

ReadZipWriteHDFS.java

```

1 package applications;
2
3 import java.io.File;
17
18 public class ReadZipWriteHDFS
19 {
20     // Expands the zip file passed as argument 1, into the
21     // directory provided in argument 2
22
23     public static String ZIP_FILE_IN = "";
24     public static String HDFS_FILE_OUT = "";
25     public static int FILES_PROCESSED = 0;
26
27     public static void main(String args[]) throws Exception
28     {
29         ReadZipWriteHDFS me = new ReadZipWriteHDFS();
30         me.go(args);
31     }
32
33     public void go(String[] args) throws FileNotFoundException, IOException
34     {
35         System.out.println("Begin execution of " + this.getClass().getName());
36
37         ZIP_FILE_IN = args[0];
38         HDFS_FILE_OUT = args[1];
39
40         System.out.println("Will read from zip file: " + ZIP_FILE_IN);
41         System.out.println("Will write to HDFS file: " + HDFS_FILE_OUT);
42
43         // see Hadoop in Action page 44
44         Configuration conf = new Configuration();
45         FileSystem hdfs = FileSystem.get(conf);
46         Path hdfsFile = new Path(HDFS_FILE_OUT);
47         FSDataOutputStream outputStream = hdfs.create(hdfsFile);
48
49         ZipFile zipFile = new ZipFile(new File(ZIP_FILE_IN));
50         InputStream is = new FileInputStream(ZIP_FILE_IN);
51         ZipInputStream inStream = new ZipInputStream(is);
52
53         recursive_extract(zipFile, inStream, outputStream);
54         outputStream.close();
55
56         System.out.println(FILES_PROCESSED + " files processed.");
57         System.out.println("End execution of " + this.getClass().getName());
58     }
59
60     public void recursive_extract(ZipFile zipFile, ZipInputStream inStream, OutputStream
61     outputStream)
62     {
63         // create a buffer to improve copy performance later.
64         byte[] buffer = new byte[2048];
65
66         try
67         {
68             // now iterate through each item in the stream. The get next
69             // entry call will return a ZipEntry for each file in the
70             // stream

```

ReadZipWriteHDFS.java

```

70     ZipEntry entry;
71     while((entry = inStream.getNextEntry())!=null)
72     {
73         /*
74         String s = String.format("Entry: %s len %d added %TD",
75                                 entry.getName(), entry.getSize(),
76                                 new Date(entry.getTime()));
77         System.out.println(s);
78         */
79
80         if (entry.getName().contains(".zip"))
81         {
82             InputStream nestedInputStream = zipFile.getInputStream(entry);
83             ZipInputStream nestedZipInputStream = new
ZipInputStream(nestedInputStream);
84             recursive_extract(zipFile, nestedZipInputStream, outputStream);
85         }
86         else
87         {
88             try
89             {
90                 int len = 0;
91                 while ((len = inStream.read(buffer)) > 0)
92                 {
93                     outputStream.write(buffer, 0, len);
94                 }
95                 FILES_PROCESSED++;
96                 if (FILES_PROCESSED % 1000 == 0)
97                 {
98                     System.out.println(FILES_PROCESSED + " files processed so far.");
99                 }
100            }
101            catch (Exception e)
102            {
103                System.out.println(" !!! Error processing file " + entry.getName());
104                // e.printStackTrace();
105            }
106        }
107    }
108 }
109 catch (Exception e)
110 {
111     e.printStackTrace();
112 }
113 }
114 }

```