

# Analyzing Data Properties using Statistical Sampling Techniques

Illustrated on Scientific File Formats and Compression Features

Julian M. Kunkel

kunkel@dkrz.de



2016-06-23

# Outline

- 1 Introduction
- 2 Exploring a Subset of Data
- 3 Statistical Sampling
- 4 Summary

# Motivation

- Understanding data characteristics is useful
  - Relation of file types to optimize relevant file formats
  - Conducting what-if analysis
    - Influence of compression, deduplication
    - Performance expectations
- Analysing large quantities of data is time consuming and costly
  - Scanning petabytes of data in  $> 100$  millions of files
  - With 50 PB of data and 5 GiB/s read, 115 node days are needed
  - The complete experiment for this paper would have cost 4000 €

⇒ Working on a representative data set reduces costs
- Conducting analysis on representative data is difficult
  - What data makes up a representative data set?
  - How can we infer knowledge for all data based on the subset?
    - Based on file numbers (i.e. a typical file is like X)
    - Based on capacity (i.e. 10% of storage capacity is like Y)
  - Many studies simply select a data set and claim it is representative

# Contribution

## Goal

- Investigation of statistical sampling to estimate file properties
  - Can we trust the results?
  - What are typical mistakes when sampling data?
- Conduct a simple study to investigate compression and file types

## Approach

- 1 Scanning a fraction of data on DKRZ file systems
  - Analyzing file types, compression ratio and speed
- 2 Investigating characteristics of the data set
- 3 Statistical simulation of sampling approaches
  - We assume the population (full data set) is the scanned subset
- 4 Discussion of the estimation error for several approaches

## 1 Introduction

## 2 Exploring a Subset of Data

- Sampling Approach
- Processing of Files
- Distribution of File Sizes
- Scientific File Formats
- Compression Ratio
- Compression Speed
- Differences Between Projects

## 3 Statistical Sampling

## 4 Summary

# Sampling of the Test Data

- DKRZ usage: 320 million files in 12 PB, 270 project dirs
- Scan of user accessible data (scan is done by a regular user)
  - Accessible data: 58 million files, 160 project dirs
- Scanned files: 380 k files (0.12%) in 53.1 TiB (0.44%) capacity
  - Discrepancy since home directories contain very small files

## Scanning Process

- 1 Run a find for each project directory, store it in a file
- 2 Select up to 10 k files from each project randomly (scan list)
- 3 Permutate the scan list
- 4 Partition the scan list into chunks (file lists)
- 5 Run multiple processes concurrently, each working on a file list
- 6 Terminate the processes after a couple of days

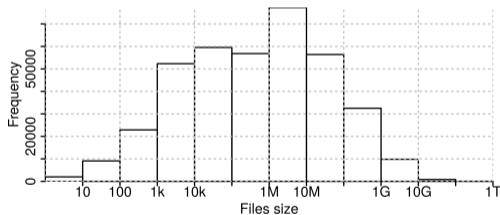
*As we will see this approach is not optimal for analyzing by capacity*

# Processing of Files

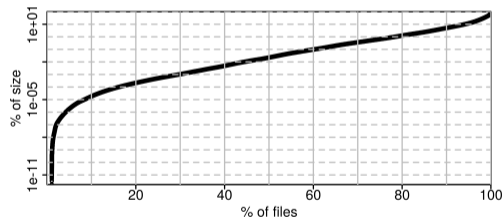
- Processing is implemented as simple shell script
  - Starter script spawns multiple scripts each working on a file list
- The process reads the file list and current progress
  - Allows to restart processes and continue processing
- Process ignores non-existing files
  - Some files are deleted after the scan
- Copy the file
  - (Best alternative, copy data into memory)
  - You would not believe how many files changed on the fly
  - Running different tools on the data falsifies results
- Run tools to investigate properties:
  - `file` to identify file type (based on magic), suboptimal
  - CDO to identify scientific file format (high accuracy)
  - Compression tools: LZMA, GZIP, BZIP2, ZIP

# Distribution of File Sizes

- For now, we analyze all scanned files!
- File size follows a heavy tailed distribution
- 90% of files consume roughly 10% capacity



(a) Histogram (logarithmic x-axis)

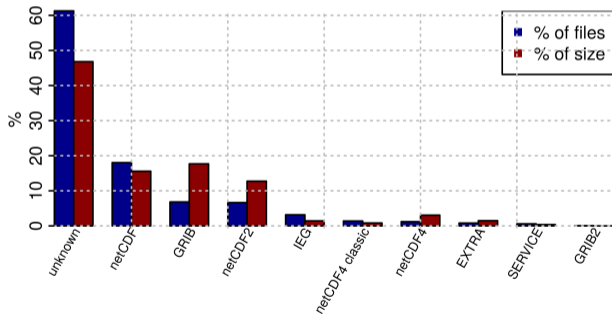


(b) Cumulative file sizes (y-axis in log scale)

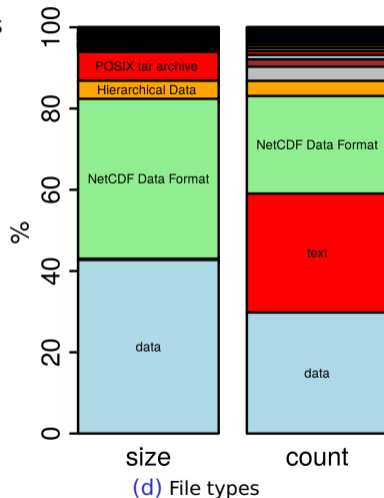


# Scientific File Formats

- The computation by file count and capacity differs
  - The heavy-tailed distribution skews analysis
- file often determines wrong formats



(c) CDO types

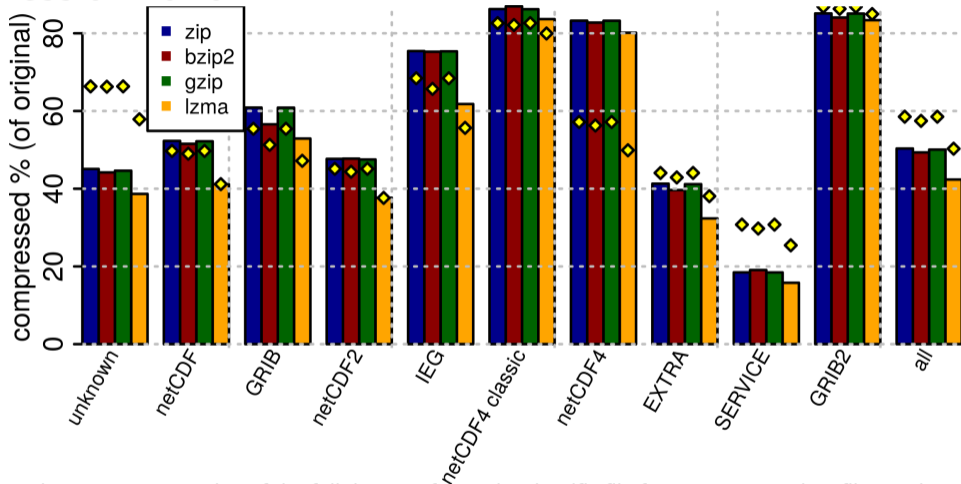


size

count

(d) File types

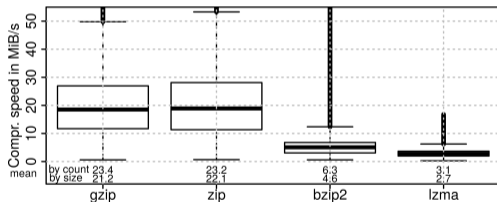
# Compression Ratio



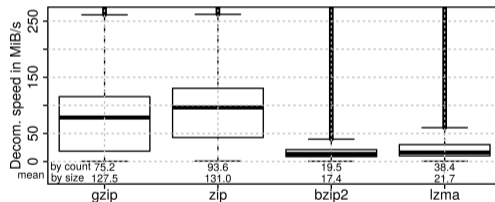
Arithmetic mean compression of the full data set for each scientific file format computed on file number. The column "all" shows the mean values for the whole data set. Yellow diamonds show compress % computed by file size

# Compression Speed

- Measured user-time for the execution of each tool
  - Ignored I/O
- Again difference between compression by size and count



(a) Compression

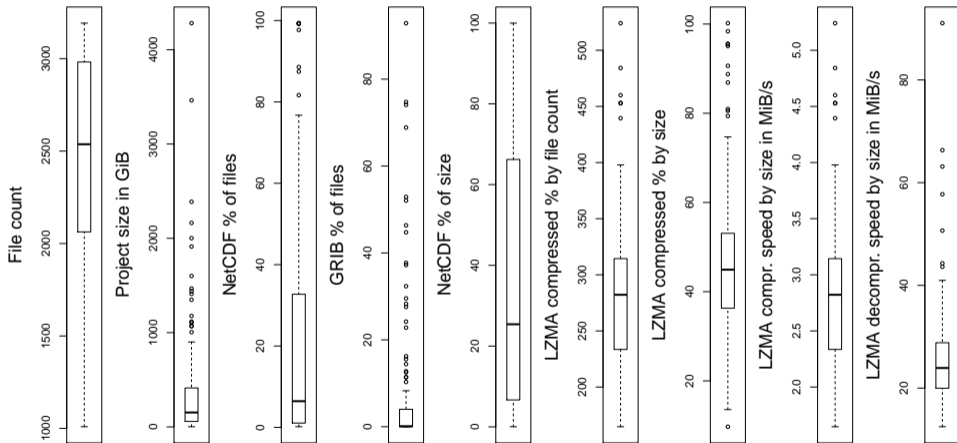


(b) Decompression

Boxplots showing compression/decompression speed per file, mean shown under the plot

# Differences Between Projects

Properties vary significantly → proper sampling requires to pick data from all projects



Analyzing 125 individual projects, each point represents the arithmetic mean value of one

- 1 Introduction
- 2 Exploring a Subset of Data
- 3 Statistical Sampling**
  - Overview
  - Demonstration of the Strategies
- 4 Summary

# Statistical Sampling

- Can we determine the error when analyzing only a fraction of data?
  - We simulate sampling by drawing samples from the totally analyzed files
- Statistics offers methods to determine confidence interval and sample size
- We analyze random variables for quantities that are continuous or proportions
- Proportions: fraction of samples for which a property holds

## Sample size and confidence intervals

- For proportions Cochran's sample size formula estimates sample size
  - (Similar number) works for extremely large population sizes
  - Error bound  $\pm 5\%$  requires 400 samples (95% confidence)
  - Error bound  $\pm 1\%$  requires 10,000 samples
- For continuous variables
  - Models require to know the distribution of the value
  - A-priori unknown, usually not Gaussian, difficult to apply → out-of-scope (here)
  - Nevertheless, we will demonstrate convergence

# Sampling Strategies

## Sampling to Compute by File Count

- 1 Enumerate all files
- 2 Create a simple random sample
  - Select a random number of files to analyze without replacement
  - For proportional variables, the number of files can be computed with Cochran's formula
  - You can use simulation to estimate the error for contiguous variables
    - You may increase sampling size if the accuracy does not suffice

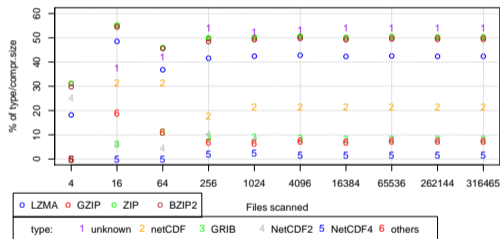
## Sampling to Compute by File Size

- 1 Enumerate all files AND determine their file size
- 2 Pick a random sample based on the probability  $\frac{filesize}{totalsize}$  with replacement
  - Large files are more likely to be chosen (even multiple times)
- 3 Create a list of unique file names and analyze them
- 4 Compute the arithmetic mean for the variables
  - If a file has been picked multiple times in Step 2., its value is used multiple times

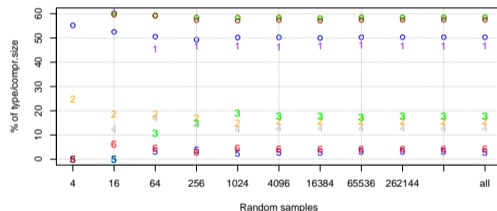
# Demonstration of the Strategies

- Apply the approach with an increasing number of samples
  - Compare true value with the estimated value

## Running one simulation for increasing sample counts



(a) Compute mean by count



(b) Compute mean by size

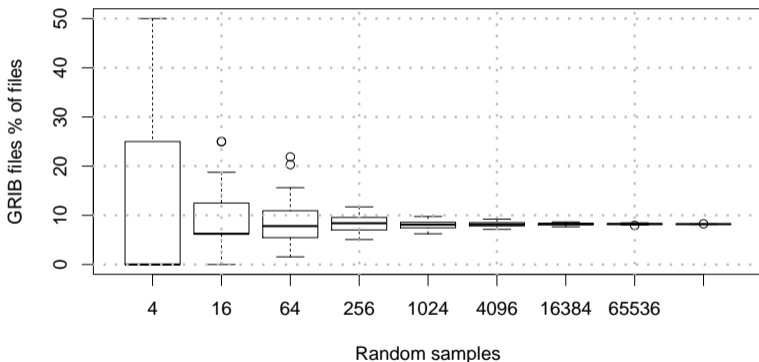
Evaluating various metrics (proportions) for an increasing number of samples

- This suggests that the results converge quickly but how trustworthy is one run?



# Investigating Robustness: Computing by File Count

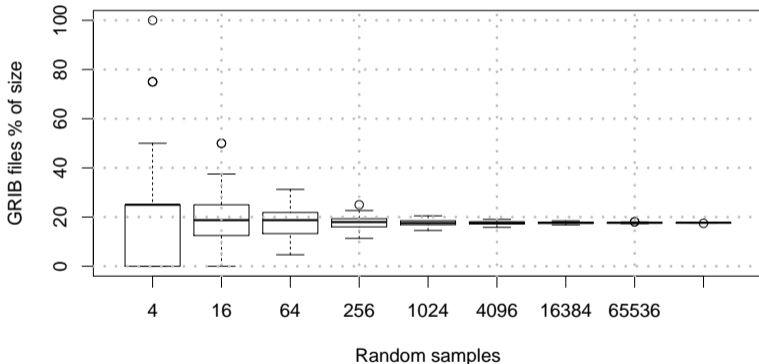
- Running the simulation 100 times to understand the variance of the estimate
- Clear convergence: thanks to Cochran's formula the total file count is irrelevant



Simulation of sampling by file count to compute compr.% by file count

# Investigating Robustness: Computing by File Size

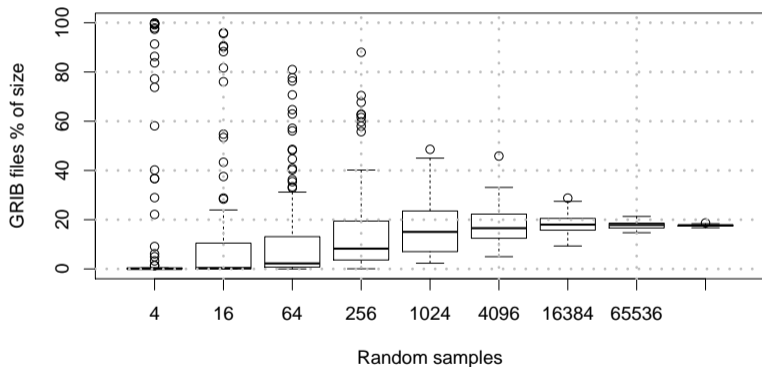
- Using the correct sampling by weighting probability with file size



Simulation of sampling to compute proportions of types by size

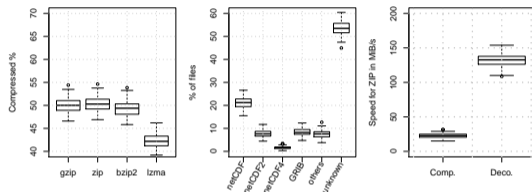
# Investigating Robustness: Computing by File Size

- Using the WRONG sampling by just picking a simple random sample
- Almost no convergence behavior; you may pick a file with 99% file size at the end

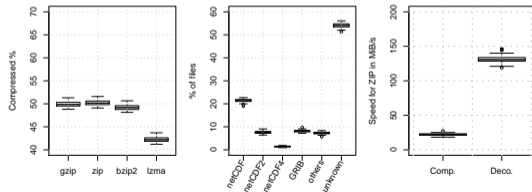


Simulation of sampling to compute proportions of types by size

# Additional Characteristics Computed by File Count

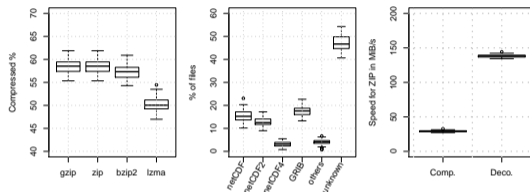


(a) Sampling 316 = 0.1% of files)

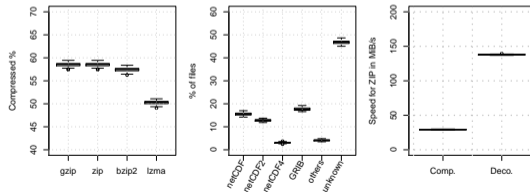


(b) Sampling 3164 = 1% of files)

# Additional Characteristics Computed by File Size



(c) Drawing 256 samples with replacement



(d) Drawing 4096 samples with replacement

# Summary

- We investigated statistical sampling to estimate data characteristics for a system
- The approach is demonstrated for analyzing scientific file formats and compression
- Several sources of error have been discussed
- Estimation of values that should be computed by file size requires proper sampling
- Statistic simulation helps to understand the error when analyzing continuous vars