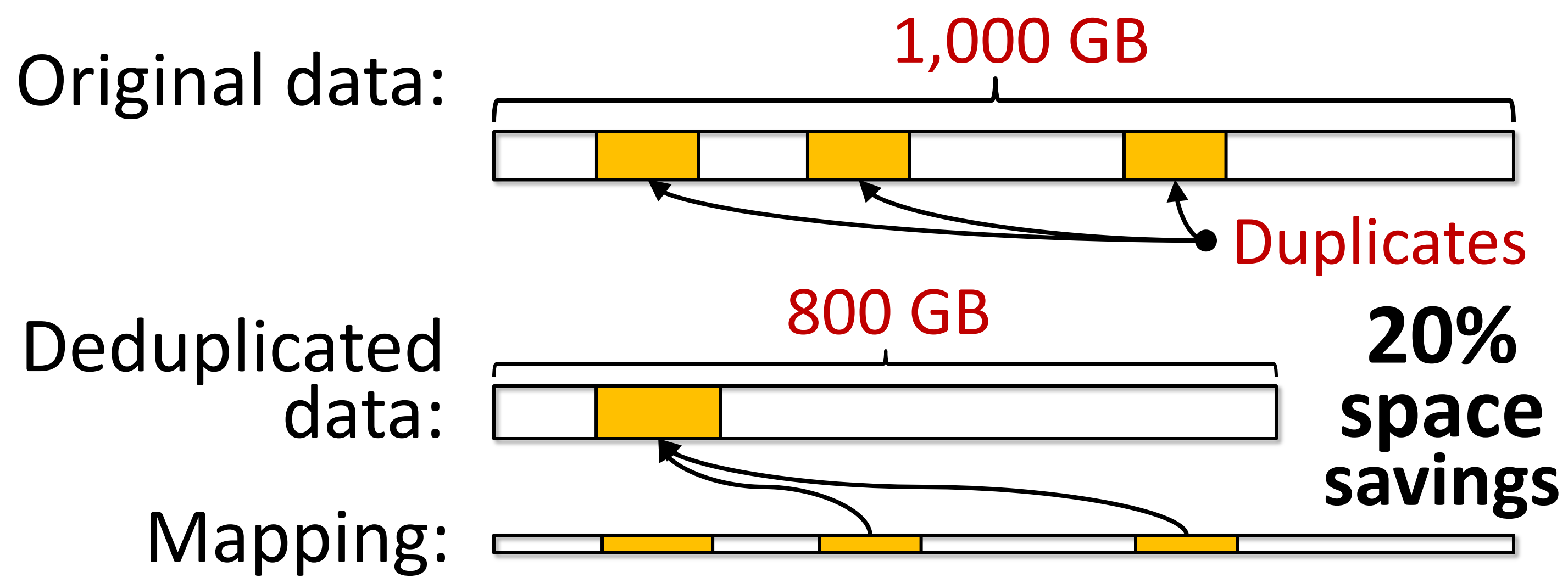


# Generating Realistic Datasets for Deduplication Analysis

Vasily Tarasov<sup>1</sup>, Amar Mudrankit<sup>1</sup>, Will Buik<sup>2</sup>, Philip Shilane<sup>3</sup>, Geoff Kuenning<sup>2</sup>

<sup>1</sup>Stony Brook University, <sup>2</sup>Harvey Mudd College, <sup>3</sup>EMC Corporation

## 1 Deduplication

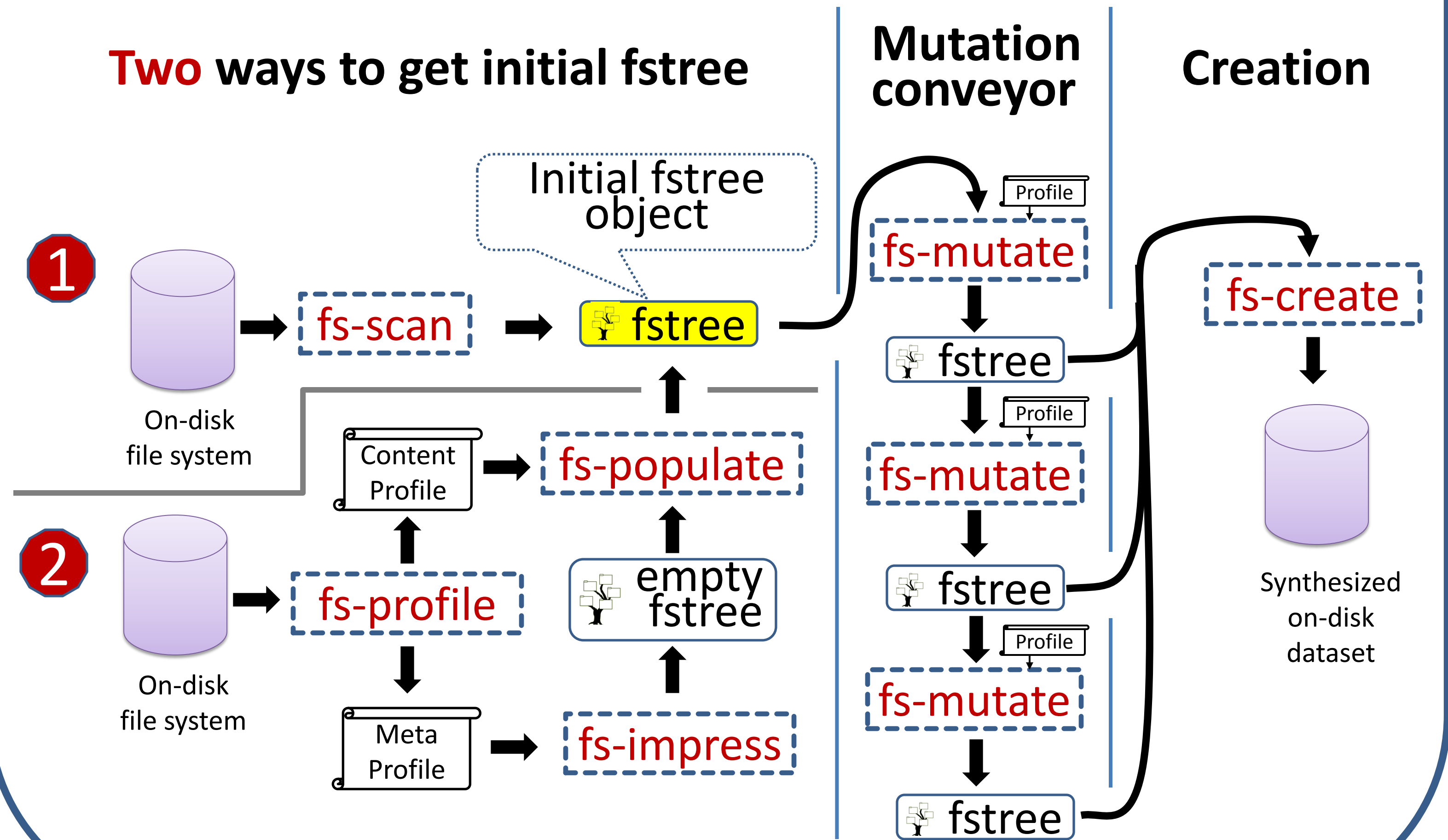


- Typical Steps:
- 1. Chunking**
    - Whole-file, Fixed, Variable size, ...
  - 2. Hashing**
    - SHA256, MD5, Tiger, ...
  - 3. Indexing**
    - Conventional databases
    - Specialized indexes
    - ....

Performance depends on dataset!

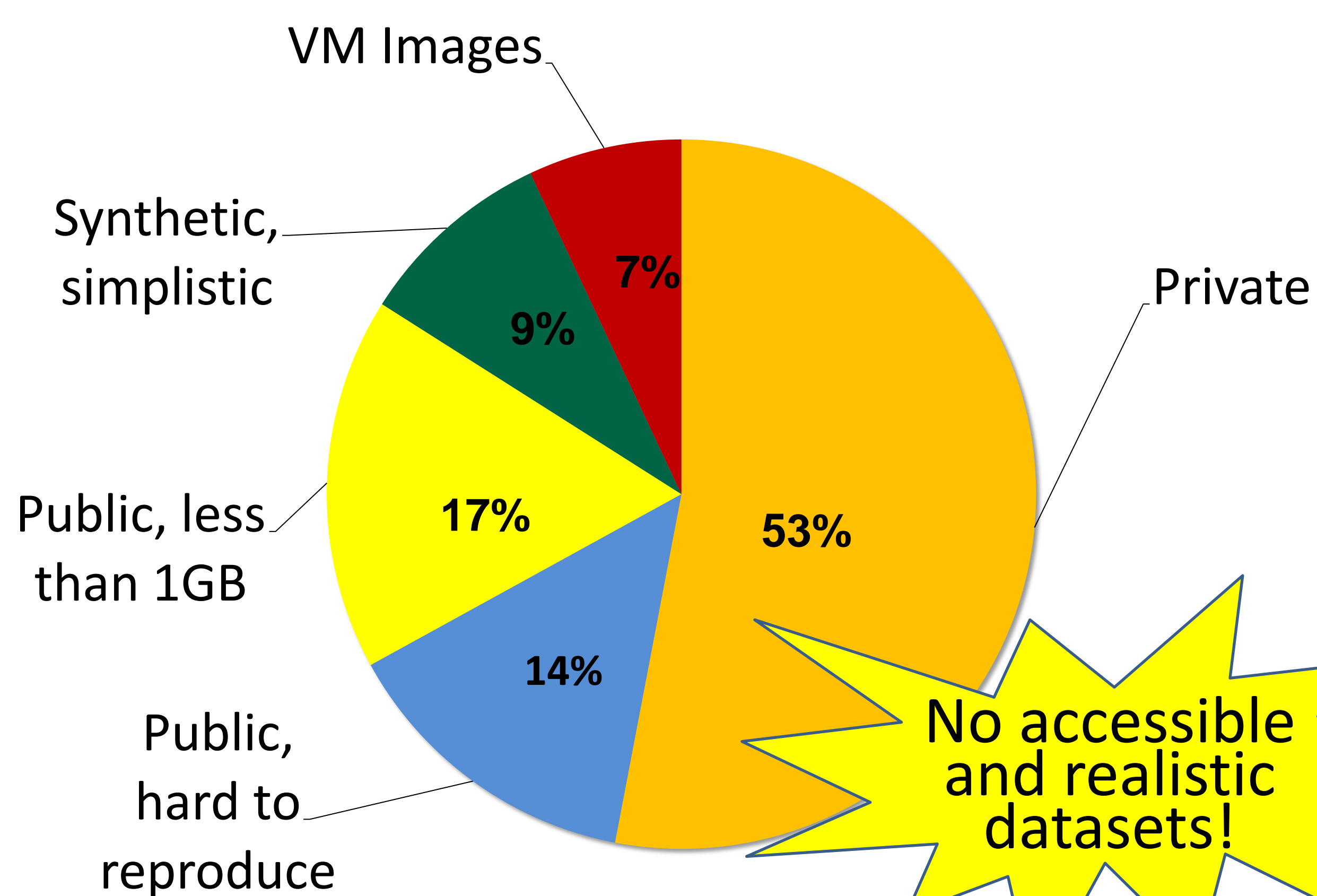
## 4 Action Modules

Modules define actions on in-memory `fstree` object



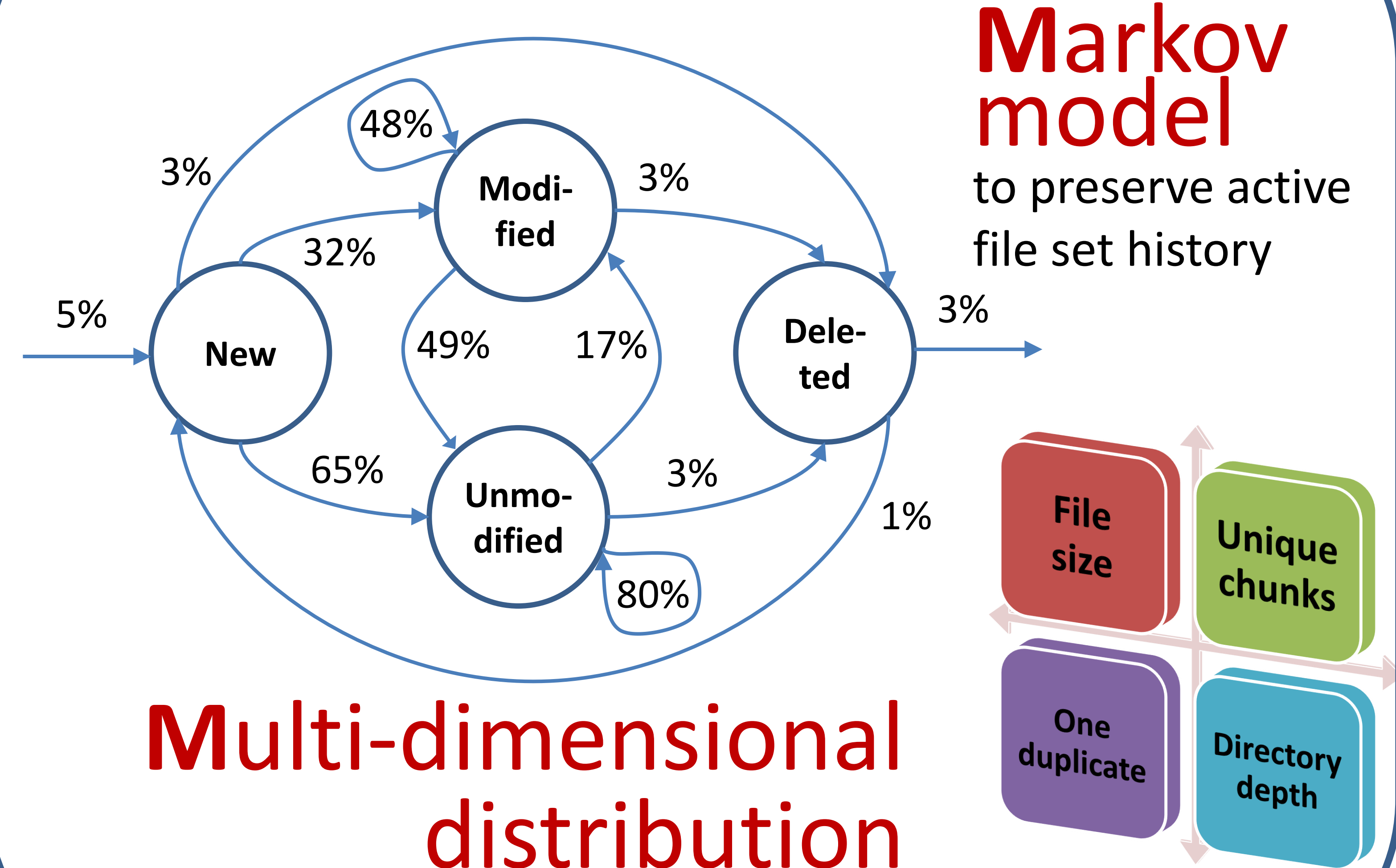
## 2 Previous Datasets

120 datasets, 29 papers from, e.g., FAST, ATC conferences

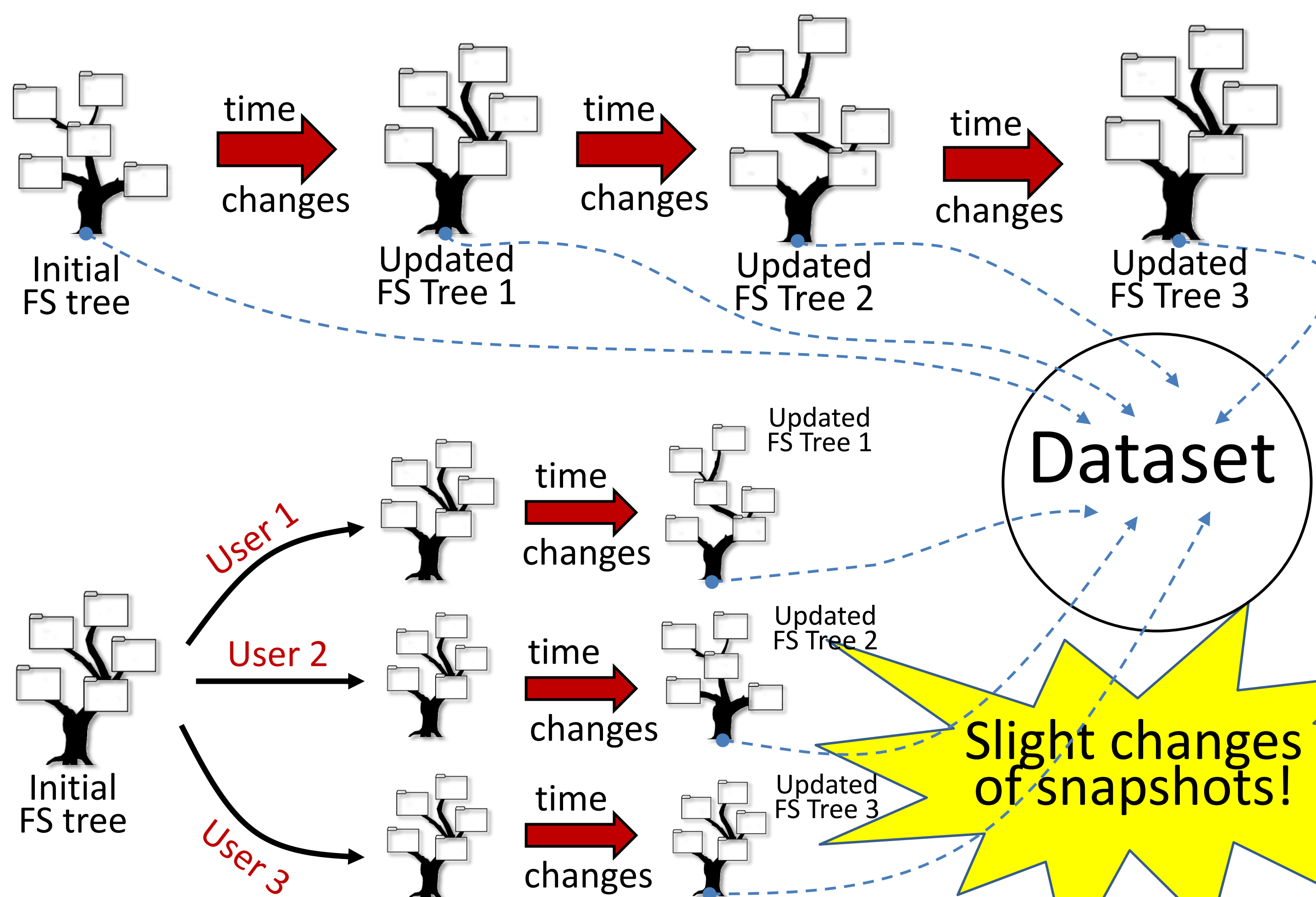


No accessible and realistic datasets!

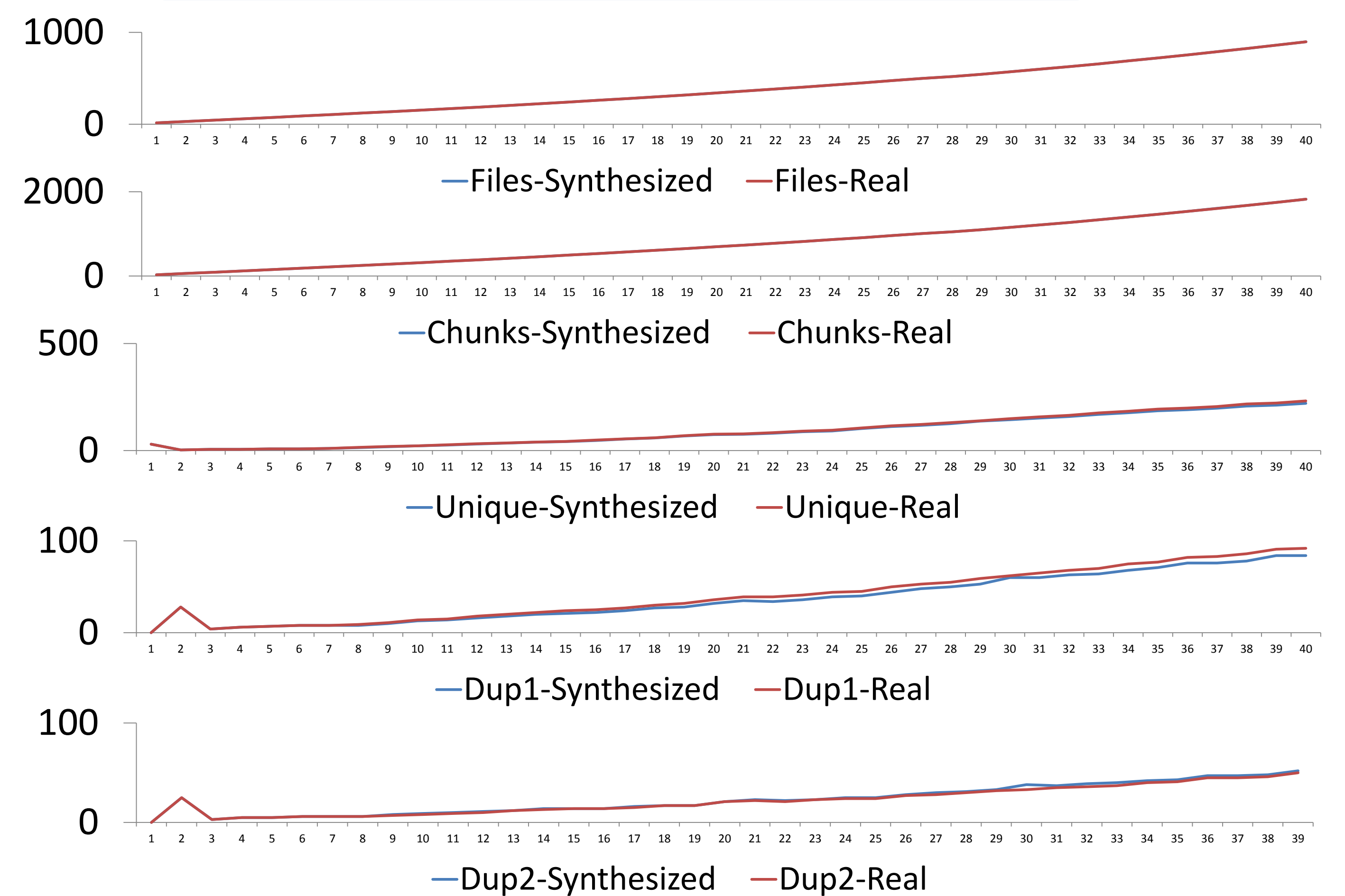
## 5 M&D Model



## 3 Dataset Formation



## 6 Preliminary Results



2.6.0 – 2.6.39 Kernels dataset: within 10% accuracy!