

## Chan-Ching Hsu

The goal is to write a program using MapReduce that can sort a large data set quickly. Sorting is a common task on large data sets, and the input is a standard data set.

The name of the data file ends with the number of records in the input, for example “gensort-out-500K” has 500 thousand input records. Since each record is 100 bytes, this file takes about 50MB.

[illegible]

The first Mapper reads in an input data with *TextInputFormat* and outputs with *SequenceFileOutputFormat* with compression. There is no reducer actually implemented for the first Mapper since we just want to convert the input into another format (*TextInputFormat*). In the second MapReduce job, we essentially apply *TotalOrderPartitioner* to partition with *RandomSampler* to attempt to evenly distribute keys among Reducers; *RandomSampler* uniformly samples keys for partitioning keys into desired groups among which keys are in order. For example, the keys assigned to the i-th group will be less than the groups after i.

The run time on of the program running on the largest dataset was roughly 295 seconds on Cystorm.

output file list

```
[cchsu@n0 lab_4]$ hdfs dfs -ls /scr/cchsu/lab4/output
Found 9 items
-rw-rw-r-- 3 cchsu 419x 0 2017-02-09 17:48 /scr/cchsu/lab4/output/_SUCCESS
-rw-rw-r-- 3 cchsu 419x 561305900 2017-02-09 17:48 /scr/cchsu/lab4/output/part-r-00000
-rw-rw-r-- 3 cchsu 419x 571102200 2017-02-09 17:48 /scr/cchsu/lab4/output/part-r-00001
-rw-rw-r-- 3 cchsu 419x 610220700 2017-02-09 17:48 /scr/cchsu/lab4/output/part-r-00002
-rw-rw-r-- 3 cchsu 419x 611212900 2017-02-09 17:48 /scr/cchsu/lab4/output/part-r-00003
-rw-rw-r-- 3 cchsu 419x 613859800 2017-02-09 17:48 /scr/cchsu/lab4/output/part-r-00004
-rw-rw-r-- 3 cchsu 419x 564000200 2017-02-09 17:48 /scr/cchsu/lab4/output/part-r-00005
-rw-rw-r-- 3 cchsu 419x 596407000 2017-02-09 17:48 /scr/cchsu/lab4/output/part-r-00006
-rw-rw-r-- 3 cchsu 419x 871891300 2017-02-09 17:48 /scr/cchsu/lab4/output/part-r-00007
```

part-r-00000 and part-r-00001

```
[cchsu@n0 lab_4]$ hdfs dfs -cat /scr/cchsu/lab4/output/part-r-00000 | head -n 10
^O!Uve 00000000000000000000000000001228D4 77778880000222444DDDDDDDEEEEO0000000CCCC7777DDDD
^3CO, 5555AAAA99999EEEE888822229999CCCCDD6666555544442222
0000000000000000000000000000158C5C5 8888BBBBDDDD1111CCCCSS556666BBB1111EEEEDDDD22229999
!&S3[/] 00000000000000000000000000002145D78 33332222FFFFBBB0000FFFAAAA666655553333DDDD3333CCCC
!,=U$,9 0000000000000000000000000000019072E3 8888AAAA11114444FFFF77773333EEEE44440000FFFF99999999
#!'cl'~ 0000000000000000000000000000EDC5C8 2222BBBB2222FFFFFFFFFFFCCC55556666666777700003333
&!epOj 000000000000000000000000000001CBBA42E FFFFFFFF66669999FFFF44446666222233330000BBBB33333333
^AJ,j;L; 000000000000000000000000000001FBDD93E 8888AAAA7777BBB7777AAAA6666DDDDDEEEE9999BBB BBBB
(Vc$SpZ 00000000000000000000000000000A9F516 5555DDDD111199998888EEEEBBB55555551111BBBBCCCC0000
*>MS1.E 00000000000000000000000000000281643F 5555EEEE888899994444FFFF1111CCCCEEEE1111EEEE6666FFFF
,K4a--v 0000000000000000000000000000001B8132
cat: Unable to write to output stream.

[cchsu@n0 lab_4]$ hdfs dfs -cat /scr/cchsu/lab4/output/part-r-00000 | tail -n 10
*q6/ ($zP 0000000000000000000000000000779761 555555555555DDDD777755599999AAAA000099995555EEEE6666
^*q6G1z~& 0000000000000000000000000000131BFAC 5555FFFFF3333AAAA22224447777222DDDD0000CCCC44445555
^*q86A*)gd 00000000000000000000000000002552351 0000DDDD7777333311116666BBBB111111113333555577776666
*q::**n*u 0000000000000000000000000000F7B410 999999999999FFFF999900008888AAAA1111DDDD9999BBB1111
^*q:@*|8Zg 0000000000000000000000000000028C028A DDDD222277770000AAAA22229999BBB777788884444FFFF7777
^*q:/ ,S1Bq 000000000000000000000000000001A0563C 4444AAAA8888555500011116666CCCC4444BBB000077775555
^*q<e?M0Y) 0000000000000000000000000000011845C0 2222BBBBDDDDDD6666DDDDDDDD8888CCCC9999333344441111
^*q?v3+Nu4 000000000000000000000000000000011437C9 33330000BBB99999AAAA66663333EEEE8888777788883333EEEE
^*qCrxf+g 000000000000000000000000000000F9701 9999555888888888882222BBBB9999AAAF777744446666
^*qCxc=9Dm 000000000000000000000000000002A26E2B BBBBAAAABBB3333EEEE9999AAAAEEEE66669999AAAF777444
[cchsu@n0 lab_4]$ hdfs dfs -cat /scr/cchsu/lab4/output/part-r-00001 | head -n 10
^*qE)r"(o 000000000000000000000000000002FD83AF BBBBB555AAAABBBBEEEE2222EEEE3333FFFF0000666655550000
^*qE*-z4Q) 00000000000000000000000000001DDOFF2 7777333322223333EEEE55556666BBBEEEE6666AAAAAAAF7777
^*qE9ao#wl 0000000000000000000000000000028DB12 5555666600001111AAAA1111BBBB1111FFFF111144440000FFFF
^*qEq-!J(J 000000000000000000000000000001C9561B 88887777AAAA22227777222266669999FFF6666AAAA00004444
^*qFN$)cHb 00000000000000000000000000000001806C8B 777788886666888855552222666622225555BBB666655554444
^*qJ;Ge4T 0000000000000000000000000000012D2C1E 222244442222111122223333BBB2222FFFF3333DDDDDDDD3333
^*qJDSr::c 0000000000000000000000000000000139F1C2 BBBBDDDD555533333333CCCCBBB0000BBBEEEE8888AAAA1111FFFF
^*qK'.wsp, 000000000000000000000000000000013263C7 EEEE0000000044440000000044441111BBBB33338888BBB8888
^*qKd:??s 00000000000000000000000000000FE6C0 111144444444DDDD0000FFFFAAAAEEEE4444AAAAEEEE44441111
^*qLE[5#05 000000000000000000000000000002DD7B1 777799990000EEEEDDDD9999111199999999AAAA888811116666
cat: Unable to write to output stream.

[cchsu@n0 lab_4]$ hdfs dfs -cat /scr/cchsu/lab4/output/part-r-00001 | tail -n 10
5Q3cnWxExD 00000000000000000000000000000015D7EDD 7777AAAAAAA66664444EEEE333355553333CCCC3333EEEEAAAA
5Q3dpMyJ7~ 000000000000000000000000000000B4DBE 8888AAAA000033336666AAAA11117777DDDDAAA0000BBB3333
5Q3eu516e& 00000000000000000000000000000D42DF5 CCCCBBBBEEEE3333BBB2222CCCC66667777BBB000099992222
5Q3f?7I)qN 0000000000000000000000000000019F29A1 3333333377776666AAAA22225555333399998888DDDDAAAA6666
5Q3fe*HL0 00000000000000000000000000000034BB6F 22226666DDDD66667777BBB9999111100002222EEEE99990000
5Q3fyQu";3 00000000000000000000000000000265606 FFFF9999AAAAADD5555EEEE1111333311111111AAAA8888BBB
5Q3k ,~bPT 0000000000000000000000000000007EC068 444422222222DDDD0000EEEE7777111133335555EEEE77779999
5Q31q!#!VT^ 000000000000000000000000000002516A1A 55551111555500008888BBB0000AAAAEEEE88883333AAAA7777
5Q3m(DupiP 0000000000000000000000000000000E30A88 22223333888833334444AAAA4444555544447777FFFFDDDD9999
5Q3nfFyAGyy 00000000000000000000000000000094BE53 BBBB11113333DDDD0000BBB3333EEEECCCCBBB2222CCCCCCCC
```

part-r-00002 and part-r-00003



```

[cchsu@n0 lab_4]$ hdfs dfs -cat /scr/cchsu/lab4/output/part-r-00004 | head -n 10
LdP'RXp',< 00000000000000000000000000000000009A434B 99992222777711110000111177774444CCCCDDDD3333DDDD4444
LdP(G)'@dh 0000000000000000000000000000000000C13170 FFFFBBBB000066668888DDDD555577770000BBBBEEEEDDDD1111
LdP2;MC_R- 00000000000000000000000000000000001BD9100 BBBB99992222FFFFEEEE5555EEEE000033338888BBBB00001111
LdP4cQJ]W 00000000000000000000000000000000001214C4E CCCC00003333CCCCCCCCCCCC9999BBBB00002222DDDD66663333
LdP5%Qe'M 0000000000000000000000000000000000795379 2222222233334444DDDDDDDDDDDDCCCC7777AAAACCCC0000EEEE
LdP6iwApM* 000000000000000000000000000000000064D4B7 2222DDDD8888FFFFAAAAABBBDDDDDEEEE7777FFFFFFFFCCCC8888
LdP7grbNaO 00000000000000000000000000000000003B5E9F AAAA99997777DDDDDEEEE00006666BBB3333DDDD999966660000
LdP=iA'Jse 00000000000000000000000000000000002EDE414 BBBB0000999966665555666666688883333888888883333DDDD
LdP>+N,;>` 0000000000000000000000000000000000140A875 6666CCCCBBBBEEEE22228888EEEEAAAA0000CCCCFFFF11112222
LdP>Bw*;3( 000000000000000000000000000000000097DBC0 CCCCBBB4444CCCCCCCCAAAAFFFFCCCC44442222555544441111
cat: Unable to write to output stream.
[cchsu@n0 lab_4]$ hdfs dfs -cat /scr/cchsu/lab4/output/part-r-00004 | tail -n 10
XCVbd1gI#4 0000000000000000000000000000000000666FAD 0000FFFFDDDD88889999000066661111EEEE8888CCCC9999AAAA
XCVd8?2VDj 00000000000000000000000000000000009957E9 999988884444EEEEBBBBEEEE222255557777CCCC11111111EEEE
XCVda2POvd 0000000000000000000000000000000000157323F 9999BBBBEEEEAAAA4444CCCC8888FFFFAAAA4444DDDDCCCC0000
XCVfP7O-w/ 000000000000000000000000000000000097DD8C 33332222AAAAABBB444499994444EEEE999955558888EEEE5555
XCVfi!PGJk 00000000000000000000000000000000009452C3 5555AAAA4444CCCCCCCC7777555577772222BBBB33335555CCCC
XCVg_F[SDf 000000000000000000000000000000000023EB102 11118888777777772222BBBB4444DDDD9999DDDD7777DDDDFFFF
XCVi+N<wy& 000000000000000000000000000000000022BEF36 0000DDDD44442222FFFFEEEE33337777AAAA0000CCCC1111BBBB
XCVi:cgzKs 00000000000000000000000000000000002A9A230 00008888DDDDDDDD6666EEEE3333FFFFDDDDAAAAABBB11111111
XCVkk'A+hr 0000000000000000000000000000000000BDF83F 55559999FFFFFFFF8888444433334444AAAA88887777CCCC0000
XCVko:Eff& 000000000000000000000000000000000047F9E2 CCCC333377778888111144449999AAAA9999999922227777FFFF
[cchsu@n0 lab_4]$ hdfs dfs -cat /scr/cchsu/lab4/output/part-r-00005 | head -n 10
XCVl!'5YH" 00000000000000000000000000000000001AF7F37 DDDD44448888EEEE2222FFFF55551111DDDDCCCC555544448888
XCVmK=*1J$ 00000000000000000000000000000000001EAEF22 FFFFDDDDFFFFBBBB55553333FFFF00006666FFFFFFFF3333FFFF
XCVpP2d_!B 00000000000000000000000000000000001CEF9FB BBBB11119999DDDD3333EEEE0000BBBB88888888000022224444
XCVq2?65}V 000000000000000000000000000000000015EE09A DDDDEEEFFFFF666611119999EEEE88886666BBB777722227777
XCVq;}X:67 0000000000000000000000000000000000164A4C7 444477774444AAAA5555444433334444BBBB99997777BBBB8888
XCVrA>80>s 00000000000000000000000000000000002C5246B 0000DDDD3333AAAAEEEE33334444FFFFBBBB55550000BBBB4444
XCVtRP>py? 000000000000000000000000000000000005C0AE4 555544442222AAAAABBB6666999977774444DDDD11112222DDDD
XCVt)+~~}q 000000000000000000000000000000000009E337 DDDD222244448888BBBBEEEE000099990000EEEE111144448888
XCVu)j([3} 000000000000000000000000000000000024B53E2 CCCC55557777666611111111888866665555222200007777FFFF
XCVx*,:;EY 000000000000000000000000000000000095F395 4444DDDD5555FFFF999966660000CCCC8888CCCCAAAAFFFF2222
cat: Unable to write to output stream.
[cchsu@n0 lab_4]$ hdfs dfs -cat /scr/cchsu/lab4/output/part-r-00005 | tail -n 10
c(J=ZjAxUJ 0000000000000000000000000000000000249BAB7 8888AAAAADDDEEEE555566666666FFFF111155559999CCCC8888
c(J>,JpY3= 000000000000000000000000000000000009CA40F EEEE1111666644446666444477779999444499995555FFFF0000
c(J>Iv n'H 00000000000000000000000000000000001D9EF43 99992222AAAAACCCCAAAAFFFF3333CCCCDDDDDDDEEEEDDDDDCCC
c(J>y#Y=OZ 00000000000000000000000000000000001340373 22223333FFFF9999CCCCBBBB44448888DDDDDEEEFFFFFAAAACCCC
c(JATY6#8K 000000000000000000000000000000000001C50536 4444EEEE1111AAAAFFFF8888AAAAABBB00002222EEEE1111BBBB
c(JFiVaG:] 000000000000000000000000000000000008342CE BBBBAAAA00004444AAAA666611119999EEEE44448888EEEE3333
c(JHX(:;z7 0000000000000000000000000000000000E7DFEB 111100001111EEEE888844445555333311117777DDDD33334444
c(JJ!-?yd{ 00000000000000000000000000000000000945FB6 FFFFCCCC333311110000CCCCBBBB22226666000055559999BBBB
c(JJN6J0Q 0000000000000000000000000000000000099CB71 44441111AAAAACCCFFFAAAA222233334444DDDDAAAA55556666
c(JKb\XEIf 00000000000000000000000000000000001F82369 9999EEEEEECCCC11117777BBBB1111FFFF111111119999EEEE

```

part-r-00006 and part-r-00007



```

[cchsu@n0 lab_4]$ hdfs dfs -cat /scr/cchsu/lab4/output/part-r-00006 | head -n 10
c(JL@Yx7ze 00000000000000000000000002B798AF 5555AAAA9999DDDDAAAA22225555EEEEFFFFF8888111155550000
c(JMO97U.r 00000000000000000000000002C5E6DD 00005555FFFF2222222233337777EEEE77771111BBBBEEEEEEAAAA
c(JMZ,M;3# 00000000000000000000000000003259 DDD22222999999994444DDDD777733330000111144442222EEEE
c(JNS5=Em` 00000000000000000000000000005EBB74 AAAA1111666644447777BBBB999944440000EEEE99995555DDDD
c(JN1~7saq 00000000000000000000000000020D0B3 4444EEEEFFFFFBBB6666CCCC555599991111CCCC66666666CCCC
c(JPS)_i"? 000000000000000000000000000226E383 111177777777FFFF3333222222227777EEEE0000EEEE9999CCCC
c(JRd$]YW 00000000000000000000000000024836D8 DDD2222444444443333FFFF22220000AAAA BBBB222244449999
c(JS&I.240 00000000000000000000000000025D5071 AAAA222222225555666633331111777744449999DDDD55556666
c(JT(4+~H7 000000000000000000000000001528FAA BBBEEEEFFFFFAAAA7777BBBB99992222CCCC7777000055557777
c(JT>yZg4# 0000000000000000000000000000801DE7 EEEE5555666644440000EEEE1111444499995555AAAA99998888
cat: Unable to write to output stream.
[cchsu@n0 lab_4]$ hdfs dfs -cat /scr/cchsu/lab4/output/part-r-00006 | tail -n 10
nHJdd^7b1* 0000000000000000000000000000130199 BBBB333366667777DDDD0000555599997777EEEEEEEEEEEEEEEE
nHJghQHP3m 0000000000000000000000000000EC0F49 666655557777AAAA4444EEEEDDDD2222FFFF5555CCCCBBBBEEEE
nHJj-,1b;G 0000000000000000000000000002A6CB2A 00000000EEEE00009999DDDD555599994444FFFF2222DDDD7777
nHJc9-OdH/ 000000000000000000000000000244EE39 11114444EEEE11118888EEEEDDDBBB6666DDDD77774444EEEE
nHJowC)JxB 0000000000000000000000000000E2F720 222233338888BBBB000088882222EEEEDDDD5555CCCC66661111
nHJqct"ceL 00000000000000000000000000020D0794 1111AAAA BBBB BBBB77779999AAAA BBBB BBBBCCCCEEEEEEEEBBBBDDDD
nHJrcG7'2P 0000000000000000000000000000C970BC 333322220000FFFFEEEE BBBB7777CCCC222244443333FFFF5555
nHJtj~ojK" 000000000000000000000000000013ACA2E CCCC AAAA FFFF BBBB FFFF2222AAAA55552222CCCC999900003333
nHJu$2H)6h 00000000000000000000000000001537DBF 6666444422228888FFFF000088880000BBBB4444222244440000
nHJx~LMTg+ 00000000000000000000000000001E253FF 33331111FFFFAAAA777788885555DDDD BBBB AAAA CCCC00000000
[cchsu@n0 lab_4]$ hdfs dfs -cat /scr/cchsu/lab4/output/part-r-00007 | head -n 10
nHJxpj\vjd 0000000000000000000000000002A75ACD 11116666DDDDDDDD3333FFFF888822226666888844447777AAAA
nHJz2)?BjC 00000000000000000000000000001032F1C FFFF333322227777DDDDDDDDAAAA DDDDDAAA BBBBEEEE11115555
nHJ{ '~x#^O 00000000000000000000000000001DF3828 6666444400001111EEEE BBBB3333999944444444FFFF BBBB9999
nHJ|_epIT< 0000000000000000000000000000024B20D 8888 BBBB44448888EEEEAAAA8888AAAAEEEE BBBB BBBB3333AAAA
nHK"TSWAS# 0000000000000000000000000000F91DE0 22226666 BBBB4444777755558888EEEE22228888AAAAAAA1111
nHK&PB^lM$ 0000000000000000000000000000DA452D 2222FFFF44443333EEEE00001111DDDD1111FFFF BBBB1111AAAA
nHK&lpm{EZ 00000000000000000000000000001E7C667 DDD22223333111177775555222211113333CCCC222211118888
nHK't.Rv=3 00000000000000000000000000001079F47 CCCC555511115555999944446666FFFF FFAA AADD DD33338888
nHK) ?X$YH/ 00000000000000000000000000008EDA0B 111188882222CCCC55559999 BBBB EEEE55556666111111114444
nHK0oQtgu{ 0000000000000000000000000002AD18FE FFFF55550000EEEE0000 BBBB AAAA CCCC66666666EEEE77773333
cat: Unable to write to output stream.
[cchsu@n0 lab_4]$ hdfs dfs -cat /scr/cchsu/lab4/output/part-r-00007 | tail -n 10
~~~r1o,;!4 0000000000000000000000000002977B8D 33337777222255556666DDDDAAAA9999444455555555BBBBAAAA
~~~ruROgur 0000000000000000000000000001E87D8A BBBBEEEE22222222BBBBAAAA55551111FFFF99995555FFFF7777
~~~s!SL~~7 0000000000000000000000000002ECCEDA 11114444FFFF BBBB11113333000000008888AAAA9999EEEE7777
~~~s/Pq,-E 00000000000000000000000000006BE930 2222DDDDDDDD77771111EEEECCCC7777BBBB4444888811111111
~~~sf&V`wv 00000000000000000000000000015B2E94 6666 BBBB1111EEEEEEEE5555CCCCDDDD55555555BBBBBBDDDD
~~~uq2k#=U 0000000000000000000000000002C06745 99991111DDDD222211110000FFFFEEEEFFFF33337777CCCC2222
~~~yK01:gE 0000000000000000000000000002048B4F CCCC11114444888822226666BBBB888855557777EEEE BBBB0000
~~~zbA Tt 00000000000000000000000000007F9F4F BBBBCCCC666655559999FFFF8888AAAA11116666AAAA BBBB0000
~~~ze0^FEg 0000000000000000000000000001E06130 4444CCCC BBBB99992222888855558888CCCCFFFF000011111111
~~~}GxjWHI 0000000000000000000000000000CA1345 777711118888AAAAAAA22221111BBBB00002222BBBBCCCC2222

```

## Resource

- <https://hadoop.apache.org/docs/r2.4.1/api/org/apache/hadoop/mapred/Partitioner.html>

## Template Code

- Chapter 9 MapReduce Features: Total Sort, Hadoop – The Definitive Guide