# dnaseq.py

```python
# Maps integer keys to a set of arbitrary values.
class Multidict:
  # Initializes a new multi-value dictionary, and adds any key-value
  # 2-tuples in the iterable sequence pairs to the data structure.
  def __init__(self, pairs=[]):
    raise Exception("Not implemented!")
  # Associates the value v with the key k.
  def put(self, k, v):
    raise Exception("Not implemented!")
  # Gets any values that have been associated with the key k; or, if
  # none have been, returns an empty sequence.
  def get(self, k):
    raise Exception("Not implemented!")

# Given a sequence of nucleotides, return all k-length subsequences
# and their hashes.  (What else do you need to know about each
# subsequence?)
def subsequenceHashes(seq, k):
  raise Exception("Not implemented!")

# Similar to subsequenceHashes(), but returns one k-length
    subsequence
# every m nucleotides.  (This will be useful when you try to use two
# whole data files.)
def intervalSubsequenceHashes(seq, k, m):
  raise Exception("Not implemented!")

# Searches for commonalities between sequences a and b by comparing
# subsequences of length k.  The sequences a and b should be
    iterators
# that return nucleotides.  The table is built by computing one hash
# every m nucleotides (for m >= k).
def getExactSubmatches(a, b, k, m):
  raise Exception("Not implemented!")

if __name__ == '__main__':
  if len(sys.argv) != 4:
    print 'Usage: {0} [file_a.fa] [file_b.fa] [output.png]'.format(
        sys.argv[0])
    sys.exit(1)

  # The arguments are, in order: 1) Your getExactSubmatches
  # function, 2) the filename to which the image should be written,
  # 3) a tuple giving the width and height of the image, 4) the
  # filename of sequence A, 5) the filename of sequence B, 6) k, the
  # subsequence size, and 7) m, the sampling interval for sequence
  # A.
  compareSequences(getExactSubmatches, sys.argv[3], (500,500), sys.
      argv[1], sys.argv[2], 8, 100)
```

## dnaseqlib.py

```python
# Produces hash values for a rolling sequence.
class RollingHash:
    def __init__(self, s):
        self.HASH_BASE = 7
        self.seqlen = len(s)
        n = self.seqlen - 1
        h = 0
        for c in s:
            h += ord(c) * (self.HASH_BASE ** n)
            n -= 1
        self.curhash = h

    # Returns the current hash value.
    def current_hash(self):
        return self.curhash

    # Updates the hash by removing previtm and adding nextitm.
        Returns the updated
    # hash value.
    def slide(self, previtm, nextitm):
        self.curhash = (self.curhash * self.HASH_BASE) + ord(nextitm
            )
        self.curhash -= ord(previtm) * (self.HASH_BASE ** self.
            seqlen)
        return self.curhash
```

```python
def compareSequences(getExactSubmatches, imgfile, imgsize, afile,
    bfile, k, m):
    a = kfasta.FastaSequence(afile)
    b = kfasta.FastaSequence(bfile)
    matches = getExactSubmatches(a, b, k, m)
    buildComparisonImage(imgfile, imgsize[0], imgsize[1],
                          kfasta.getSequenceLength(afile),
                          kfasta.getSequenceLength(bfile), matches)
```

## kfasta.py

```
1   # An iterator that returns the nucleotide sequence stored in the
         given FASTA file.
2   class FastaSequence:
3       def __init__(self, filename):
4           self.f = open(filename, 'r')
5           self.buf = ''
6           self.info = self.f.readline()
7           self.pos = 0
8       def __iter__(self):
9           return self
10      def next(self):
11          while '' == self.buf:
12              self.buf = self.f.readline()
13              if '' == self.buf:
14                  self.f.close()
15                  raise StopIteration
16              self.buf = self.buf.strip()
17          nextchar = self.buf[0]
18          self.buf = self.buf[1:]
19          self.pos += 1
20          return nextchar
```

## Iterators vs Generators

```
1  class Reverse:
2      """Iterator for looping over a sequence backwards."""
3      def __init__(self, data):
4          self.data = data
5          self.index = len(data)
6      def __iter__(self):
7          return self
8      def next(self):
9          if self.index == 0:
10             raise StopIteration
11         self.index = self.index - 1
12         return self.data[self.index]
13
14 # >>> rev = Reverse('spam')
15 # >>> iter(rev)
16 # <__main__.Reverse object at 0x00A1DB50>
17 # >>> for char in rev:
18 # ...     print char
19 # ...
20 # m
21 # a
22 # p
23 # s
```

```
1  def reverse(data):
2      for index in range(len(data)-1, -1, -1):
3          yield data[index]
4
5  # >>> for char in reverse('golf'):
6  # ...     print char
7  # ...
8  # f
9  # l
10 # o
11 # g
```

```
1  >>> data = 'golf'
2  >>> list(data[i] for i in range(len(data)-1,-1,-1))
3  ['f', 'l', 'o', 'g']
4
5  >>> sum(i*i for i in range(10))              # sum of squares
6  285
7
8  >>> xvec = [10, 20, 30]
9  >>> yvec = [7, 5, 3]
10 >>> sum(x*y for x,y in zip(xvec, yvec))      # dot product
11 260
```

6.006 Introduction to Algorithms
Fall 2011