

2 (Unix)

In []:

```
# Extract first 250,000 lines into "training_set_tweets_250000.txt"
!head -500000 edges.csv > edges_500000.csv
```

In [18]:

```
# Swap order if userID is larger than followerID and store the result into "edges_500000_dup.csv"
!awk -F "," '{if($1<$2) printf("%d,%d\n", $1,$2);if($1>$2) printf("%d,%d\n", $2,$1)}' edges_500000.csv > edges_500000_dup.csv
```

In [19]:

```
# Find pairs that appear twice (reciprocal follower) and store it into "output.csv"
!sort edges_500000_dup.csv | uniq --count --repeated > output.csv
```

In [34]:

```
# Report reciprocal followers  
!grep -E -o " [0-9]+,[0-9]+$" output.csv | awk -F "," '{printf("%d,%d\n%d,%d\n",  
$1,$2, $2,$1)}' > result_reciprocalFollowers.txt
```

100591,100721
100721,100591
102898,122546
122546,102898
13232,18205
18205,13232
13232,63255
63255,13232
134409,134410
134410,134409
135546,135684
135684,135546
15574,15926
15926,15574
192865,192899
192899,192865
19628,19821
19821,19628
19628,20033
20033,19628
201063,40997
40997,201063
201078,201607
201607,201078
22196,76473
76473,22196
23503,41422
41422,23503
31866,32002
32002,31866
32173,32452
32452,32173
33099,62167
62167,33099
33884,34046
34046,33884
33884,34101
34101,33884
3682,5276
5276,3682
40704,40997
40997,40704
40704,41039
41039,40704
40997,41039
41039,40997
40997,62623
62623,40997
58783,58875
58875,58783
60887,70696
70696,60887
63255,65435
65435,63255
65411,65435
65435,65411
65435,93260
93260,65435
70696,70772
70772,70696
78182,78464

```
78464,78182
80092,80096
80096,80092
89222,89350
89350,89222
93260,93427
93427,93260
```

In [106]:

```
# Number of reciprocal followers: 34 * 2
!grep -E -o " [0-9]+,[0-9]+$" output.csv | awk -F "," '{printf("%d,%d\n",%d,%d\n",
$1,$2, $2,$1)}' | wc -l
```

68

In [115]:

```
# Shell script for 2
!cat 2.sh
```

```
#!/bin/sh
awk -F "," '{if($1<$2) printf("%d,%d\n", $1,$2);if($1>$2) printf("%d,%d\n", $2,$1)}' edges_500000.csv > edges_500000_dup.csv
sort edges_500000_dup.csv | uniq --count --repeated > output.csv
grep -E -o " [0-9]+,[0-9]+$" output.csv | awk -F "," '{printf("%d,%d\t%d,%d\t", $1,$2, $2,$1)}'
echo "\n"
```

In [116]:

```
# Runtime of 2 using Unix command
!time bash 2.sh
```

```
100591,100721    100721,100591    102898,122546    122546,102898    132
32,18205        18205,13232     13232,63255     63255,13232     134
409,134410      134410,134409   135546,135684   135684,135546   155
74,15926        15926,15574     192865,192899   192899,192865   196
28,19821        19821,19628     19628,20033     20033,19628     201
063,40997       40997,201063    201078,201607   201607,201078   221
96,76473        76473,22196     23503,41422     41422,23503     318
66,32002        32002,31866     32173,32452     32452,32173     330
99,62167        62167,33099     33884,34046     34046,33884     338
84,34101        34101,33884     3682,5276       5276,3682       407
04,40997        40997,40704     40704,41039     41039,40704     409
97,41039        41039,40997     40997,62623     62623,40997     587
83,58875        58875,58783     60887,70696     70696,60887     632
55,65435        65435,63255     65411,65435     65435,65411     654
35,93260        93260,65435     70696,70772     70772,70696     781
82,78464        78464,78182     80092,80096     80096,80092     892
22,89350        89350,89222     93260,93427     93427,93260     \n
1.53user 0.03system 0:00.71elapsed 219%CPU (0avgtext+0avgdata 55792
maxresident)k
0inputs+13344outputs (0major+14257minor)pagefaults 0swaps
```