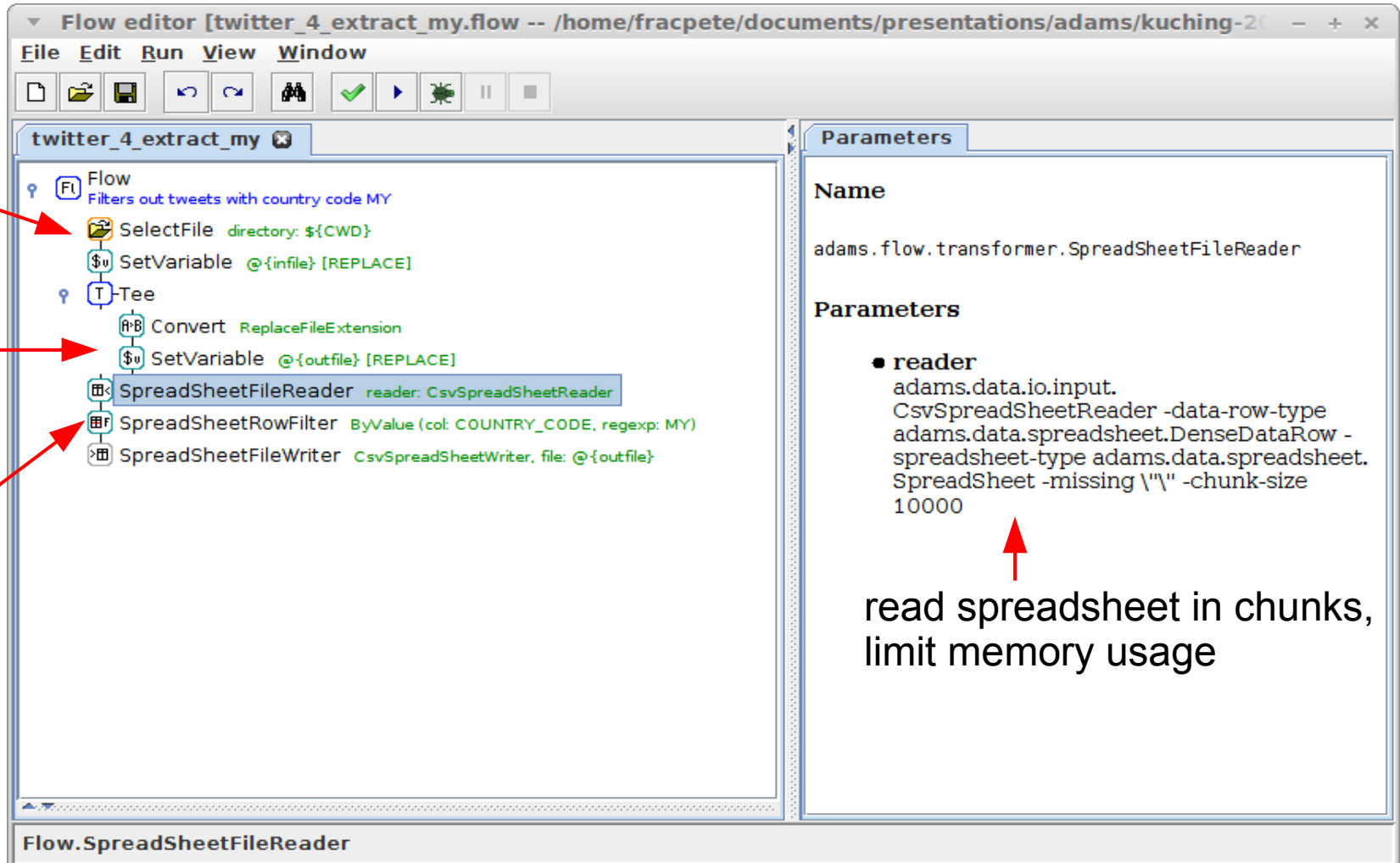
collect as much  
meta-data as possible

# Filtering archives

---

- At time of writing ~3 million tweets per day
- Experiments can operate on subsets to speed things up considerably
  - pre-filter spreadsheet archives
- Actors to use
  -  SpreadSheetFileReader
  -  SpreadSheetRowFilter
  -  SpreadSheetFileWriter

# Filtering by country



The screenshot shows the Flow editor window titled "Flow editor [twitter\_4\_extract\_my.flow -- /home/fracpete/documents/presentations/adams/kuching-20...]. The interface includes a menu bar (File, Edit, Run, View, Window) and a toolbar with icons for file operations, flow control, and execution. The main canvas displays a flow diagram for "twitter\_4\_extract\_my" with the following steps:

- Flow**: Filters out tweets with country code MY
- SelectFile**: directory: \${CWD}
- SetVariable**: @-{infile} [REPLACE]
- Tee**: A junction point for the flow.
- Convert**: ReplaceFileExtension
- SetVariable**: @-{outfile} [REPLACE]
- SpreadSheetFileReader**: reader: CsvSpreadSheetReader (highlighted with a red arrow)
- SpreadSheetRowFilter**: ByValue (col: COUNTRY\_CODE, regexp: MY) (highlighted with a red arrow)
- SpreadSheetFileWriter**: CsvSpreadSheetWriter, file: @-{outfile}

On the right side, the **Parameters** panel is visible, showing the configuration for the selected **SpreadSheetFileReader** widget:

- Name**: adams.flow.transformer.SpreadSheetFileReader
- Parameters**:
  - reader**: adams.data.io.input.CsvSpreadSheetReader -data-row-type adams.data.spreadsheet.DenseDataRow -spreadsheet-type adams.data.spreadsheet.SpreadSheet -missing \"\" -chunk-size 10000

A red arrow points to the **chunk-size** parameter, with the annotation: "read spreadsheet in chunks, limit memory usage".

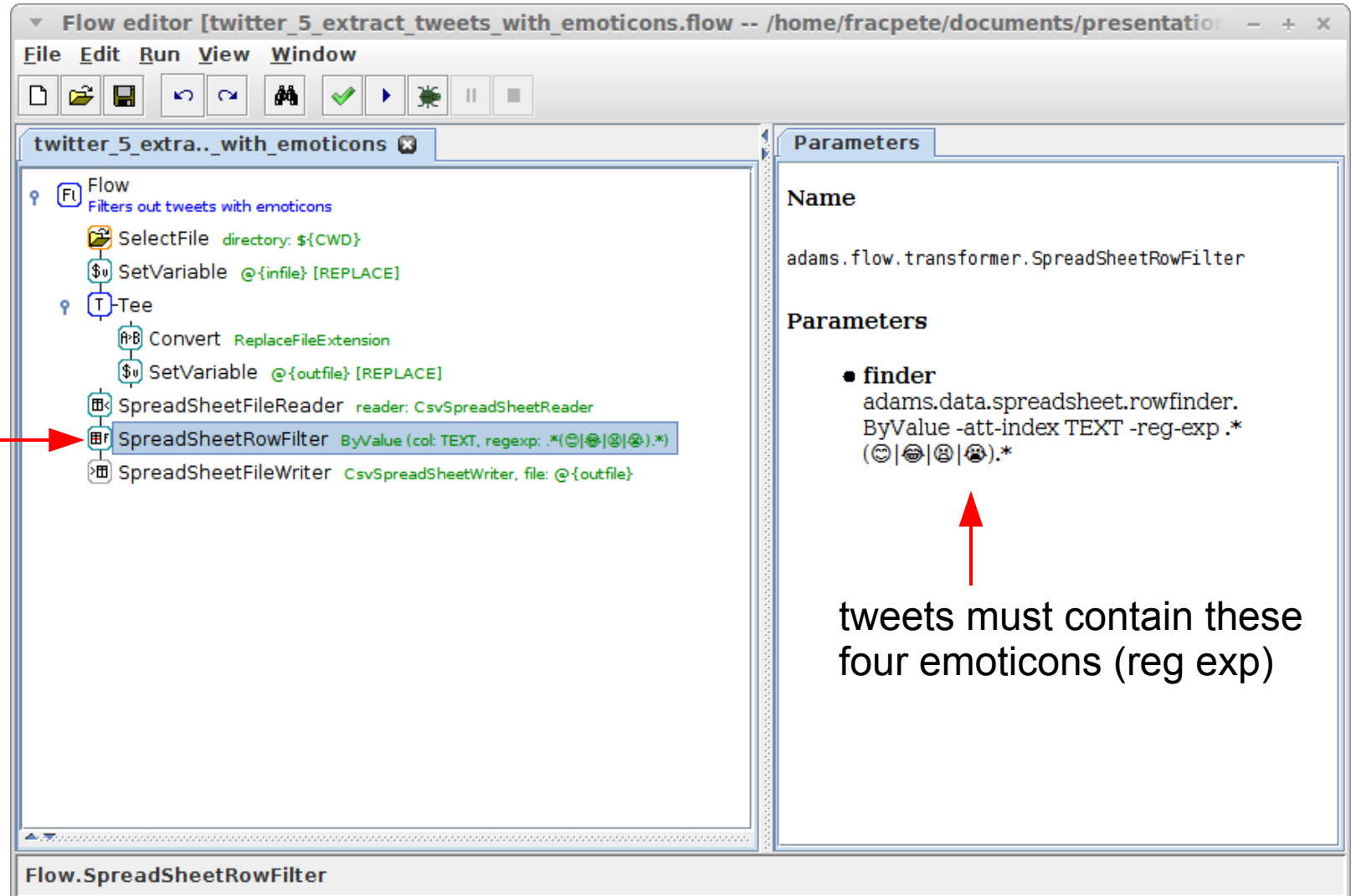
Annotations on the left side of the image, with red arrows pointing to the corresponding flow steps, include:

- "select archive files to filter" (points to **SelectFile**)
- "create output filename based on input file" (points to **SetVariable @-{outfile} [REPLACE]**)
- "leave only rows with value 'MY' in column 'COUNTRY'" (points to **SpreadSheetRowFilter**)

The status bar at the bottom of the window displays "Flow.SpreadSheetFileReader".

# Filtering by content

filter column  
"TEXT" using  
a regular  
expression



Flow editor [twitter\_5\_extract\_tweets\_with\_emoticons.flow -- /home/fracpete/documents/presentation] - + x

File Edit Run View Window

twitter\_5\_extra..\_with\_emoticons

Flow  
Filters out tweets with emoticons

- SelectFile directory: \${CWD}
- SetVariable @-{infile} [REPLACE]
- Tee
  - Convert ReplaceFileExtension
  - SetVariable @-{outfile} [REPLACE]
  - SpreadsheetFileReader reader: CsvSpreadSheetReader
  - SpreadsheetRowFilter ByValue (col: TEXT, regexp: .\* (☺|☹|😄|😞).\*)**
  - SpreadsheetFileWriter CsvSpreadSheetWriter, file: @-{outfile}

Parameters

Name

adams.flow.transformer.SpreadSheetRowFilter

Parameters


- finder  
adams.data.spreadsheet.rowfinder.  
ByValue -att-index TEXT -reg-exp.\*  
(☺|☹|😄|😞).\*

Flow.SpreadSheetRowFilter

tweets must contain these  
four emoticons (reg exp)

# Replaying archives

---

- Replaying archives is excellent for repeatable experiments
- Use cases
  - further filtering
  - feature extraction
  - plotting
- Actor
  -  TweetReplay



# New Years “Happiness”

plots happy vs sad

select archives to replay

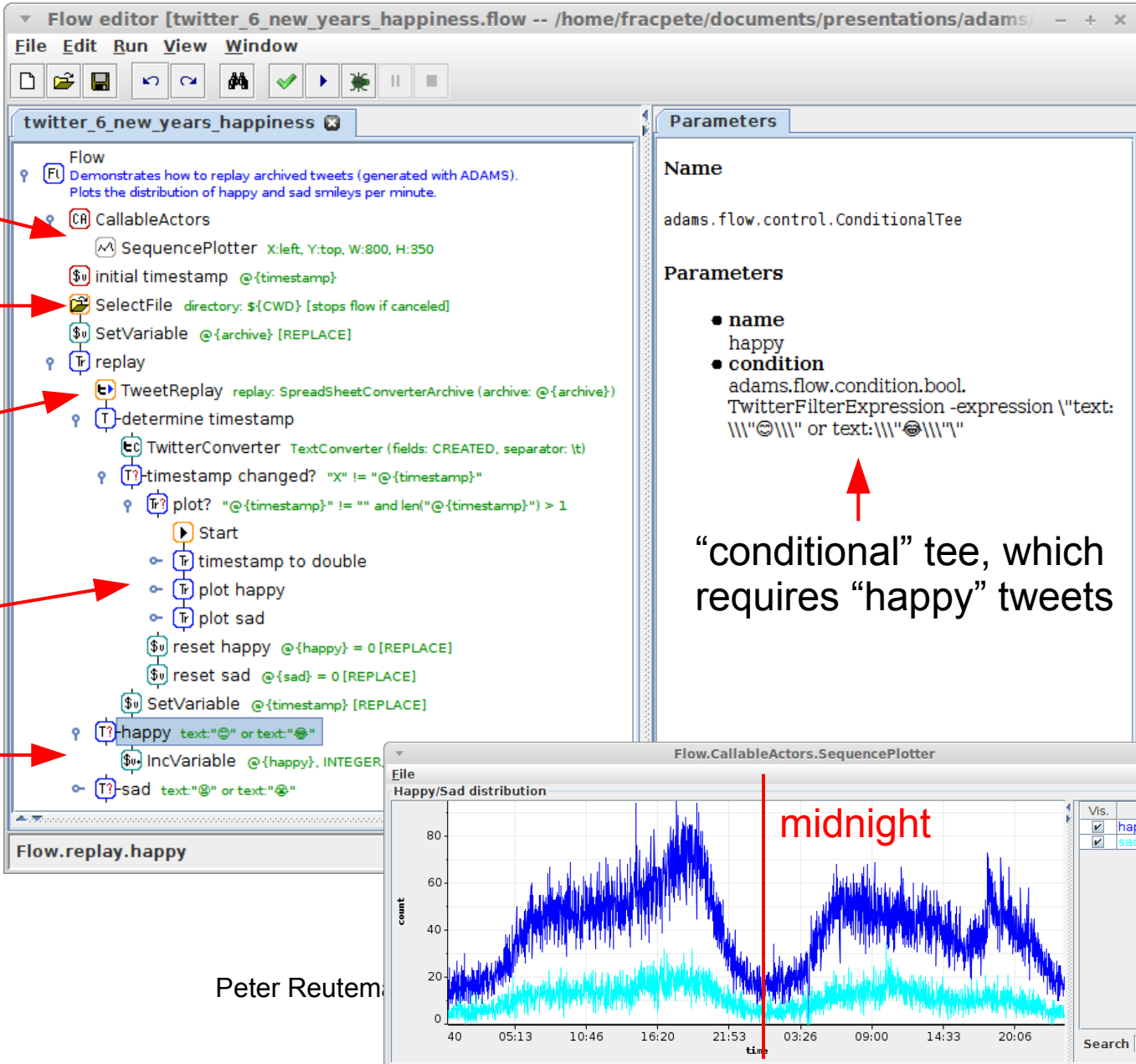
replay current archive, creating “fake” tweets from archived data

create plot containers each minute

count “happy” and “sad” tweets

“conditional” tee, which requires “happy” tweets

midnight



The screenshot shows a Flow editor window titled "Flow editor [twitter\_6\_new\_years\_happiness.flow -- /home/fracpete/documents/presentations/adams]". The flow is named "twitter\_6\_new\_years\_happiness" and contains the following steps:

- Flow: Demonstrates how to replay archived tweets (generated with ADAMS). Plots the distribution of happy and sad smileys per minute.
- CallableActors: SequencePlotter (X:left, Y:top, W:800, H:350)
- initial timestamp: @timestamp
- SelectFile: directory: \${CWD} (stops flow if canceled)
- SetVariable: @archive [REPLACE]
- replay: TweetReplay (replay: SpreadsheetConverterArchive (archive: @archive))
- determine timestamp: TwitterConverter (TextConverter (fields: CREATED, separator: \t))
- timestamp changed? "X" != "@timestamp"
- plot? "@timestamp" != "" and len("@timestamp") > 1
- Start
- timestamp to double
- plot happy
- plot sad
- reset happy @happy = 0 [REPLACE]
- reset sad @sad = 0 [REPLACE]
- SetVariable @timestamp [REPLACE]
- happy text:"😊" or text:"😄"
- IncVariable @happy, INTEGER
- sad text:"😞" or text:"😡"


The Parameters panel on the right shows the following details for the "SequencePlotter" actor:

- Name: adams.flow.control.ConditionalTee
- Parameters:
  - name: happy
  - condition: adams.flow.condition.bool. TwitterFilterExpression -expression "\"text: \\\"😊\\\" or text: \\\"😄\\\""

The bottom right shows a plot titled "Happy/Sad distribution" with a y-axis labeled "count" (0 to 80) and an x-axis labeled "time" (40 to 20:06). The plot shows two data series: "happy" (blue) and "sad" (cyan). A vertical red line marks "midnight" at 21:53. The "happy" series shows a significant peak around 16:20 and another peak around 09:00, while the "sad" series remains relatively low.

# Create dataset

---

- Build dataset from tweets as basis for predictive model
- Use happy vs sad as the “state” of the tweet (= class attribute)
- Add class attribute column using
  -  SpreadSheetInsertColumn

boolean tweet  
filter expression

only let happy/sad  
tweets pass through

insert column  
with “state”

Pet

# Create features

---

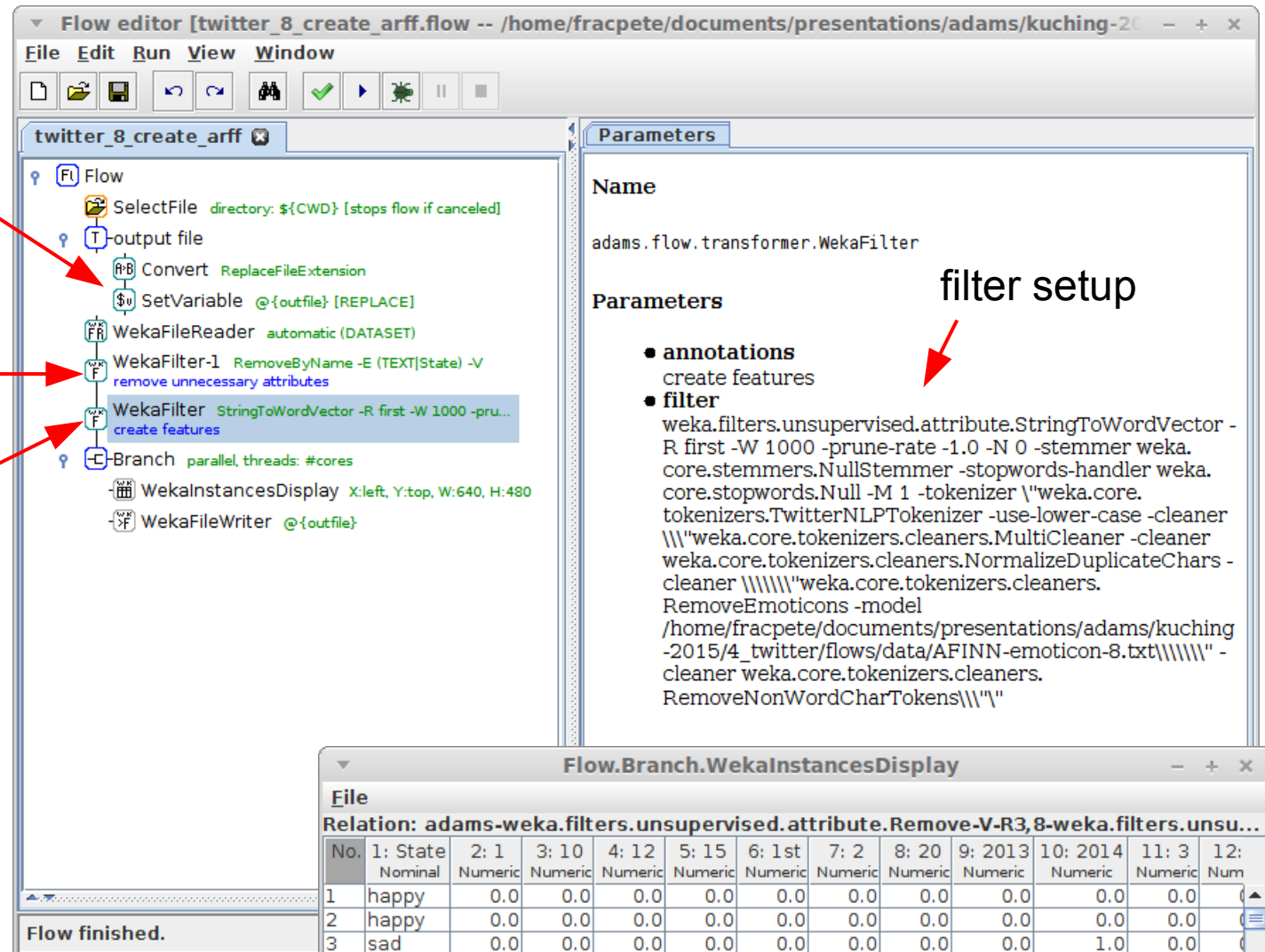
- Using WEKA's StringToWordVector to create features from tweet text
- Apply token cleaners to remove unwanted content
  - users
  - URLs
  - hashtag
  - emoticons

# Create features

create output filename

remove some attribute

generate features from tweet text



**Flow editor [twitter\_8\_create\_arff.flow -- /home/fracpete/documents/presentations/adams/kuching-20 -- + x]**

File Edit Run View Window

twitter\_8\_create\_arff

Flow

- SelectFile directory: \${CWD} [stops flow if canceled]
- output file
- Convert ReplaceFileExtension
- SetVariable @-{outfile} [REPLACE]
- WekaFileReader automatic (DATASET)
- WekaFilter-1 RemoveByName -E (TEXT)State -V remove unnecessary attributes
- WekaFilter StringToWordVector -R first -W 1000 -pru... create features
- Branch parallel, threads: #cores
- WekaInstancesDisplay X:left, Y:top, W:640, H:480
- WekaFileWriter @-{outfile}

Flow finished.

**Parameters**

Name

adams.flow.transformer.WekaFilter

Parameters

- annotations  
create features
- filter  
weka.filters.unsupervised.attribute.StringToWordVector -R first -W 1000 -prune-rate -1.0 -N 0 -stemmer weka.core.stemmers.NullStemmer -stopwords-handler weka.core.stopwords.Null -M 1 -tokenizer "\"weka.core.tokenizers.cleaners.MultiCleaner -cleaner weka.core.tokenizers.cleaners.NormalizeDuplicateChars -cleaner \"\"weka.core.tokenizers.cleaners.RemoveEmoticons -model /home/fracpete/documents/presentations/adams/kuching-2015/4\_twitter/flows/data/AFINN-emoticon-8.txt\"\" -cleaner weka.core.tokenizers.cleaners.RemoveNonWordCharTokens\"\"\"

**Flow.Branch.WekaInstancesDisplay**

File

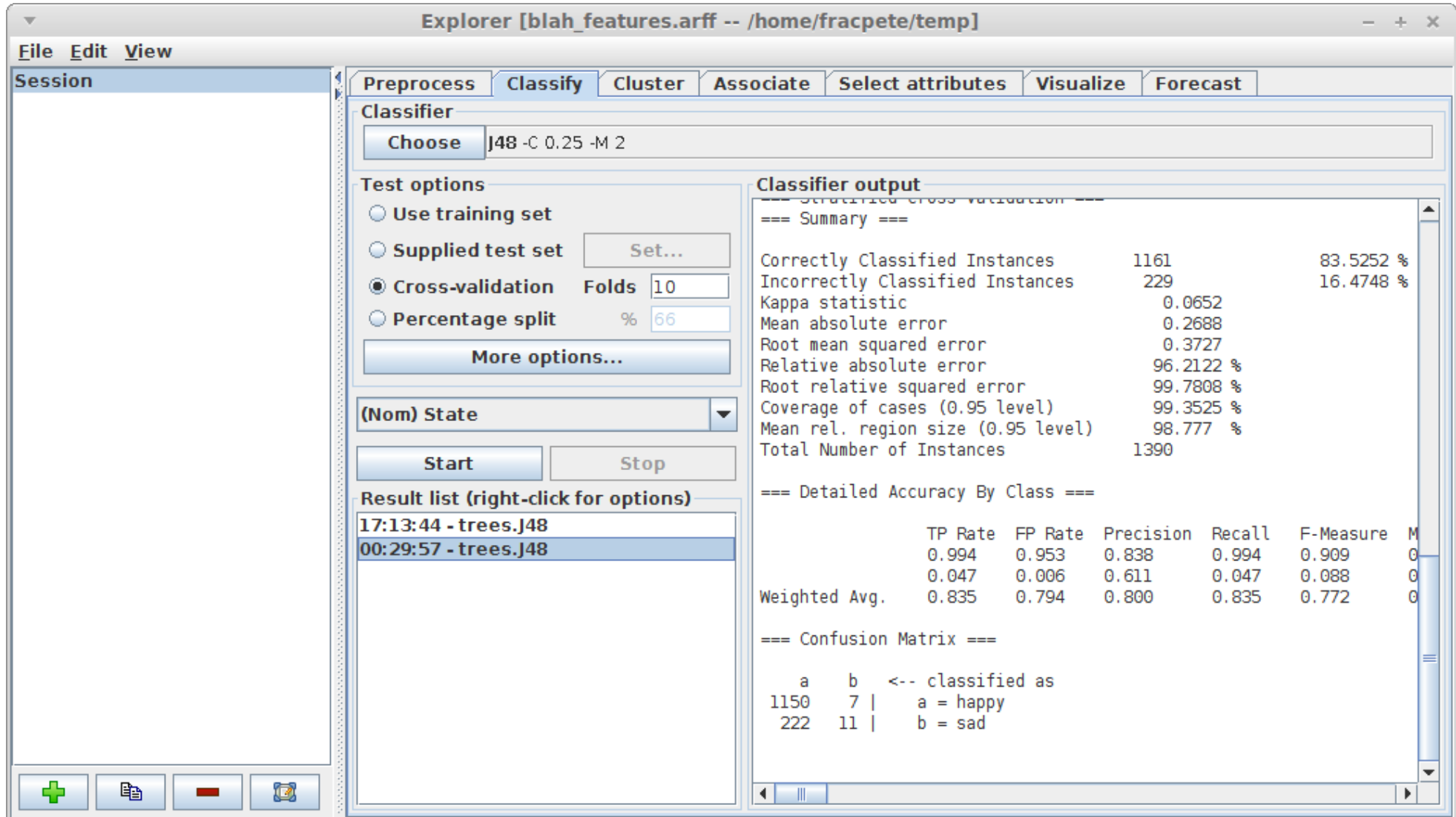
Relation: adams-weka.filters.unsupervised.attribute.Remove-V-R3,8-weka.filters.unsu...

No.	1: State Nominal	2: 1 Numeric	3: 10 Numeric	4: 12 Numeric	5: 15 Numeric	6: 1st Numeric	7: 2 Numeric	8: 20 Numeric	9: 2013 Numeric	10: 2014 Numeric	11: 3 Numeric	12: Num
1	happy	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
2	happy	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
3	sad	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	
4	happy	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
5	sad	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
6	sad	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
7	happy	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
8	happy	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
9	happy	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
10	happy	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
11	happy	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
12	happy	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	

Pet

# Create features

- using J48 on 1% percent sub-sample



The screenshot shows the Explorer software interface with the following components:

- File Edit View** menu bar.
- Session** pane on the left.
- Preprocess Classify Cluster Associate Select attributes Visualize Forecast** tabs.
- Classifier** section with a **Choose** button and the text **J48 -C 0.25 -M 2**.
- Test options** section with radio buttons for **Use training set**, **Supplied test set**, **Cross-validation** (selected), and **Percentage split**. The **Folds** field is set to **10** and the **%** field is set to **66**. A **More options...** button is also present.
- (Nom) State** dropdown menu.
- Start** and **Stop** buttons.
- Result list (right-click for options)** section showing two entries: **17:13:44 - trees.J48** and **00:29:57 - trees.J48** (selected).
- Classifier output** section displaying the following text:



```
=== Stratified cross validation ===
=== Summary ===
Correctly Classified Instances      1161      83.5252 %
Incorrectly Classified Instances    229      16.4748 %
Kappa statistic                    0.0652
Mean absolute error                0.2688
Root mean squared error            0.3727
Relative absolute error            96.2122 %
Root relative squared error        99.7808 %
Coverage of cases (0.95 level)    99.3525 %
Mean rel. region size (0.95 level) 98.777 %
Total Number of Instances         1390

=== Detailed Accuracy By Class ===
               TP Rate  FP Rate  Precision  Recall  F-Measure  M
0.994    0.953    0.838    0.994    0.909    0
0.047    0.006    0.611    0.047    0.088    0
Weighted Avg.    0.835    0.794    0.800    0.835    0.772    0

=== Confusion Matrix ===
      a    b    <-- classified as
1150    7 |    a = happy
 222   11 |    b = sad
```

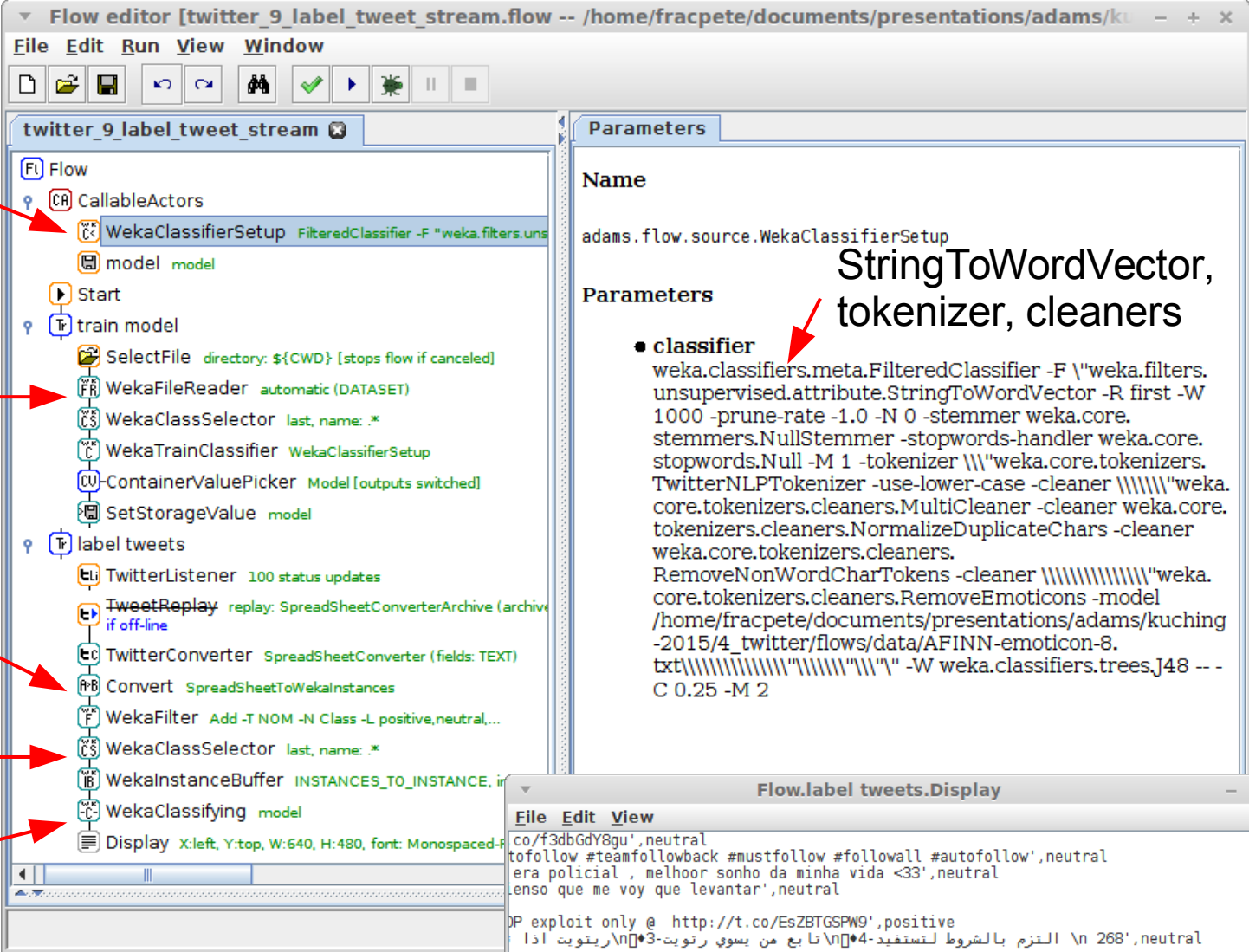
# Label tweet stream

---

- Build model on training set
  - FilteredClassifier: StringToWordVector and J48
- Predict labels for new tweets coming through
- Actors to use
  -  WekaTrainClassifier
  -  WekaClassifying



# Label tweet stream



setup generates features from tweet content

train model on training set

create Weka dataset and add class attribute

set class attribute and output row by row

apply trained model output row with filled in class value

**Parameters**

Name  
adams.flow.source.WekaClassifierSetup

Parameters

- classifier**  
weka.classifiers.meta.FilteredClassifier -F "\"weka.filters.unsupervised.attribute.StringToWordVector -R first -W 1000 -prune-rate -1.0 -N 0 -stemmer weka.core.stemmers.NullStemmer -stopwords-handler weka.core.stopwords.Null -M 1 -tokenizer '\"weka.core.tokenizers.TwitterNLPTokenizer -use-lower-case -cleaner '\"weka.core.tokenizers.cleaners.MultiCleaner -cleaner weka.core.tokenizers.cleaners.NormalizeDuplicateChars -cleaner weka.core.tokenizers.cleaners.RemoveNonWordCharTokens -cleaner '\"weka.core.tokenizers.cleaners.RemoveEmoticons -model /home/fracpete/documents/presentations/adams/kuching-2015/4\_twitter/flows/data/AFINN-emoticon-8.txt '\"weka.classifiers.trees.J48 -- C 0.25 -M 2

**Flow.label tweets.Display**

```
co/f3dbGdY8gu',neutral
tofollow #teamfollowback #mustfollow #followall #autofollow',neutral
era policial , melhor sonho da minha vida <33',neutral
enso que me voy que levantar',neutral

OP exploit only @ http://t.co/EsZBTGSPW9',positive
التزم بالشروط لتستفيد-4*تابع من يسوي رتوت-3*رتوت ادا
',neutral
stop you from being the best person you can be.',neutral
',neutral
AGg3',positive
```



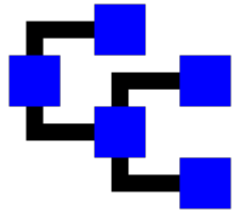
# Competition: Twitter

---

- Build a good classification model
- Use the “sanders.arff” dataset
- Perform cross-validation
- Beat ZeroR: 67.1409 %
- Hints
  - Use FilteredClassifier with StringToWordVector
  - Use TwitterNLPTokenizer
  - Use token cleaners

# Questions?

---



<https://adams.cms.waikato.ac.nz/>

@TheAdamsFlow