

COMP26120 Lab 2 Report

Louise A. Dennis

August 28, 2022

1 Experiment 1 – Sorting Performance

1.1 Theoretical Best Case

Hypothesis The behaviour for insertion sort on sorted input is $O(n)$.

Theoretically, the best case for insertion sort is $O(n)$ when the input is sorted. The algorithm loops over the input and inserts each element at the start of the sorted list.

Experimental Design To test the hypothesis, random sorted input dictionaries were produced. 5 dictionaries of sizes 10K, 20K, 30K, 40K and 50K were generated. The spell checking program was then run on each dictionary with an input file containing a single word that was not in the dictionary in order that the look-up time should have as little impact on the comparative performance on each dictionary as possible. The time for the program to execute was measured using the UNIX `time` command summing the `user` and `sys` values output by the command.

The process of generating dictionaries and computing the time was automated using the shell scripts shown in Appendix B.

Results The results were then plotted using `gnuplot` and `gnuplot`'s `fit` functionality was used to fit a line $f(x) = m \times x + q$ calculating values for m and q in the process. The results are shown in Figure 1 and the raw data can be seen in Appendix A. As can be seen from the graph the line $f(x) = mx + q$ has a good fit and the values of 0.000007 and 0.0104 have been computed for m and q respectively. This confirms the hypothesis about the behaviour of the algorithm on sorted input, and allows us to predict the time taken to sort a dictionary of size x as $0.000007x + 0.0104$.

1.2 Theoretical Worst Case

Hypothesis The behaviour for insertion sort on reverse sorted input is $O(n^2)$.

Theoretically, the worst case for insertion sort is $O(n^2)$ when the input is reverse sorted. The algorithm loops over the input and inserts each element at the end of the sorted list.

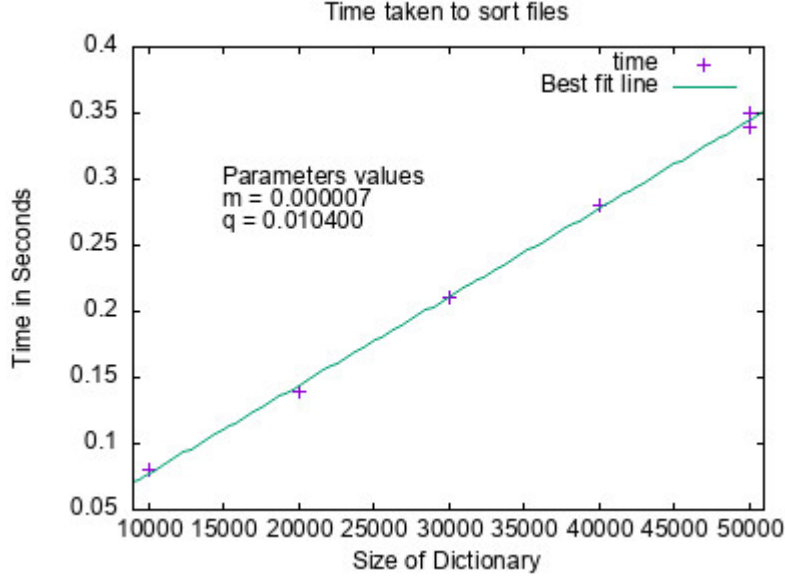


Figure 1: Time taken to look up a word in a sorted dictionary, with best fit line $f(x) = mx + q$ shown and values for the parameters m and q

Experimental Design To test the hypothesis, random reverse sorted input dictionaries were produced. 5 dictionaries of sizes 10K, 20K, 30K, 40K and 50K were generated. The spell checking program was then run on each dictionary with an input file containing a single word that was not in the dictionary in order that the look-up time should have as little impact on the comparative performance on each dictionary as possible. The time for the program to execute was measured using the UNIX `time` command summing the `user` and `sys` values output by the command.

The process of generating dictionaries and computing the time was automated using the shell scripts shown in Appendix B.

Results The results were then plotted using `gnuplot` and `gnuplot`'s `fit` functionality was used to fit a line $f(x) = m \times x^2 + q$ calculating values for m and q in the process. Note that q was initially set to 0.1 based on viewing the data points on the graph, and the fitting process didn't vary this value¹. The results are shown in Figure 2 and the raw data can be seen in Appendix A. As can be seen from the graph the line $f(x) = mx^2 + q$ has a good fit to the data and the values 0.006×10^{-7} and 0.1 have been computed for m and q respectively. This

¹This seems to be a bug in the `gnuplot` fitting algorithm where it doesn't vary q for quadratic functions. For a two hour lab setting an initial value was fine, for a scientific paper this wouldn't be acceptable without some justification.

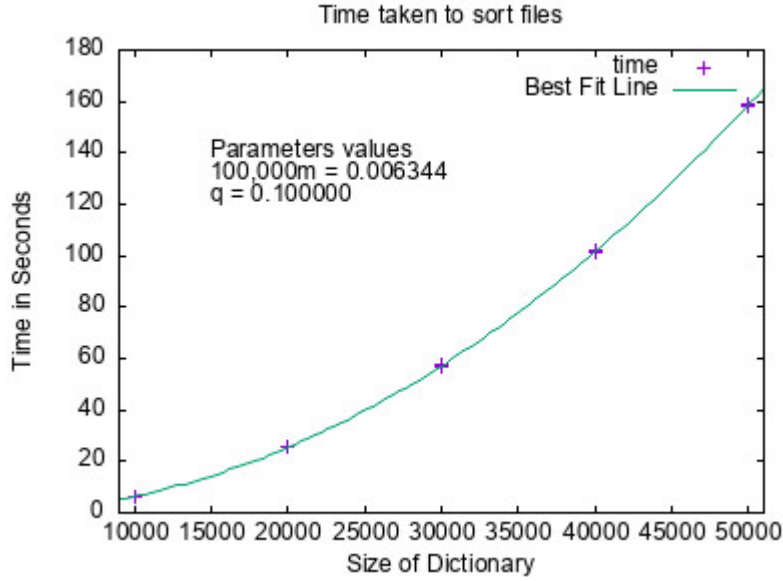


Figure 2: Time taken to look up a word in a reverse sorted dictionary, with best fit line $f(x) = mx^2 + q$ shown and values for the parameters m and q

confirms our hypothesis about the behaviour of the algorithm on reverse sorted input, and allows us to predict the time taken to sort a dictionary of size x as $0.006 \times 10^{-7}x^2 + 0.1$.

1.3 Average Case

Hypothesis The behaviour for insertion sort on random input is somewhere between the theoretical best case $O(n)$ and the theoretical worst case $O(n^2)$.

If we've identified the best and worst cases correctly, then average case performance should lie between these two. However we don't know whether the behaviour will be linear (if slower than best case) or quadratic (if faster than worst case). Plotting the data and comparing to best and worst cases should help us determine this.

Experimental Design To test the hypothesis, random input dictionaries were produced. 5 dictionaries of sizes 10K, 20K, 30K, 40K and 50K were generated. The spell checking program was then run on each dictionary with an input file containing a single word that was not in the dictionary in order that the look-up time should have as little impact on the comparative performance on each dictionary as possible. The time for the program to execute was measured using the UNIX `time` command summing the `user` and `sys` values output by

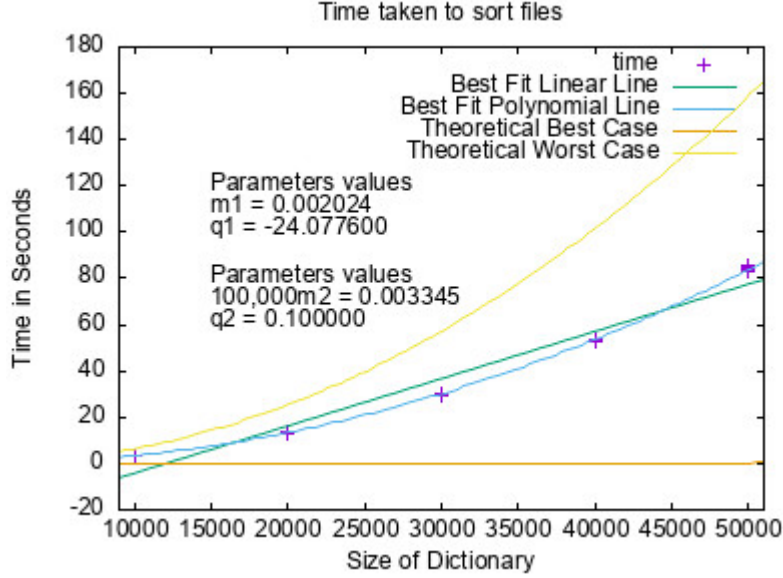


Figure 3: Time taken to look up words in a random dictionary, with best fit lines $f1(x) = m1x + q1$ and $f2(x) = m2x^2 + q2$ shown and values for the parameters $m1$, $m2$, $q1$ and $q2$

the command.

The process of generating dictionaries and computing the time was automated using the shell scripts shown in Appendix B.

Results The results were then plotted using **gnuplot** and **gnuplot**'s **fit** functionality was used to fit two lines, $f1(x) = m1 \times x + q1$ (linear) and $f2(x) = m2 \times x^2 + q2$ (quadratic) calculating values for $m1$, $m2$, $q1$ and $q2$ in the process. As above $q2$ was initially set to 0.1. The results are shown in Figure 3 which also shows the computed lines for best and worst case performance. The raw data can be seen in Appendix A. As can be seen from the graph the line both lines $f1(x)$ and $f2(x)$ can be fitted to the data. However $f2(x) = m2x^2 + q2$ is a better fit than $f1$. For $f2(x)$ the values 0.003×10^{-8} and 0.1 have been computed for $m2$ and $q2$ respectively. This confirms our hypothesis that the behaviour of the average case lies between that of the best and worst cases and lets us predict that on average the time taken to look up an entry in a dictionary of size x will be $0.003 \times 10^8 x^2 + 0.1$ seconds.

1.4 Validation and Discussion

While the average case lies between that of the best and worst cases. The fact that the line $f2(x) = m2x^2 + q2$ is a better fit for the data than a linear fit means that, in complexity terms, the performance in the average case of insertion sort is polynomial rather than linear. So average case performance of insertion sort is more comparable to worst case performance than to best case performance.

2 Experiment 2

Hypothesis Given a dictionary of size, k , and n queries. The point where it becomes quicker to sort the dictionary using insertion sort and then use binary search to check the query, as opposed to using linear search on the unsorted dictionary falls somewhere between $\frac{n}{10}$ queries and $10n$ queries.

Given a dictionary of size k and n queries. The time taken to look up the queries in the dictionary using linear search will be around nk and the time taken to sort the dictionary (average case) and then perform the look up using binary search will be around $k^2 + n \times (\log k)$. If we assume that, as the numbers become larger, the addition of $n \times (\log k)$ becomes negligible then we would expect the crossover point to occur around the moment where n becomes larger than k (so nk becomes larger than k^2). For the purpose of our hypothesis we are assuming this is in the range where n and k have the same order of magnitude – so $\frac{k}{10} \leq k \leq 10k$.

Experimental Design To test the hypothesis, random input dictionaries were produced. 5 dictionaries of sizes 10K, 20K, 30K, 40K and 50K were generated. For each dictionary five input files of queries were also produced where, if k is the size of the dictionary the query files were of the size $\frac{k}{10}$, $\frac{k}{2}$, k , $5k$ and $10k$. The spell checking program was then run on each dictionary with each input file in two modes – the first mode performed only linear search and the second mode performed insertion sort followed by binary search. The time for the program to execute was measured using the UNIX `time` command summing the `user` and `sys` values output by the command.

The process of generating dictionaries, query input files and computing the time was automated using the shell scripts shown in Appendix D.

Results The results were then plotted using `gnuplot` and `gnuplot`'s `fit` functionality was used to fit the line, $f(x) = m \times x + q$ for each value of k (Note that if k is fixed then both approaches should be linear in the size of n). The results are shown in Figure 4. The raw data can be seen in Appendix C. As can be seen from the graphs the lines for $f1(x)$ and $f2(x)$ cross at around 2.5 for all sizes of dictionary. This confirms our hypothesis that it becomes quicker to sort the dictionary of size n using insertion sort and then use binary search to check the query, as opposed to using linear search on the unsorted dictionary

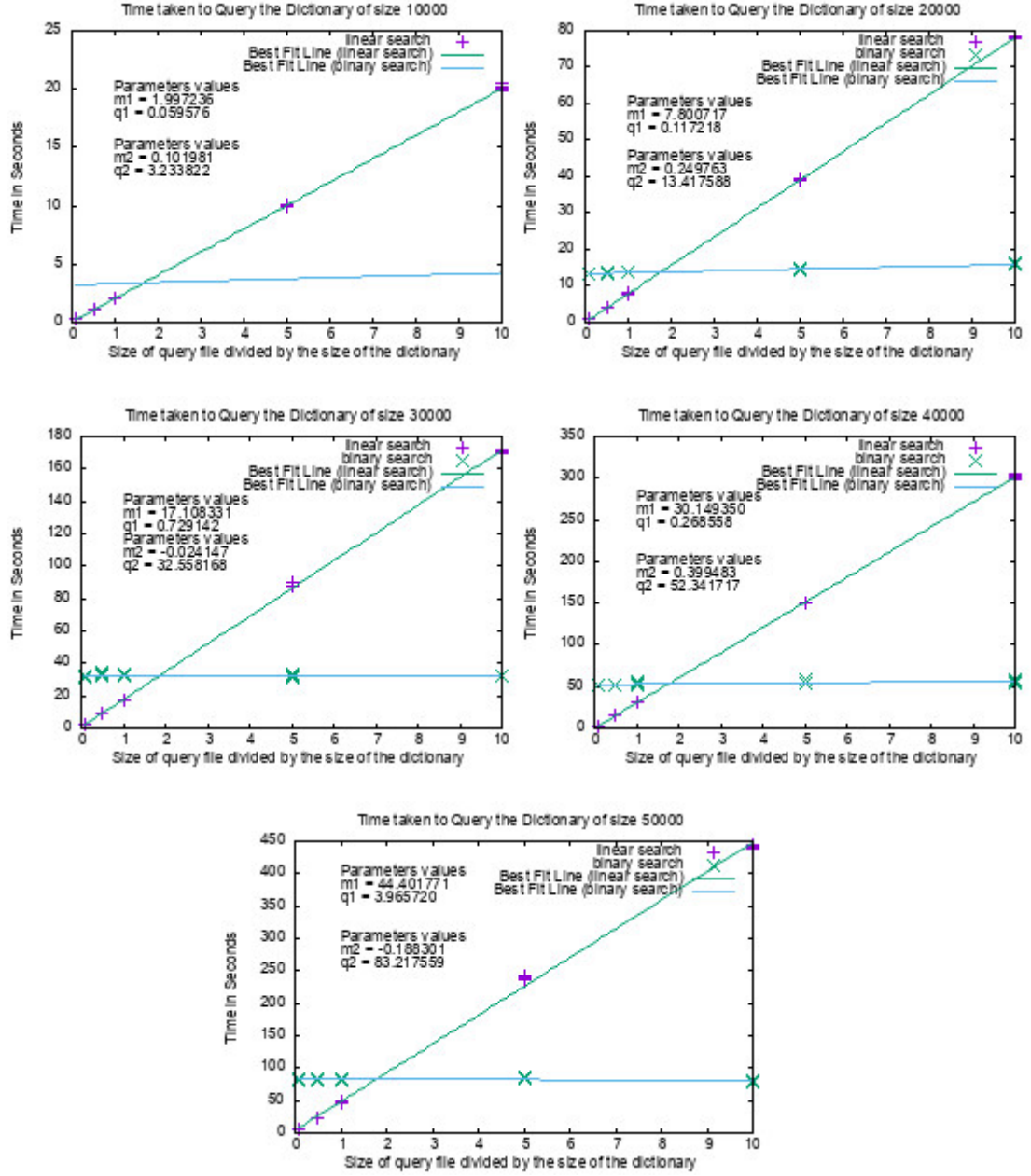


Figure 4: Time taken to look up a word in a random dictionaries of sizes 10000, 20000, 30000, 40000 and 50000 with best fit lines $f_1(x) = m_1x + q_1$ (linear search) and $f_2(x) = m_2x^2 + q_2$ (insertion sort plus binary search) shown and values for the parameters m_1, m_2, q_1 and q_2

somewhere between $\frac{n}{10}$ queries and $10n$ queries. In fact it tells us that this point occurs somewhere around $2.5n$ queries.

A Raw Data for Experiment 1

Sorted Data		Reverse Sorted Data		Random Data	
Size	Time (s)	Size	Time (s)	Size	Time (s)
10000	.08	10000	6.43	10000	3.18
10000	.08	10000	6.40	10000	3.20
10000	.08	10000	6.40	10000	3.20
10000	.08	10000	6.38	10000	3.20
10000	.08	10000	6.39	10000	3.26
20000	.14	20000	25.80	20000	13.12
20000	.14	20000	25.97	20000	13.03
20000	.14	20000	25.60	20000	12.94
20000	.14	20000	25.48	20000	12.92
20000	.14	20000	25.62	20000	13.41
30000	.21	30000	57.49	30000	29.19
30000	.21	30000	57.46	30000	29.65
30000	.21	30000	57.40	30000	30.01
30000	.21	30000	57.30	30000	29.45
30000	.21	30000	56.88	30000	29.36
40000	.28	40000	101.80	40000	52.39
40000	.28	40000	101.79	40000	52.40
40000	.28	40000	101.66	40000	52.57
40000	.28	40000	101.64	40000	53.69
40000	.28	40000	101.64	40000	53.49
50000	.34	50000	158.16	50000	83.03
50000	.35	50000	159.24	50000	84.15
50000	.34	50000	158.28	50000	84.41
50000	.34	50000	158.52	50000	86.00
50000	.35	50000	158.69	50000	84.93

B Shell Scripts for Experiment 1

B.1 Generating Dictionaries

```
STATES="sorted reverse none"
SIZES="10000 20000 30000 40000 50000 60000 70000 80000 90000 100000"

for STATE in $STATES
do

for SIZE in $SIZES
```

```

do

echo $SIZE

for COUNT in 1 2 3 4 5
do

    python3 generate.py random dict_${SIZE}_${STATE}_${COUNT} query_${SIZE}_${STATE}_${COUNT} $S

done

done

done

```

B.2 Script for Computing Run Times

```

STATES="sorted reverse none"
SIZES="10000 20000 30000 40000 50000"

rm data/sorted_data_python.dat
rm data/reverse_data_python.dat
rm data/none_data_python.dat
rm data/sorted_data_python.csv
rm data/reverse_data_python.csv
rm data/none_data_python.cvs

for STATE in $STATES
do

for SIZE in $SIZES
do

for COUNT in 1 2 3 4 5
do

    # Debugging statement to check program calls working as expected
    python3 ../../search_and_sort_lab/python/speller_darray.py -d ../dictionaries_and_queries

    ALL_TIME='(time -p python3 ../../search_and_sort_lab/python/speller_darray.py -d ../dictionaries_and_queries

    RUNTIME=0
    for i in $ALL_TIME;
    do RUNTIME='echo $RUNTIME + $i|bc';
    done
    echo $SIZE $RUNTIME >> data/${STATE}_data_python.dat

```



```
echo $SIZE, $RUNTIME >> data/${STATE}_data_python.csv
```

done

done

done

C Raw Data for Experiment 2

Dictionary Size:10K, Linear Search		Dictionary Size:10K, Binary Search	
$\frac{k}{n}$	Time (s)	$\frac{k}{n}$	Time (s)
0.1	.27	0.1	3.25
0.1	.27	0.1	3.25
0.1	.27	0.1	3.24
0.1	.27	0.1	3.24
0.1	.27	0.1	3.25
0.5	1.05	0.5	3.28
0.5	1.05	0.5	3.27
0.5	1.07	0.5	3.28
0.5	1.05	0.5	3.29
0.5	1.05	0.5	3.29
1	2.06	1	3.35
1	2.05	1	3.33
1	2.05	1	3.35
1	2.11	1	3.34
1	2.06	1	3.34
5	10.05	5	3.73
5	10.16	5	3.73
5	9.94	5	3.74
5	9.99	5	3.75
5	9.94	5	3.72
10	19.84	10	4.22
10	19.82	10	4.21
10	19.98	10	4.27
10	20.15	10	4.25
10	20.44	10	4.34

Dictionary Size:20K, Linear Search		Dictionary Size:20K, Binary Search	
$\frac{k}{n}$	Time (s)	$\frac{k}{n}$	Time (s)
0.1	.93	0.1	13.50
0.1	.93	0.1	13.42
0.1	.92	0.1	13.46
0.1	.91	0.1	13.35
0.1	.93	0.1	13.42
0.5	4.02	0.5	13.54
0.5	4.02	0.5	13.59
0.5	4.05	0.5	13.41
0.5	3.95	0.5	13.54
0.5	4.00	0.5	13.54
1	7.96	1	13.73
1	7.90	1	13.73
1	7.96	1	13.77
1	8.13	1	13.78
1	7.94	1	13.80
5	38.96	5	14.73
5	39.37	5	14.51
5	38.74	5	14.34
5	38.81	5	14.46
5	39.01	5	14.74
10	78.54	10	16.25
10	78.07	10	15.86
10	78.24	10	16.02
10	77.73	10	15.95
10	78.37	10	15.73

Dictionary Size:30K, Linear Search		Dictionary Size:30K, Binary Search	
$\frac{k}{n}$	Time (s)	$\frac{k}{n}$	Time (s)
0.1	1.93	0.1	30.77
0.1	1.94	0.1	32.13
0.1	1.96	0.1	32.12
0.1	1.96	0.1	31.91
0.1	1.94	0.1	32.36
0.5	8.94	0.5	32.11
0.5	9.01	0.5	33.11
0.5	9.00	0.5	32.26
0.5	8.99	0.5	32.22
0.5	9.09	0.5	34.35
1	17.60	1	32.84
1	17.78	1	33.42
1	17.60	1	33.04
1	17.76	1	33.30
1	17.63	1	32.48
5	87.92	5	33.29
5	87.36	5	31.97
5	87.46	5	33.86
5	87.70	5	32.33
5	89.91	5	30.80
10	169.84	10	32.35
10	171.04	10	32.00
10	171.03	10	32.42
10	171.01	10	32.24
10	171.82	10	32.27

Dictionary Size:40K, Linear Search		Dictionary Size:40K, Binary Search	
$\frac{k}{n}$	Time (s)	$\frac{k}{n}$	Time (s)
0.1	3.27	0.1	52.43
0.1	3.30	0.1	52.71
0.1	3.26	0.1	52.40
0.1	3.30	0.1	51.89
0.1	3.27	0.1	51.73
0.5	15.37	0.5	52.48
0.5	15.31	0.5	52.31
0.5	15.42	0.5	51.95
0.5	15.58	0.5	52.16
0.5	15.61	0.5	51.22
1	30.99	1	51.36
1	30.51	1	52.11
1	30.50	1	53.30
1	30.94	1	57.01
1	30.60	1	54.80
5	151.42	5	57.45
5	150.14	5	53.14
5	149.35	5	52.92
5	151.06	5	52.93
5	149.44	5	52.88
10	299.36	10	54.77
10	300.39	10	56.61
10	302.89	10	56.45
10	303.19	10	57.56
10	304.64	10	57.13

Dictionary Size:50K, Linear Search		Dictionary Size:50K, Binary Search	
$\frac{k}{n}$	Time (s)	$\frac{k}{n}$	Time (s)
0.1	5.04	0.1	82.93
0.1	5.09	0.1	82.73
0.1	5.07	0.1	82.71
0.1	5.10	0.1	82.09
0.1	5.08	0.1	82.57
0.5	24.04	0.5	83.71
0.5	24.19	0.5	84.22
0.5	24.21	0.5	80.27
0.5	23.45	0.5	80.75
0.5	23.38	0.5	80.57
1	46.81	1	80.64
1	46.67	1	83.44
1	48.23	1	83.83
1	47.83	1	83.65
1	47.93	1	83.76
5	239.90	5	86.15
5	236.81	5	85.79
5	238.46	5	86.55
5	241.26	5	84.96
5	236.24	5	85.07
10	441.12	10	79.45
10	443.57	10	80.05
10	444.69	10	80.09
10	441.90	10	79.61
10	438.42	10	79.22

D Shell Scripts for Experiment 2

D.1 Generating Dictionaries and Queries

```

SIZES="10000 20000 30000 40000 50000"

QUERY_SIZES="0.1 0.5 1 5 10"

for SIZE in $SIZES
do

for QSIZE in $QUERY_SIZES
do

QUERY_SIZE=$(( echo "$SIZE * $QSIZE/1" | bc))

echo $QUERY_SIZE

```

```

for COUNT in 1 2 3 4 5
do

    python3 generate.py random dict_${SIZE}_${QUERY_SIZE}_${COUNT} query_${SIZE}_${QUERY_SIZE}

done

done

done

```

D.2 Computing Run Times

```

SIZES="10000 20000 30000 40000 50000"
QUERY_SIZES="0.1 0.5 1 5 10"

for SIZE in $SIZES
do

    rm data/${SIZE}_amortised_data_linear_python.dat
    rm data/${SIZE}_amortised_data_binary_python.dat
    rm data/${SIZE}_amortised_data_linear_python.csv
    rm data/${SIZE}_amortised_data_binary_python.csv

    for QSIZE in $QUERY_SIZES
    do

        QUERY_SIZE=$(( echo "$SIZE * $QSIZE/1" | bc))

        for COUNT in 1 2 3 4 5
        do

            # Debugging command to check program call actually works
            python3 ../../search_and_sort_lab/python/speller_darray.py -d ../dictionaries_and_queries

            ALL_TIME='(time -p python3 ../../search_and_sort_lab/python/speller_darray.py -d ../dictionaries_and_queries

            RUNTIME=0
            for i in $ALL_TIME;
            do RUNTIME='echo $RUNTIME + $i|bc';
            done

```

```

echo $QSIZE $RUNTIME >> data/${SIZE}_amortised_data_linear_python.dat
echo $QSIZE, $RUNTIME >> data/${SIZE}_amortised_data_linear_python.csv

# Debugging command to check program call actually works
python3 ../../search_and_sort_lab/python/speller_darray.py -d ../dictionaries_and_queries

ALL_TIME2='(time -p python3 ../../search_and_sort_lab/python/speller_darray.py -d ../dictionaries_and_queries

RUNTIME2=0
for i in $ALL_TIME2;
do RUNTIME2='echo $RUNTIME2 + $i|bc';
done

echo $QSIZE $RUNTIME2 >> data/${SIZE}_amortised_data_binary_python.dat
echo $QSIZE, $RUNTIME2 >> data/${SIZE}_amortised_data_binary_python.csv

done

done

done

```