



**CALIFORNIA STATE SCIENCE FAIR  
2012 PROJECT SUMMARY**

<b>Name(s)</b> <b>Samyukta Yagati</b>	<b>Project Number</b> <b>S1429</b>
<b>Project Title</b> <b>Detecting Duplicate Content in Text Documents Using N-Gram Indexing</b>	
<p style="text-align: center;"><b>Abstract</b></p> <p><b>Objectives/Goals</b> The goal was to develop a novel test for accurately identifying duplicated content in a large text corpus and build an efficient program to detect duplicate content based on the test.</p> <p><b>Methods/Materials</b> The basic algorithm is to break each document in the collection into n-grams (which are consecutive runs of n words), compile an index of these n-grams, and search for clusters of documents that share a large number of n-grams. Then, a second, smaller paragraph-level index is assembled from these clusters to ascertain the proximity of shared n-grams. I experimented with different n-gram sizes, document types, and document sizes to maximize the effectiveness of the program.</p> <p>The program was implemented in the Java language using DrJava IDE and JDK 6.0 on an iMac and a MacBookPro.</p> <p><b>Results</b> The program successfully identified duplicate paragraphs in large documents (70 to 150 pages in length). It also pinpointed duplicate content in hundreds of news articles from the web and identified duplicate content within a single, large document through the paragraph indexing analysis. Through my experiments, I discovered that an n-gram size of three provides the best balance between storage space and accuracy.</p> <p><b>Conclusions/Discussion</b> I developed a simple, effective criterion for finding duplicate content in document sets of moderate size, which I implemented into a fast, easy-to-use, accurate stand-alone program that allows the user to check for duplicate content in a group of documents.</p>	
<b>Summary Statement</b> I built a duplicate content detector for large text document corpora using a simple trigram overlap test.	
<b>Help Received</b> My father helped me to learn to use the IDE and find relevant information about Java language constructs on the Internet.	