

# Text File Similarity Measurement Report

School of Computer Science and Information Technology

COSC2468 Scripting Language Programming

Semester 1 2018 – Assignment 1 – Perl

Jenna Afarian

S3511312

The intent of this report is to explain the choices made in Part D to measure the similarity of two text files. The following interpretation of the specification was made:

1. Similarity measurement should be based on a word frequency count (or how many times words overlap with each other between files).
2. If files contain exactly the same words, they match 100%.
3. If files contain no duplicate words, they match 0%.

Using Text::Similarity::Overlaps I was able to count the number of times a word overlapped between files. I then chose to get a percentage of overlaps based on the number of words in the largest file by counting the words in both files.

To determine the percentage I used the following calculation:

$\text{<number of overlaps> / <number of words in largest file> * 100}$

#### Example of a 100% match

<b>File 1 contents</b>	I am a test file.
<b>File 2 contents</b>	I am a test file.
<b>Number of overlaps</b>	5
<b>Word count of largest file</b>	5
<b>Calculation</b>	$5 / 5 * 100$
<b>Result</b>	100%

#### Example of a 0% match

<b>File 1 contents</b>	I am a test file.
<b>File 2 contents</b>	This is not your lunch, Mr. Potter.
<b>Number of overlaps</b>	0
<b>Word count of largest file</b>	7
<b>Calculation</b>	$0 / 7 * 100$
<b>Result</b>	0%