

Transformations for Compressing Quality Scores of Next Generation Sequencing Data

Raymond Wan

`r-wan@cb.k.u-tokyo.ac.jp`

University of Tokyo and CBRC, AIST

August 2, 2011

FASTQ Data

FASTQ

Compression

Next Generation

Sequence Data

Quality Scores

Compression

Compression

Background

Quality Scores

Compression

Results

Conclusion

Appendix

FASTQ Data

FASTQ Data

FASTQ
Compression

Next Generation
Sequence Data
Quality Scores
Compression

Compression
Background

Quality Scores
Compression

Results

Conclusion

Appendix

Motivation for current research:

- Amount of next generation sequence data steadily increasing.
- Offers no new **insight** into the data; but focus on storing the data – somewhat of a hidden issue.
- Compression aims to improve disk storage + transmission **costs**.

FASTQ Data

FASTQ
Compression

Next Generation
Sequence Data
Quality Scores
Compression

Compression
Background

Quality Scores
Compression

Results

Conclusion

Appendix

Overview of work by others:

- Relatively less work done in this area so far.
- Most work to date is on DNA sequence compression (since 1999).
- In the last 1-2 years, some work on read compression:
 - ◆ Most describe **software** with pre-selected parameters.
 - ◆ Little explanation about what parameters worked and what did not.

FASTQ format

```
@SRR032209.3000 length=36
AATTAACATGTAAAGGCAGTTTGCTTCCATGCTACC
+SRR032209.3000 length=36
BA@?@???@?@;A<9?>@>?(54423894+0+0038;
...
```

[FASTQ Data](#)

[FASTQ](#)

[Compression](#)

[Next Generation
Sequence Data](#)

[Quality Scores](#)

[Compression](#)

[Compression](#)

[Background](#)

[Quality Scores](#)

[Compression](#)

[Results](#)

[Conclusion](#)

[Appendix](#)

FASTQ format

```
→ @SRR032209.3000 length=36
   AATTAACATGTAAAGGCAGTTTGCTTCCATGCTACC
→ +SRR032209.3000 length=36
   BA@?@???@?@;A<9?>@>?(54423894+0+0038;
   ...
```

FASTQ Data

FASTQ

Compression

Next Generation
Sequence Data

Quality Scores
Compression

Compression
Background

Quality Scores
Compression

Results

Conclusion

Appendix

FASTQ format

```
@SRR032209.3000 length=36
→ AATTAACATGTAAAGGCAGTTTGCTTCCATGCTACC
+SRR032209.3000 length=36
BA@?@???@?@;A<9?>@>?(54423894+0+0038;
...
```

FASTQ Data

FASTQ

Compression

Next Generation
Sequence Data

Quality Scores
Compression

Compression
Background

Quality Scores
Compression

Results

Conclusion

Appendix

FASTQ format

```
@SRR032209.3000 length=36
```

```
AATTAACATGTAAAGGCAGTTTGCTTCCATGCTACC
```

```
+SRR032209.3000 length=36
```

```
→ BA@?@???@?@;A<9?>@>?(54423894+0+0038;
```

```
...
```

FASTQ quality scores range from “!” to “~”, or ASCII code 33 to 126 (the range of printable ASCII characters).

FASTQ Data

FASTQ

Compression

Next Generation
Sequence Data

Quality Scores
Compression

Compression
Background

Quality Scores
Compression

Results

Conclusion

Appendix

FASTQ Data

FASTQ

Compression

Next Generation
Sequence Data

Quality Scores
Compression

Compression
Background

Quality Scores
Compression

Results

Conclusion

Appendix

FASTQ format

```
@SRR032209.3000 length=36
AATTAACATGTAAAGGCAGTTTGCTTCCATGCTACC
+SRR032209.3000 length=36
BA@?@???@?@;A<9?>@>?(54423894+0+0038;
...
```

FASTQ quality scores range from “!” to “~”, or ASCII code 33 to 126 (the range of printable ASCII characters).

e.g.	“!”	1.000
	“+”	0.100
	“5”	0.010
	“?”	0.001

The number of possible characters is **94**.

[FASTQ Data](#)

[FASTQ](#)

[Compression](#)

[Next Generation](#)

[Sequence Data](#)

[Quality Scores
Compression](#)

[Compression](#)

[Background](#)

[Quality Scores
Compression](#)

[Results](#)

[Conclusion](#)

[Appendix](#)

Reasons for compressing quality scores:

- Perhaps the least important, but also has attracted the least attraction (so far).
- Systems are being developed to make use of quality scores – mapping, assembly, and read trimming.

FASTQ Data

**Compression
Background**

Compression
Model

Coder
Example: Gamma
Code

Model + Coder

Quality Scores
Compression

Results

Conclusion

Appendix

Compression Background

FASTQ Data

Compression Background

Compression

Model

Coder

Example: Gamma Code

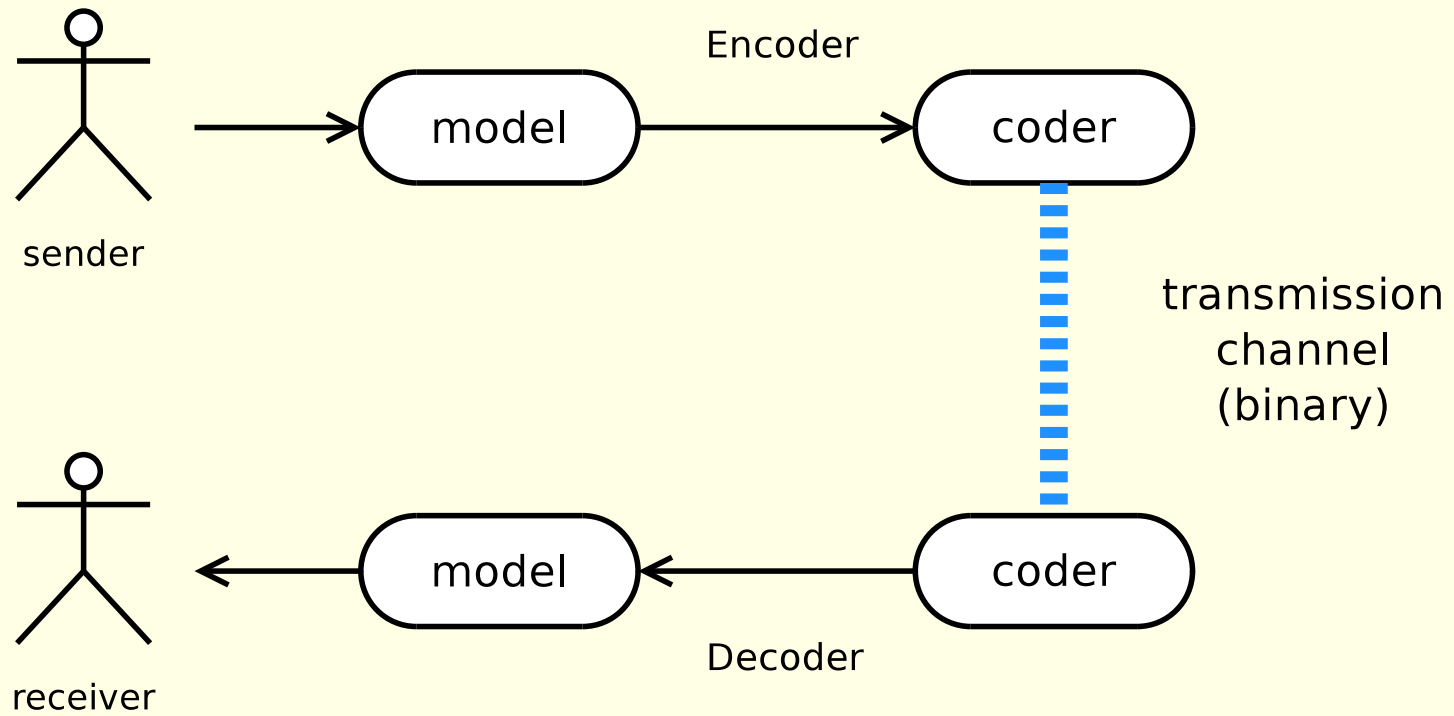
Model + Coder

Quality Scores Compression

Results

Conclusion

Appendix



FASTQ Data

Compression
Background

Compression

Model

Coder
Example: Gamma
Code

Model + Coder

Quality Scores
Compression

Results

Conclusion

Appendix

Model:

- What if we **remove** unimportant information from the data?
 - ◆
- What if we use “surrounding” data to create a model?
 - ◆
- What if the data is similar to data that we **just** saw?
 - ◆

FASTQ Data

Compression

Background

Compression

Model

Coder

Example: Gamma

Code

Model + Coder

Quality Scores

Compression

Results

Conclusion

Appendix

Model:

- What if we **remove** unimportant information from the data?
 - ◆ **MP3 compression for music**
- What if we use “surrounding” data to create a model?
 - ◆ **Image compression; also facsimile transmission**
- What if the data is similar to data that we **just** saw?
 - ◆ **gzip system**

→ We select/develop a model that depends on some prior knowledge of the data.

FASTQ Data

Compression
Background

Compression
Model

Coder
Example: Gamma
Code

Model + Coder

Quality Scores
Compression

Results

Conclusion

Appendix

Humans communicate in characters (Latin alphabet, kanji, hiragana, etc.) but computers use a binary alphabet $\{0, 1\}$.

Storing or transmitting data means that a conversion to binary digits (**bits**) is required.

That is the **coder's job** using the output of the model.

FASTQ Data

Compression
Background

Compression
Model

Coder

Example: Gamma
Code

Model + Coder

Quality Scores
Compression

Results

Conclusion

Appendix

ASCII maps the Latin alphabet, punctuation, and numbers to bits.

i.e., “I” → ASCII code 73 → 01001001

This code assumes all 256 characters in ASCII are **equally-probable**.

FASTQ Data

Compression
Background

Compression
Model

Coder

Example: Gamma
Code

Model + Coder

Quality Scores
Compression

Results

Conclusion

Appendix

- What if we encode a number x in two parts – (1) the largest logarithmically-sized **bin** and then (2) position in the bin?



- What if the **bins** are equal in size and determined by some parameter?



- What if we give short codes to frequently occurring symbols?



FASTQ Data

Compression
Background

Compression
Model

Coder

Example: Gamma
Code

Model + Coder

Quality Scores
Compression

Results

Conclusion

Appendix

- What if we encode a number x in two parts – (1) the largest logarithmically-sized **bin** and then (2) position in the bin?
 - ◆ **Gamma code**

- What if the **bins** are equal in size and determined by some parameter?
 - ◆ **Golomb code**

- What if we give short codes to frequently occurring symbols?
 - ◆ **Huffman coding**

→ Ideally, select a coder that fits the output of the model.

Create logarithmically-sized bins (base-2): 1, 2, 4, 8, 16, 32, 64, ...

x	Bin	Position in bin	Length
1	0		1
2	10	0	3
3	10	1	3
4	110	00	5
5	110	01	5
7	110	11	5
8	1110	000	7
...			
15	1110	111	7
16	11110	0000	9
...			
31	11110	1111	9
32	111110	00000	11

gzip uses the LZ77 algorithm to represent new data as a back-pointer into a window of recently-seen data.

[FASTQ Data](#)

[Compression Background](#)

[Compression Model](#)

[Coder Example: Gamma Code](#)

[Model + Coder](#)

[Quality Scores Compression](#)

[Results](#)

[Conclusion](#)

[Appendix](#)

gzip uses the LZ77 algorithm to represent new data as a back-pointer into a window of recently-seen data.

	Tuple	Text
(a)	$\langle 0, 0, w \rangle$	w <u>o</u> odchuck_lchuck_lif

FASTQ Data

Compression
Background

Compression
Model

Coder
Example: Gamma
Code

Model + Coder

Quality Scores
Compression

Results

Conclusion

Appendix

gzip uses the LZ77 algorithm to represent new data as a back-pointer into a window of recently-seen data.

	Tuple	Text
(a)	$\langle 0, 0, w \rangle$	woodchuck_chuck_if
(b)	$\langle 0, 0, o \rangle$	woodchuck_chuck_if

FASTQ Data

Compression
Background

Compression
Model

Coder
Example: Gamma
Code

Model + Coder

Quality Scores
Compression

Results

Conclusion

Appendix

gzip uses the LZ77 algorithm to represent new data as a back-pointer into a window of recently-seen data.

	Tuple	Text
(a)	$\langle 0, 0, w \rangle$	w <u>o</u> odchuck_chuck_if
(b)	$\langle 0, 0, o \rangle$	w <u>o</u> odchuck_chuck_if
(c)	$\langle 1, 1, d \rangle$	w <u>o</u> odchuck_chuck_if

FASTQ Data

Compression
Background

Compression
Model

Coder
Example: Gamma
Code

Model + Coder

Quality Scores
Compression

Results

Conclusion

Appendix

gzip uses the LZ77 algorithm to represent new data as a back-pointer into a window of recently-seen data.

	Tuple	Text
(a)	$\langle 0, 0, w \rangle$	w oodchuck_chuck_if
(b)	$\langle 0, 0, o \rangle$	w o odchuck_chuck_if
(c)	$\langle 1, 1, d \rangle$	w o o d chuck_chuck_if
(d)	$\langle 0, 0, c \rangle$	w o o d c huck_chuck_if

FASTQ Data

Compression
Background

Compression

Model

Coder

Example: Gamma
Code

Model + Coder

Quality Scores
Compression

Results

Conclusion

Appendix

gzip uses the LZ77 algorithm to represent new data as a back-pointer into a window of recently-seen data.

	Tuple	Text
(a)	$\langle 0, 0, w \rangle$	w oodchuck_chuck_if
(b)	$\langle 0, 0, o \rangle$	w o odchuck_chuck_if
(c)	$\langle 1, 1, d \rangle$	w o o d chuck_chuck_if
(d)	$\langle 0, 0, c \rangle$	w o o d c chuck_chuck_if
(e)	$\langle 0, 0, h \rangle$	w o o d c h chuck_chuck_if

FASTQ Data

Compression
Background

Compression
Model

Coder
Example: Gamma
Code

Model + Coder

Quality Scores
Compression

Results

Conclusion

Appendix

gzip uses the LZ77 algorithm to represent new data as a back-pointer into a window of recently-seen data.

	Tuple	Text
(a)	$\langle 0, 0, w \rangle$	w oodchuck_chuck_if
(b)	$\langle 0, 0, o \rangle$	w o odchuck_chuck_if
(c)	$\langle 1, 1, d \rangle$	w o d chuck_chuck_if
(d)	$\langle 0, 0, c \rangle$	w o o d chuck_chuck_if
(e)	$\langle 0, 0, h \rangle$	w o o d c h uck_chuck_if
(f)	$\langle 0, 0, u \rangle$	w o o d c h u ck_chuck_if

gzip uses the LZ77 algorithm to represent new data as a back-pointer into a window of recently-seen data.

	Tuple	Text
(a)	$\langle 0, 0, w \rangle$	w oodchuck_chuck_if
(b)	$\langle 0, 0, o \rangle$	w o odchuck_chuck_if
(c)	$\langle 1, 1, d \rangle$	w o d odchuck_chuck_if
(d)	$\langle 0, 0, c \rangle$	w o d c huck_chuck_if
(e)	$\langle 0, 0, h \rangle$	w o d c h huck_chuck_if
(f)	$\langle 0, 0, u \rangle$	w o d c h u ck_chuck_if
(g)	$\langle 3, 1, k \rangle$	w o d c h u c k _chuck_if

FASTQ Data

Compression Background

Compression

Model

Coder

Example: Gamma Code

Model + Coder

Quality Scores Compression

Results

Conclusion

Appendix

gzip uses the LZ77 algorithm to represent new data as a back-pointer into a window of recently-seen data.

	Tuple	Text
(a)	$\langle 0, 0, w \rangle$	w oodchuck_lchuck_lif
(b)	$\langle 0, 0, o \rangle$	w o odchuck_lchuck_lif
(c)	$\langle 1, 1, d \rangle$	w o o d chuck_lchuck_lif
(d)	$\langle 0, 0, c \rangle$	w o o d c huck_lchuck_lif
(e)	$\langle 0, 0, h \rangle$	w o o d c h huck_lchuck_lif
(f)	$\langle 0, 0, u \rangle$	w o o d c h u ck_lchuck_lif
(g)	$\langle 3, 1, k \rangle$	w o o d c h u c k _lchuck_lif
(h)	$\langle 0, 0, _ \rangle$	w o o d c h u c k _ _lchuck_lif

gzip uses the LZ77 algorithm to represent new data as a back-pointer into a window of recently-seen data.

	Tuple	Text
(a)	$\langle 0, 0, w \rangle$	w oodchuck_lchuck_lif
(b)	$\langle 0, 0, o \rangle$	w o odchuck_lchuck_lif
(c)	$\langle 1, 1, d \rangle$	w o o dchuck_lchuck_lif
(d)	$\langle 0, 0, c \rangle$	w o o d c huck_lchuck_lif
(e)	$\langle 0, 0, h \rangle$	w o o d c h uck_lchuck_lif
(f)	$\langle 0, 0, u \rangle$	w o o d c h u ck_lchuck_lif
(g)	$\langle 3, 1, k \rangle$	w o o d c h u c k _lchuck_lif
(h)	$\langle 0, 0, _ \rangle$	w o o d c h u c k _ _lchuck_lif
(i)	$\langle 6, 6, i \rangle$	w o o d c h u c k _ _c h u c k _ _l i f

FASTQ Data

Compression Background

Compression

Model

Coder

Example: Gamma Code

Model + Coder

Quality Scores Compression

Results

Conclusion

Appendix

gzip uses the LZ77 algorithm to represent new data as a back-pointer into a window of recently-seen data.

	Tuple	Text
(a)	$\langle 0, 0, w \rangle$	w oodchuck_lchuck_lif
(b)	$\langle 0, 0, o \rangle$	w o odchuck_lchuck_lif
(c)	$\langle 1, 1, d \rangle$	w o o dchuck_lchuck_lif
(d)	$\langle 0, 0, c \rangle$	w o o d c huck_lchuck_lif
(e)	$\langle 0, 0, h \rangle$	w o o d c h uck_lchuck_lif
(f)	$\langle 0, 0, u \rangle$	w o o d c h u ck_lchuck_lif
(g)	$\langle 3, 1, k \rangle$	w o o d c h u c k _lchuck_lif
(h)	$\langle 0, 0, _ \rangle$	w o o d c h u c k _ _lchuck_lif
(i)	$\langle 6, 6, i \rangle$	w o o d c h u c k _ _c h u c k _ _l i f
(j)	$\langle 0, 0, f \rangle$	w o o d c h u c k _ _c h u c k _ _l i f

FASTQ Data

Compression Background

Compression

Model

Coder

Example: Gamma Code

Model + Coder

Quality Scores Compression

Results

Conclusion

Appendix

gzip uses the LZ77 algorithm to represent new data as a back-pointer into a window of recently-seen data.

	Tuple	Text
(a)	$\langle 0, 0, w \rangle$	w oodchuck_lchuck_lif
(b)	$\langle 0, 0, o \rangle$	w o odchuck_lchuck_lif
(c)	$\langle 1, 1, d \rangle$	w o o dchuck_lchuck_lif
(d)	$\langle 0, 0, c \rangle$	w o o d c huck_lchuck_lif
(e)	$\langle 0, 0, h \rangle$	w o o d c h uck_lchuck_lif
(f)	$\langle 0, 0, u \rangle$	w o o d c h u ck_lchuck_lif
(g)	$\langle 3, 1, k \rangle$	w o o d c h u c k _lchuck_lif
(h)	$\langle 0, 0, _ \rangle$	w o o d c h u c k _ _lchuck_lif
(i)	$\langle 6, 6, i \rangle$	w o o d c h u c k _ _c h u c k _ _l i f
(j)	$\langle 0, 0, f \rangle$	w o o d c h u c k _ _c h u c k _ _l i f

Also called the “sliding-window” algorithm.

Source: Wan [2003, pg. 21]

FASTQ Data

Compression Background

Compression

Model

Coder

Example: Gamma Code

Model + Coder

Quality Scores Compression

Results

Conclusion

Appendix

gzip uses the LZ77 algorithm to represent new data as a back-pointer into a window of recently-seen data.

	Tuple	Text
(a)	$\langle 0, 0, w \rangle$	w oodchuck_chuck_if
(b)	$\langle 0, 0, o \rangle$	w o odchuck_chuck_if
(c)	$\langle 1, 1, d \rangle$	w o o d chuck_chuck_if
(d)	$\langle 0, 0, c \rangle$	w o o d c huck_chuck_if
(e)	$\langle 0, 0, h \rangle$	w o o d c h uck_chuck_if
(f)	$\langle 0, 0, u \rangle$	w o o d c h u ck_chuck_if
(g)	$\langle 3, 1, k \rangle$	w o o d c h u c k _chuck_if
(h)	$\langle 0, 0, _ \rangle$	w o o d c h u c k _ _chuck_if
(i)	$\langle 6, 6, i \rangle$	w o o d c h u c k _ _c h u c k _ _i f
(j)	$\langle 0, 0, f \rangle$	w o o d c h u c k _ _c h u c k _ _i f

Also called the “sliding-window” algorithm.

Source: Wan [2003, pg. 21]

FASTQ Data

Compression Background

Compression

Model

Coder

Example: Gamma Code

Model + Coder

Quality Scores Compression

Results

Conclusion

Appendix

FASTQ Data

Compression
Background

**Quality Scores
Compression**

First Steps

Lossless
Transformations

Lossy
Transformations
Compression

Results

Conclusion

Appendix

Quality Scores Compression

FASTQ Data

Compression
Background

Quality Scores
Compression

First Steps

Lossless
Transformations
Lossy
Transformations
Compression

Results

Conclusion

Appendix

Our first step is to examine what real data looks like – something that others have already done:

Accession	Species	# reads ($\times 10^6$)	Length	Size (MiB)
SRR032209	<i>Mus musculus</i>	18.8	36	662.8
SRR070788_1	<i>Homo sapiens</i>	24.8	100	2,393.2
SRR089526	<i>H. sapiens</i>	23.9	48	1,115.7

FASTQ Data

Compression
Background

Quality Scores
Compression

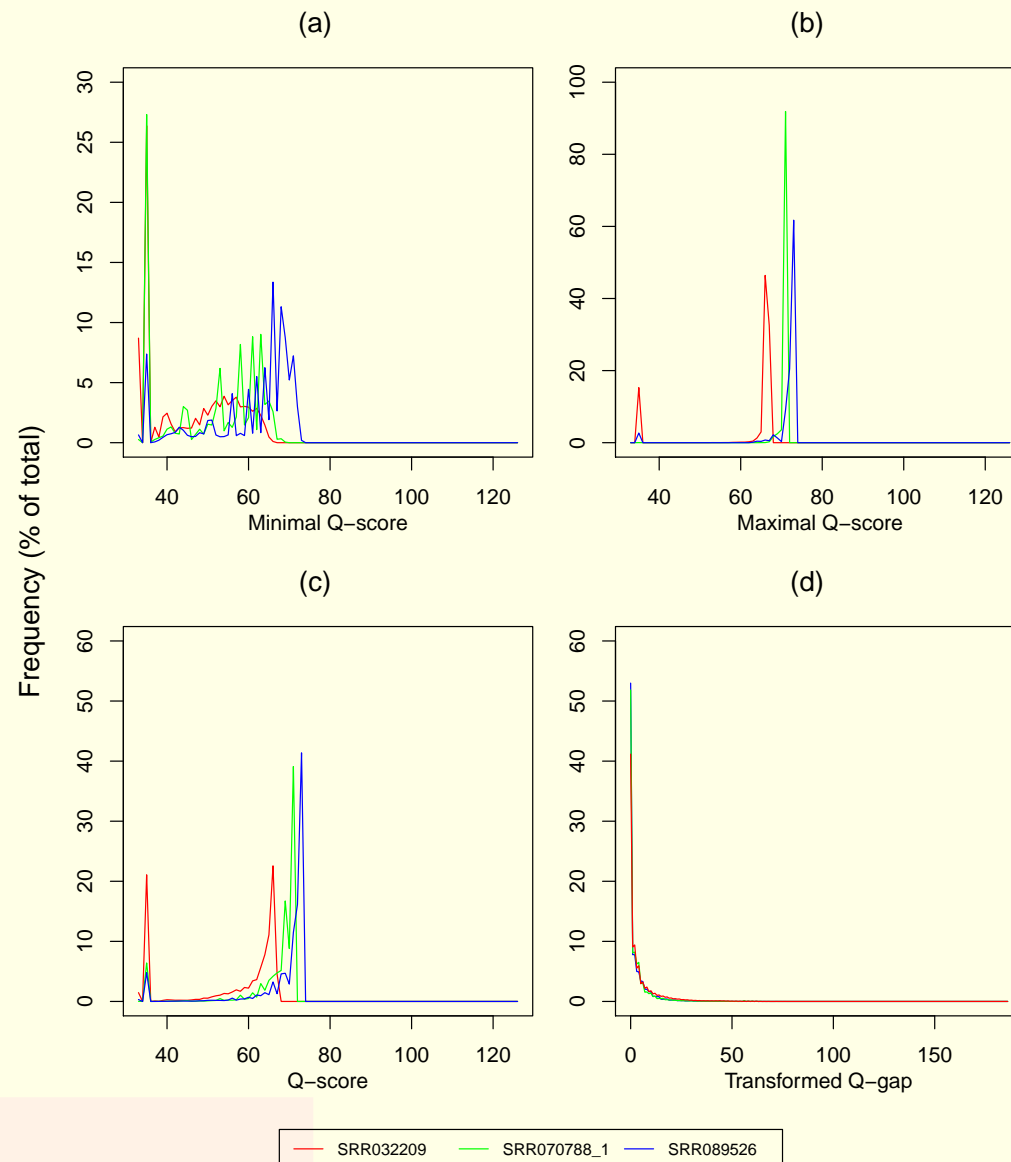
First Steps

Lossless
Transformations
Lossy
Transformations
Compression

Results

Conclusion

Appendix



FASTQ Data

Compression
Background

Quality Scores
Compression

First Steps

Lossless
Transformations
Lossy
Transformations
Compression

Results

Conclusion

Appendix

Findings:

- About half the range of possible quality scores is unused in real data.
- Non-uniform distribution of quality scores.
- The absolute differences between adjacent quality scores is close to 0.

FASTQ Data

Compression
Background

Quality Scores
Compression

First Steps

Lossless
Transformations

Lossy
Transformations
Compression

Results

Conclusion

Appendix

ReadBlocking Create blocks of reads so that each block is processed independently.

MinShifting Shift all of the quality scores in the block so that the smallest value is 1.

FreqOrdering Reorder quality scores by frequency so that the most frequent ones have the smallest value of 1.

GapTranslating Convert quality scores in a block into differences.

FASTQ Data

Compression
Background

Quality Scores
Compression

First Steps

Lossless
Transformations

Lossy
Transformations

Compression

Results

Conclusion

Appendix

Example: ADBCCCEE \rightarrow {65, 68, 66, 67, 67, 67, 69, 69}

MinShifting {1, 4, 2, 3, 3, 3, 5, 5}

FreqOrdering {3, 5, 4, 1, 1, 1, 2, 2}

GapTranslating {65, 3, -2, 1, 0, 0, 2, 0}

\rightarrow {65, 5, 4, 1, 0, 0, 3, 0}

\rightarrow {65, 6, 5, 2, 1, 1, 4, 1}

The aim is to change the distribution of symbols so that small (integer values) appear more often.

FASTQ Data

Compression
Background

Quality Scores
Compression

First Steps

Lossless
Transformations

Lossy
Transformations

Compression

Results

Conclusion

Appendix

Truncating Truncate the range of quality scores so that values greater than x are changed to x .

LogBinning PHRED scores places quality scores in logarithmically-sized bins. Continue this practice but use only some of the bins. So, with a parameter of 3, we have: $\{33, 34, 35, 36, 37, 38, 39, \dots\} \rightarrow \{33, 33, 33, 36, 36, 36, 39, \dots\}$.

UniBinning Uniformly re-bin quality scores so that each probability range maps to a bin. So, with a parameter of 5, the first bin has the probabilities $[0, 0.2)$. The next one is $[0.2, 0.4)$, and so on.

We **unify** the parameters to the new alphabet size for easier comparison.

FASTQ Data

Compression
Background

Quality Scores
Compression

First Steps

Lossless
Transformations

Lossy
Transformations

Compression

Results

Conclusion

Appendix

Static codes Encodes a value x the same way regardless of neighboring symbols.

Parameterized codes Encodes a value x after examining the entire block to determine overall statistics.

Entropy codes Encodes the values in a block close to their “zero-order” self-information. So, frequent symbols get short codes.

Compression systems Employ commonly-used systems that already combine a model and a coder.

FASTQ Data

Compression
Background

Quality Scores
Compression

Results

Overview

Compression Results

Mapping

Bringing it All
Together

Conclusion

Appendix

Results

FASTQ Data

Compression Background

Quality Scores Compression

Results

Overview

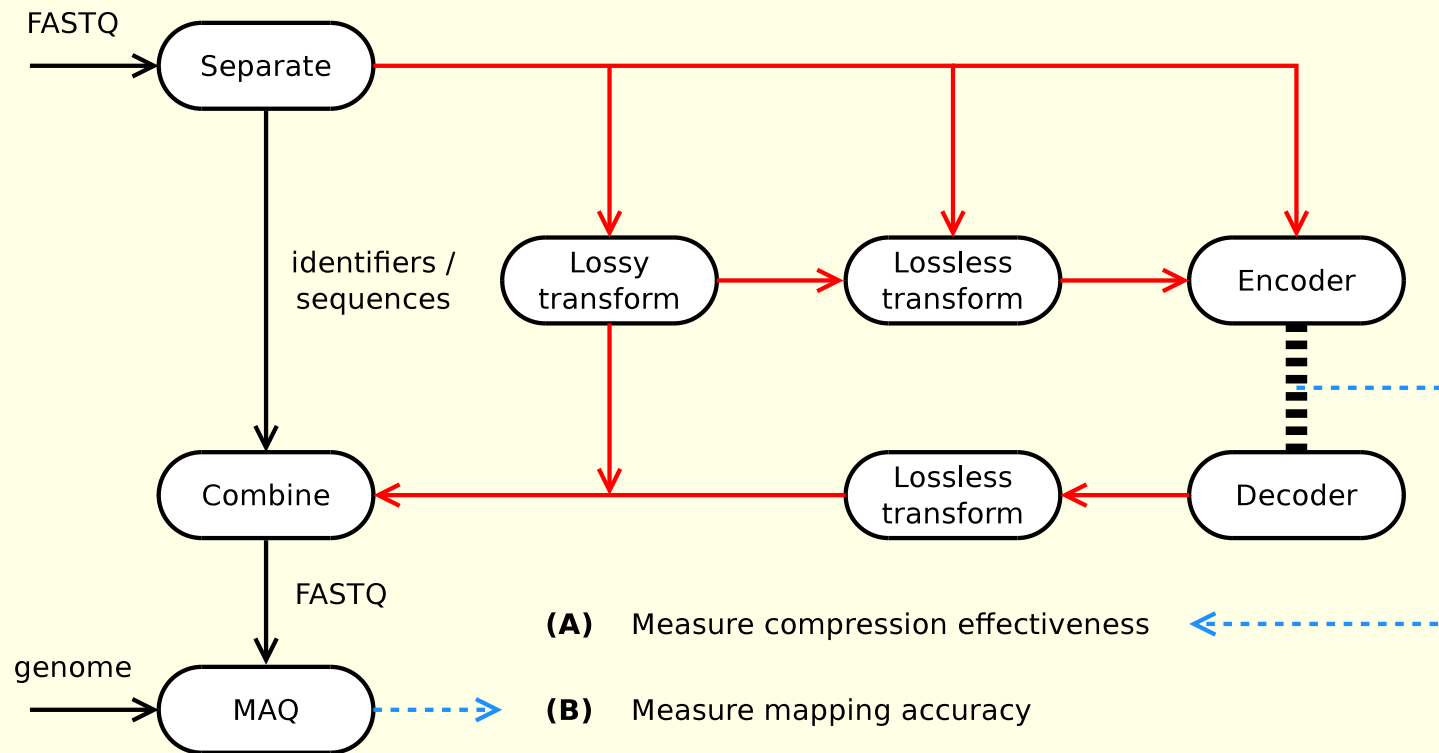
Compression Results

Mapping

Bringing it All Together

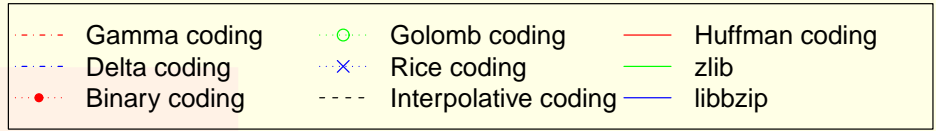
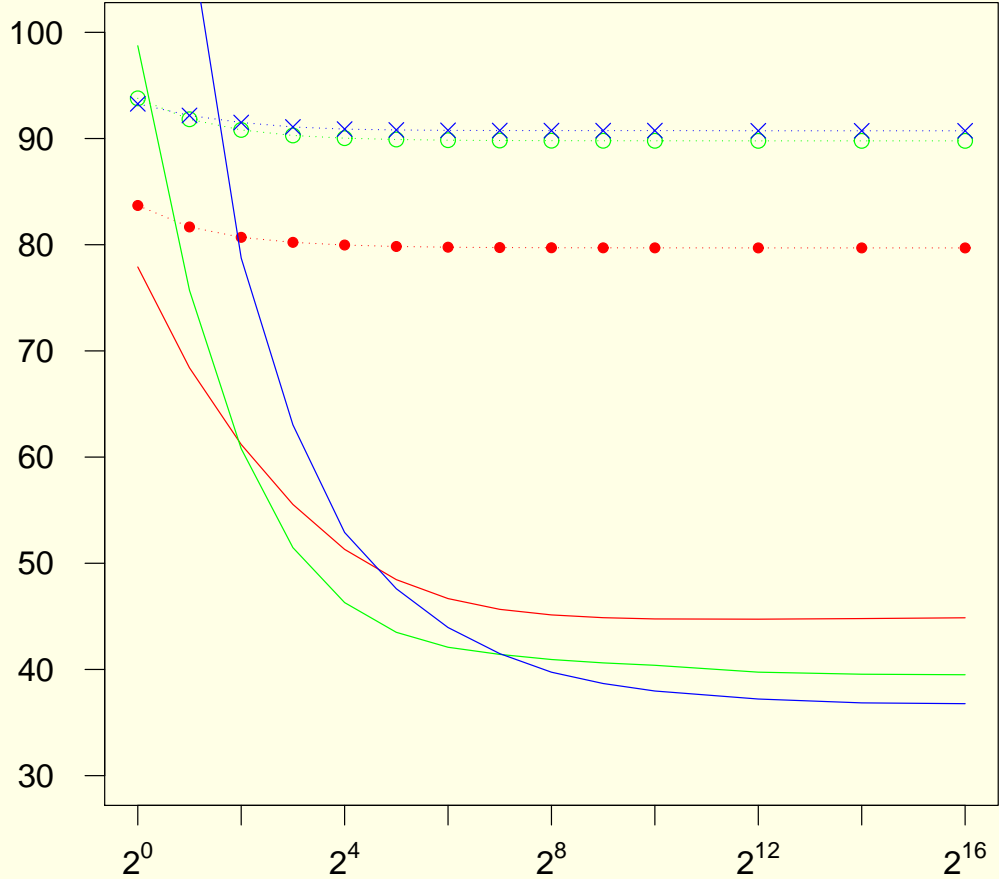
Conclusion

Appendix



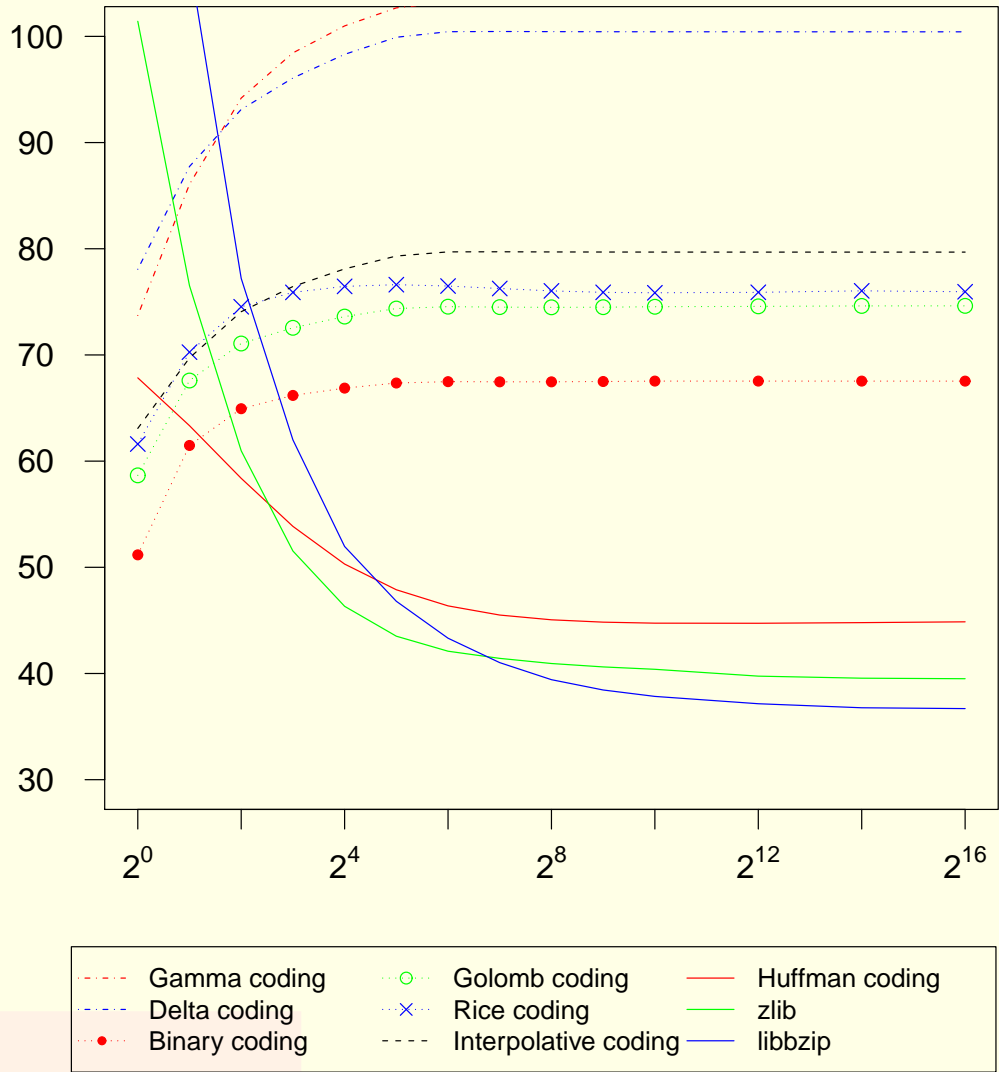
- FASTQ Data
- Compression Background
- Quality Scores Compression
- Results
- Overview
- Compression Results**
- Mapping
- Bringing it All Together
- Conclusion
- Appendix

No Transformation:



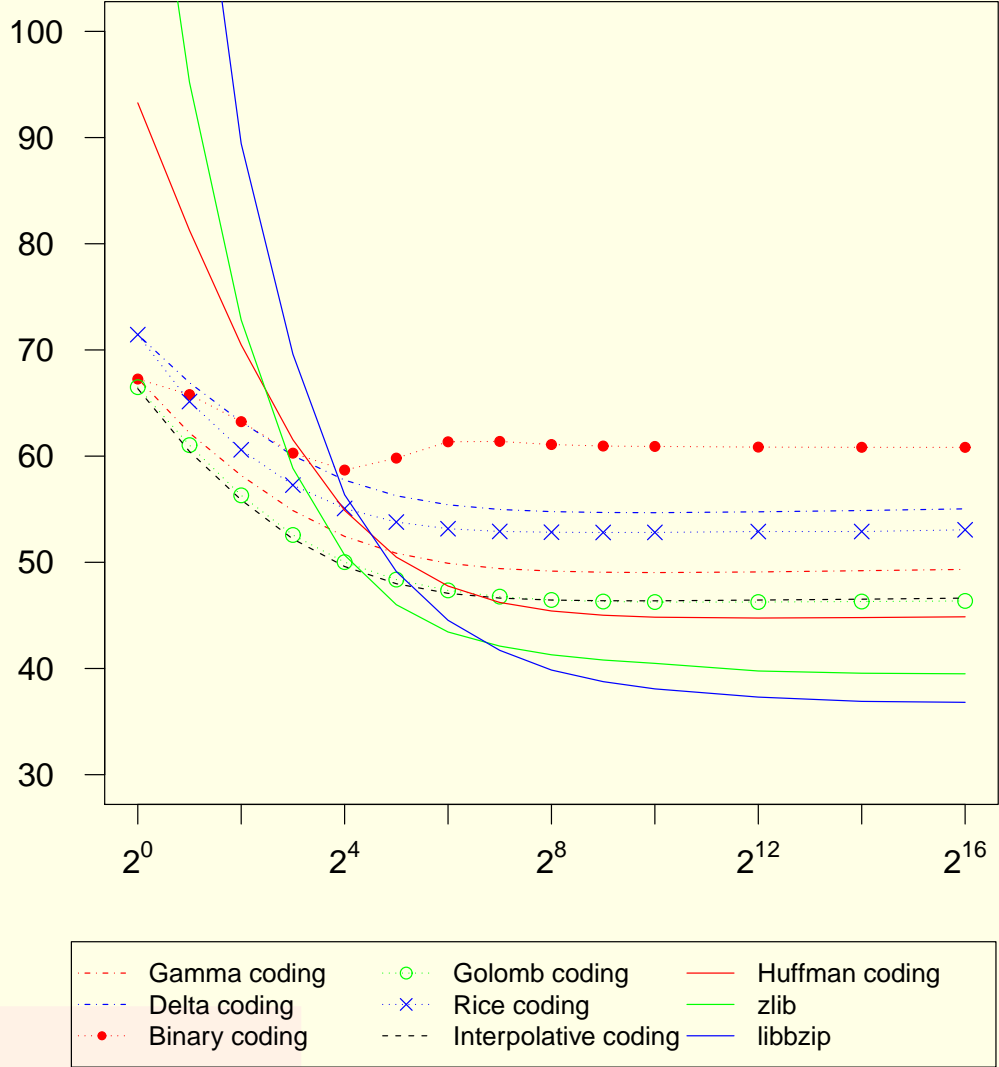
- FASTQ Data
- Compression Background
- Quality Scores Compression
- Results Overview
- Compression Results**
- Mapping
- Bringing it All Together
- Conclusion
- Appendix

MinShifting:



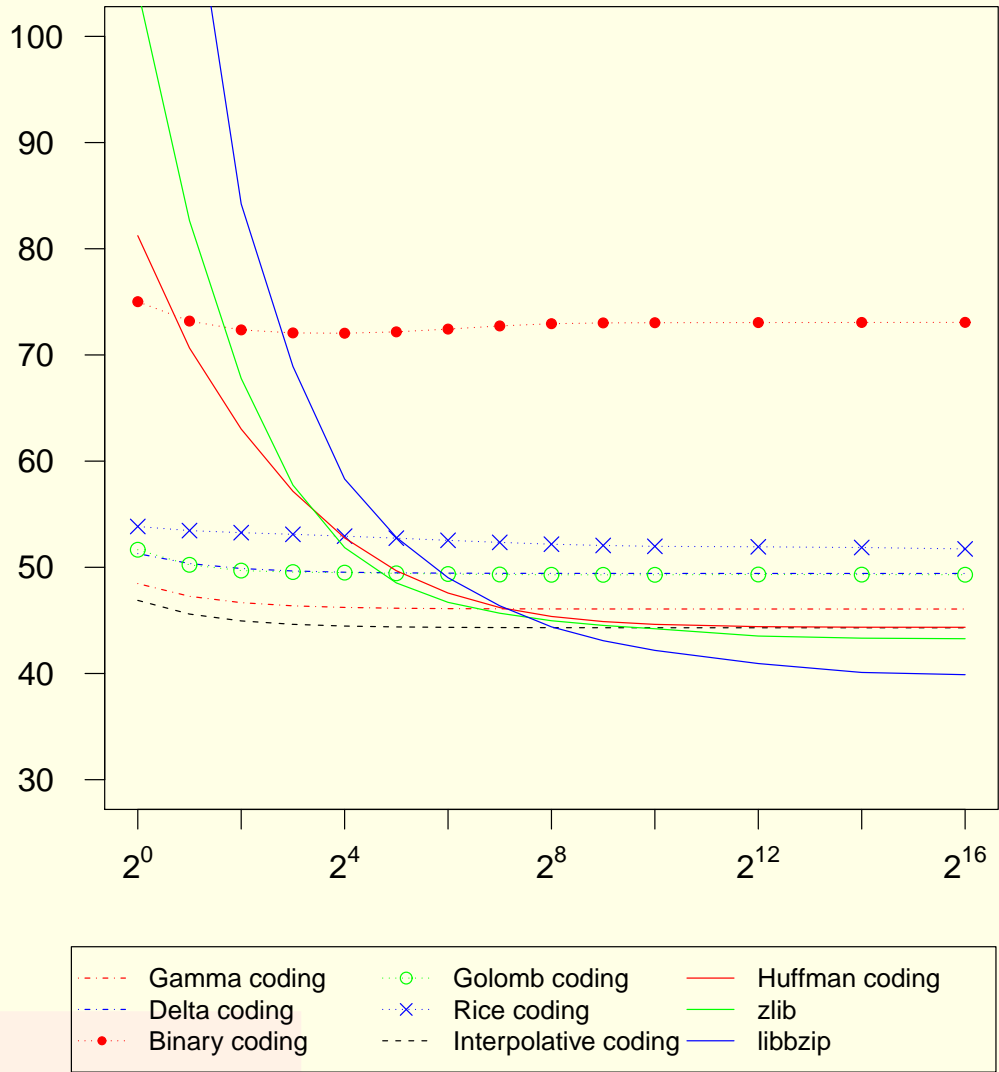
- FASTQ Data
- Compression Background
- Quality Scores Compression
- Results
- Overview
- Compression Results**
- Mapping
- Bringing it All Together
- Conclusion
- Appendix

FreqOrdering:



- FASTQ Data
- Compression Background
- Quality Scores Compression
- Results
- Compression Results**
- Mapping
- Bringing it All Together
- Conclusion
- Appendix

GapTranslating:



Compression Results

FASTQ Data

Compression Background

Quality Scores Compression

Results

Overview

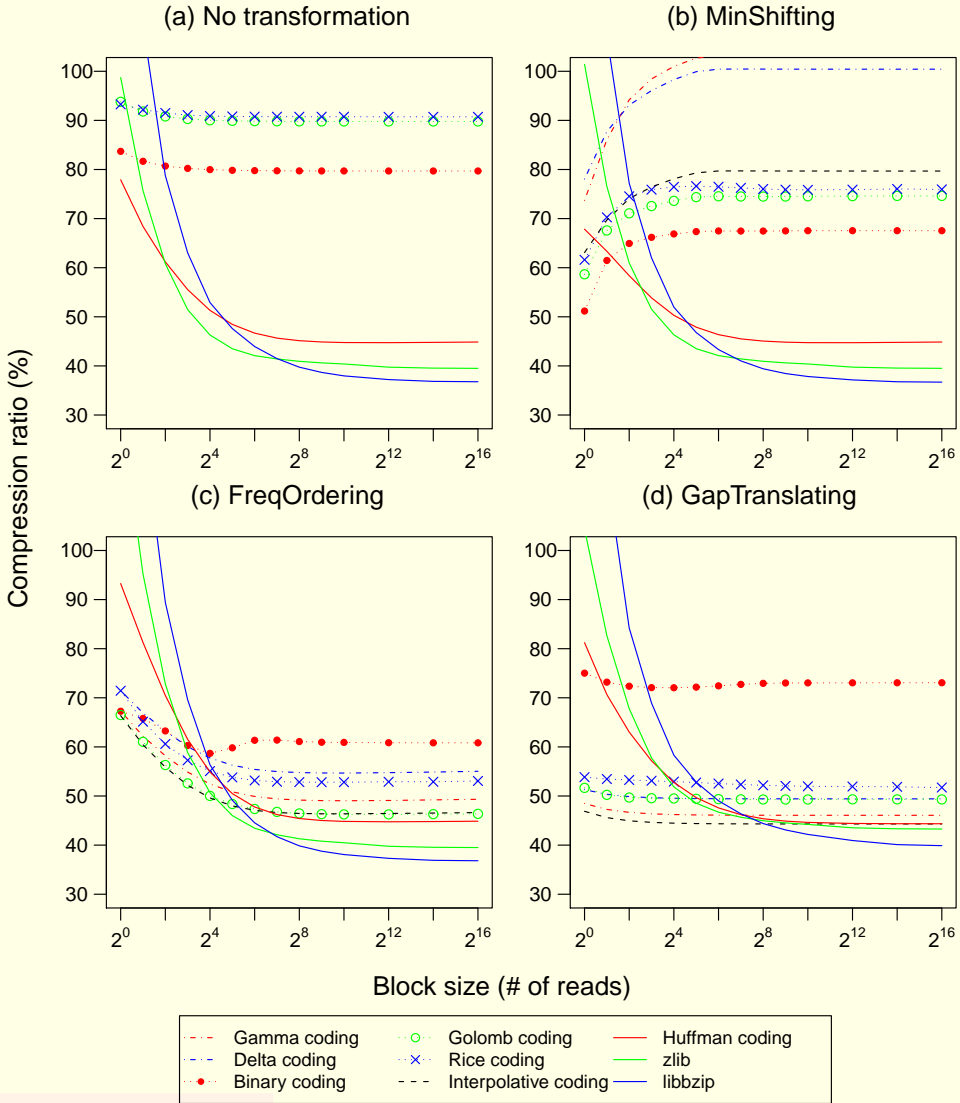
Compression Results

Mapping

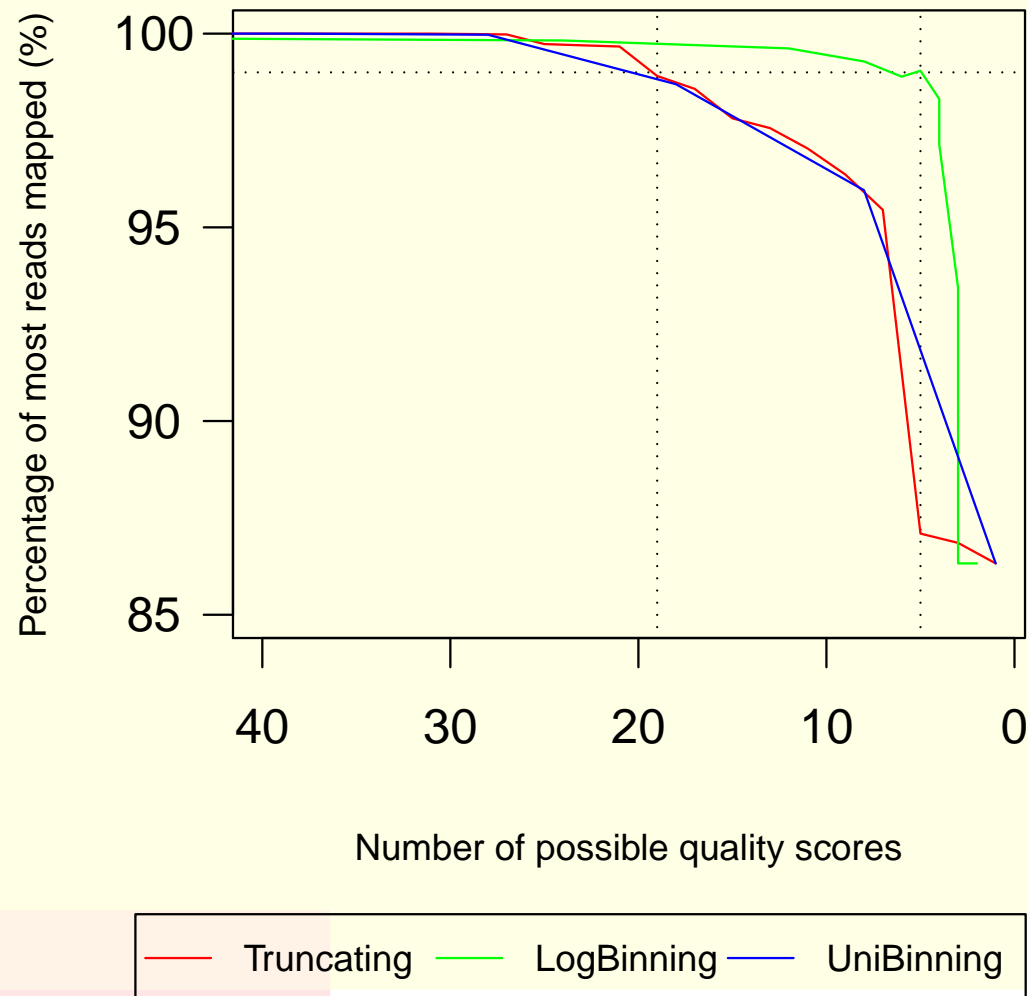
Bringing it All Together

Conclusion

Appendix



Of the 18.8×10^6 reads, 12.1×10^6 were unambiguously mapped (about 64.5%).



Bringing it All Together

FASTQ Data

Compression Background

Quality Scores Compression

Results

Overview

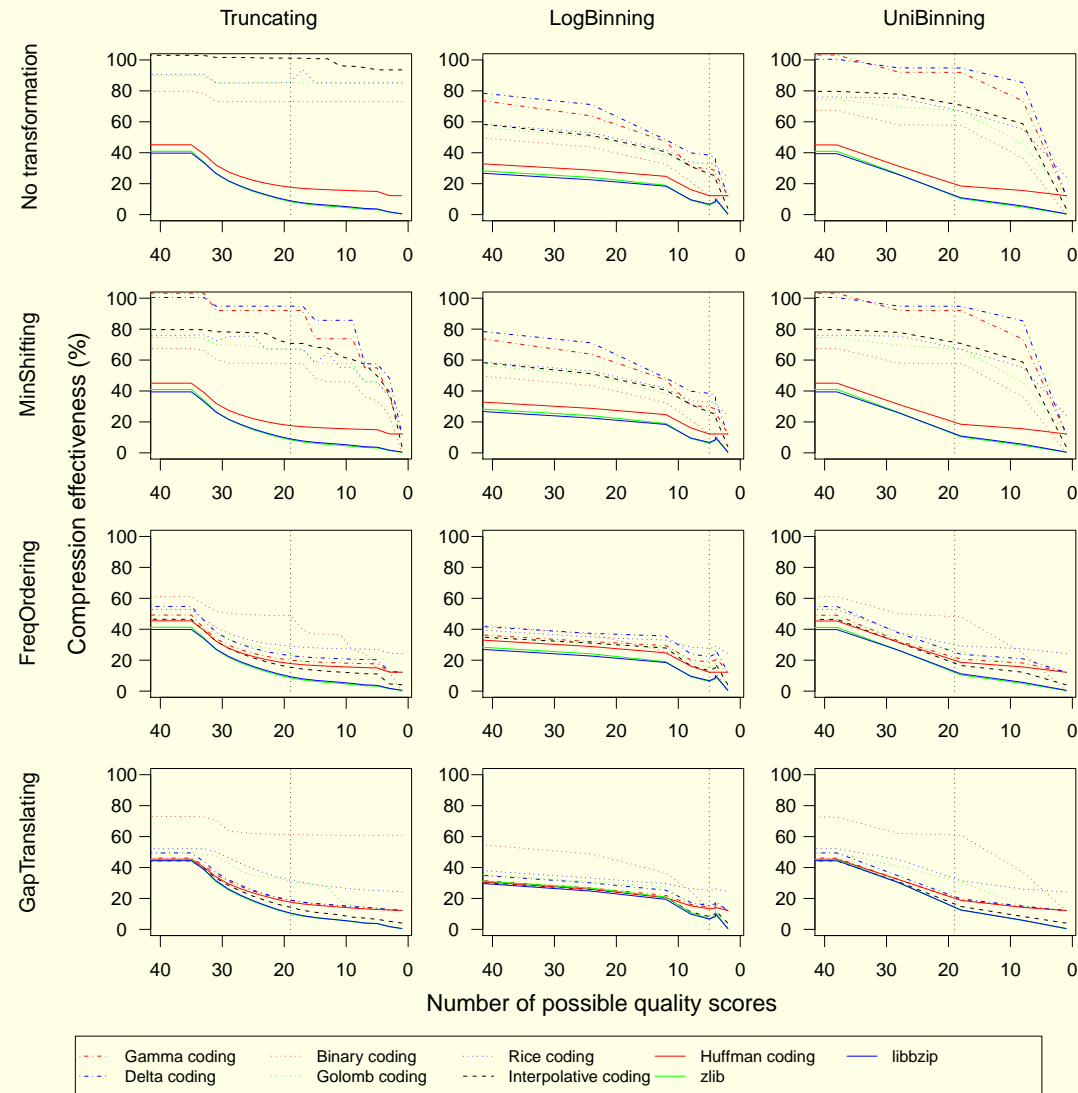
Compression Results

Mapping

Bringing it All Together

Conclusion

Appendix



FASTQ Data

Compression
Background

Quality Scores
Compression

Results

Conclusion

Advantages

Disadvantages

Questions?

Appendix

Conclusion

FASTQ Data

Compression
Background

Quality Scores
Compression

Results

Conclusion

Advantages

Disadvantages

Questions?

Appendix

- Simple coding schemes perform as well as more complicated compression systems.
- To attain this performance, lossless transformations are necessary.
- Small block sizes reduce memory usage (the amount of data in memory), minimize the effects from data corruption, and allow random access.

-
- Their compression time is less than `zlib` and `libbzip2`.
 - Compression ratios and compression/uncompression time of these “non-compression systems” for `FreqOrdering` and `GapTranslating` are **unaffected** by block size, essentially removing one parameter.

FASTQ Data

Compression
Background

Quality Scores
Compression

Results

Conclusion

Advantages

Disadvantages

Questions?

Appendix

- zlib and libbzip2 is well-known to general users; the other methods are more popular in the index compression field.
 - zlib still outperforms all other methods, if a large block size is chosen.
-
- zlib and libbzip2 is fast at decoding. This is partly due to their underlying methods, but due to their age, we expect them to be well-tuned for performance.

The most likely questions in your head...

FASTQ Data

Compression
Background

Quality Scores
Compression

Results

Conclusion

Advantages

Disadvantages

Questions?

Appendix

FASTQ Data

Compression
Background

Quality Scores
Compression

Results

Conclusion

Advantages

Disadvantages

Questions?

Appendix

The most likely questions in your head...

1. “Huh???”

FASTQ Data

Compression
Background

Quality Scores
Compression

Results

Conclusion

Advantages

Disadvantages

Questions?

Appendix

The most likely questions in your head...

1. “Huh???”
2. “Yes, I kind of see – but there is already BAM...what is the point of this work?”

FASTQ Data

Compression
Background

Quality Scores
Compression

Results

Conclusion

Advantages

Disadvantages

Questions?

Appendix

“Yes, I see – but there is already BAM...what is the point of this?”

- SAM/BAM stores mapped reads, including location information, “control information” to indicate changes to the sequence, etc.
- **Can** be used to store FASTQ data, but some overhead is added on top.
- BAM uses an “algorithm” called BGZF – essentially variable-length **gzipped** blocks.

So, our comparative results is roughly equivalent to BGZF, and thus BAM. Of course, our gzip results are not fully equivalent because of the above differences.

FASTQ Data

Compression
Background

Quality Scores
Compression

Results

Conclusion

Advantages

Disadvantages

Questions?

Appendix

- Dr. Vo Ngoc Anh [University of Melbourne]
- Prof. Kiyoshi Asai [University of Tokyo and CBRC, AIST]

FASTQ Data

Compression
Background

Quality Scores
Compression

Results

Conclusion

Appendix

References

Appendix

FASTQ Data

Compression
Background

Quality Scores
Compression

Results

Conclusion

Appendix

References

R. Wan. *Browsing and Searching Compressed Documents*. PhD thesis, University of Melbourne, Australia, December 2003