

# CPR E 419 - Lab 3

Anjana Deva Prasad

## Experiment 1:

“Sort the data in “input-5m” by keys, using TotalOrderPartitioner.  
Use up to 10 reducers.”

Command:

```
cpre419@cpre419-VirtualBox:~/hadoop/sbin$ hadoop jar ../../Downloads/lab3.jar sorting
```

Class: sorting.java

### NO Command line Arguments

Sampling parameters:

```
InputSampler.RandomSampler(double freq,  
                             int numSamples,  
                             int maxSplitsSampled)
```

freq: 0.1

numSamples: 1000

maxSplitsSampled: 100

The following output files were generated. The size is indicative of how well Hadoop's TotalPartitioner was able to perform load balancing.

```
cpre419@cpre419-VirtualBox:~/Desktop/lab3_output/exp1/5m/output$ ls -ltr  
total 278340  
-rw-r--r-- 1 cpre419 cpre419 0 Mar 6 19:52 _SUCCESS  
-rw-r--r-- 1 cpre419 cpre419 30201165 Mar 6 19:52 part-r-00000  
-rw-r--r-- 1 cpre419 cpre419 27744750 Mar 6 19:52 part-r-00001  
-rw-r--r-- 1 cpre419 cpre419 27040515 Mar 6 19:52 part-r-00002  
-rw-r--r-- 1 cpre419 cpre419 25479684 Mar 6 19:52 part-r-00003  
-rw-r--r-- 1 cpre419 cpre419 31994100 Mar 6 19:52 part-r-00004  
-rw-r--r-- 1 cpre419 cpre419 28690608 Mar 6 19:52 part-r-00005  
-rw-r--r-- 1 cpre419 cpre419 29331402 Mar 6 19:52 part-r-00006  
-rw-r--r-- 1 cpre419 cpre419 28824216 Mar 6 19:52 part-r-00007  
-rw-r--r-- 1 cpre419 cpre419 24942972 Mar 6 19:52 part-r-00008  
-rw-r--r-- 1 cpre419 cpre419 30750588 Mar 6 19:52 part-r-00009
```

The first and last 5 lines for each part file.

I generated this by doing a cat followed by head -n5/tail -n5 for each file.

*We can see that each file has been sorted based on their ASCII files.*

*Moreover, we can also see that none of the constraints have been violated. The biggest key in  $R_i$  is still smaller than the smallest key in  $R_{i+1}$  where  $i$  can take values from 0 to 8.*

0

0001MeE7B7JPCWj Hslz3NRPzxugmQIH9s1UW1ku3jPc79T73HdtFxGJ  
0001lcFyJCuKi4l vRlpXC3AuprUAMabMLgsleSZxURRHKS0mi9VprSX  
00023bbrDdDCyzn Un47nRNdoCLyDoh9irbhbRDuxlRgxB4HctA7Wubd  
0006rvR9W2eYJL5 87fUZNdKD6C12lFoDagEcrJxZXce4sPTn4X57Qas  
0008P8NCnnLVhJh qmRAduMXpSRodW8PCiPDtCqc5oJqLCuJaix9YFLp

6SMW0BH0ld80FGK Ik3qoeSlaEBczrhIALOWJ8FcoGEM9bTpi985BEys  
6SMeXV3Dvo8LrIY 02voIVCBECACHa01SsGtNryvuY5NFKKU5UQ2b0bl  
6SMf9LKeRXJts66 XmontLhpUiontK6oxeKR83L3X2hdtFbg3VextOSM  
6SMmRrS8rece4AJ r3NrE4oSAveBVVT5NtrZmcqgiCiPKi8URMcaKq2v  
6SMn8WMqLVyemyl YLpV3N4qbHQbs5ByDmkJvu9TpoCb6n5funkch6Js

1

6SMnVeGiyb2vmIJ hdkMub5RUkzkENMZkR15AFm1uEBSAI3rq8BpoZzs  
6SMp39HEUKvBiei eTR5SPdqHZ9jvcptpbQVGk4iYmpqF53PpECrTNmS  
6SMpZK6dBBI88uh 9oqkOcs36bKf3VPAEBxFiH1dNiEolaLcOOc7O3Bd  
6SMsHm7WUOPWOZp Y90FCff7O0SHXZhL2NtHCUr08IYbMM5860UR33uf  
6SMtcP7cu6o675s 5FPjOsOzRUfd5qe3PkF5YWUSGhdemCgMRv8x1zRu

COz3tLH4KgpVArH e6JD95p1YW0eFOpeBCI3BY9O2bl5As5iJz8dSFxg  
COz3zea2AfrPKMZ ziiJD2TD22YBnWNCJf6cNdOIC3emtFxrW28Q25UY  
COz5512IKQWxO38 IFTjmuVMifiqRKKqbFsNazPQlCAMSdNrN6po1UWf  
COz6FY33MgPKr81 WtMhKLato40kno3zFniBbX00dYp9SAFU4WcWlrVO  
COz6et8kUcm56i8 SLAxFPH396IRuik8RNsTJoUrm1mlAso11N4TQ1Ev

2

COzCqeLDSxAVOcS BO0BWDrntWjpxgbyuJHTV7gHalNJXyPbNkmlsq4h  
COzD1PIZj3Ru2xv 1j0SNp0qnXSt7ZakAnpU6XGP1DxuXAvB0drB45Qj  
COzG9A27I9XnbVa a2mbpeVQziJ6lBXPgGLlZIHJKKydvXx3L15s11F  
COzHJ2oFZKeN7rj tozkalYiorLy6ayvTAssbHQJBZDoOm5Y7RXOsBE8  
COzJjVN4ydv9eCj GeGba65UbnhcKompB9xJhA0iNXCnTAjoll9APAUg

IBVPAusR12JW3UO aj28Zul5an11VmAXDUqy9AOhvy6IEzNu8W3qqqfh  
IBVVFNBgGVf96SJ 53YEIrEbCK1RQKC2Xnq98ClapyyKCfUm01gWXOMk

IBVXejuJcCV2TTk kvsp9ryxJKSuku5MBqyxJ7NtM1kz1NUMBBKkSmCu  
IBVYkYd6FxSWBMA byl6ARuYzB2Ilf0cjbnynaA2PgcYIQ6T5pWS3nEI  
IBVbKKxqLPt9KCb lqHqigAmaYM9LF0mHyCT5Sj6X8INAIVYJPB60SCN

3

IBVdOH4UeLqLUxh 4Fh4RWbkiaC2HUHBjsrYptlGdKFkvefSuuUSPT6L  
IBVeg1TmhhJxr24 ALd6uhqo1PdpE8YLCXuhM1JZ1ZTDFXeS5uValr60  
IBVqtDTHxzKE6FR oKjG2MiTYmQMpnxb99CvMnlhBoq2UZoPHQfmbbsPT  
IBVx1tZESzbNIAV cNypmUVKTOWNdgYGF0lhXRxlG5KK4pDB2s6aliJ3  
IBW0UZIqXSI9XuL vRjjO4OJ1zuBDGQUWj9drWagvSGCOpLZhd5zUMx0

NdEVSe0FUyq9Vf9 J32WVPQHfXnRnYtGmtDYM1S6o1SugzecbJMe3pjA  
NdEVUCcF7JAQGH2 5OIxDYGF2zpLv3GK8VyPuMz1N3dvHrX0noMONLEE  
NdEbd6crCm2OuCb tdRNOICzxFNmqdAv0H448nHiNUUNLcFbtzImEUD3  
NdEfBAgzONRTV5W 9Yln4BNztnAV149stIAOR61O7jAd72MOUyPiDc2i  
NdEjDn4srGL6lOi L2Ji6Q48jFL4Cb1to7OXPX61H49rrVE9gkMrkDto

4

NdEosVBxjfWcnRq X6kKb6P6q25qaZrekpX9t2dCWV2zLRjkuNxDGOdC  
NdEp8vFVFxjxMb jfai0vXGUsxirbXg3G99Sijrn3zBh0rN98pzGz3v  
NdEsunaETyk12ZR 25HT0IepscCxxQuooA68z18ZAEjk4z79uXWNouCe  
NdEyx7MioCuUqdP behe6QtE0GCF4euE5WGZJ5h8s5QFSadWPWu6Z64H  
NdEyzA0OvHBllj4 4dt60lyjbbqsSkKnVp0Qzi1dMFK6tx2GYTXmzhnV

UULop9d9FU73YQS mKkfV9oWJABHN3AvPKfQpUgQs2S8WPYjhmEHHDzF  
UULoqfLmO7hkDAk uOCX6QTrlu4k6Az9nWNh4kkUMjjfDqgrkL29yCKD  
UULpgD8fyhq46kd tubCQtVjdoTz9A8BPS3SxxnjWPWbX6G23aW1jyHK  
UULtetMKOzQUtFN N5D1i3cGsob8HcNflVSf75lcCeTO637nki4gPBh  
UULtqSEG6IW8B5Z lI7HCOO5Dzlt5ipyT1TJO7QAcZ3Rk2Aj0dFHmUum

5

UULuNX9d1O8b0qc aL0LAXdXbb9N3pT9BvVY048TaW1uVkpYYsUTn2SC  
UULvZfffsQSuJDy GhqZ4mCNIP3yAN9YflTjTj0EljxAr9cPigLb6sYJ  
UULyP4Xzil7zxIO 21iYKUZIldoXOrRQox8IVpkQkdrqeVg2DpbQyj5cr  
UULyeuRccfnqMgC ENUsFsKtMS9DiNbaSg3N2MTai6rlNdJ7NSMz1blu  
UUM2tNGXmXKaUp9 rOr10lnkbFRpAZaRZJfKKI28lliaZ4Ex5ujmpX5

ac6fafHsuetb6lA Z0nOy27Eza0ml6Rovb6gpLsFTPj7QER22jpdzne6  
ac6frbyddjpPRIC 1xnOYu5VLZJr8k83AK4BQTJOinJvt4YjiYSgBFFO  
ac6h994fmNPQDOG qCt7b6McfelKhWvoVplgzGln3udWy0n07yWxt90m  
ac6kRI8neCm2MSu AXKRf0HvkE9fOukeqUJRIJJC5uOM1LO9pmWjLQ7y

ac6qaQNV5lI8WCT ZVXFMzq2yT1Dg9nCt6t6WdFlp244bMt8Cj77ifyd

6

ac6xFQ2bLlaaOYq Lrk7yAZeNh2Rgvtuv0C3oaddvRSGgXTYBqxSFQfK  
ac6xnjeLqcnb5F9 X8UbmYcxzxH85Z7TYm1HGnLfMpSL9zGtXOQ7prFJ  
ac70g0URJ0L31Uj Sa1ttqnFy9SrRQFAK8IKWS01UFcxWeuoZAN4TI0M  
ac73V2Ik6334jbF XhjxjOfUUHIMkLJpa55hCUDAaVQCuVQIn4hHczzg  
ac74dgPdJkLyy7X WtHulPqj40Q7Zj8HVVqZoTZkzrlW5ksnc85jbUL5

gsgncMKUuL8jM85 lryj3Zz3atuYs2BcvDkq6jnJYluXEylrEeBkgbrl  
gsgqGnEZHQt7v94 NIDUOZZeZINt5TAfl0B6scTmOrgOWpNqe2rOfEM7  
gsgt3Hvc2hEAg3q 21Yu3loabqfVsf5t2dDfuqYUoFBaC8v4SXygCWuG  
gsgt5BHFALvrMc3 GRGWc6JlZbiDfPNJMLMW5nHAXqNLe8G9euCzoNd  
gsgv36uxSARG6Wp xbh9ZpCKuVcaoXfMIB7z8HQEs0vfWkyOSX6A2JOI

7

gsgyKPqabMuNrGJ Q5CSzMqbUKedAzpTGBVIHbffgd9AiKSGKfFsUCFB  
gsh4VFm6GApm2my bj8BIW6zCiZyT4kl4VE4CUuYoCWMAAeMjhmH1Guk  
gsh5WxnIBfGuD8t k0xPuGoJjWE0TFhL4Wgv8MV7Tt6UQMyAAKMNLbWq  
gsh5hDol8cXXI7s kbjqkVLmVTJYMeFFtKoSyzGHTslb1WHgEjtxKch5  
gsh6XWZsnNiCz3D ycd3md7bQJWU6iLpqzoENCxP6jXEXzk7XU2qFFC6

n5K3iBNHFInmYOv e2ZzR7CTyff6QKTvcNRJJWCryGs32iXy6ohRjNiC  
n5K4sBnfVnglvJn FiHPeYvlh2PuEyzJGzWyGZUtV42Z20ckUWlfpTJH  
n5K4ub7Xs3SHTPa MA8sKEhtlx99XYEBjr2FOeNWql0fxkGLLM5Pqul  
n5K71e7g0BK5K17 aquoVdmmMHZDzZOtlUzUAllvRUev7rTGphKKkPZd  
n5K9DRYO7Im0IOI kfiBybi4dfEEoXIX1WKGzG7x9AdTpCmKS5XH4rhx

8

n5KBvtTZKogooqD BsuBS8eSCeo21oYtqaxNILld987BA5y9hYg2nig5  
n5KEzdGmW6WgUam x3RRaV8zx9YCVrtk4lfjNH3NGxzI7Jf1t56IKWBV  
n5KGNiBvWuC3ml8 aZtYRKKcT2mcU8qQkt0ePrpqBAmeOfsH7YAAtK03  
n5KHFXxaf0Fc6xp WftxbTIJY9OOqarKtpuTWLdrxNtqnjiK0AXIk12U  
n5KIkJlGU9MuWJY seXW1kQSjSdtjxTrXtUYy8u4z3hUsVyvPj7rNyOT

sQHebmjZ9DsM8tM mBrs5Q3k0DMLgN6EK9jVhqlkFhOdMghWrD66teAH  
sQHIZItIAdPUFp WlaJm0WZXkHTu5xgSxlnVtWAiplaGDJH5LdCA5AN  
sQHs9tacSJ5lgYD 8BQPTGcEVeimohvGyxfsnvS5faiSUSXRh9Z4aVbq  
sQHxnFqKAiZ7vg8 8avHFhx7DofVSFICcmpAWrNe6vJokf0L8t2p8bYu  
sQHxzmxB4Yl8nvQ 5amjvUeRoQrrmfijJiUCv9cWCccHRuqlSQtdqWq

9

sQl07EIPvz7KT1t rt9JeuSLfBFn7YITH2CghWndEa0KBxUsy3qPzS81  
sQl59icVtJFn867 f8LnTI51ArVuJHzGovtompv9fyeli1pPhgy7InGM  
sQlJlkD9hiW9Z3N Rth8unZBWngqcn3XWSpyJgU4rleExxtNRAQX2h4e  
sQITMUCxSgjmFBT bQ63abNutv1THvs99nkSWpLBTZyIVYKVMBRF2V34  
sQlY717C4Gg1PqW UT6DZmN2Ro96h5ZID6hoZx3ehPHsqink0q3EfaJl

zzzlsNEhLA6XXr4 358YpGiful8Wvb0f917CLkUTmclMXn9UFpWFAiNm  
zzznxiU8di17hBk zMz4cONqHlhp8PMnHaLtU7VB2pDlrS5Oo5kKYbip  
zzzp2osLN7Didmk NA759V0qkWGK2Nfd8EvPc5FL9FgAIUbIX52JR4o3  
zzzuoUL6N4qMAFu JzmJyyjISIKaSWPi5ObPcZoT9hEtHg3VDYQPTlco  
zzzz0yxLW36DFE3 U32iJXn1rJyHPFOCvXEObCDsLT2pgv1ldMqBaeFI

## Experiment 2:

**“Write your own “MyPartitioner” class and use it to sort the data in input-5 00k. You MUST NOT use InputSampler class. Explain your algorithm/strategy to partition in the submission. Use up to 10 reducers.”**

## Command:

```
cpred419@cpred419-VirtualBox:~/hadoop/sbin$ hadoop jar ../../Downloads/lab3.jar customSort /lab3/input-500k /lab3/exp2/500k/output/
```

Class: customSort.java

Command line arguments: src, dest

Src path(Dataset location): /lab3/input-500k

Dest path: /lab3/exp2/500k/output

## Sampling Strategy:

I choose samples based on ideas borrowed from sampling based on tossing a coin.

This is done in the Reducer phase of the first Map-Reduce phase. Using this strategy, I collect 10% samples from the dataset which is further sorted to compute the boundary for partitioning the data.

### *Coin Toss - Random Sampling*

We randomly toss a coin for each sample (programmatically, this is equivalent to calling `Math.random()` and storing the result in a *toss* variable). Now, based on the outcome of the toss(the value of the variable *toss* in our case, this can be a double value between 0 and 1) we can decide whether the sample can be chosen or not depending if its head or tails(programmatically, this is comparing the value of the variable *toss* with a threshold(say 0.5) and discarding samples based on some criteria(for ex: value <

threshold)). Thus, as long as we've not collected enough samples, we call `Math.random()`, compare the outcome with the defined threshold and either choose or discard based on whether the defined criteria is satisfied or not. Once enough samples have been collected, we can stop looking at the remaining samples. The number of samples chosen, threshold value and criteria decide whether the performed sampling is representative of the data or not. In this case, picking 50000 samples for a dataset of size 500000 and discarding toss values that are greater than 0.5 gives good results.

### Load Balanced output:

The size of each file is indicative of good load-balancing. The input file is 29MB and each sorted part is approx 2.9M(29/10)

```
cpre419@cpre419-VirtualBox:~/Desktop/lab3_output/exp2/500k/output$ ls -ltr
total 27856
-rw-r--r-- 1 cpre419 cpre419      0 Mar  6 19:52 _SUCCESS
-rw-r--r-- 1 cpre419 cpre419 2811240 Mar  6 19:52 part-r-00000
-rw-r--r-- 1 cpre419 cpre419 2811468 Mar  6 19:52 part-r-00001
-rw-r--r-- 1 cpre419 cpre419 2871261 Mar  6 19:52 part-r-00002
-rw-r--r-- 1 cpre419 cpre419 2823438 Mar  6 19:52 part-r-00003
-rw-r--r-- 1 cpre419 cpre419 2846979 Mar  6 19:52 part-r-00004
-rw-r--r-- 1 cpre419 cpre419 2883972 Mar  6 19:52 part-r-00005
-rw-r--r-- 1 cpre419 cpre419 2907285 Mar  6 19:52 part-r-00006
-rw-r--r-- 1 cpre419 cpre419 2826858 Mar  6 19:52 part-r-00007
-rw-r--r-- 1 cpre419 cpre419 2855529 Mar  6 19:52 part-r-00008
-rw-r--r-- 1 cpre419 cpre419 2861970 Mar  6 19:52 part-r-00009
cpre419@cpre419-VirtualBox:~/Desktop/lab3_output/exp2/500k/output$
```

### High-level end to end flow:

In the first Map-Reduce phase, we compute the boundary values after sampling. The results are written to a temporary file. This temp file has 10 lines where each line gives the upper limit boundary. Once the first MR is completed, we read the temp file in the main function, and store the results in `configuration(conf.set(key,val))` which is used by the Mapper Phase in the second MR for partitioning the data.

The first and last 5 lines for each part file.

I generated this by doing a `cat` followed by `head -n5/tail -n5` for each file.

*We can see that each file has been sorted based on their ASCII files.*

*Moreover, we can also see that none of the constraints have been violated. The biggest key in  $R_i$  is still smaller than the smallest key in  $R_{i+1}$  where  $i$  can take values from 0 to 8.*

0

000ckhLmtS7AeAE 4D2NVCzvxJxrXCeX0NIXDrVx2gxNQp3knjt4rVzv  
000dJJ1TKvD9JY3 qyElVmbqdmMQi27871R6OncsDoJRAMxp0MjfEGtW  
001SZ6pXMRQSVxz vm3eElvShsy3lkQ6MaiP7UejGmq1Rsu2dQl2VdmM  
002CGxZhgFpf49N 5MVBj8BvxZ0FehM35RIG6ZdL6D0TVeps4smgONDS  
002DLQWKhfVtQPz ItqHFG2F0Het6RZUFODMRj2s9MYhXJFgj3Qu9Oqc

60DhfQr3ZvU3zoF nuaxeubxssmOFnxVuOZhXOIUO9G5JiAIMOpPiXHG  
60DvUS1J20ZngNF 38ApED5TtSg1CKIz48774zaEL9qTkvYW7iRDlgNJ  
60E8MvMFGPcmqhe NV0pjPOJSs6llsA6YlraMXyZb1rFH3n5AKTTClA3  
60EBEXcVX3JWm87 ojc0S6jxbQB2KZ1Q6fGoKIT0ldb5RCUzbXi7eg2Q  
60EHSVQL6h2aljf ZboWosup1smsJi2kNOYjf3tm2NCIGja4oHa6qRfu

1

60ETr3LZtpsbtWTz 1fryGjRcKFFbc6aSlZy4c8RUFtcExAJaEEp9MFzF  
60F050xtCQQVGLz XtyODCnXGkJIW99ldvMINSF7AXCrWG13TLdM8ckE  
60G4ZaiZp10UBdi iGiXR5MQmS1oTnS4W7G8gkb77GXGASS6Uol9bHau  
60HFnL3almjf3ZU 8mzttY5j3lqUKVc5B1pJNq9LmCREEBEcnYvg3lv3  
60HKznb9TSo2yCE rusMb061Rb4S7HHy7PziWiQHmNUoFUB1pm3pB5aV

C0eJ2H9zOH70ECZ lmbTSWrQKINLpGI1OVer5xfn5sCDF1dGEReKRocP  
C0fByQit8WjeTM0 pQhBLILgNTyHsCn2GlyaPoM5A19CDNIrfMIT2A7h  
C0fnBPJv3azCO5M BRWt98i60sbAvXObdWOrDzvS5tPdxjR5idccsjz3  
C0fnQUYMQ1CGrGd bAiHep6mJ6sVvrQr6G2s1jv6kDlmaWLIUxyyh384  
C0g0RddcP6fa7CA A0y7XTHx8vo5PIToanHqnn3VMFM1zO0MUHgzrtr

2

C0gQMXpMQjXaOgZ 3gE4mBrqSWUtNo4VU1CuKEfhJDQX8FsT1YoKoHd3  
C0hQT4VfjksZ1Ye UeNBen1oqiBWPFdf2VYGsDRPcnSAf4qhySV4UUvN  
C0hVZWIfHYVM0k2 T6q4WygXpzS3EPoCsNZlczexLI7iOoY1sAPVWtjR  
C0hgJQ9TFR89G9t W4kimFQA0DbADL0ifx5o0Yq6y5OLmUJTnx4ctHSC  
C0iomRXAskIxCcE MnZvlfBZebNMkvaJCv1ATVyl0DG61FzZBNjYk2rY

ICdA6FTN4dT6WT1 H6ex5aMzqUTBa5SHOi1xxqr0JI5BH0Ap7eC4sU3n  
ICePythBPrAaDOJ sPcWnUtQky8iYjg9Z1Al1oEK8afHOlihmU1FVbFB  
ICeevGDv21fOiQh x03QbcNUqRsNuCC8ajPF2oDfaHACgB9HLVfN9iKH  
ICgGSOZUgci8vCd 14sLafZ6xjhrCXE9V4tMCT9qEaC4ja8nuTJyc0rm  
ICgJt4n0NiiQgbQ sli6zubZ0hM9Q84NMQ1WMTQzLPrsckKN73q5lqTJ

3

ICgLTJjn6QJdxuv 0bYhVRSQSCXQ99V3NtKS2vgDiCLgayYH5NYWqiOL  
ICgmJLNxfhYJE3Z Vlyz6MY4v6xdBHWiGK9fHtaFChdh92I7u0EF99TV  
ICiQ2vxaQOsT3BR 10oKaXjHPhxPfUED1WgNyLo7qKOUSXO692uO31Qq  
ICj4LKRpJVCnOfT 6BXLNgGMkbb5IY3IVH5YqyY4i9ST41fMmgFstSnU  
ICjnciKGgV1QqJZ bhj3JMOIAFVa7X7cqPBtpvocn26fkLVaL8jOm0WZ

OFFR5VvxE4JikVq eEX0dOr7HqfMp5pYaqvEMSzZLkvEzRbvC0vWo1Jc  
OFFeYcJTlazFixj eOqyHCJXgdF2fH1nSAvaj6AcsX939u7N2d2W8OGM  
OFG7vdEObZEFWOf aAZEHT6GxHJJ9skCO9gU3eWShsMu0aessJz16nQW  
OFGx8qXGoJvqDxS H5ymzdmypbJoS1r44JkvDxtfK7U5pRACvKZ0QkFa  
OFHLsd019tPh4ZM x1AZ41B8Wsot6UaH3KURrFqHt0vxA7iYRE7kSKWm

4

OFHbrrZgdyKHBAn Xs6TNL1SdGJAK15gFlyCq91S8fdBuVOLJYEtmRk  
OFIzncsDWfiX8SE 2X8fWZOdkq0v8cePIkTheBnnfp7NI69Mnl9BNvrA  
OFJ0EUU8hD4SM4k SmmkPSGQLK1s5HI44NPtycYlladc3Xyxxt9EcY2K  
OFJHstFHhxYIUuX R9nsQKkriMn3GN7oAAc3pxODT5yGcNPacob26Y06  
OFK6E2lkYczslDs BiPuqOPAEXiTSfNyexxFzQG3vAZWmfd0PeGWCHLE

ULUaOfI0p62PQIP ALWyu1ui0HxHN5rX0Y40LQrGNYOzPhZHAMuzdmll  
ULUqJIEgbQqIYRh 0Mi3Le69VbedpjtPL655571D5A8zpDENIGkG7GQ1  
ULUuMOLYQTIcmsL OkHnbc5Kz8ovcCy5QvykiaslppryozuzyfvmxW  
ULVLiGg1yKX1kql rqMo4ODOY0Fj4sLq4zEbEntqTYZtAla3o7BVdfQS  
ULVbFnRse3JKt9E zbnKRdfYkR2RVc1gdCbSsn5tOH3WLfyjkmMmtyzD

5

ULVetiGyYLIcGRV 1fTeYCp3xBkJJmcVooLp24IQzdEdkRPR7ulhvWeD  
ULVgbZvG00Qqk6F hVEF5truOXZDMvfP13sRtWFFVFKN0PHBNXeZlydV  
ULVhDa9JUxLqVgg qeq7axUf57sAZsUeZYmGGXzgnNvm9keISLOKd4fF  
ULWGCvSc5NRV8Fa 05nY3x1r5LYQkFoBo3OAR6qjlv4A3PMM0fq0Ggal  
ULXHD5177DrnqAb ax3ngLceA6Ft51Wo42R6nrX5RnmmeQSpteoAUiol

aXVWZZDklIqmj5q x54fJ9cahsJo98ka0fE5SZ2LrG4fjcPNKp0fBuKF  
aXVmTg2x8FMFrVG tEGKpN2t3z1HnocGjSDovVUL6ZaGIQQpvEWDR3Fo  
aXW0DLnFuVxY7I8 stnru3FVYyBqUPjpE8c4zlnJ2v6fKqBYuQrAGX7s



aXWMI6JunCA0TXP UHCQRRaNLp8aE7Or71it3DROi7iXfAbzXZz5kuot  
aXXT1gJnVhTAjEj UJblaFr6FAQNEusbcGVifTREzt4RZ60Om2vLMSOJ

6

aXXeJAVbvKjKp7V 3zTZVZfsJGYeM9y9O4u6K7GtNIGNZdZI5Sij6nXv  
aXYXaleOVdhxUes DSHrTPJSagn0xPmudAuN03YZDJBZcNOHcX7zSdxQ  
aXZiaMOgyJa8bUk qlW8Z8K8RWkdMsgjJvG0Tq3afZZsuKUvCNnLhspS  
aXZiuAfqZFuf6N1 GGSD2QPM7NYPB3QWRbZQfLYpchYLaTxSiJTikFCf  
aXZljylxH02qxpX xlsGnFfsmAK8h11vMpZm29irQfR9sCQM6TC9QUFP

giYyMdOpS96drpS nmJvLyKShhNdx18qlU9DqqUFWjgUeoAr5cCpY4Tj  
giZFtNZ1urs61zg 4xoikRul5f2DYNxhDUnrR0Eld1q3NB3pbKYVUzMV  
giZixRZbRfs2ex4 EnWLBc1DfmTRexgPgFSqW3rYKgpvqEX1YU5OjQNu  
giaqjn129BWYou0 dmieuSMajt08zjXSIEJX9FNWI1GJBY4FDH5vS5ht  
giarVuCNtiEMO8y XOAbEHSZX0La5GaCx0AmHWFnae7ju5cU4hqK4FB2

7

gibbqPBXABJLclX sDZyo8Pq0kyhiK8D1042gh5NMTJYefVFvsPgCH96  
gibjHzT85yPJ8q1 QD6FyWB7zmS9EIV4Wo1qboClZuBmoiqYKy3HQbau  
gidSZpvJpGMA2TI Sgu2oaXfE4yD3JnAbVNubnDNs9lljxBbQqy3bQUE  
gidVyK7JJ710g88 Y70suuq6yDGrjDyc8fR3yGWVcjmBt3a0frJTl8y7  
gie7fXm7FG5GKs0 sjdj0aeePC17mVbrafWcC4ed95a2ji2ZrNbj8cS2

mmDKfUFnKEOjLrl MrVdSZhglCCWrYQeVMko8YfFDCI2js4hxmKMnRfj  
mmDotpkLKyCcoRo FTIMSmGsMhASYR7FJMWNABDGKqyzzmWESiT6oWPK  
mmDx9T5A94dxF5O jnPvfQ9ltJVnxsOv7xjAULReVVNf44E275gVragM  
mmECyaCyE4Tsl7B 224Vlo21GzB89MdWyhVHsKEMze4ASgAnKlv0W46V  
mmEGqLOPW9oDCfq GPhj59XVvY2GiYQ6c30DEMeJGjQXAj70mx2Biesd

8

mmEfndnfJMCQyNP dCGeAgPVy5Ppenq0aZx3UfGk4YeRcnElklbbin5b  
mmEmFtKnmNat3h9 MDEBPoDjGz3voLdZYQ5V4qV4dH3lvMXQ7HpNr45r  
mmEn1COZ2913vOV 3TZ7q4Jad6JSsLcNECqUq2PxXMc1aSLu52W7gsZW  
mmEvYacVsR4HEJ7 RnZ4nbYigb1y96zlkvlHsl3q9PbDtDQz125Vmh2O  
mmG2CKpmoPNm014 U1S5oH5XttzpETrnT9l3FUGLDI7yoQAVVIO90D48

stOXO42519OK3Jz Mm2v2atrEvtHBIEEfMdqyFO9oDPrHBONF9gOJq7A

stPNrJpLe3hJ5l2 vPNYeOIUJ0h7iiUJuvCulEPqpGpQMj6zXTYp8KAK  
stPP5aEEvb5uoWJ ZOJ1EsoZVVJSz4HPgcR0sFLuOlqOioS2MkNzVxxN  
stQECjK9WQV1FcL NFyCC8LEf1YuCfnzlpPMdpYjPhCyCoFGH8RcRrmV  
stQPNHdWZk0QLGV Xg2LgyOJxxdARfVGEuSxGfSISuZxbPSXrEKb9ma4

9

stQmlsgneskgg25 W5NOyl9ErMYOYs95LT1b73lJbv44QImocYUDo2Ss  
stR20E6vgRnHgHa reafevjUtgFm0TGxs5PQgZfclzrlm9MMJ0FyHJkJ  
stR786q5xeVQFcx cgSSOFVxC4pc2Ju58KzleGxP1IHBdlFJuhXt63ok  
stR9ifOBkTRjJKy 5r91M3nsdd2G7V4iRIqE6gbWpsIIKYXXYtbDAAti  
stRCblqk3yzoWnB jaeaeAnmY84FHHS6rk04N2lxQlhEBKK3c6gbyMVj

zzxZrATXhf7FU0B sulCYFi1sUFqLbENYAx8aGzZd4A3xtcPkGm3nJvW  
zzxe4LDV1mQeUWp 2q5PF7MzFMCVksSz1EGkWYngpWtOBLkUH4TeXZzDC  
zzxuUEUGmDrflNI DZEWS7p6RI0Fq72rVbP1HlXr6h7mhEbNJzi2nU8Q  
zzyzcFdkM8StXh1 Vg2N3sG7kNcvNMu693xGZTV5BMNPI5vLUTCyyN8P  
zzzLyRRmSPdcdgF ql7TSuKJy216BRbO7R7n5M4Mn3FC6TlciqttycTp

### Experiment 3:

**“Question: We have a list of datasets (see below), sorted by their size. Try TotalOrderPartitioner and your own partitioner (MyPartitioner) on each dataset. Do you get different results? What makes it different? Which is the largest dataset your solution can sort? You can get the bonus if you can sort 5m dataset by using only up to 10 reducers and your own partitioner.”**

I tried my own partitioner on input-5k, input-500k and input-5m. It seems to work really well for all cases. In fact, I'm able to sort input-5m pretty well both in terms of speed and load balancing efficiency.

As far as the results are concerned, the overall data is still sorted in the same way if all the parts are combined and the entire data is represented as one huge file.

However, the data that is stored in each part is different. This happens since we are using a different sampling strategy which affects the boundary values that are chosen. The parameters for Hadoop's TotalPartitioner and the criteria which we use to sample the data decides the bucket size and boundary for comparison which isn't the same.

```

cpre419@cpre419-VirtualBox:~/Desktop/lab3_output/exp2/5m/output$ ls -ltr
total 278340
-rw-r--r-- 1 cpre419 cpre419 0 Mar 6 19:52 _SUCCESS
-rw-r--r-- 1 cpre419 cpre419 28341426 Mar 6 19:52 part-r-00000
-rw-r--r-- 1 cpre419 cpre419 28790871 Mar 6 19:52 part-r-00001
-rw-r--r-- 1 cpre419 cpre419 28466256 Mar 6 19:52 part-r-00002
-rw-r--r-- 1 cpre419 cpre419 28372776 Mar 6 19:52 part-r-00003
-rw-r--r-- 1 cpre419 cpre419 28013619 Mar 6 19:52 part-r-00004
-rw-r--r-- 1 cpre419 cpre419 28835673 Mar 6 19:52 part-r-00005
-rw-r--r-- 1 cpre419 cpre419 29221449 Mar 6 19:52 part-r-00006
-rw-r--r-- 1 cpre419 cpre419 28581852 Mar 6 19:52 part-r-00007
-rw-r--r-- 1 cpre419 cpre419 27946701 Mar 6 19:52 part-r-00008
-rw-r--r-- 1 cpre419 cpre419 28429377 Mar 6 19:52 part-r-00009

```

As shown in the above figure, each part has approximately 1/10th of the data. This was done by considering 50000 samples for input sampling for 5m samples. This is just 1% not even 10% of the data and yet it does really well!!

#### First and Last 5 lines for custom parishioner on input-5m

```

cpre419@cpre419-VirtualBox:~/Desktop/lab3_output/exp2/5m/output$ ./print.sh
0

```

```

0001MeE7B7JPCWj Hslz3NRPzxugmQIH9s1UW1ku3jPc79T73HdtFxGJ
0001lcFyJCuKI4l vRlpXC3AuprUAMabMLgsleSZxURRHKS0mi9VprSX
00023bbrDdDCyzn Un47nRNdoCLyDoh9irbhbRDuxlRgxB4HctA7Wubd
0006rvR9W2eYJL5 87fUZNdKD6C12lFoDagEcrJxZXce4sPTn4X57Qas
0008P8NCnnLVhJh qmRAduMXpSRodW8PCiPDtCqc5oJqLCuJaix9YFLp

```

```

63pVFg6M0SZeMj0 LdFrJXtkYbeBKjp2l5pz3AjSj3EGPJOtCtsLFnSV
63pVRUZmuzsniSF 7PKv4s3eSuDu7WJQsNQoeTio7EhHNWJNc9EXe4Mh
63pWJBOVN3Mrthi WN8yxHbbgiZyKKnWz2BIPeneOAoltB0TPiWhtlC6
63pZASnOh6MR3QV AU6x3FhkVRXT1pcrWOgeiVak32FLNI7iMlfPWCHn
63pc4fyYdKekJmF voLXIsKVHHeDaacNLm4ICZmpY349WmtOB4bZMfL7

```

1

```

63pmPxM5Jpg1mfL cLSVLIEu76MGUNIoP8msy47BPxslTalR8rVEIQhB
63pplU5S15MyhAC tdIX3DzgExZ6XnHOAf4kWHD9t242stO21L66B8N3
63pre8CmnmduZ3g GI0hgNyZOtRdpLTNIXIM5ESABQaG2uQRdjg5LM3V
63q0rsbVCPNPZSG XEcnO1k75Jfsj4m7AIVtx9PdYsRpt0LdFfU6mEqp
63q4rjQE0TQc0Ui FnEQ68jSPMYn8ST8lOtFzR0oObxpQ3J8MWqJ7vRA

```

```

CELZpcmTzDaCjmE 9QPXqUyYylljSOKUkLkVkRUyHjetKhVyAEcpn1Hz
CELaTPRpLQf0lgN ZmhuS6iAWkBSMWSceL2P35JEU8Jt9MQI3jEOZ9fc

```

CELezpbKWrWzpNo 44eSuKRoboJ1EtZCzfRzC7cfrkla0Vxl5kgzZ6m  
CELhRcKmbdao9jz EyjNRgy2olcNk5xZsykGKIKCZMvAoyK7yhvpGPs3  
CELmsZ27c7Bp0DE ZIEHn3LXNHRKp6zVnulcBcErYmftzGMvTCS4Tqdi

2

CELoLDsen7WEyEs FIHE1C7piz3VlnMqqH4KvntfQo0iUfoCWtm53lpq  
CELSMRImMz099zj aGgc1msLzrRea97gU9DselHNslHnu9JcOR7UE66P  
CELSYReSYCSNJLM RNjJm2Oby0Ou6CxJCjocz5GBBb9TvEqcNOLEOXOm  
CELvIDJ7KhiutyU FZriWdO31xhJyNq9qZ2TBMWQppPoxU9TmXKtRa6L  
CELyT1RKTlc3250 zkrPsNGmvXnqjl6lQXYAU8XnebT8zKhMjpYEXEZt

IJREmthA8A53iC Tc687EKj0LHI9mCaPpeKF14K6DXkgCbWxtFxn8fA  
IJRNBYIfEgjObaL EXgV8q5QDumF4CyiPq2E3GmfCTkKCqWzP65mLLeo  
IJRO3nGL6UcZAfk Ea5cihrpQ5oxe6y4sWroaSJ8LnpXtusLOOXE3dls  
IJRP33UXNvNY7W0 PFKgr1HiNcFVF0D04vD4bCkjX69MXI9oVfV6R0xQ  
IJRRmW2MGyc1OR MIW7xIBQ5mzBdt3XMXqFUkju1cPslK1gHgK9F2b

3

IJRSY5VI1krIk9e fH77oSDsqHPs3kTyETQ1Do0mbLsBGIVuSCol7YkF  
IJRSYxaAxTjO9cc lqjxHIPYyyuP5SHlmtPBP7IB5vDTME5Eg31NZeUf  
IJRURAdB1Vt2NEg pSXcUtYgEOu2tHxjMPSDgjc3qlU7T4zE57o4ujh  
IJRXXCk9LvXmgeN xXseRZCzNXvCysdqaE3xkuWakPdUpeVmaMyj9iyg  
IJRXxC4ZxJDXLju UO9xBp2oeFYXxmvF6pbjsINGNkNCCCExpafad7nt

OO8AFRj96OVcf85 osfLqoQMHJ152ujjTA3ncoa63qqLlr4t1GMR0zDO  
OO8F0y9YoU3eikO oTI8T26ZXICbcuUFm1CqIL7S9d3BiahOyuPSqr3q  
OO8IVJol8MZOJnQ szKKYq4oGxGZ5pjdHEq4eGZm2HjLdPzoHX0R6ejv  
OO8JxeZ9XJ6GBtJ xosLZITcpD1kbiTuaoQ05hSYD40h6JiFqqOnAC6FE  
OO8Ka7zJ7mj1fmx e3cVzD7b5zcEs6Y7mCo4TuEXQrOOh9pXLdhY9UqX

4

OO8KaHV3cD4vYWG 5xYKNfqIOfz225KitplzkHNDIFDN3J6zmIE9Yxs  
OO8KqZ7hl2xW07N HGaBnQWLKtKjhviGckYCffEZmaoN7IfBRZ98xx3b  
OO8U8eRxVsHv87F qHr8M0RRfvVHhHj6RguRuTKi44A9NPL7LcOqeQUN  
OO8WP3h9nN7Sguy sqnaeVqy0eBNkEvsA1P3lOtm9rTSMUXt5EUvdR6S  
OO8Y3n1SeylvEzy eQq3ahBXnRT9fpaxsQP8qZiBIIAyTHsyqJufGpcs

UOAPCyfauH7v76T QobrsoxbsyJ0qdpsCavxqAsnCWqtKayG7eiDisjt  
UOAUE9OtCIhoKBI WNp4fRiXD62glZfioXDmAVe1jmd4voiUmCzOvhFM  
UOAYZLt19Pz6kCa bQ7kndfLVcZTOdCoSWVKzvr1eBqCdNI6Bcs95oKO  
UOAemWZdaKMPDE2 guPiut8li42RcBkBbqetaDWE6AtDLK5BJerFLMCI  
UOAetLJelf2sZ1R SChHCThdZU2yKjJDJYKGGIM1WOpL45OdVMexFJDE

5

UOAfPum4BvVF4Nr gDWnJxKNkXUgRIET92aWRxxiNufmUdYVSRa2qA2c  
UOAztT2c9bdRD5V WfTRy3VLFm6QcEhvUuVlvDkUni1t2iaCttW4sOFT  
UOAz4HYbkmmMIGF q7Qvl3itZ3iZ9jdH2VDuKJODA56ZBrEn0EoHJOVj  
UOAz5vXTzC74i7T 9Fj46Yz3D5bkfg7WEsCgNFBjzIX7eYsjCJ72V0vH  
UOAzl2EP9Bsl4yC iJLCEO2PMMAnFRn5hSnEBjWcIP0O9dFLR4LCkyBu

aXokdH8jFMqU3OC SrHZybKC9Qszeg4KahMdDJUHveOP1TLQQO09M1e9  
aXot12tKF7gUn5b 2ZCbXqEftOIGgDhxDIYzZF2v9MU8hnWuh5H1a8yp  
aXoucZEJQX63IIX Ws9Rc58ruU6Uu2WevPIJ0Cr8U3xoCTYpamhlgBUM  
aXoy35dUNIceFnT iLT5KkO3US77IUVURdhnXMvuKVe2Ed5hVTzXqf2d  
aXp0v0uXE0UEOFN 6YSLVUGxErRQORjVTHubNgWeRuiarSjvCMq5zgQj

6

aXp4nfgtyocy46p bmXC5pbN1GQclDiuEGXbi22LJxK8bInjryVjfdS  
aXp5mMFT84sYT0u T47SeJAJ9z0N8Ru1eibhvE89NYA9n9O403tc380g  
aXp6ctT2PdBZ7dS J3dQcTDrlx4MAP1IPdqMyuDLReFyJAMOUg33U0p6  
aXp9S3LNPAOObkR fBAQJByr1obFd55zfmBKJqZ2eBk8ta2bicvCF2WX  
aXpAoJIIMoaRIzx Yfm5UXR1uafTbJNY5HG3nkEu2HzFM24slfmUriNx

gn2llklYeUnYFkq WcdccWMf3yVjJjCslgFkWhlnL9GiDbymVbTRVXDj  
gn2mB7gCRVm38ur gQ0nPEfL4UQgDGe808Rg9Yz015UsDDRvZrK2RtVt  
gn2qLaaqAhS29Ph OzfTRh1lyAsxVifqq2liG7gVvx31id8caac6IASa  
gn2rITilfhmvl70 x3PgqFPPT0mxCF3CPY9o6oQWDLJxlZt8y9gtqdd  
gn2s5Y83Qv2L0Jc xiFa1A8idQNIyPABd1hl3PSrJhBqxZSPX6LUcSTb

7

gn2sTsH2ldbjp9O vZkOT86NoCUdEYag1dILKd623N8xtyfMfY56hIDk  
gn2ukS0WgCZIM2a iWahH7ngYRd61d3JB9nrXtk5ggIWDncNJK1PG51L  
gn2uyrnalpYzzLq EXc1JKAG85itombTWBzXmYjKJbGKFchj0dnve5EF  
gn308UaZKb7kHKv 0XNeukfLLe8bBAQZMevfALbLxjLE3EOGHXuferm

gn32OAKiBoq1LHJ zhyy4Ggu7paoa1SKmmexkGhmWHITHBqKkTYa11QU

mvN2euixkl9LjQO sHzbYqOGKr7OajmWPB8lGaWiOk5fl2DOIANLmBi7  
mvN3ftmal59bcVG eYg6jZ40jFpZJ4VyfZ97DyODgeR0zmH00DhhhBDs  
mvN3iMtZ4crsoUX xr3m4VyPy2zChEPcsDSj380xORxN7a99P0FSMmO7  
mvN8Ax2uHD9qtEe SseT3rS7H5QqH7qb5emonJktGCCep0enbbxoOWBy  
mvN9mt6UAVYWHAK NOMaJ6loA5mMFKVK9gDXnsyyiyk4XnbixkJxbA6m

8

mvNEnEghCUEADO6 T9t32olNhlqOjAhcHAEe5WMF2JZaZ4L3NBeyF37C  
mvNF9rnhIMt9buq zeolnu0BmYjtkZJdESqkEtUZpHLqjvNgYH1GQlTe  
mvNFTGNyU4d3FqW lGbRKW7UndFWHj5SLzcLIXQ6dU3TCRbsaHnX0ZN8  
mvNGZ7BCdAK9QgV 4l0fFqrJGBlXkGF8iOzBR3fgp9PbGfMqGfpWhz0f  
mvNHlm0BUy4Snrh moJgQREmJxvqqL4aC1stlqdWPOCMEJKQ31coQWLM

suWuEPARi3TU6xP RfcHZucB64laLWVrhbybmsFV9NEQAL38E7sP7o4  
suWv3KDKl64TsTA pHfavzpNgN0U1y0FY176H86qEE0223nKSc50LbB3  
suWytESX1eRjYae cQW3Gr4bck4KXzXOFb4HRvlj0Fp3xd5zFzTIFx4l  
suWzDJBpdMJv6vy bvJUc68nloiDcCL625gfOBzsdoPTKXAm2XpY2KQN  
suX0eFqhBl8sOOr sr68Vaif2glP1T3oxzbOtLDFN6OI0ml5jkTM2kie

9

suX7qL1spb4XEaH q6q8DMmMiQjBnrf61Vp1TaIVRQIJTPxVQEGMWZfE  
suX9crKYPxl2yBy Ukcr8IBZ7leNjomPxJphhr7ifdk8RkdCrgMBpb8B  
suXCMfkiFGgUU34 4MPtDM53AAydi8H1CXizG8pJXUKUNPv3E4O9USRQ  
suXCjm9dxZXiAl0 HKphEXpWO3fKZ4bFV6iXNR5P9iur7PFEDdbHRz0V  
suXELCKfeL4bSdq FPvn6vUUAxigxcequYQNxZEzyaV9aQDx2dvvYQDM

zzzlsNEhLA6XXr4 358YpGiful8Wvb0f917CLkUTmclMXn9UFpWFAiNm  
zzznxiU8di17hBk zMz4cONqHlhp8PMnHaLtU7VB2pDlrS5Oo5kKYbip  
zzzp2osLN7Didmk NA759V0qkWgK2Nfd8EvPc5FL9FgAlUbIX52JR4o3  
zzzuOUL6N4qMAFu JzmJyyjISIKaSWPi5ObPcZoT9hEtHg3VDYQPTlco  
zzzz0yxLW36DFE3 U32iJXn1rJyHPFOCVXEObCDsLT2pgv1ldMqBaeFI