**Saloni Shah**
**NUID: 001617242**

# CS6240 : Assignment 3

## Design Discussion

**Pre-Processing Job:**

The pre-processing follows the steps mentioned (in assignment). The pre-processing occurs in the following manner:

1. Input file is read line by line, and each line is sent to Bz2WikiParser.
2. The parser filters out relevant links, and strips URLs to only the page name. It discards any duplicate pages occurring inside a page's linked pages list, eliminates self referenced links, and pages containing ~ as well as links containing these characters.
3. It then creates an adjacency list representation for the page name.
4. It emits the pagename and its adjacencyList
5. Then for each page in the adjacencyList, it emits that page with empty adjacency list (this is to handle the condition where a page is present in adjacency list, but does not exist in the database, and hence has to be considered as dangling node)
6. Reducer then collects all the adjacency list, combine them, and forms a node representation : (pagename, pagerank, adjacencyList) (initial pagerank is set to -1).

**PageRank Job:**

The **PageRank algorithm** uses the **exact pseudo code** mentioned in the **Module 6, 3.3**.

To **handle the dangling nodes problem and for calculation of** $\delta$, I use the following approach: 'Merging computation of $\delta$ into previous Reduce phase'. For each iteration i, $\delta$ is calculated in Reduce phase, and used by Mapper in iteration (i+1) to correct the page rank obtained in ith iteration. To achieve this, a global counter is used that is updated by all Reduce calls for dangling nodes (This is the **Solution 2 mentioned in module for Computing $\delta$ (Reference: Module 6, 3.6)**). I chose this version since $\delta$ is calculated in single job, instead of adding a separate job just to calculate $\delta$. This saves us the time and memory of running an entire Map-Reduce job just to calculate $\delta$. The problem with Order inversion approach is that we need to pass value of dangling nodes to all the Reducers, so that they can calculate $\delta$ to be used in that iteration. Hence, we need to emit all dangling contribution for each node to all the Reducers, increasing unnecessary data transfer from Mappers to Reducers. The only additional step that needs to be taken in the approach I have taken, is to include an extra Map job to correct the pagerank of each node after the 10th iteration. So, the final correct pagerank values are generated in 11th iteration (where Reduce task is set to 0).

Estimation of Convergence:

The value of alpha chosen in 0.15. Alpha can't be kept very low since it will increase the probability of random jumps, while setting alpha to a higher value, allows navigation from one page to another through links only. Hence, if we keep it very low or very high, it will never converge to 1. So, it should be set such that it converges to 1.

**Top-k Job:**

The **TopK algorithm** uses **exact pseudo** code mentioned in **Module 5, 2.13** (approach using local top-K and global top-K). This approach is best suited for small K, since we scan the input just once, and store only the required top_K in the memory. While, if we use sorting approach for small K, we would need to store entire input data into memory, which is inefficient in terms of space. Since value of K is just 100 here, using local top-K and global top-K approach is better. [Note: If any 2 pages have the same pagerank, both of them are considered for top 100 selection (i.e. they are not eliminated)]

**Amount of Data Transfer:**

(All data transfer is shown in bytes)

| Iteration No. | 6 m4.large machines | | 11 m4.large machines | |
|---|---|---|---|---|
| | Mapper -> Reducer | Reducer -> S3 | Mapper -> Reducer | Reducer -> S3 |
| 1 | 1191750288 | 1128436083 | 1235298193 | 1128436348 |
| 2 | 1290786548 | 1128436646 | 1335725820 | 1128439561 |
| 3 | 1291188567 | 1128439085 | 1335970177 | 1128437803 |
| 4 | 1291347215 | 1128424731 | 1336301945 | 1128423233 |
| 5 | 1291366313 | 1128422100 | 1336307907 | 1128419171 |
| 6 | 1291176684 | 1128420073 | 1336293443 | 1128418257 |
| 7 | 1291297111 | 1128420268 | 1336361791 | 1128417323 |
| 8 | 1291231886 | 1128414957 | 1336421072 | 1128421708 |
| 9 | 1291548675 | 1128413883 | 1336506066 | 1128423248 |
| 10 | 1291342422 | 1128425099 | 1336617100 | 1128423683 |

(Here, in last (11th) iteration there is no data transfer between Mappers and Reducers, since only Map job is needed, and Mappers will directly write to S3. Hence I have not show it in the table above)

In any system configuration, in each iteration, the amount of data transferred from Mappers to Reducers and from Reducers to S3 does change slightly. This is because, in each iteration, the pagerank of each node changes, (pagename and adjacency list remain constant) and depending on this change, the number of bytes that are transferred or written changes. After each iteration, the pagerank usually becomes smaller and smaller, and hence to represent it correctly, exponential form and precision of floating points is needed, which causes few more or less bytes to be transferred from Mappers to Reducers, and from Reducers to S3.

## Performance Comparison
**Running Times:**

|  | **6 m4.large machines** | **11 m4.large machines** |
|---|---|---|
| **Pre-Processing Job** | 42 minutes 34 seconds | 21 minutes 6 seconds |
| **PageRank Job (10 iteration)** | 27 minutes 55 seconds | 16 minutes 26 seconds |
| **Top-100 Job** | 38 seconds | 27 seconds |

Each phase in Run 2 (11 m4.large machines) shows a good speedup compared to Run 1 (6 m4.large machines). This is because in Run 2, there are 5 more m4.large machines working on the same dataset given to Run 1. Due to this, the work on each machine becomes half of what it is on in Run 1. And hence, each phase in Run 2 completes the work in almost half the time taken in Run 1 (Running times above can do prove this observation). Here, though Top-100 job in Run 2 completes faster than Run 1, it does not show much of a good speedup. This is because, in Top-K job, there might be many mappers, but there is only 1 reducer working on the top 100 from each mapper. Due to this, reducer in Run 1 and Run 2 will get about 100 * number_of_mappers records, so if number of mapper machine increases in last job, reducers workload also increases. Due to this reason, even though there are more machines in Run 2, the speed up for Top-100 is not so good.

**Top-100 Pages:**
**Simple Dataset:**
United_States_09d4    0.0051890090002740434

Wikimedia_Commons_7b57   0.00480676647470988
Country   0.003940284687713574
England   0.0027524814361112155
Water   0.0026878096234471574
Animal   0.0025540875651497643
City   0.0025108240807830287
United_Kingdom_5ad7   0.002358647093612773
Germany   0.002350401697711995
Earth   0.0023247348599551684
France   0.0023236079471426027
Europe   0.002038097037168201
Wiktionary   0.0017538842142764614
English_language   0.0017496771217548222
Government   0.0017323446521037042
Computer   0.001716840484713746
India   0.0017131709183853
Money   0.0016673836980231798
Japan   0.0015516905685357793
Plant   0.0015235595093602682
Italy   0.001507433090498333
Canada   0.0014814073434532187
Spain   0.0014711236922238576
Food   0.0014246868489679767
Human   0.0014120970062699617
China   0.0013967150612732362
People   0.0013822485250560876
Australia   0.0013298542407507953
Asia   0.0012844361711364049
Capital_(city)   0.0012742684212522326
Television   0.0012649972257606518
Sun   0.0012602100811783014
Number   0.0012432362289291035
State   0.0012403756814549144
Sound   0.0012352116672222275
Science   0.0012325431753597168
Mathematics   0.0012310566392958523
Metal   0.001192304623749709
Year   0.0011770925835108761
2004   0.001173357313768757

Language    0.001150165884858011
Russia    0.0011461817792128453
Wikipedia    0.001123330280988467
Religion    0.0010985666999662946
19th_century    0.0010965391417803436
Music    0.0010874313232146736
Scotland    0.0010548007350065563
20th_century    0.0010537049832591268
Greece    0.0010492227329348632
Latin    0.0010298606131876865
London    0.00102735544285155
Greek_language    0.001004357256650529
Energy    9.990118103796386E-4
World    9.863508479979037E-4
Centuries    9.759058651368076E-4
Culture    9.452039652115251E-4
History    9.364696034256512E-4
Liquid    9.145230968002311E-4
Netherlands    9.057245076491723E-4
Planet    9.049322622392159E-4
Light    9.016763526865974E-4
Society    9.014920621454229E-4
Atom    8.900226406531608E-4
Wikimedia_Foundation_83d9    8.88440070776325E-4
Scientist    8.883836105737015E-4
Image    8.876884860222222E-4
Law    8.862908055986277E-4
Geography    8.788451614551093E-4
List_of_decades    8.785742942839124E-4
Uniform_Resource_Locator_1b4e    8.618845063634374E-4
Africa    8.605699671526503E-4
Turkey    8.448863678892099E-4
Inhabitant    8.30479488232508E-4
Capital_city    8.230488140439364E-4
Plural    8.215155955104328E-4
Electricity    8.137230016666818E-4
Poland    7.972379043155155E-4
Building    7.971238925722246E-4
Car    7.946540606240864E-4

Sweden    7.917125562342923E-4
Book    7.914884705321319E-4
Biology    7.869328964315926E-4
War    7.708172945482264E-4
Chemical_element    7.681607959198563E-4
God    7.609357218915576E-4
North_America_e7c4    7.562868644168624E-4
September_7    7.547781812642647E-4
Website    7.462973500605942E-4
Nation    7.426671526407832E-4
Politics    7.397103787590738E-4
2006    7.332900172260957E-4
Fish    7.322371112911346E-4
Species    7.308711176294948E-4
Mammal    7.216744135950795E-4
Island    7.178090203037469E-4
Portugal    7.171070596607501E-4
Gas    7.155515366540768E-4
River    7.115777513010706E-4
Switzerland    7.061075074386641E-4
World_War_II_d045    7.020304931583214E-4

**Full Dataset:**
United_States_09d4    0.002622883307725724
2006    0.0012284974115401603
United_Kingdom_5ad7    0.0012031345232478765
Biography    9.820750030583663E-4
2005    9.170453114331424E-4
England    8.802045052385164E-4
Canada    8.559019243189323E-4
Geographic_coordinate_system    7.716537557510497E-4
France    7.250155425564715E-4
2004    7.198917516046923E-4
Australia    6.804752357198294E-4
Germany    6.543395104727504E-4
2003    5.873910170218375E-4
India    5.834188603062393E-4
Japan    5.828499867966542E-4
Internet_Movie_Database_7ea7    5.335068278947029E-4

Europe    5.092684279282765E-4
Record_label    4.914575092040242E-4
2001    4.8700951198761414E-4
2002    4.8287569488536823E-4
World_War_II_d045    4.7805172711679826E-4
Population_density    4.703435073017509E-4
Music_genre    4.6719637178231063E-4
2000    4.646639470823794E-4
Italy    4.458079830035117E-4
Wiktionary    4.362093187146297E-4
Wikimedia_Commons_7b57    4.352977195224375E-4
London    4.3479475608461675E-4
English_language    4.184924190124008E-4
1999    4.0593676886523377E-4
Spain    3.6292229527105577E-4
1998    3.563095348985902E-4
Russia    3.438958027851477E-4
1997    3.3728506998715403E-4
Television    3.3629707612170177E-4
New_York_City_1428    3.3462856024990344E-4
Football_(soccer)    3.26148648392111E-4
1996    3.236267727634881E-4
Census    3.235551257749954E-4
Scotland    3.22189805812045E-4
1995    3.1015498593562127E-4
China    3.086407053476629E-4
Population    3.043214375168833E-4
Square_mile    3.04056159848861E-4
Scientific_classification    3.0401129926075406E-4
California    3.0166613242840735E-4
1994    2.9069059165481116E-4
Sweden    2.876209953787776E-4
Public_domain    2.8741664930924404E-4
Film    2.8626953981236556E-4
Record_producer    2.8411279243647825E-4
New_Zealand_2311    2.8310101842408004E-4
New_York_3da4    2.7888558279744717E-4
Netherlands    2.76671181070038E-4
Marriage    2.758133039378725E-4

1993    2.748027246452099E-4
United_States_Census_Bureau_2c85    2.7466711649185965E-4
1991    2.718970189676913E-4
1990    2.683246782500269E-4
1992    2.663656156472363E-4
Politician    2.6489459038802444E-4
Album    2.605577884155138E-4
Latin    2.6045696116246966E-4
Actor    2.583393632505134E-4
Ireland    2.5810098404018743E-4
Per_capita_income    2.5564270352658393E-4
Studio_album    2.5185786280951093E-4
Poverty_line    2.511650008893579E-4
Km²    2.4950708971558256E-4
1989    2.4688974587744404E-4
Norway    2.4086685269665328E-4
Website    2.3901474110413337E-4
1980    2.3532256907970485E-4
Animal    2.2937819007781048E-4
Area    2.292130433722194E-4
1986    2.2703360707975189E-4
Personal_name    2.2624086525437702E-4
Poland    2.261199647608192E-4
Brazil    2.256619988669503E-4
1985    2.2402853548642287E-4
1987    2.233052142740763E-4
1983    2.2175551866755638E-4
1982    2.21097659767572E-4
French_language    2.193810555473214E-4
1981    2.1934770408862716E-4
1979    2.193298954042148E-4
1984    2.1878974281640544E-4
World_War_I_9429    2.1869361511968075E-4
1988    2.185763275043908E-4
Paris    2.180114096060794E-4
1974    2.179757176312975E-4
Mexico    2.156691801773946E-4
19th_century    2.118571806277182E-4
1970    2.1132376508534002E-4

January_1    2.1086786200188968E-4
USA_f75d    2.1070856929063453E-4
1975    2.0860252359153428E-4
1976    2.084679274023311E-4
Africa    2.0779879925956986E-4
South_Africa_1287    2.0736014983858958E-4

Some of the pages in top 100 do seem reasonable, and might have important information contained inside. But some pages like Personal_name, Website, Km², Record_label, Car, Plural, etc. do not seem to be so important to be present in top 100. This generally happens because multiple pages in Wikipedia keep referencing the same article when needed. For example, article on Audi might have reference to Car, and BMW might also have reference to same link of Car, increasing pagerank of Car, even though it is not an important page.