# CONCURRENCY II

Instructors: Crista Lopes

# Threads

- Cute and furry beasts

# Threads

□ Must be restrained

# Restrained concurrency models

- Actors
  - Good for independent tasks
  - Good for discriminate producers/consumers of data
- Tuple spaces
  - Good for indiscriminate producers/consumers of data
- Map-reduce
  - Good for data-intensive, parallelizable situations

# Actors

# Actor model

- Letterbox style (Ch 11) + Threads
- Actor = Object with its own thread
  - Aka "active object"
- Actors send messages to each other
  - Avoid shared memory
- Messages are placed in actors' queues
  - Queues must be "thread-safe"
  - Sender places message and moves on
    - Asynchronous request

# Active Object (Python)

```python
7  class ActiveWFObject(Thread):
8      def __init__(self):
9          Thread.__init__(self)
10         self.name = str(type(self))
11         self.queue = Queue()          ←── Thread-safe queue
12         self._stop = False
13         self.start()
14
15     def run(self):
16         while not self._stop:
17             message = self.queue.get()    Block until there
18             self._dispatch(message)       is a message
19             if message[0] == 'die':       Message loop
20                 self._stop = True
21
22 def send(receiver, message):
23     receiver.queue.put(message)       Utility (could be a method)
24
```
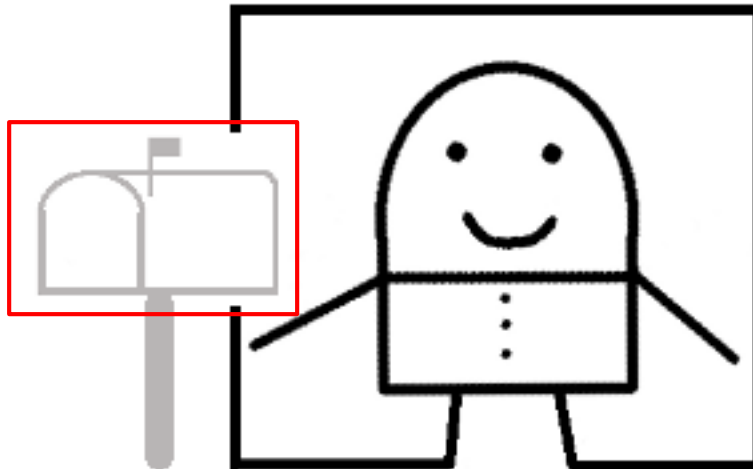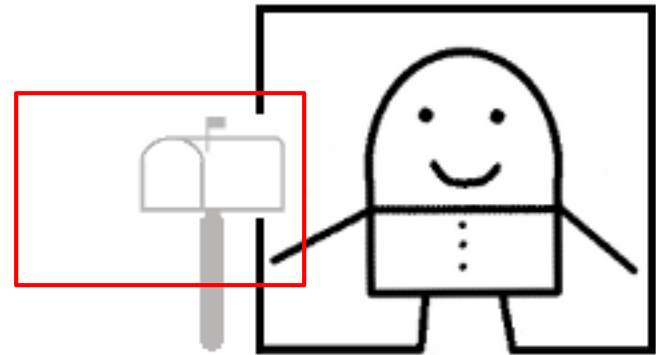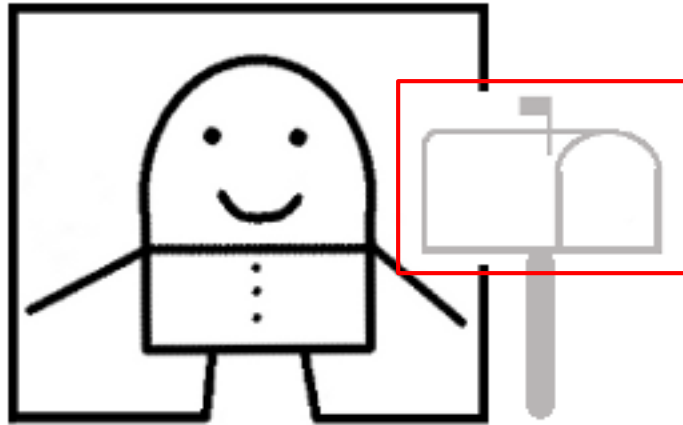
# Active Object Queue

# Queues

- Put / Enqueue / Send
- Get / Dequeue / Receive


- Operations must be thread safe
  - No items can be lost

# Thread-safe queues

- Java: ArrayBlockingQueue

- C#: ConcurrentQueue

- C++ / Boost: message_queue

- Other langs: search for it or do it yourself

# Actor example

superclass

```
78  class WordFrequencyManager(ActiveWFObject):
79      """ Keeps the word frequency data """
80      _word_freqs = {}
81
82      def _dispatch(self, message):
83          if message[0] == 'word':
84              self._increment_count(message[1:])
85          elif message[0] == 'top25':
86              self._top25(message[1:])
87
88      def _increment_count(self, message):
89          word = message[0]
90          if word in self._word_freqs:
91              self._word_freqs[word] += 1
92          else:
93              self._word_freqs[word] = 1
94
95      def _top25(self, message):
96          recipient = message[0]
97          freqs_sorted = sorted(self._word_freqs.iteritems(), key=
98              operator.itemgetter(1), reverse=True)
            send(recipient, ['top25', freqs_sorted])
```

dispatch messages

Send messages to other actors

# Actor model

- Concurrency constrained by
  - Associating [certain] objects with threads
  - Using message queues in each actor
  - Having threads on a loop
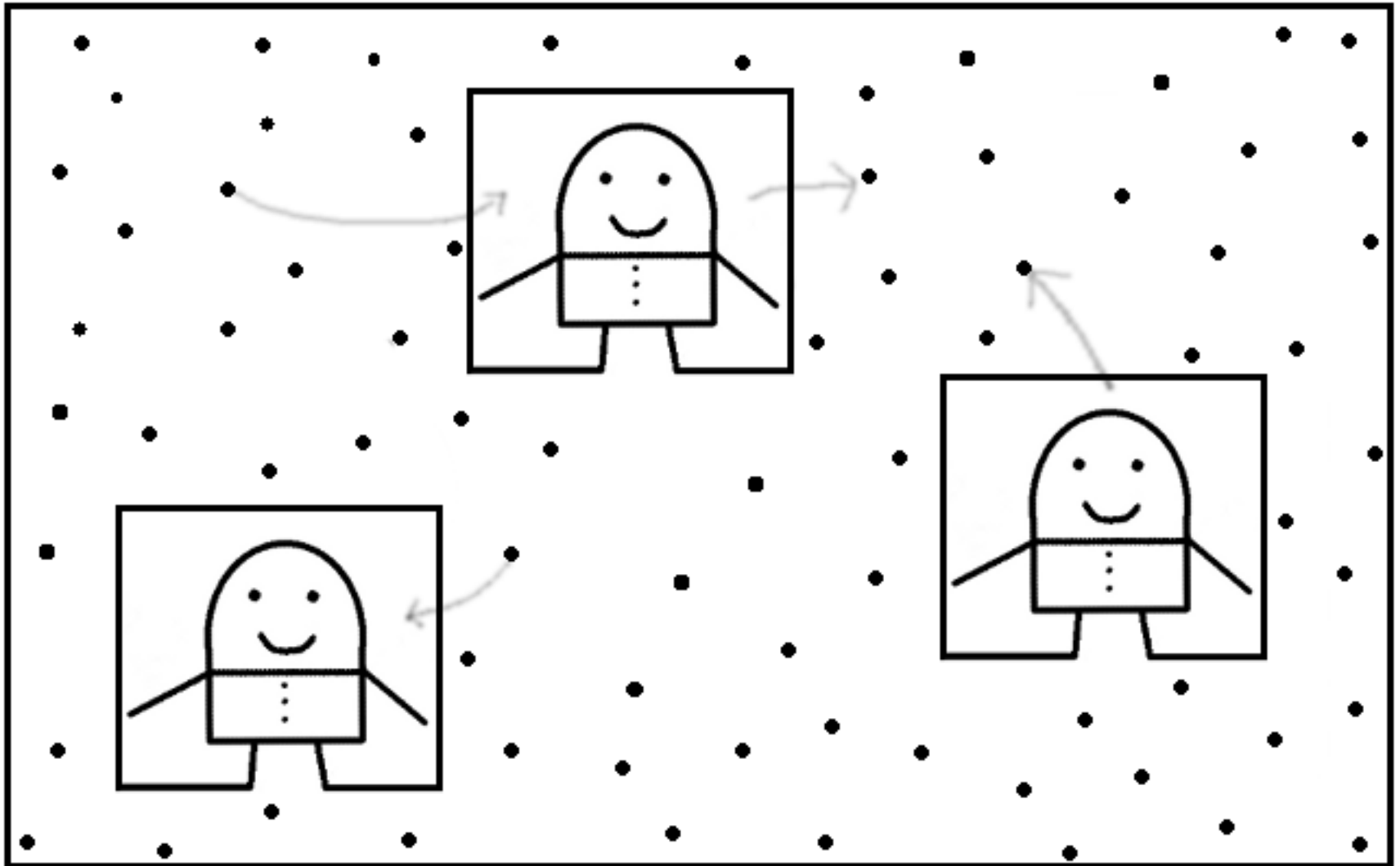- Programmer needs to refrain from passing shared mutable objects around or else...

# Tuple Spaces

# Tuple space model

- Concurrent threads
  - Consumers and producers of data items
- Shared data structures (queues, lists, trees, etc.)
  - Must be "thread-safe"
- Producers add items and move on
    - Asynchronous deposit
- Consumers take items and process them

- Similar to previous model, but where the queues are outside the objects/functions, and may not be queues

# Tuple space model

# TF Tuple spaces

```
4  # Two data spaces
5  word_space = Queue.Queue()
6  freq_space = Queue.Queue()
```

# TF Producer

```
26  # Let's have this thread populate the word space
27  for word in re.findall('[a-z]{2,}', open(sys.argv[1]).read().lower
        ()):
28      word_space.put(word)
```

# TF Consumer / Producer worker

```python
10  # Worker function that consumes words from the word space
11  # and sends partial results to the frequency space
12  def process_words():
13      word_freqs = {}
14      while True:
15          try:
16              word = word_space.get(timeout=1)
17          except Queue.Empty:
18              break
19          if not word in stopwords:
20              if word in word_freqs:
21                  word_freqs[word] += 1
22              else:
23                  word_freqs[word] = 1
24      freq_space.put(word_freqs)
```

# Starting workers

```python
30 # Let's create the workers and launch them at their jobs
31 workers = []
32 for i in range(5):
33     workers.append(threading.Thread(target = process_words))
34 [t.start() for t in workers]
```

(functional style of creating threads in Python)

# Tuple space model

☐ Can be functional or OOP style

☐ OOP style: worker functions are threaded objects

☐ Best fit: data processing parallelization

Actors example decomposition:
- DataStorageManager
- StopWordManager
- WordFrequencyManager
- WordFrequencyController

Tuple space example decomposition:
- Producers of words
- Consumers of words / producers of word frequencies
- Consumers of word frequencies

# Tuple space model

- Concurrency constrained by
  - Having shared, thread-safe collections of items
  - Having producers/consumers of items in those collections
  - No further communication between threaded code
- Programmer needs to refrain from passing shared mutable objects around or else...
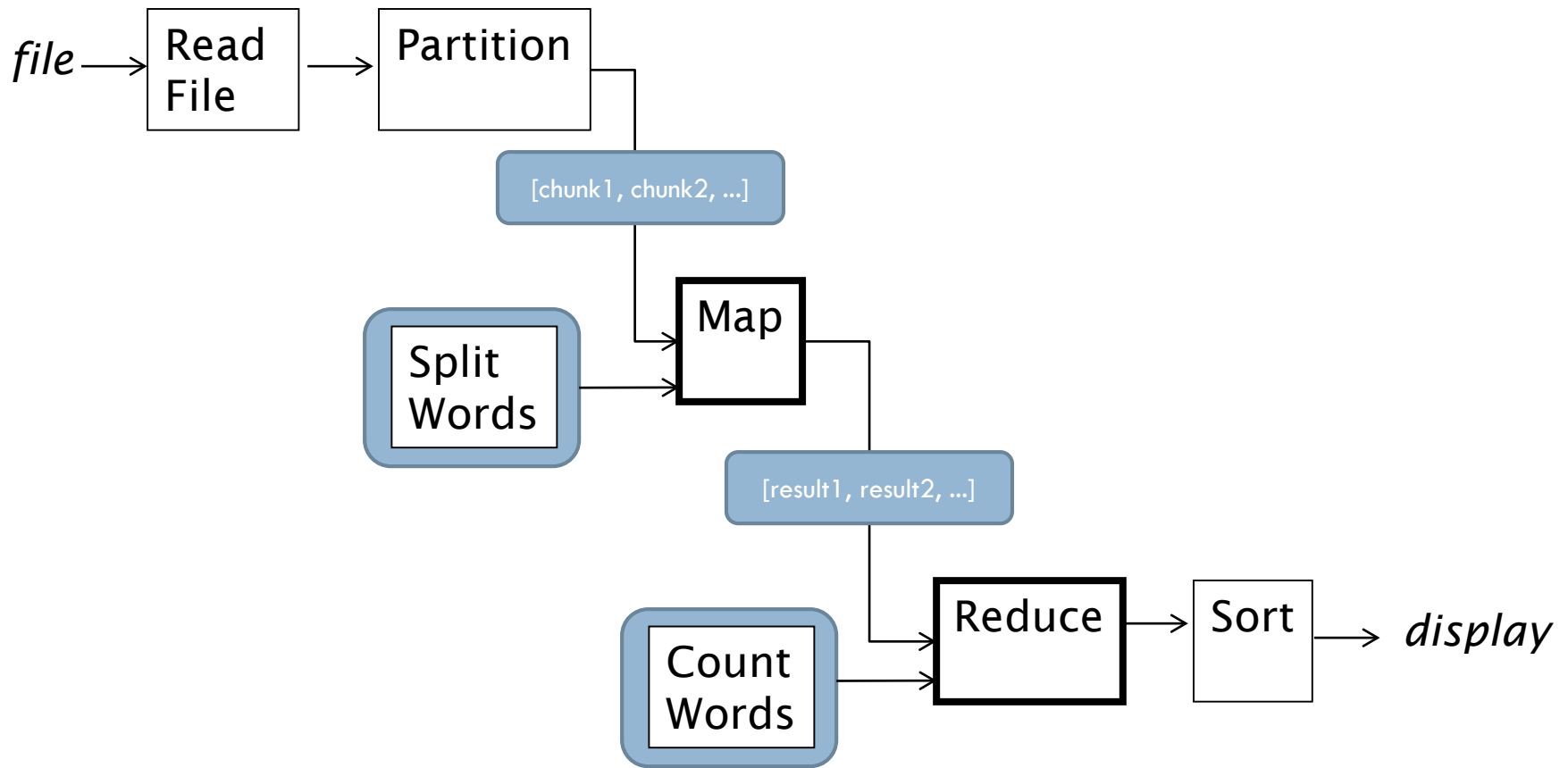
# Map-Reduce

# Map-Reduce model

- Big data situations
  - Problem at hand must be data-parallelizable
- Data is split into chunks
- Chunks are processed independently, produce partial results
  - A function is "mapped" to the chunks of data, <u>potentially in parallel</u>
- Partial results are then "reduced" to final result
  - This step is sequential

```
splits = map(split_words,partition(read_file(sys.argv[1]),200))
splits.insert(0, []) # normalize input to reduce
word_freqs = sort(reduce(count_words, splits))
```

# Data partitioning

```python
7  def partition(data_str, nlines):
8      """
9      Partitions the input data_str (a big string)
10     into chunks of nlines.
11     """
12     lines = data_str.split('\n')
13     for i in xrange(0, len(lines), nlines):
14         yield '\n'.join(lines[i:i+nlines])
```

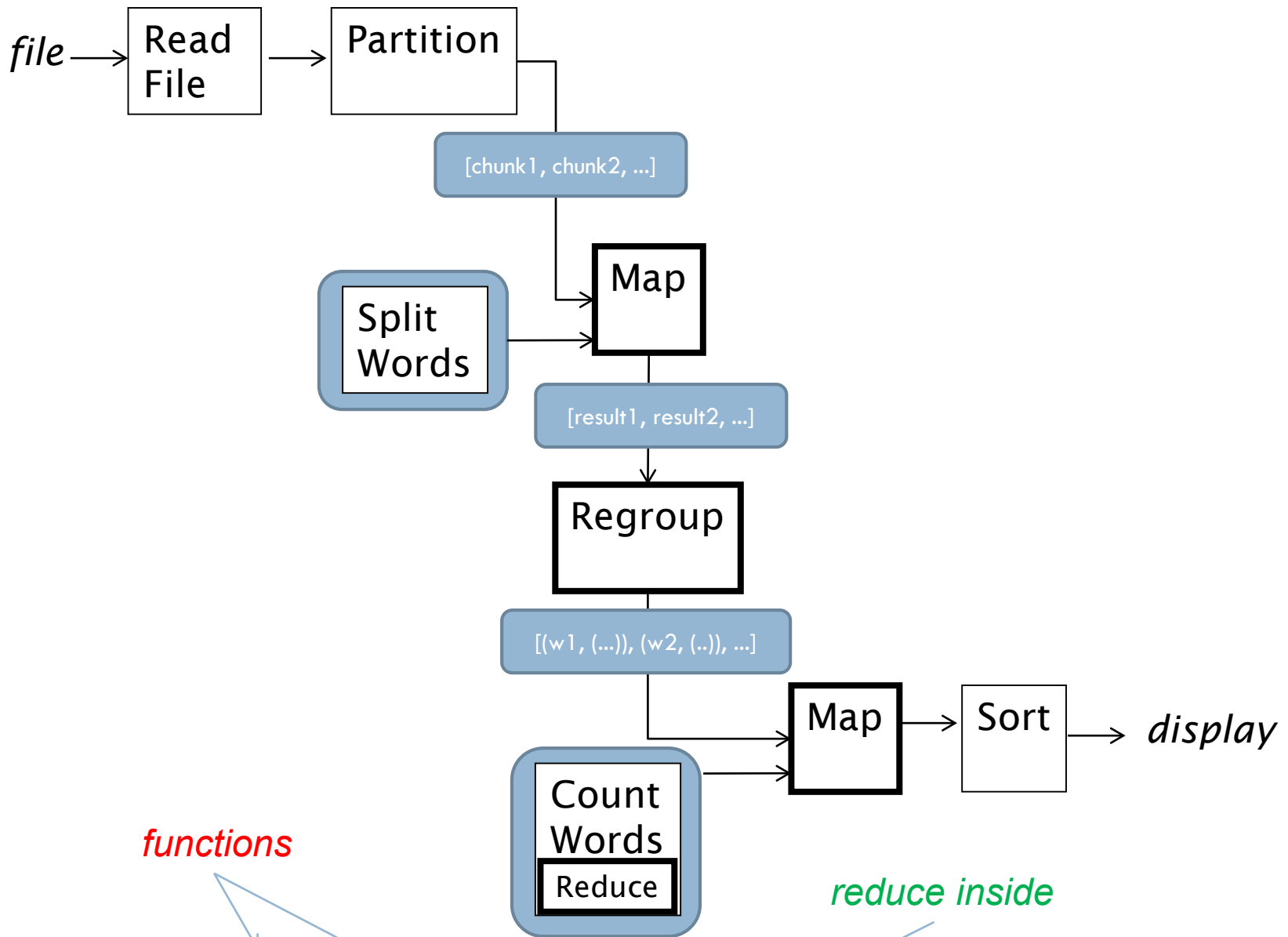# Mapper – parsing words – emit

```python
16  def split_words(data_str):
17      """
18      Takes a string,  returns a list of pairs (word, 1),
19      one for each word in the input, so
20      [(w1, 1), (w2, 1), ..., (wn, 1)]
21      """
22      def _scan(str_data):
23          pattern = re.compile('[\W_]+')
24          return pattern.sub(' ', str_data).lower().split()
25
26      def _remove_stop_words(word_list):
27          with open('../stop_words.txt') as f:
28              stop_words = f.read().split(',')
29          stop_words.extend(list(string.ascii_lowercase))
30          return [w for w in word_list if not w in stop_words]
31
32      # The actual work of splitting the input into words
33      result = []
34      words = _remove_stop_words(_scan(data_str))
35      for w in words:
36          result.append((w, 1))
37      return result
```

# Reducer – counting words

```python
39  def count_words(pairs_list_1, pairs_list_2):
40      """
41      Takes a two lists of pairs of the form
42      [(w1, 1), ...]
43      and returns a list of pairs [(w1, frequency), ...],
44      where frequency is the sum of all the reported occurrences
45      """
46      mapping = dict((k, v) for k, v in pairs_list_1)
47      for p in pairs_list_2:
48          if p[0] in mapping:
49              mapping[p[0]] += p[1]
50          else:
51              mapping[p[0]] = 1
52      return mapping.items()
```

# Map-Reduce, Hadoop

- The previous style allows for parallelization of the map step, but requires serialization of the reduce step. Google map-reduce and Hadoop use a slight variation that makes the reduce step also potentially parallelizable. The main idea is to regroup, or reshuffle, the list of results from the map step so that the regroupings are amenable to further mapping of a reducible function.

```
file → Read File → Partition
                         [chunk1, chunk2, ...]

Split Words → Map
                    [result1, result2, ...]

                    Regroup
                    [(w1, (...)), (w2, (..)), ...]

                              Map → Sort → display
                    Count
                    Words
                    Reduce
```

*functions*

*reduce inside*

```
splits = map(split_words,partition(read_file(sys.argv[1]),200))
splits_per_word = regroup(splits)
word_freqs = sort(map(count_words, splits_per_word.items()))
```

# Regroup

```python
def regroup(pairs_list):
    """
    Takes a list of lists of pairs of the form
    [[(w1, 1), (w2, 1), ..., (wn, 1)],
     [(w1, 1), (w2, 1), ..., (wn, 1)],
     ...]
    and returns a dictionary mapping each unique word to the
    corresponding list of pairs, so
    { w1 : [(w1, 1), (w1, 1)...],
      w2 : [(w2, 1), (w2, 1)...],
      ...}
    """
    mapping = {}
    for pairs in pairs_list:
        for p in pairs:
            if p[0] in mapping:
                mapping[p[0]].append(p)
            else:
                mapping[p[0]] = [p]
    return mapping
```

# Map-reduce

- Java: try [Functional Java](#) library
  - Or do "mapper" and "reducer" classes yourself

- PHP: try [this](#)

# Map-Reduce model

- Concurrency constrained by
  - Having worker threads work on mutually exclusive chunks of data
  - No communication between threaded code