

Compression IOI

by Uğur Çayoğlu

13. February 2020

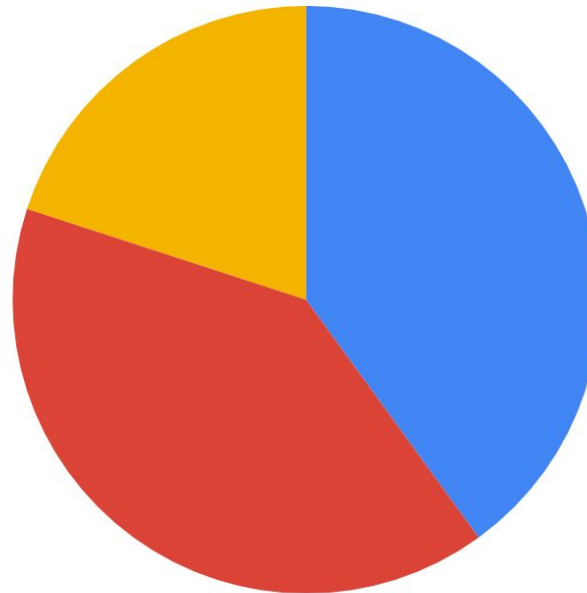
STEINBUCH CENTRE FOR COMPUTING (SCC) &
INSTITUTE FOR METEOROLOGY AND CLIMATE RESEARCH (IMK-ASF)

Me, myself and I



- Studied Information Management and Engineering @ KIT

● Economics
● Computer Science
● Law



Me, myself and I



- Studied Information Management and Engineering @ KIT
- Bachelor thesis:
 - *Improve Collaboration Between **Marketing and Sales** and Analyze the Impact on the Quality of **Customer Data** in the Enterprise*



Me, myself and I



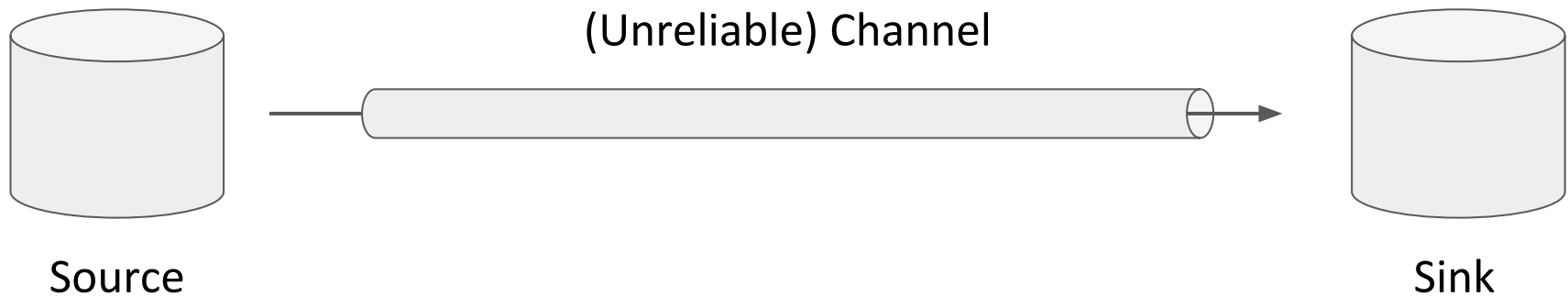
- Master thesis:
 - *Conceptualization and Prototypical Implementation of a Tool to Determine the Optimal **XML Similarity Measurement***
- Doctoral thesis:
 - ***Compression Methods** for Structured Floating Point Data and their Application in Climate Research*



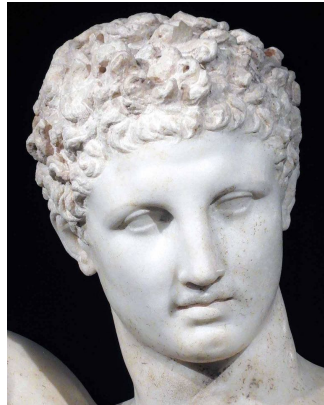
History



- Compression has its origin in communication technology
- Hence compression == source coding

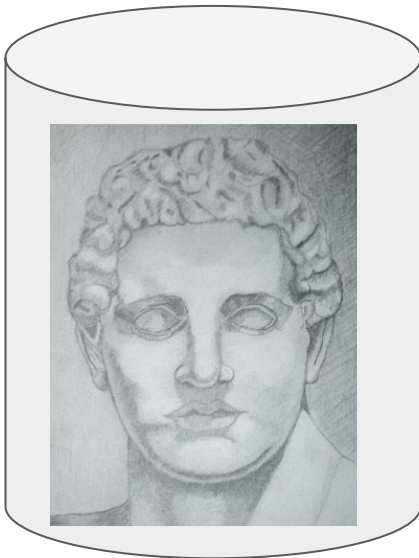
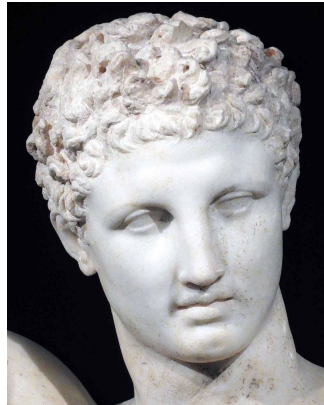


Compression IOI



Source [A]

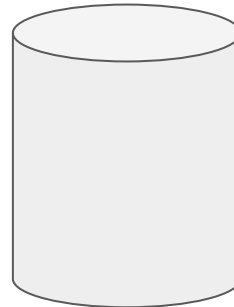
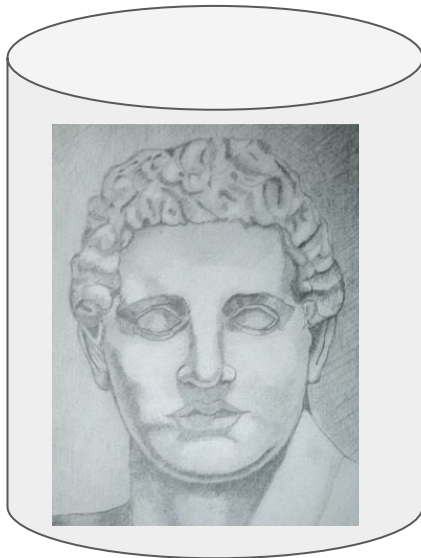
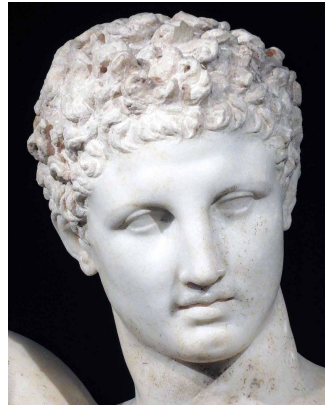
Compression IOI



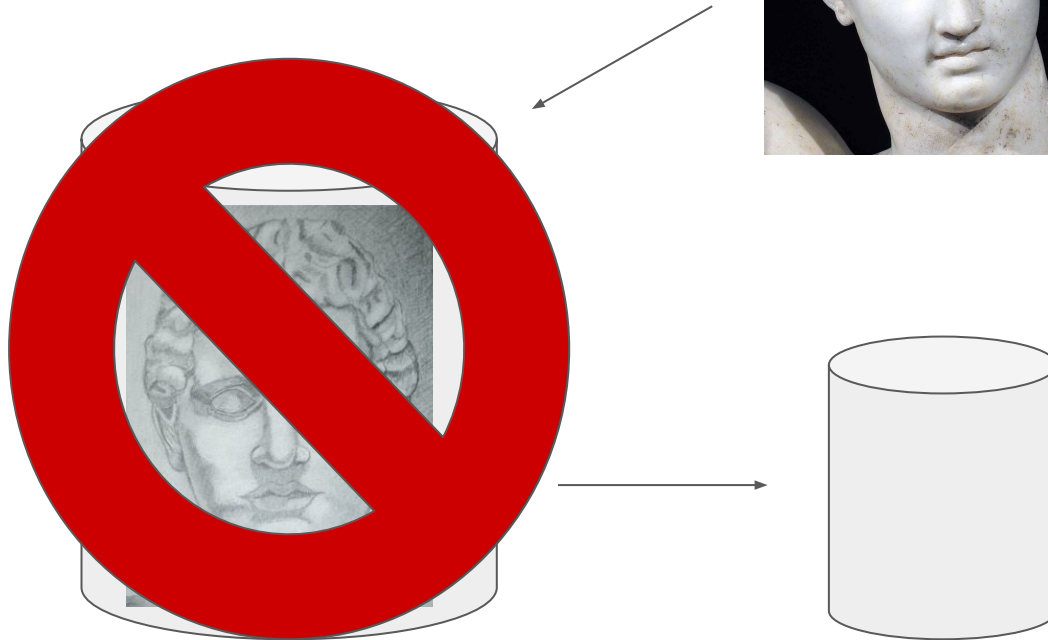
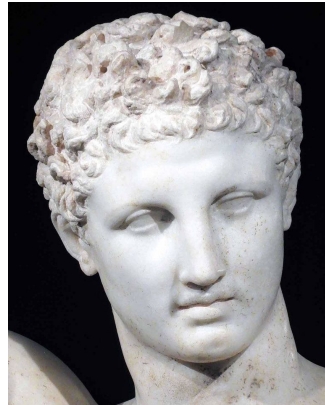
KIT
Karlsruhe Institute of Technology

Source [B]

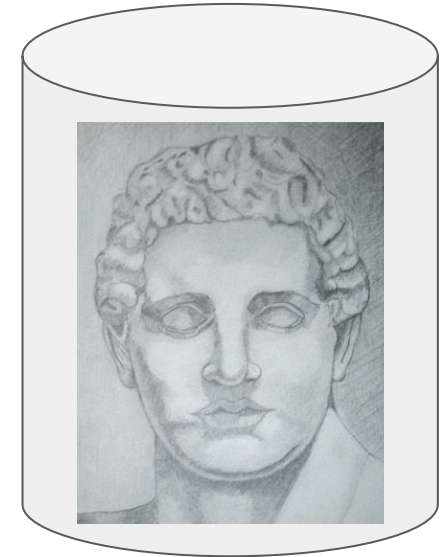
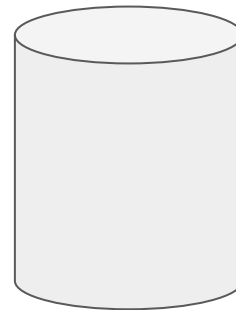
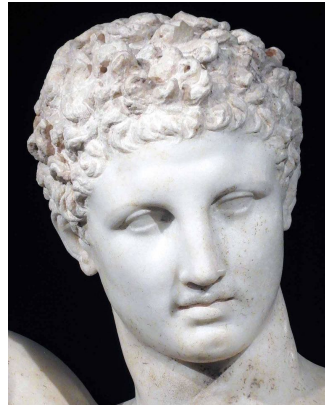
Compression IOI



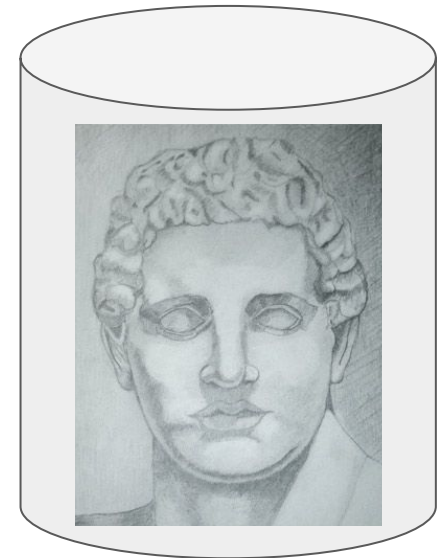
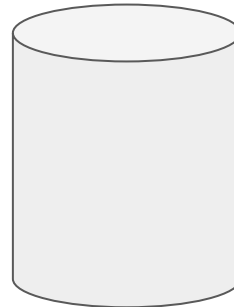
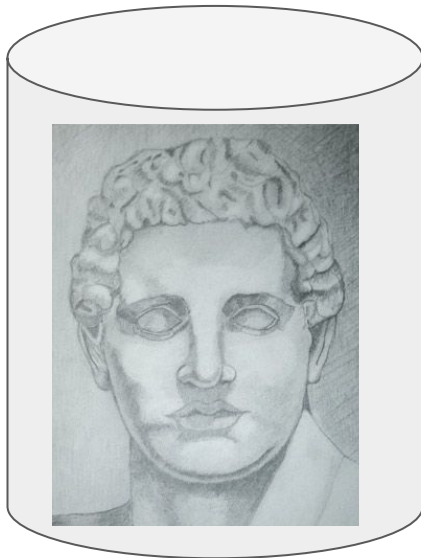
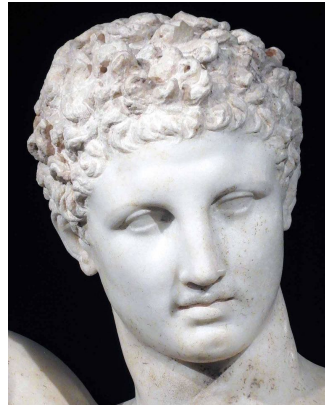
Compression IOI



Compression IOI

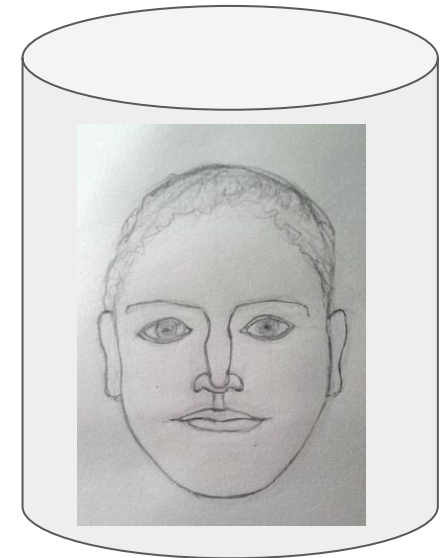
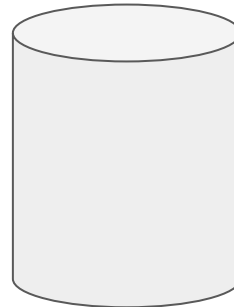
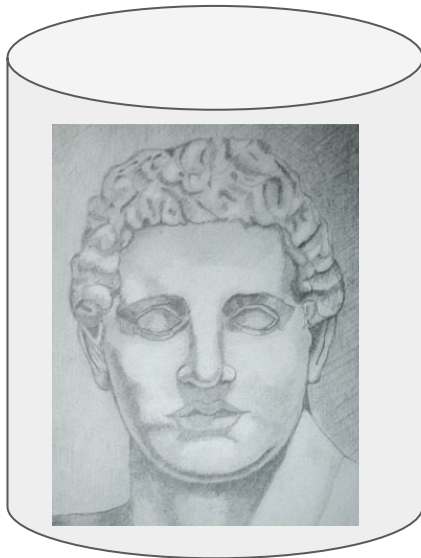
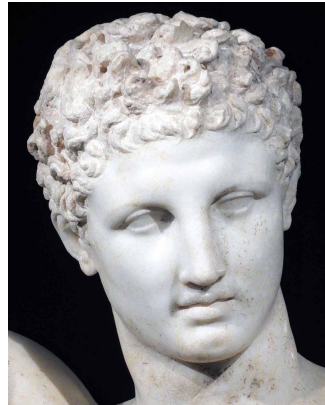


Compression IOI



Lossless compression

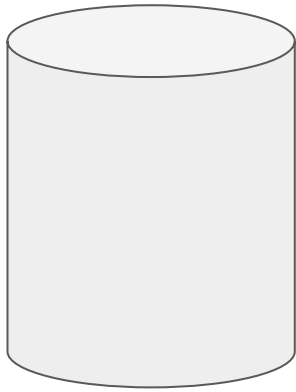
Compression IOI



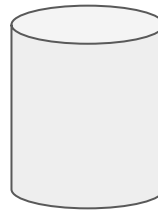
Lossy compression

Source [B]

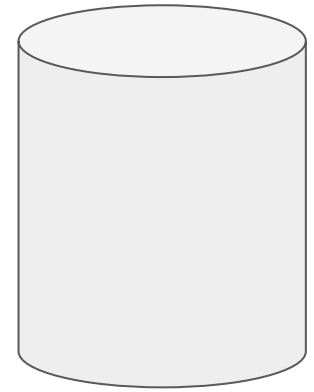
Nomenclature



Original data



Encoded data

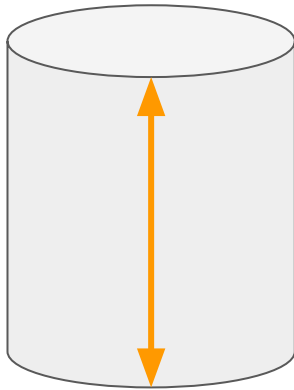


Decoded data

Nomenclature

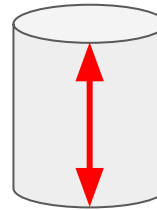


Size of original data

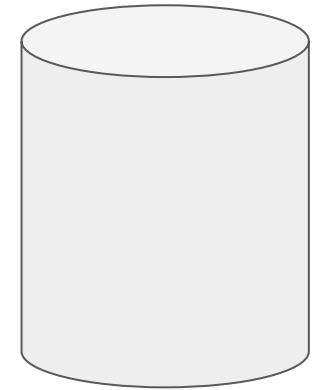


Original data

Size of encoded data



Encoded data

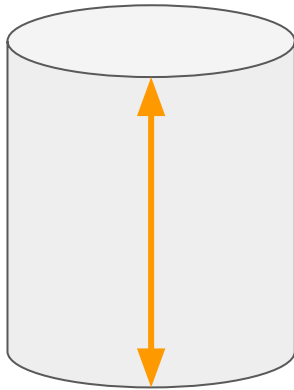


Decoded data

Nomenclature

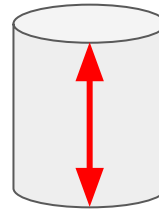


Size of original data

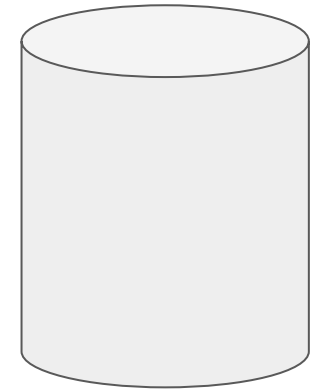


Original data

Size of encoded data



Encoded data



Decoded data

Size of encoded data

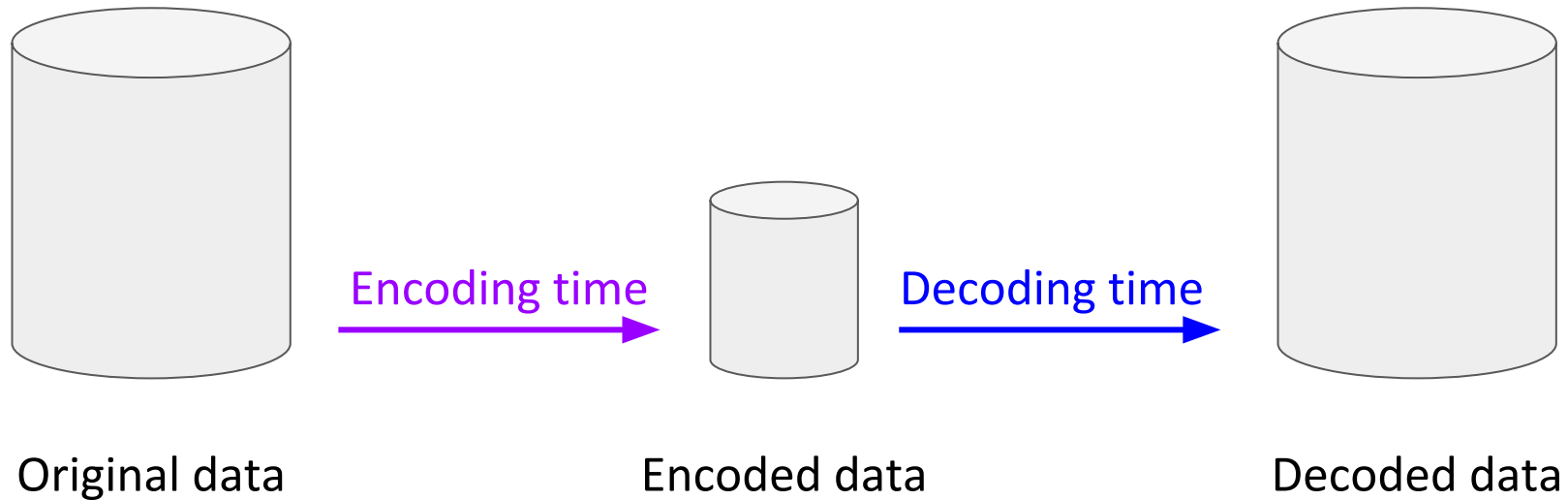


Size of original data



Compression Factor
[no dimension]

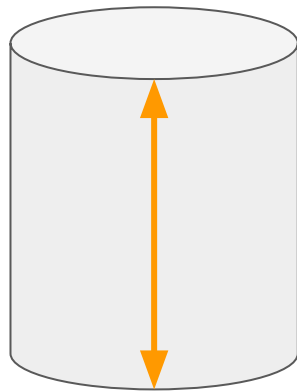
Nomenclature



Nomenclature



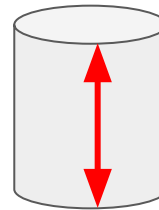
Size of original data



Original data

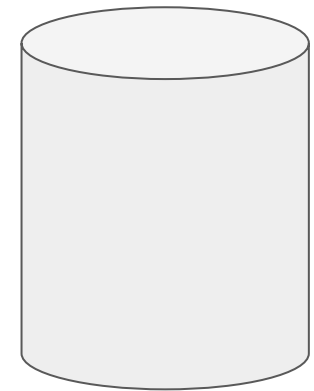
Encoding time

Size of encoded data



Encoded data

Decoding time



Decoded data

Encoding time



Size of original data



**Encoding
Throughput
[Bytes/sec]**

Decoding time



Size of encoded data

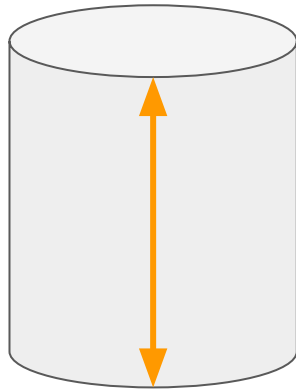


**Decoding
Throughput
[Bytes/sec]**

Nomenclature

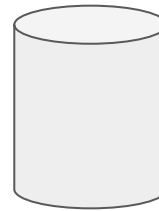


Size of original data

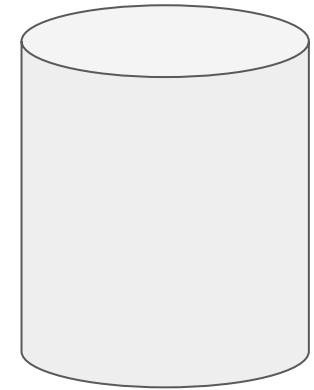


Original data

Encoding time
→



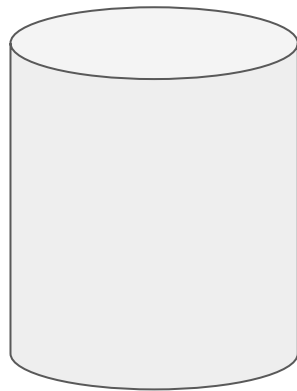
Encoded data



Decoded data

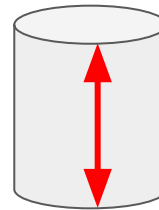


Nomenclature



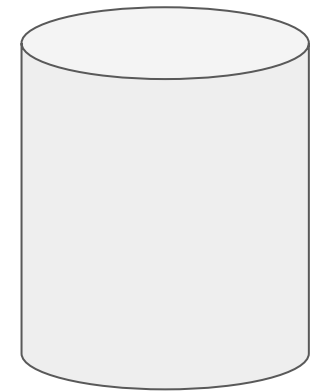
Original data

Size of encoded data



Encoded data

Decoding time
→



Decoded data

Decoding time
→

Size of encoded data

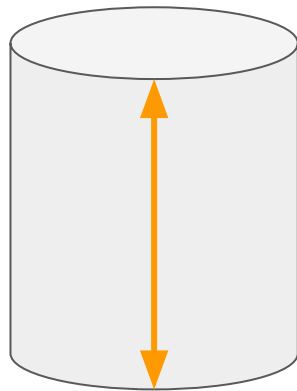


**Decoding
Throughput
[Bytes/sec]**

Nomenclature



Size of original data

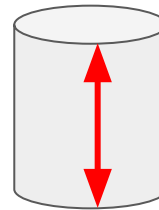


Original data

Encoding time

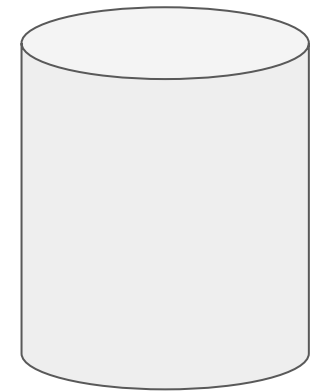


Size of encoded data



Encoded data

Decoding time

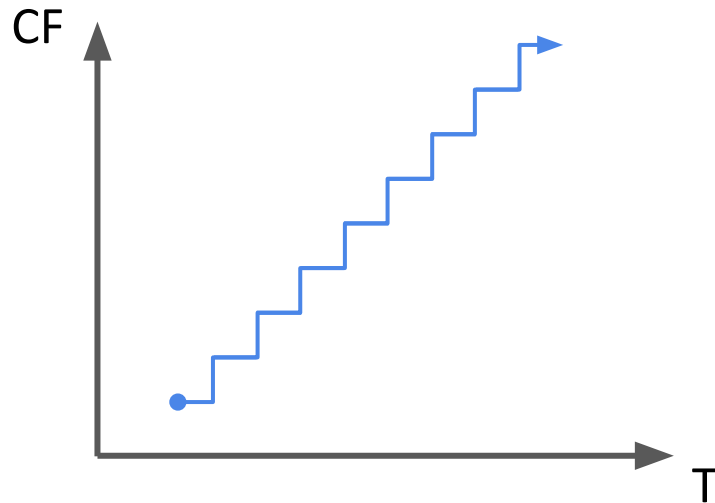


Decoded data

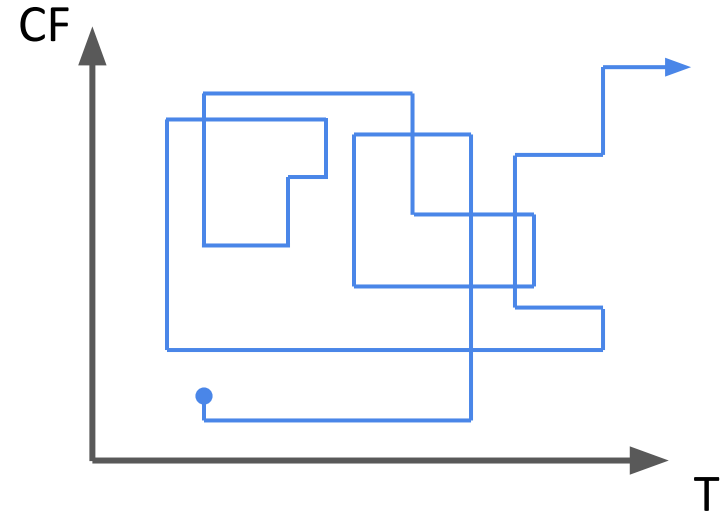
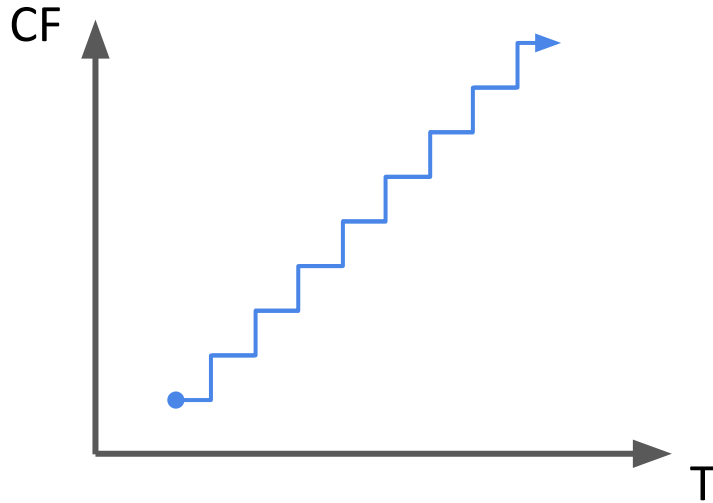
Compression Factor (CF)
[no dimension]

Throughput (T)
[Bytes/sec]

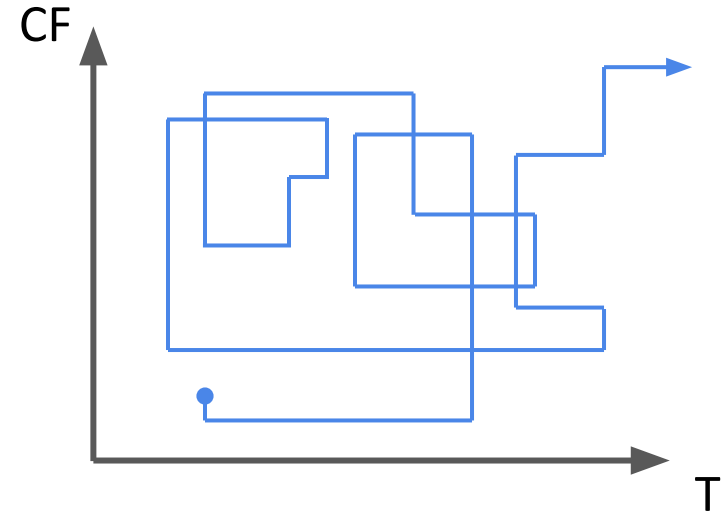
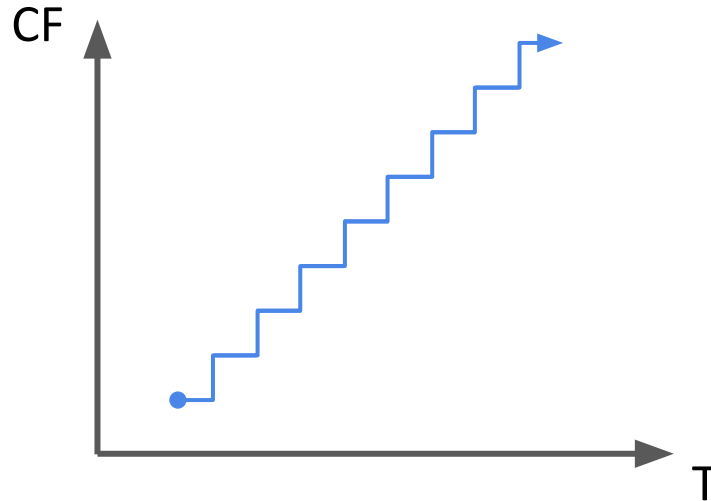
Iterative Process



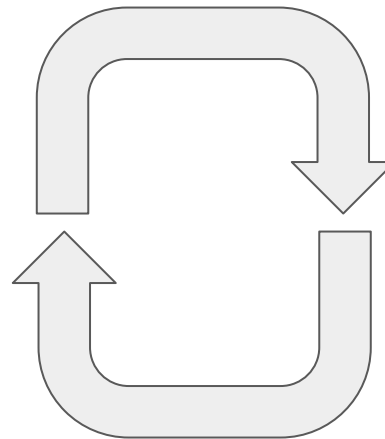
Iterative Process



Iterative Process



Compression Factor (CF)
[no dimension]



Throughput (T)
[Bytes/sec]

Timeline



Past

Compression Factor

Understanding the information and intrinsic structure of the data

Custom scientific field

Present

Throughput

Research of data structures and algorithms

Applied computer science

Future

Compression Factor

AI/ML methods for understanding of interdependencies

Custom + related scientific field

Far Future

Throughput

...

...

Timeline



Past

Present

Future

Far Future

Compression Factor

Throughput

Compression Factor

Throughput

Understanding the
information and
intrinsic structure
of the data

Research of data
structures and
algorithms

AI/ML methods for
understanding of
interdependencies

...

Domain
knowledge

Applied computer
science

Custom + related
scientific field

...

Timeline

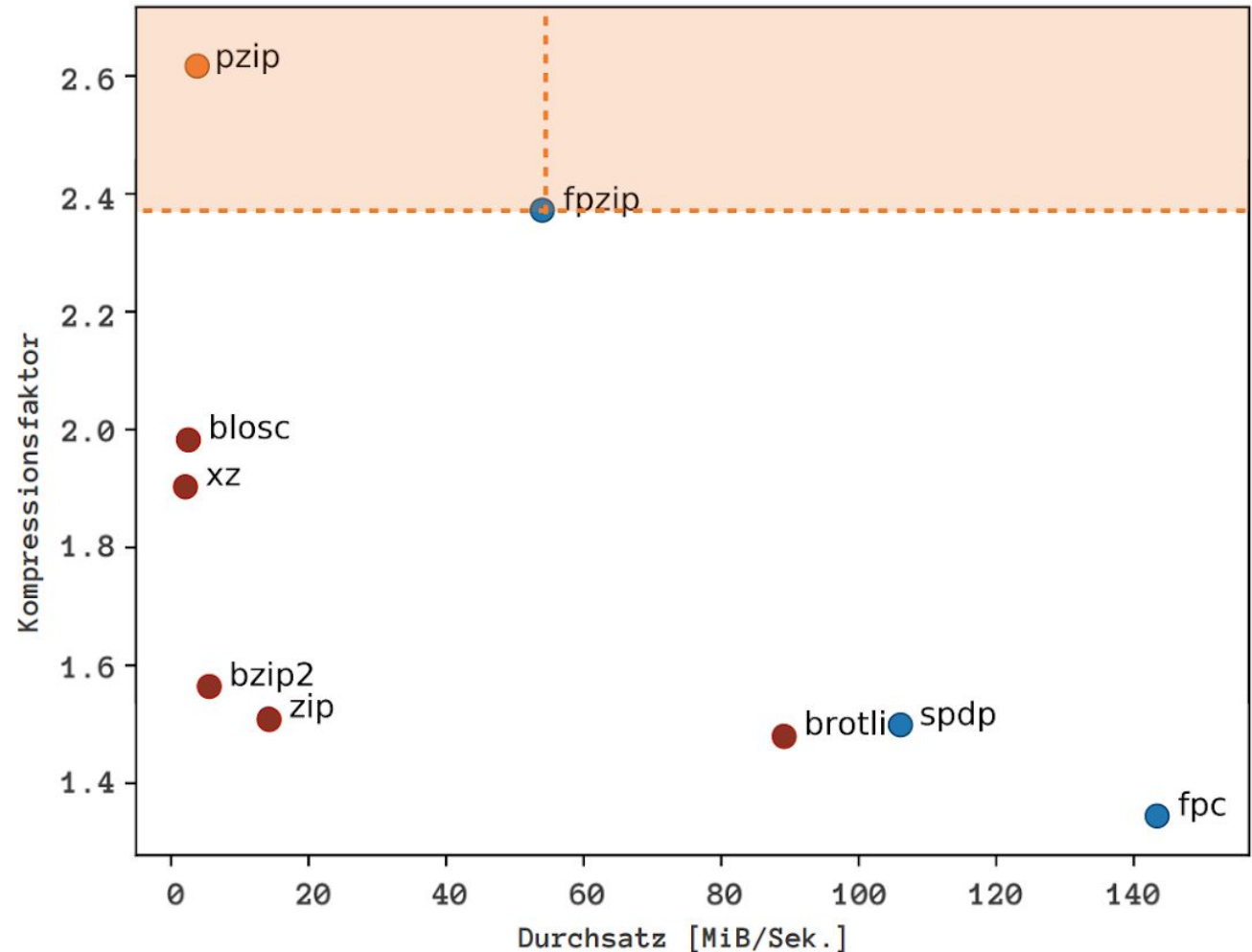


Past

Compression Factor

Understanding the
information and
intrinsic structure
of the data

Domain
knowledge



Timeline



Past

Present

Future

Far Future

Compression Factor

Throughput

Compression Factor

Throughput

Understanding the information and intrinsic structure of the data

Research of data structures and algorithms

AI/ML methods for understanding of interdependencies

...

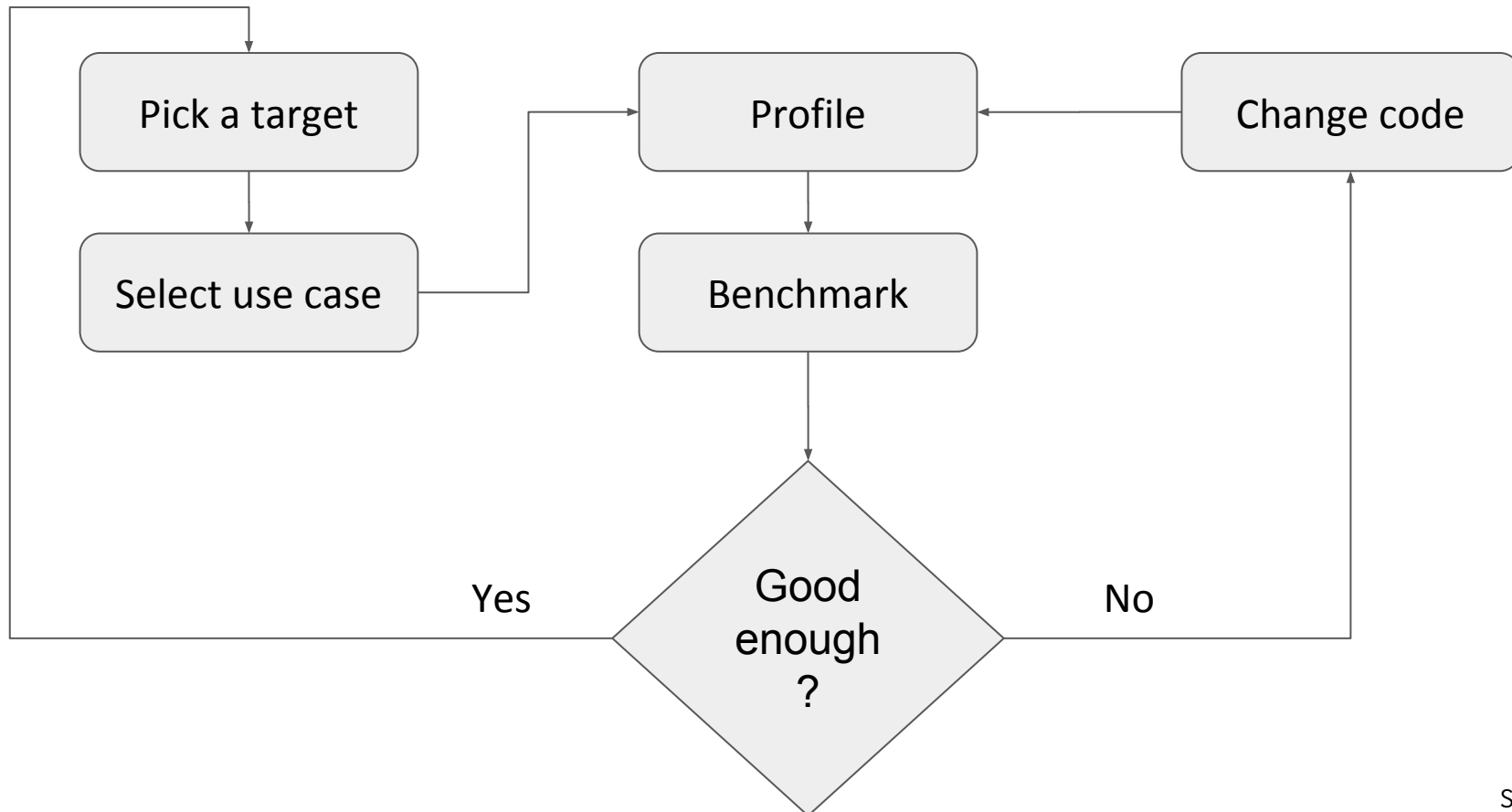
Domain knowledge

Applied computer science

Custom + related scientific field

...

Workflow for optimizing throughput or compression factor



Source [C]

Profiling and Benchmarking in this context

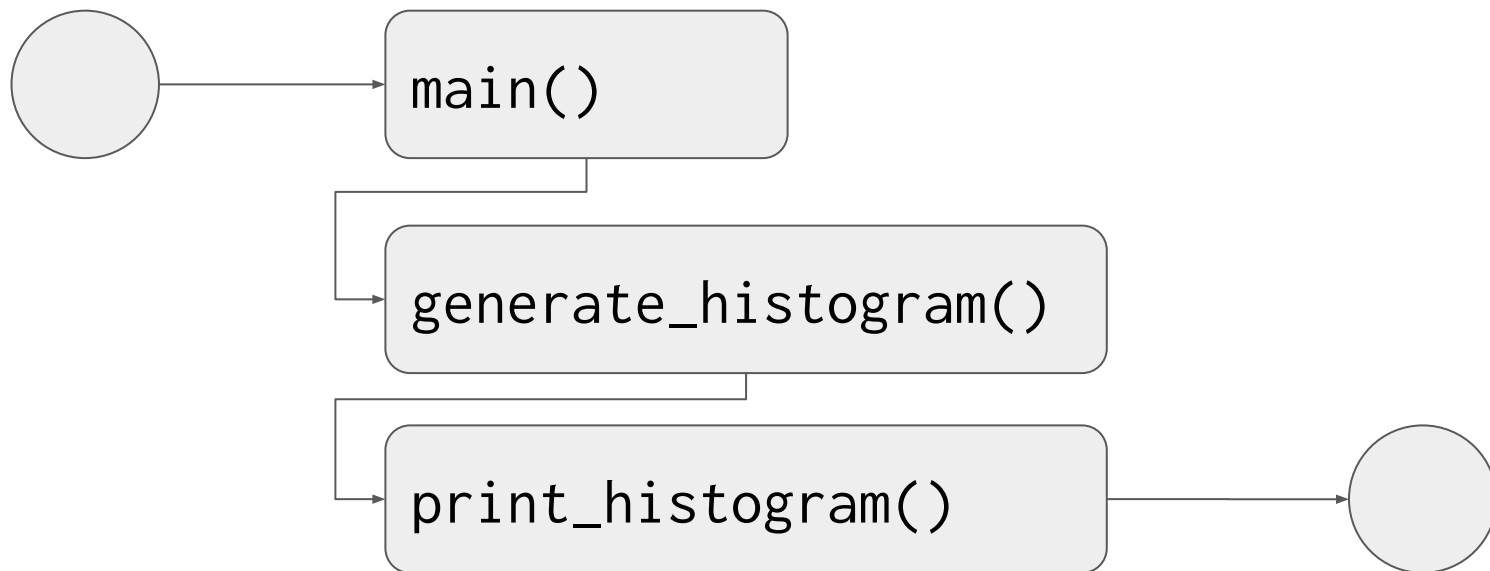


- Profiling: Analysis of the **program as a whole** (might be split down to function level)
- Benchmarking: Comparison of **evolutions/versions** of functions

Profiling and Benchmarking in this context



- Profiling: Analysis of the **program as a whole** (might be split down to function level)
- Benchmarking: Comparison of **evolutions/versions** of functions



Profiling and Benchmarking in this context



- Profiling: Analysis of the **program as a whole** (might be split down to function level)
- Benchmarking: Comparison of **evolutions/versions** of functions

`generate_histogram()`

`fast_histogram()`

`instant_histogram()`

Profiling and Benchmarking in this context



- Profiling: Analysis of the **program as a whole**
(might be split down to function level)
 - [hyperfine](#) (time)
 - [valgrind](#) (time + memory)
 - [gnu time](#) (time + memory)
- Definitely worth looking into ...
 - [flamegraph](#)
 - [perf](#)
 - *[not-perf](#)
 - *[cargo-instruments](#)

based on [C]

Profiling and Benchmarking in this context



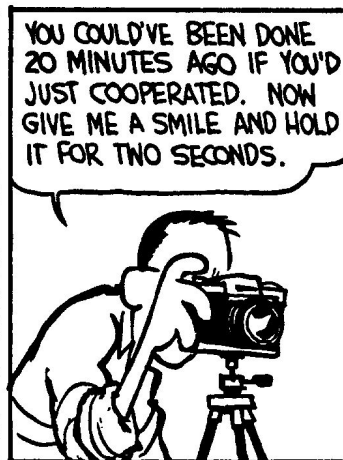
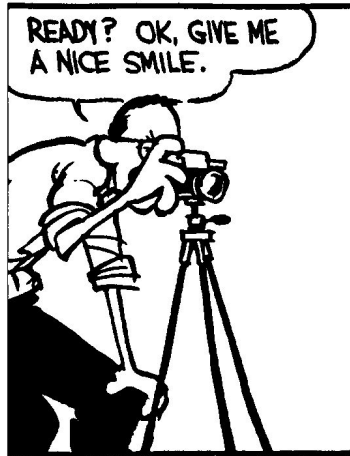
- Benchmarking: Comparison of **evolutions/versions** of functions
 - Very specific for each programming language
 - No general recommendation possible
 - Rust: [criterion](#)

Demo time



- Profile I/O of my laptop
- First generate some random data
- Copy this data from a to b
- Profile different copy commands
 - cp
 - dd
 - rsync
- Applications we will use
 - hyperfine (**time**)
 - valgrind (time + **memory**)
 - gnu time (time + **memory**)

That's all folks



Source [Z]

Sources



- [Z] [Calvin Photoshoot](#)
- [A] [Greece Olympia Marble Statue](#)
- [B] <https://drawingacademy.com/drawing-progress>
- [C] https://fosdem.org/2020/schedule/event/rust_optimizing_rav1e/