Lecture 32: Sorting Big Data: Sorting on the Disk

```
1    def split(file, size):
2        open_files = []
3        with open(file) as fin:
4            tmp = tempfile.TemporaryFile('w+t')
5            for line in fin:
6                if os.fstat(tmp.fileno()).st_size < size:
7                    print(line,file=tmp,end="")
8                else:
9                    tmp.flush()
10                   tmp.seek(0)
11                   open_files.append(tmp)
12                   tmp = tempfile.TemporaryFile('w+t')
13                   print(line,file=tmp,end="")
14
15           tmp.flush()
16           tmp.seek(0)
17           open_files.append(tmp)
18
19       return open_files
```

```python
 1   def sort_files(files):
 2
 3       sorted_files = []
 4
 5       for file in files:
 6           contents = [line for line in file]
 7           contents.sort()
 8           tmp = tempfile.TemporaryFile('w+t')
 9           for line in contents:
10               print(line, file=tmp, end="")
11           tmp.flush()
12           tmp.seek(0)
13           sorted_files.append(tmp)
14           file.close()
15
16       return sorted_files
```

```
1   def merge_files(files, final):
2
3       tmp = files[0]
4       for file in files[1:]:
5           tmp2 = tempfile.TemporaryFile('w+t')
6           for line in merge_iter(tmp, file):
7               print(line,file=tmp2,end="")
8           tmp.close()
9           tmp = tmp2
10          tmp.flush()
11          tmp.seek(0)
12
13      with open(final, 'w+t') as fout:
14          for line in tmp:
15              print(line,end='',file=fout)
```

```
 1   def merge_iter(iter1, iter2):
 2       try:
 3           val1 = next(iter1)
 4           val2 = next(iter2)
 5           while True:
 6               if val1 < val2:
 7                   yield val1
 8                   val1 = next(iter1)
 9               else:
10                   yield val2
11                   val2 = next(iter2)
12
13       except StopIteration:
14           # one of the two iterators is empty, but we don't know which, so
15           # just yield all the remaining values in both (the one without
16           # any remaining values won't yield anything
17           for val in iter1:
18               yield val
19           for val in iter2:
20               yield val
```

```
1   def bigsort(input, output, size=2*20):
2       merge_files(sort_files(split(input, size)), output)
```