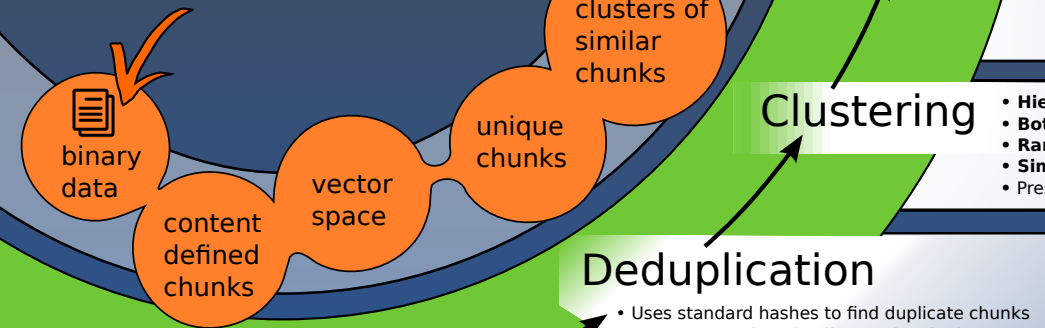




Czech Technical University in Prague  
Faculty of Information Technology  
Department of Theoretical Computer Science

# Incremental Clustering-Based Compression

Student: Luboš Krčál  
Supervisor: Jan Holub



## Chunking

- Binary data is split into **content-defined chunks**
- Resistant to character inserts or shifts
- Uses **partial pattern matches** to mark delimiters

Input data:  
A: rRLUt gBFzKS9MUCyOVQj I I SSu0bFbxQqYj 9cE  
B: rRLUt gBFzKS9MUCyOVQj I I SSu0bFbxQq669cE  
C: rRLUt gBFzKS9MUCyOVQj I I SSu0bFbxQqYj 9cE

Fixed size (blocks):

A: rRLUt gBFzKS9MUCyOVQj I I SSu0bFbxQqYj 9cE  
B: rRLUt gBFzKS9MUCyOVQj I I SSu0bFbxQq666 9cE  
C: rRLUt gBFzKS9MUCyOVQj I I SSu0bFbxQqYj 9cE

Inserting single character shifts all the blocks

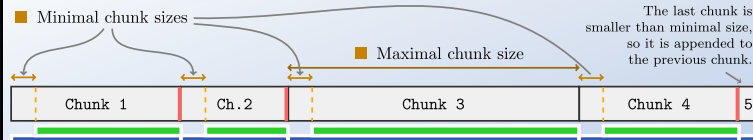
Content-defined delimiters are resistant to this

Variable size (chunks):

A: rRLUt gBFzKS9MUCyOVQj I I SSu0bFbxQqYj 9cE  
B: rRLUt gBFzKS9MUCyOVQj I I SSu0bFbxQq66 9cE  
C: rRLUt gBFzKS9MUCyOVQj I I SSu0bFbxQqYj 9cE  
D: rRLUt gBFzKS9MUCyOVQj I I SSu0bFbxQqYj 9cE

## Similarity Hashing

- **Features:** **n-gram**, compression features, etc. Every feature is hashed
- Standard simhash: merges feature hashes based on majority value per hash position (vector over  $Z_2$ , same length as feature hashes)
- **Extended simhash:** feature hashes are combined together into a vector and normalized (vector over  $Z_m$  of arbitrary length)
- Vector norms are used as distance measures



Example of binary simhash for the string "SWISS MISS".

Note that the length of the simhash determines the dimension of the vector space and the width determines the cardinality of the space.

Individual 3-gram hashes are computed. These are then combined into a vector.

hash("ISS") = 11101000110110100000101111011011  
hash("SWI") = 10000101000111111001001100010111  
hash("WIS") = 11000011101000011100100100110110  
hash("SS ") = 10000101000111111001001100010111  
hash("S M") = 10000101000111111001001100010111  
hash("S MI") = 000010111101110010011101000111  
hash("MIS") = 01011000101000111011110000010111  
simhash("SWISS MISS") = 10000001000111111000101100010111

## Archiver

- It is possible to store more information about the deduplication, simhashing, clustering and grouping - the **archivedata**
- This allows for convenient **CRUD** operations over files in the archive

## Compression

- Compression groups are **compressed using standard compressors**, such as Deflate, BZip2, etc.
- **Metadata** is then added, representing all the chunks and compression groups

## Grouping & Reordering

- Clustering is gradually processed and **compression groups** are determined
- Compression groups are represented with a single cluster

## Clustering

- **Hierarchical clustering** - binary tree, leaves represent chunks
- **Bottom-up** - new node inserted from the bottom
- **Randomized KD-trees** used to find the nearest neighbor
- **Simhash calculated for parent nodes** - as linear combination of children
- Preserves heap property on inter-cluster distances (with simhash distances)

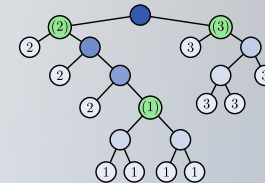
Example of metadata implementation:

```
// magic number
0000: 0000 0123

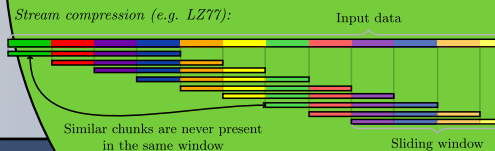
// number of compression
// groups and chunks
0004: 0000 1C28 002D 40AC

// original data size
0014: 0000 0000 06C1 229F

// offsets of the 1. group to
// compressed and orig data
001C: 0003 DBB2 0013 6002
.....
// offset of the 1. chunk
1C3C: 00E3 2421
.....
```



And, the standard compression techniques fail on data with distant similarities?



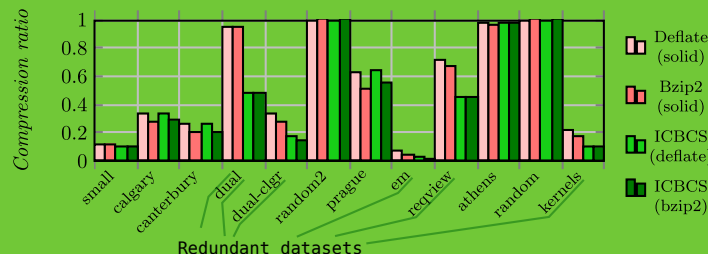
We want to achieve an excellent and efficient compression over a large dataset. With the possibility to read, add, modify and remove data conveniently.

## Results

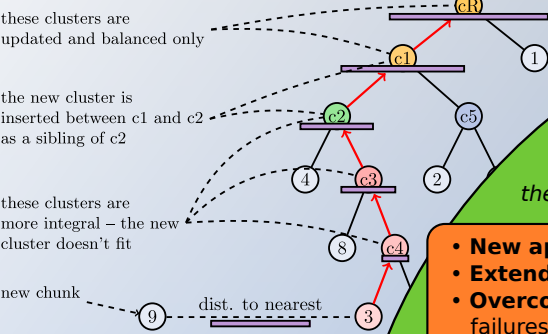
One of the major contributions was to design a compression and archival system based on clustering.

This was successfully completed with the following results:

- **New approach** to compression
- **Extends** current near-deduplication systems
- **Overcomes locality-based** redundancy removal failures of standard compression algorithms
- Has the potential of being an **ultimate archiver**



- Compression ratio significantly improved for very redundant datasets
- For non-redundant dataset, the ratio was almost the same. This is a success, since the files were now compressed in very small batches
- Using weaker compression on appropriately preprocessed data results in both faster and better compression than using a strong compressor as it is



- Optional extra balancing (in scenarios where heap property is not inherent)
- Optional deep distance (use multiple levels of descendants)
- Optional representatives (use only selected descendants)

Disbalanced clustering:

