# EPL646 – Advanced Topics in Databases
# Advanced Hadoop

http://www.cs.ucy.ac.cy/~dzeina/courses/epl646/labs/lab.html

# Calculate time

- How to calculate the time

```java
long begin = System.currentTimeMillis();
job.waitForCompletion(true);
long end = System.currentTimeMillis();
long second = (end - begin) / 1000;
System.err.println(job.getJobName() + " takes " +
second + " seconds");
```

# Task1: N-Gram

- Change the code in WordCount so that it counts how many times each set of five consecutive words appears

- You can find the code of WordCount from the solution of the previous lab

- If you don't have the datasets you can download them from the previous lab

# Task1: N-Gram

- Function *map* will have as input:
  - key = line offset (we can ignore it)
  - value = a whole line from one of the input files
- Function *map* will have as output:
  - key = five words
  - value = 1
- Function *reduce* will have as input:
  - key = five words
  - value = [a list of number 1]
- The list will as many 1 as there are appearances of the five consecutive words in our data
- Function *Reduce* will have as final output:
  - key = five words
  - value = the sum of all 1 (i.e. the same as WordCount)

University of Cyprus

# Task2: Anagram

- An anagram is a word that can be created by the movement of the letters of another word

- E.g.
  - Refills➔fillers
  - Relayed➔layered
  - Rentals➔antlers
  - Rebuild➔builder

- You must find the anagrams in a huge input file. How would you do it?

  **public static boolean isAnagram(String first, String second) {**
      // Checks that the two inputs are anagrams, by checking they have all the same characters.
      // Left as exercise for the user...
    }

University of Cyprus

# Task2: Anagram

- **Hadoopifying…**
  - (**input**) <k1, v1> →
  - map → <k2, v2> →
  - combine → <k2, v2>→
  - reduce → <k3, v3> (**output**)

- Download
  - /usr/share/dict/words or /usr/dict/words

University of Cyprus

# Jar file configuration

You need to set the jar by class parameter:

```
Configuration conf = new Configuration();
Job job = Job.getInstance(conf, "word count");
job.setJarByClass(WordCount.class);
```

**Else you will get:**

java.lang.RuntimeException: java.lang.ClassNotFoundException:

# Export the .jar file

# Run the jar file



- If the node is unhealthy you may need to execute the following command and then restart hadoop:

  chown -r epl-646:epl-646  /app/hadoop/

# Questions?

http://www.cs.ucy.ac.cy/~dzeina/courses/epl646/labs/lab.html