

# Assignment 2: Implementing MapReduce

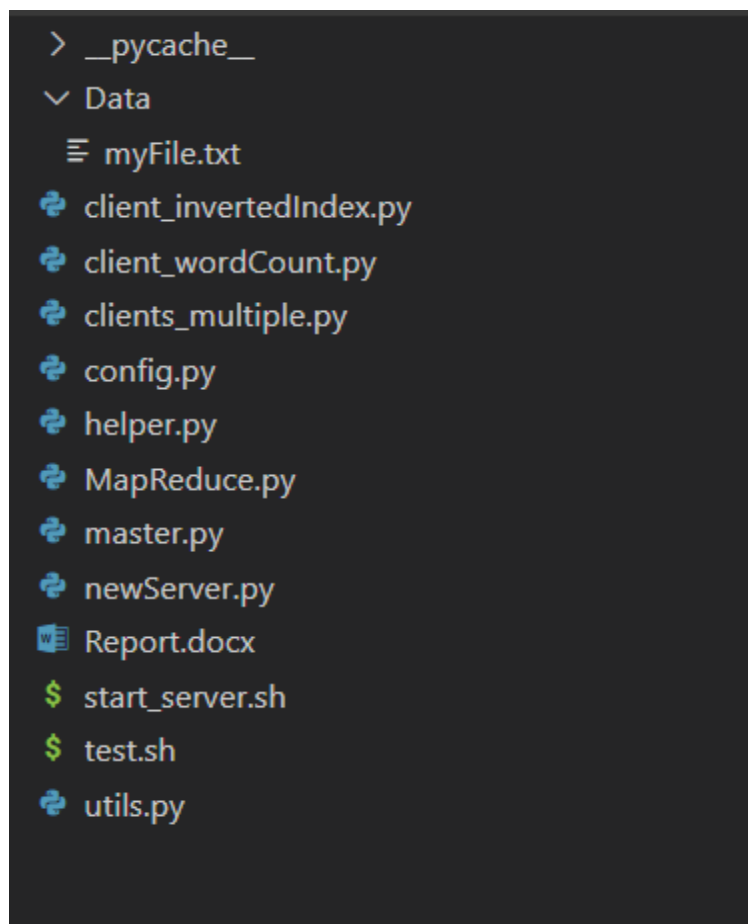
## INTRODUCTION:

MapReduce is a programming methodology for processing huge data collections. A map function analyzes a key/value pair to yield a collection of intermediate key/value pairs, whereas a reduce function merges all intermediate values associated with the same intermediate key. This model may be used to represent a wide range of real-world operations.

In this assignment, the two main applications are:

1. Word Count
2. Inverted Index

## ARCHITECTURE:



The implementation has the following files:

SUBMITTED BY:  
PRASHASTI KARLEKAR  
prkarl@iu.edu

# Assignment 2: Implementing MapReduce

## 1. **SERVER:**

A TCP server written in Python 3 that receives client requests and performs map reduce operations based on the provided input file. This server is capable of handling connections from several clients at the same time (the response time majorly depends on the capability of the hardware its running on). The server creates and calls an object of main, which is located in the Master implementation. The result of the main is returned to the client by the server.

## 2. **MASTER:**

This acts as the master node. Based upon the request by the client, master spawns the respective map and reduce functions on the input data located at the input location provided by the client. If the map function is “word\_count”, the master will create an object of class WordCount which then handles further operations. If it receives “inverted\_index”, it will create an object of InvertedIndex class.

## 3. **MAPREDUCE:**

This contains the implementations of WordCount and InvertedIndex classes.

### ➤ WordCount class implementation:

The mapper tasks reads the input file corresponding to the respective file split. It then emits key-value pairs from the file and passes this to the combiner. The combiner reformats the produced list by the mapper function into key, value pairs that complements as an input to the reduced function by combining all the values corresponding to each key. It returns tuples of (word, [# of times the word appears]). This combined data is saved on the local file system. The reduce function then reduces the computed results by each worker into the final output by summing all the values.

### ➤ InvertedIndex class implementation:

The first task performed on the input is the split\_input operation where the file is divided into chunks equal to the number of mappers. The mappers run parallelly with each mapper getting a chunk and the output of each mapper is combined at the end and stored as MergedData. After this, reducer function gets a unique key and performs reduce operation on it. Once all the reducers have been executed, we combine the results and store it as an output file at the output location provided.

## 4. **UTILS:**

This file contains the helper functions required during the map and reduce operations. These include split\_input, merge\_files and join\_files.

## 5. **HELPER :**

This includes the map and reduce functions corresponding to WordCount and InvertedIndex objects.

## 6. **CONFIG:**

This file contains Host and port information which can be changed accordingly.

## Assignment 2: Implementing MapReduce

Below are the output files for both the applications:

word\_count\_output - Notepad

File Edit View

[('spake', 182), ('zarathustra', 627), ('book', 16), ('all', 895), ('none', 15), ('friedrich', 3), ('nietzsche', 110), ('translated', 1), ('thomas', 1), ('common', 15), ('pg', 2), ('note', 31), ('archaic', 1), ('spelling', 1), ('punctuation', 1), ('usages', 1), ('have', 675), ('not', 1013), ('been', 153), ('changed', 12), ('particular', 7), ('quotations', 1), ('often', 45), ('closed', 7), ('several', 4), ('paragraphs', 8), ('dw', 1), ('contents', 2), ('introduction', 5), ('mrs', 2), ('forster', 3), ('thus', 404), ('first', 116), ('part', 52), ('prologue', 9), ('discourses', 7), ('three', 34), ('metamorphoses', 5), ('i', 7), ('academic', 7), ('chairs', 6), ('virtue', 133), ('iii', 6), ('backworldsmen', 8), ('iv', 6), ('despisers', 17), ('body', 69), ('joys', 12), ('passions', 8), ('v', 2), ('pale', 11), ('criminal', 5), ('vii', 2), ('reading', 7), ('writing', 11), ('viii', 2), ('tree', 34), ('hill', 5), ('ix', 4), ('preachers', 20), ('death', 73), ('war', 24), ('warriors', 6), ('xi', 2), ('new', 157), ('idol', 10), ('xii', 2), ('files', 11), ('market', 21), ('place', 62), ('xiii', 2), ('chastity', 10), ('xiv', 2), ('friend', 58), ('xv', 3), ('thousand', 23), ('one', 725), ('goals', 5), ('xvi', 2), ('neighbour', 30), ('love', 297), ('xvii', 2), ('way', 117), ('creating', 45), ('xviii', 4), ('old', 173), ('young', 26), ('women', 28), ('xix', 2), ('bite', 16), ('adder', 5), ('xx', 2), ('child', 52), ('marriage', 35), ('xoi', 3), ('voluntary', 18), ('xxii', 4), ('bestowing', 17), ('second', 28), ('xxiii', 3), ('mirror', 14), ('xxiv', 3), ('happy', 25), ('isles', 15), ('xxv', 2), ('pitiful', 12), ('xxvi', 2), ('priests', 9), ('xxvii', 2), ('virtuous', 17), ('xxviii', 2), ('rabble', 22), ('xxix', 3), ('tarantulas', 8), ('xxxi', 3), ('famous', 12), ('wise', 49), ('ones', 225), ('xxx', 2), ('night', 74), ('song', 36), ('xxxi', 2), ('dance', 39), ('xxxiii', 3), ('grave', 18), ('xxxiv', 3), ('self', 66), ('surpassing', 5), ('xxxv', 3), ('sublime', 14), ('xxxvi', 5), ('land', 22), ('culture', 6), ('xxxvii', 3), ('immaculate', 6), ('perception', 8), ('xxxviii', 3), ('scholars', 7), ('xxxix', 3), ('poets', 26), ('xl', 4), ('great', 239), ('events', 8), ('xli', 3), ('soothsayer', 35), ('xlii', 4), ('redemption', 9), ('xliii', 3), ('manly', 7), ('prudence', 11), ('xliv', 3), ('stillest', 12), ('hour', 62), ('third', 16), ('xlv', 2), ('wanderer', 16), ('xvi', 4), ('vision', 8), ('enigma', 7), ('xlvii', 3), ('involutionary', 6), ('bliss', 17), ('xlviii', 3), ('before', 115), ('sunrise', 6), ('xlix', 4), ('bedwarfing', 6), ('olive', 6), ('mount', 14), ('li', 4), ('passing', 11), ('lii', 3), ('apostates', 4), ('liii', 4), ('return', 38), ('home', 46), ('liv', 4), ('evil', 115), ('things', 168), ('lv', 3), ('spirit', 170), ('gravity', 19), ('lvi', 7), ('tables', 27), ('lvii', 10), ('convalscent', 11), ('lviii', 2), ('longing', 56), ('lix', 2), ('lx', 3), ('seven', 17), ('seals', 3), ('fourth', 6), ('last', 162), ('lxi', 3), ('honey', 28), ('sacifice', 18), ('lxii', 3), ('cry', 45), ('distress', 23), ('lxiii', 5), ('talk', 33), ('kings', 33), ('lxiv', 3), ('leech', 11), ('lvv', 4), ('magician', 34), ('lxvi', 3), ('out', 255), ('service', 14), ('lxvii', 3), ('ugliest', 26), ('man', 380), ('lxviii', 3), ('beggar', 15), ('lxix', 3), ('shadow', 38), ('lxx', 4), ('noon', 4), ('tide', 3), ('lxxi', 4), ('greeting', 7), ('lxii', 3), ('supper', 3), ('higher', 145), ('lxix', 3), ('melancholy', 27), ('lxxv', 3), ('science', 10), ('lxxv', 3), ('among', 69), ('daughters', 6), ('desert', 12), ('lxxvii', 4), ('awakening', 6), ('lxxviii', 3), ('pass', 44), ('festival', 16), ('lxxix', 3), ('drunken', 15), ('lxxx', 3), ('sign', 16), ('appendix', 2), ('notes', 17), ('anthony', 3), ('ludovici', 3), ('how', 254), ('came', 87), ('into', 256), ('being', 48), ('my', 976), ('most', 150), ('personal', 5), ('work', 83), ('history', 13), ('his', 898), ('individual', 12), ('experiences', 4), ('friendships', 4), ('ideals', 6), ('raptures', 1), ('bitterest', 1), ('disappointments', 5), ('sorrows', 1), ('above', 57), ('however', 373), ('there', 485), ('soars', 1), ('transfiguring', 1), ('image', 11), ('greatest', 50), ('hopes', 11), ('remotest', 7), ('aims', 5), ('brother', 33), ('had', 262), ('figure', 9), ('mind', 12), ('from', 448), ('very', 78), ('earliest', 2), ('youth', 30), ('he', 1089), ('once', 155), ('told', 13), ('me', 825), ('even', 253), ('dreamt', 5), ('him', 406), ('at', 509), ('different', 19), ('periods', 4), ('life', 186), ('would', 260), ('call', 101), ('this', 671), ('haunter', 1), ('dreams', 16), ('names', 17), ('end', 45), ('declares', 4), ('subject', 9), ('do', 519), ('persian', 3), ('honour', 33), ('...')]

### Output for Word Count

inverted\_index\_output - Notepad

File Edit View

("abandon": [[3, 7]], "abandoned": [[1, 35], [3, 7]], "abasement": [[0, 7]], "abash": [[0, 21]], "abashed": [[0, 7]], "abet": [[1, 7]], "abide": [[1, 7], [2, 21], [3, 7]], "ability": [[1, 7], [2, 7], [3, 14]], "abjures": [[3, 7]], "able": [[0, 49], [1, 7], [2, 14], [3, 63]], "abnegation": [[2, 7]], "abode": [[0, 7]], "abodes": [[0, 7]], "about": [[0, 259], [1, 259], [2, 273], [3, 140]], "above": [[0, 105], [1, 133], [2, 77], [3, 91]], "abroad": [[0, 7], [1, 21]], "abrupt": [[0, 7]], "absent": [[1, 7]], "absolute": [[0, 14], [2, 28], [3, 7]], "absolutely": [[0, 7]], "absoluteness": [[0, 7]], "absolved": [[2, 7]], "absorbed": [[0, 7], [2, 14], [3, 7]], "abstain": [[2, 7]], "abstruse": [[3, 14]], "absurd": [[0, 7]], "abundance": [[0, 21], [1, 14], [2, 21]], "abuse": [[1, 7]], "abused": [[1, 14]], "abusive": [[3, 7]], "abut": [[1, 7]], "abyssmal": [[1, 21], [2, 14], [3, 7]], "abyss": [[0, 14], [1, 56], [2, 56]], "abysses": [[1, 35], [2, 14]], "academic": [[0, 35], [1, 7], [3, 7]], "accept": [[0, 14], [1, 7], [3, 21]], "accepted": [[3, 21]], "accepting": [[0, 14], [3, 7]], "access": [[1, 7], [3, 70]], "accessed": [[3, 7]], "accessible": [[3, 7]], "accident": [[3, 7]], "accidents": [[0, 28]], "accommodate": [[1, 7]], "accompanied": [[0, 7], [2, 7]], "accompany": [[3, 7]], "accomplish": [[0, 7]], "accomplished": [[1, 7]], "accord": [[0, 21], [1, 21]], "accordance": [[2, 7], [3, 14]], "according": [[0, 21], [1, 14], [2, 21], [3, 35]], "account": [[1, 118], [1, 112], [2, 112], [3, 21]], "accredit": [[0, 7]], "accumulate": [[0, 7]], "accumulated": [[2, 7]], "accuracy": [[0, 7]], "accusation": [[2, 7], [3, 7]], "accuse": [[2, 7]], "accuser": [[2, 7]], "accusers": [[2, 7]], "accuses": [[3, 7]], "accusing": [[2, 7]], "achieve": [[3, 7]], "acknowledge": [[0, 7], [2, 7]], "acknowledged": [[2, 7]], "acorns": [[0, 7]], "acquaintances": [[0, 7]], "acquainted": [[3, 21]], "acquire": [[0, 14], [1, 7], [3, 14]], "acquired": [[1, 7]], "acquires": [[0, 7]], "acquiring": [[3, 7]], "acquisition": [[3, 7]], "acquit": [[0, 7]], "across": [[0, 21], [1, 21], [2, 14], [3, 14]], "act": [[0, 7], [2, 7], [3, 14]], "acted": [[1, 7]], "acting": [[2, 7]], "action": [[0, 14], [1, 7], [3, 14]], "actions": [[3, 7]], "active": [[3, 42]], "activity": [[0, 7], [3, 49]], "actor": [[0, 7]], "actors": [[0, 35], [1, 42]], "facts": [[3, 14]], "actual": [[1, 7], [2, 7], [3, 28]], "actually": [[0, 28], [1, 14], [2, 14], [3, 56]], "actuated": [[3, 7]], "actumen": [[3, 7]], "actust": [[2, 7]], "adapt": [[3, 7]], "adaptation": [[3, 14]], "adapted": [[2, 7]], "added": [[3, 7]], "addend": [[1, 7]], "adder": [[0, 35]], "addition": [[1, 7], [3, 7]], "additional": [[3, 28]], "additions": [[3, 7]], "address": [[3, 21]], "addresses": [[3, 14]], "adequate": [[0, 7]], "adjective": [[0, 7]], "adjustment": [[3, 14]], "admirably": [[0, 7], [2, 7]], "admiration": [[2, 7], [3, 14]], "admire": [[2, 7]], "admiring": [[2, 7]], "admits": [[3, 7]], "admonished": [[1, 7]], "admonishing": [[1, 7]], "ado": [[0, 7], [2, 7]], "adopted": [[0, 7], [3, 14]], "adopting": [[3, 7]], "adopts": [[3, 14]], "adoration": [[1, 14]], "adorations": [[1, 7]], "adore": [[0, 7], [1, 7], [2, 21]], "adored": [[2, 7]], "adorer": [[1, 7]], "adores": [[1, 7]], "adoring": [[1, 7]], "adorn": [[0, 7]], "adornment": [[1, 7]], "adown": [[2, 14]], "adult": [[3, 7]], "adulterated": [[1, 7]], "adultery": [[0, 7]], "advance": [[0, 21], [1, 14], [2, 7]], "advanced": [[1, 7], [2, 7]], "advances": [[1, 7]], "advantage": [[0, 14], [1, 7], [2, 7]], "advent": [[0, 7]], "adventures": [[1, 14]], "adventurers": [[0, 7]], "adversary": [[1, 7], [2, 7]], "adverse": [[1, 21]], "advice": [[3, 14]], "advise": [[0, 63]], "advocacy": [[3, 7]], "advocate": [[1, 7], [2, 28], [3, 14]], "advocates": [[0, 7], [1, 7], [3, 7]], "afar": [[1, 14]], "affair": [[1, 7]], "affairs": [[3, 14]], "affect": [[0, 7], [1, 7]], "affectation": [[2, 7], [3, 7]], "affirmative": [[1, 7]], "afflicted": [[0, 14], [2, 14]], "affliction": [[0, 42], [1, 14], [2, 42]], "afflictions": [[1, 21], [2, 14]], "afford": [[3, 14]], "afraid": [[1, 28], [2, 7]], "afresh": [[1, 7]], "afric": [[2, 7]], "after": [[0, 154], [1, 91], [2, 24], [3, 91]], "afterglows": [[1, 7]], "afternoon": [[0, 14], [1, 42], [2, 21], [3, 7]], "afterward": [[2, 7]], "afterwards": [[2, 21]], "again": [[0, 217], [1, 343], [2, 392], [3, 168]], "against": [[0, 154], [1, 84], [2, 140], [3, 133]], "agape": [[2, 7]], "age": [[0, 14], [1, 7], [2, 7], [3, 42]], "aged": [[3, 7]], "agency": [[2, 7]], "agent": [[3, 7]], "ages": [[0, 14], [3, 14]], "ageist": [[0, 7]], "agitation": [[2, 7]], "agitators": [[3, 7]], "agnosticism": [[3, 7]], "ago": [[0, 14], [1, 7], [2, 7], [3, 14]], "agree": [[3, 63]], "agreeable": [[0, 14]], "agreed": [[3, 14]], "agreement": [[3, 126]], "aha": [[2, 7]], "aim": [[0, 28], [1, 7], [3, 28]], "aimlessness": [[3, 7]], "aims": [[0, 14], [3, 21]], "air": [[0, 35], [1, 63], [2, 175], [3, 7]], "aired": [[1, 7]], "alarm": [[1, 14], [2, 21]], "alarmed": [[0, 7]], "alarm": [[1, 7]], "alarming": [[2, 7]], "alas": [[0, 91], [1, 119], [2, 112]], "alien": [[2, 7]], "alight": [[0, 7], [1, 14]], "alighted": [[2, 7]], "alike": [[0, 14], [1, 28], [2, 28], [3, 14]], "alive": [[0, 7], [1, 14], [2, 7], [3, 7]], "all": [[0, 1505], [1, 2198], [2, 1694], [3, 1001]], "allayed": [[1, 7]], "allegorical": [[3, 7]], "allegories": [[3, 7]], "allievates": [[3, 7]], "allevation": [[0, 7]], "allotted": [[0, 7]], "allow": [[0, 7], [1, 14], [2, 14], [3, 21]], "allowed": [[0, 7], [1, 14], [3, 7]], "allowing": [[3, 7]], "allure": [[0, 21], [2, 14]], "allured": [[1, 7], [2, 14]], "allureth": [[2, 7]], "alluringly": [[3, 7]], "allusion": [[2, 7], [3, 14]], "almight": [[3, 7]], "almost": [[0, 49], [1, 21], [2, 49], [3, 7]], "alms": [[0, 21], [2, 7]], "almshouse": [[1, 7]], "alms": [[0, 35], [1, 49], [2, 77], [3, 7]], "alone": [[0, 112], [1, 119], [2, 119], [3, 168]], "aloneness": [[0, 14]], "along": [[0, 7], [1, 7], [2, 21], [3, 7]], "alongside": [[0, 28], [3, 7]], "aloud": [[1, 7], [2, 14], [3, 7]], "alpha": [[2, 14]], "alredy": [[0, 84], [1, 126], [2, 259], [3, 84]], "also": [[0, 630], [1, 875], [2, 499], [3, 168]], "altars": [[1, 7]], "alter": [[3, 7]], "alteration": [[0, 7], [2, 7]], "altered": [[0, 7]]

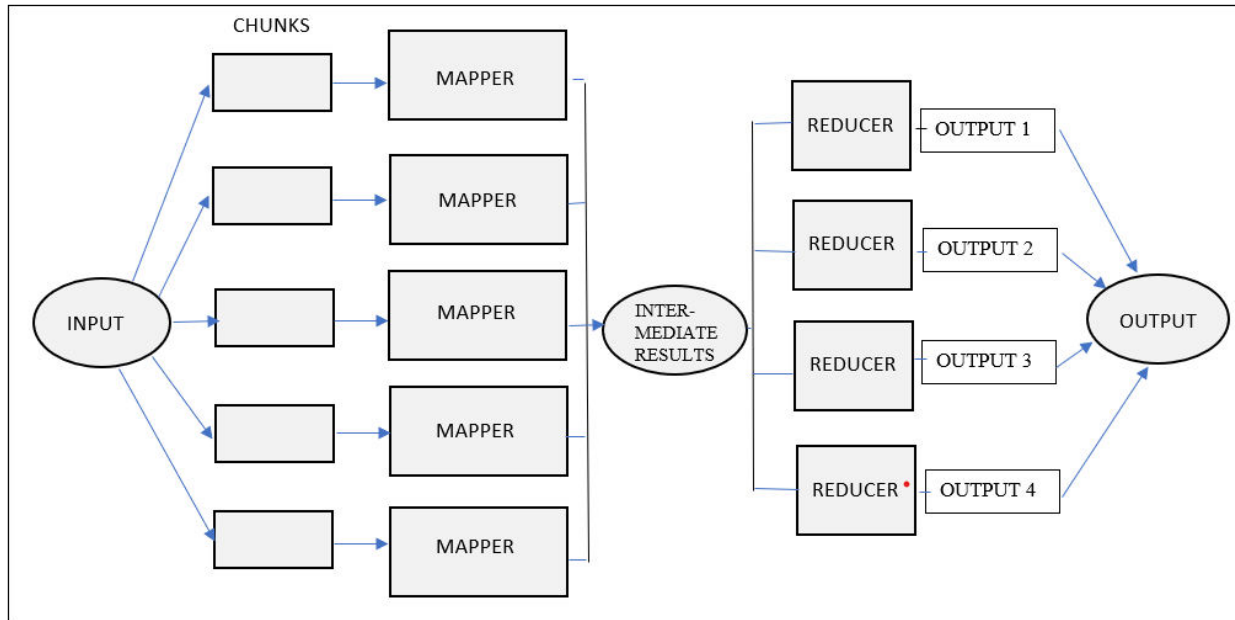
### Output for Inverted Index

### FLOW CHART:

This flow chart describes the working of the Map Reduce model that has been implemented.

SUBMITTED BY:  
PRASHASTI KARLEKAR  
prkarl@iu.edu

# Assignment 2: Implementing MapReduce



## BONUS COMPONENTS:

- Protocol Buffers are implemented as a bonus component, since all the operations adhere to the same protocol with which intermediate data is passed to reducers. This leads to an efficient implementation of reducers.
- Fault Tolerance has been implemented as a bonus component. If the master fails, the client can resend the request again. If the map function or the reducer function fails, they get restarted again. In case of reducer function, it will simply read the stored intermediate data.

## TESTS:

To start the server, run :

```
./start_server.sh
```

To run the test files, run the below script as follows:

```
./test.sh
```

There are three test cases in these files:

1. `client_invertedIndex`: This file passes input to run inverted index application on the input file.

```
-----TEST CASE TO RUN INVERTED INDEX APPLICATION-----
```

```
Response from the server
Server -- Task completed.
```

SUBMITTED BY:  
PRASHASTI KARLEKAR  
prkarl@iu.edu

# Assignment 2: Implementing MapReduce

2. client\_word\_count: This file passes input to run word count application on the input file.

```
-----TEST CASE TO RUN WORD COUNT APPLICATION-----  
  
Server response:  
Server -- Task completed.
```

3. clients\_multiple: This file tests if the server is able to handle multiple clients parallelly.

```
-----TEST CASE TO CONNECT MULTIPLE CLIENTS TO THE SERVER-----  
  
Connecting to Client 1  
Connecting to Client 2  
Response from Client 1:  
Server -- Task completed.  
Response from Client 2:  
Server -- Task completed.
```

4. Below is the result of running the ./start\_server.sh and ./test.sh command:

```
prashasti@DESKTOP-CD2K870:/mnt/c/STUDY/EDS/OldMapReduce4.0$ ./test.sh  
-----TEST CASE TO RUN WORD COUNT APPLICATION-----  
  
Server response:  
Server -- Task completed.  
-----TEST CASE TO RUN INVERTED INDEX APPLICATION-----  
  
Response from the server  
Server -- Task completed.  
-----TEST CASE TO CONNECT MULTIPLE CLIENTS TO THE SERVER-----  
  
Connecting to Client 1  
Connecting to Client 2  
Response from Client 1:  
Server -- Task completed.  
-----Execution of Client 1 request is completed.-----  
Response from Client 2:  
Server -- Task completed.  
-----Execution of Client 2 request is completed.-----  
prashasti@DESKTOP-CD2K870:/mnt/c/STUDY/EDS/OldMapReduce4.0$
```

Output for ./test.sh

SUBMITTED BY:  
PRASHASTI KARLEKAR  
prkarl@iu.edu

# Assignment 2: Implementing MapReduce

```
prashasti@DESKTOP-CD2K870:/mnt/c/STUDY/EDS/OldMapReduce4.0$ ./start_server.sh
Waiting for a Connection..
Connected to: 127.0.0.1:53518
Thread Number: 1
SERVER -- PROCESSING THREAD
SERVER ----Task completed
Connected to: 127.0.0.1:53520
Thread Number: 2
SERVER -- PROCESSING THREAD
SERVER ----Task completed
Connected to: 127.0.0.1:53522
Thread Number: 3
SERVER -- PROCESSING THREAD
Connected to: 127.0.0.1:53524
Thread Number: 4
SERVER -- PROCESSING THREAD
SERVER ----Task completed
SERVER ----Task completed
```

Output for ./start\_server.sh

## LIMITATIONS:

- The application does not any other applications of map reduce and only includes Word Count and Inverted Index.
- The application uses local file system as a datastore and does not work on a distributed file system.
- The number of clients connected to the server at the same time is restricted by the hardware settings on which the server will be executed.

## FUTURE SCOPE:

- Some nodes are slow in terms of execution. A correct model with load balancing implementation can be applied.
- A distributed datastore can be implemented within the application which will allow it to run in a more scalable environment and hence improving the performance.

SUBMITTED BY:  
PRASHASTI KARLEKAR  
prkarl@iu.edu