# Brute Force and Indexed Approaches to Pairwise Document Similarity Comparisons with MapReduce

## Jimmy Lin

Dandy Fenz
Thomas Berger

December, 15th 2009

# First results



- Given situation:
  - Database with wikipedia abstracts
  - Task: Show „more-like-this" articles

- Implementation of following tasks:
  - Preparation of input data ✓
  - Brute Force algorithm ✓
  - Inverted index creation ✓
  - Postings Cartesian Product (PCP) algorithm ✓
  - Prove PCP efficiency
  - PCP approxima

# Our collection: frequency of terms

# Brute Force (BF)

$$\sum_{t \in V} w_{t,d_i} \cdot w_{t,d_j}$$

# Postings Cartesian Product (PCP)



$$w_{t,d_i} \cdot w_{t,d_j}$$

# Comments on Hadoop

- Framework relatively easy to configure and understand

- Documentation slightly outdated
  - Plenty of example code now marked as deprecated
  - Hints for current version rare

- Complex data types implementations
  - Unfortunately not provided by Hadoop Framework
  - Found useful implementation on author's website:
    - http://www.umiacs.umd.edu/~jimmylin/cloud9/docs/index.html

# Next steps

- **General optimizations**
  - Efficient usage of data structures

- **Use approximations**
  - Implementation of proposed approx. for PCP
  - Research for alternative aproaches

- **Generation of statistics**
  - Compare efficiency of different algorithms
  - Analysis of results
    - ◇ Is comparison of wikipedia abstracts reasonable?