

Compression Options In Hadoop – A Tale of Tradeoffs

Govind Kamat, Sumeet Singh

Hadoop Summit (San Jose), June 27, 2013



YAHOO![®]

Introduction

Govind Kamat

Technical Yahoo!, Hadoop
Cloud Engineering Group



701 First Avenue
Sunnyvale, CA 94089 USA



- Member of Technical Staff in the Hadoop Services team at Yahoo!
- Focuses on HBase and Hadoop performance
- Worked with the Performance Engineering Group on improving the performance and scalability of several Yahoo! applications
- Experience includes development of large-scale software systems, microprocessor architecture, instruction-set simulators, compiler technology and electronic design

Sumeet Singh

Director of Products, Hadoop
Cloud Engineering Group



701 First Avenue
Sunnyvale, CA 94089 USA



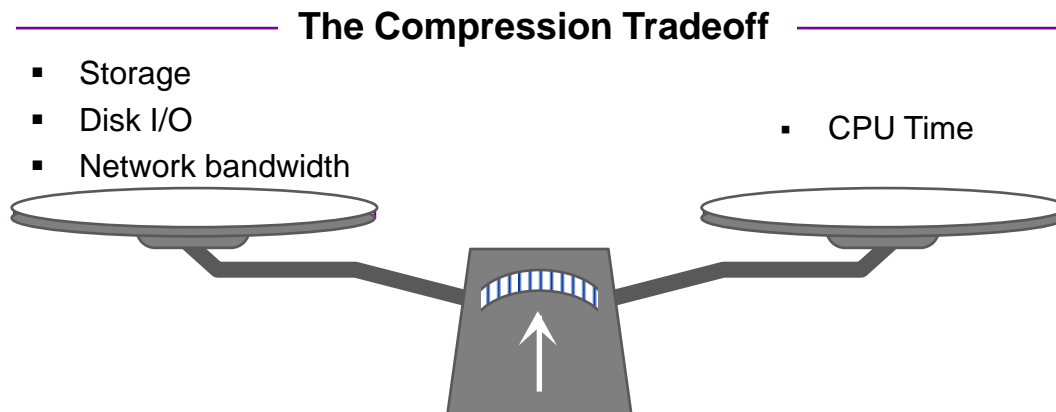
- Leads Hadoop products team at Yahoo!
- Responsible for Product Management, Customer Engagements, Evangelism, and Program Management
- Prior to this role, led Strategy functions for the Cloud Platform Group at Yahoo!

Agenda

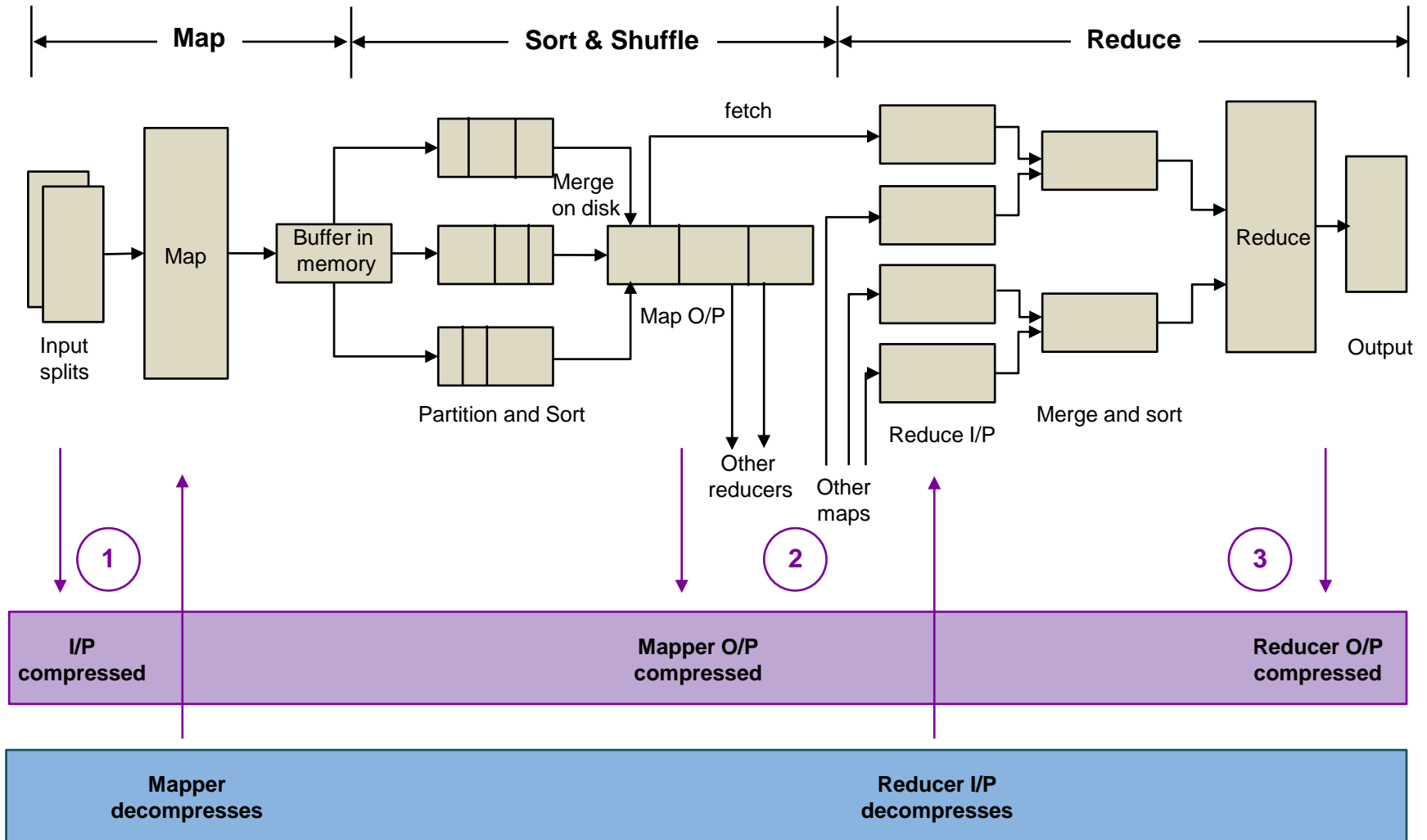
- 1 Data Compression in Hadoop
- 2 Available Compression Options
- 3 Understanding and Working with Compression Options
- 4 Problems Faced at Yahoo! with Large Data Sets
- 5 Performance Evaluations, Native Bzip2, and IPP Libraries
- 6 Wrap-up and Future Work

Compression Needs and Tradeoffs in Hadoop

- Hadoop jobs are data-intensive, compressing data can **speed up the I/O operations**
 - MapReduce jobs are almost always I/O bound
- Compressed data can **save storage space** and **speed up data transfers** across the network
 - Capital allocation for hardware can go further
- Reduced I/O and network load can bring **significant performance improvements**
 - MapReduce jobs can finish faster overall
- On the other hand, **CPU utilization and processing time increases** during compression and decompression
 - Understanding the tradeoffs is important for MapReduce pipeline's overall performance



Data Compression in Hadoop's MR Pipeline



Source: Hadoop: The Definitive Guide, Tom White



Compression Options in Hadoop (1/2)

| Format | Algorithm | Strategy | Emphasis | Comments |
|---------------|--|--|--|---|
| zlib | Uses DEFLATE (LZ77 and Huffman coding) | Dictionary-based, API | Compression ratio | Default codec |
| gzip | Wrapper around zlib | Dictionary-based, standard compression utility | Same as zlib, codec operates on and produces standard gzip files | For data interchange on and off Hadoop |
| bzip2 | Burrows-Wheeler transform | Transform-based, block-oriented | Higher compression ratios than zlib | Common for Pig |
| LZO | Variant of LZ77 | Dictionary-based, block-oriented, API | High compression speeds | Common for intermediate compression, HBase tables |
| LZ4 | Simplified variant of LZ77 | Fast scan, API | Very high compression speeds | Available in newer Hadoop distributions |
| Snappy | LZ77 | Block-oriented, API | Very high compression speeds | Came out of Google, previously known as Zippy |

Compression Options in Hadoop (2/2)

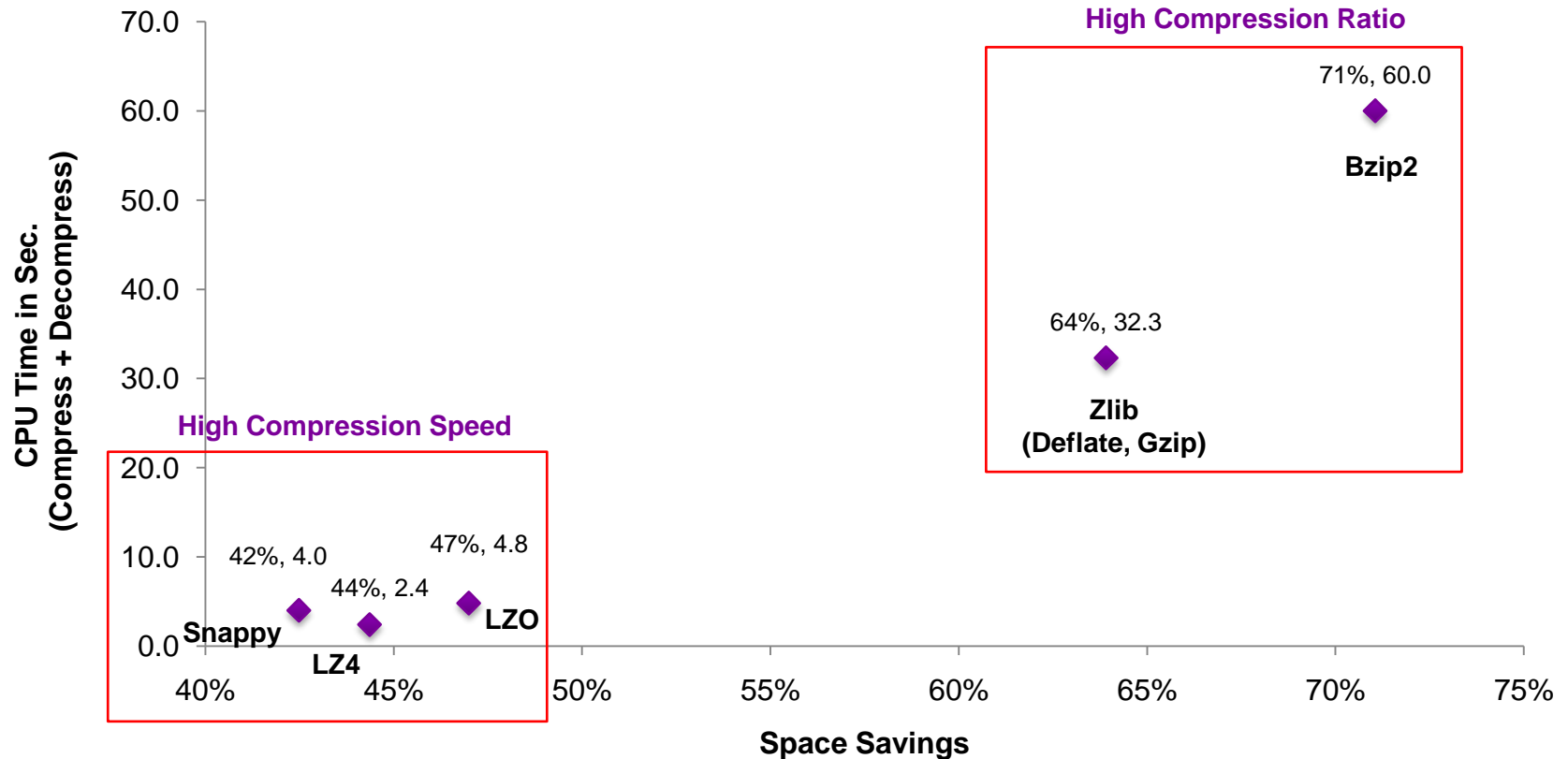
| Format | Codec (Defined in io.compression.codecs) | File Extn. | Splittable | Java/ Native |
|---------------------------------|--|------------|------------|-----------------|
| zlib/ DEFLATE (default) | org.apache.hadoop.io.compress.DefaultCodec | .deflate | N | Y/ Y |
| gzip | org.apache.hadoop.io.compress.GzipCodec | .gz | N | Y/ Y |
| bzip2 | org.apache.hadoop.io.compress.BZip2Codec | .bz2 | Y | Y/ Y |
| LZO (download separately) | com.hadoop.compression.lzo.LzoCodec | .lzo | N | N/ Y |
| LZ4 | org.apache.hadoop.io.compress.Lz4Codec | .lz4 | N | N/ Y |
| Snappy | org.apache.hadoop.io.compress.SnappyCodec | .snappy | N | N/ Y |

NOTES:

- **Splittability** – Bzip2 is “splittable”, can be decompressed in parallel by multiple MapReduce tasks. Other algorithms require all blocks together for decompression with a single MapReduce task.
- **LZO** – Removed from Hadoop because the LZO libraries are licensed under the GNU GPL. LZO format is still supported and the codec can be downloaded separately and enabled manually.
- **Native bzip2 codec** – added by Yahoo! as part of this work in Hadoop 0.23

Space-Time Tradeoff of Compression Options

Codec Performance on the Wikipedia Text Corpus



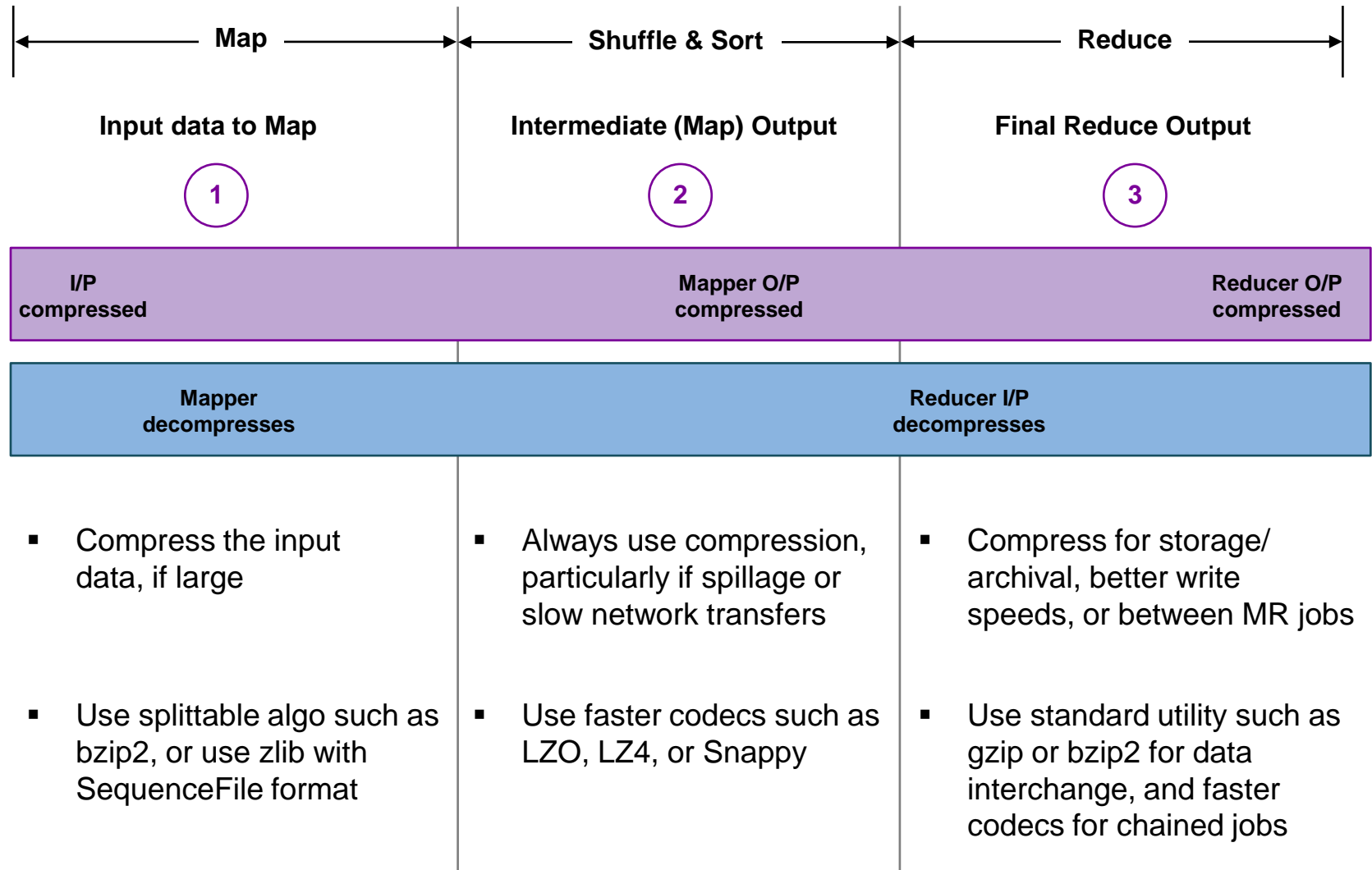
Note:

A 265 MB corpus from Wikipedia was used for the performance comparisons.
Space savings is defined as $[1 - (\text{Compressed} / \text{Uncompressed})]$

Using Data Compression in Hadoop




| Phase in MR Pipeline | Config | Values |
|-----------------------------|---|---|
| 1 Input data to Map | File extension recognized automatically for decompression | File extensions for supported formats <u>Note</u> : For SequenceFile, headers have the information [compression (boolean), block compression (boolean), and compression codec] |
| | One of the supported codecs | one defined in io.compression.codecs |
| 2 Intermediate (Map) Output | mapreduce.map.output.compress | false (default), true |
| | mapreduce.map.output.compress.codec | one defined in io.compression.codecs |
| 3 Final (Reduce) Output | mapreduce.output.fileoutputformat. compress | false (default), true |
| | mapreduce.output.fileoutputformat. compress.codec | one defined in io.compression.codecs |
| | mapreduce.output.fileoutputformat. compress.type | Type of compression to use for SequenceFile outputs: NONE, RECORD (default), BLOCK |

When to Use Compression and Which Codec

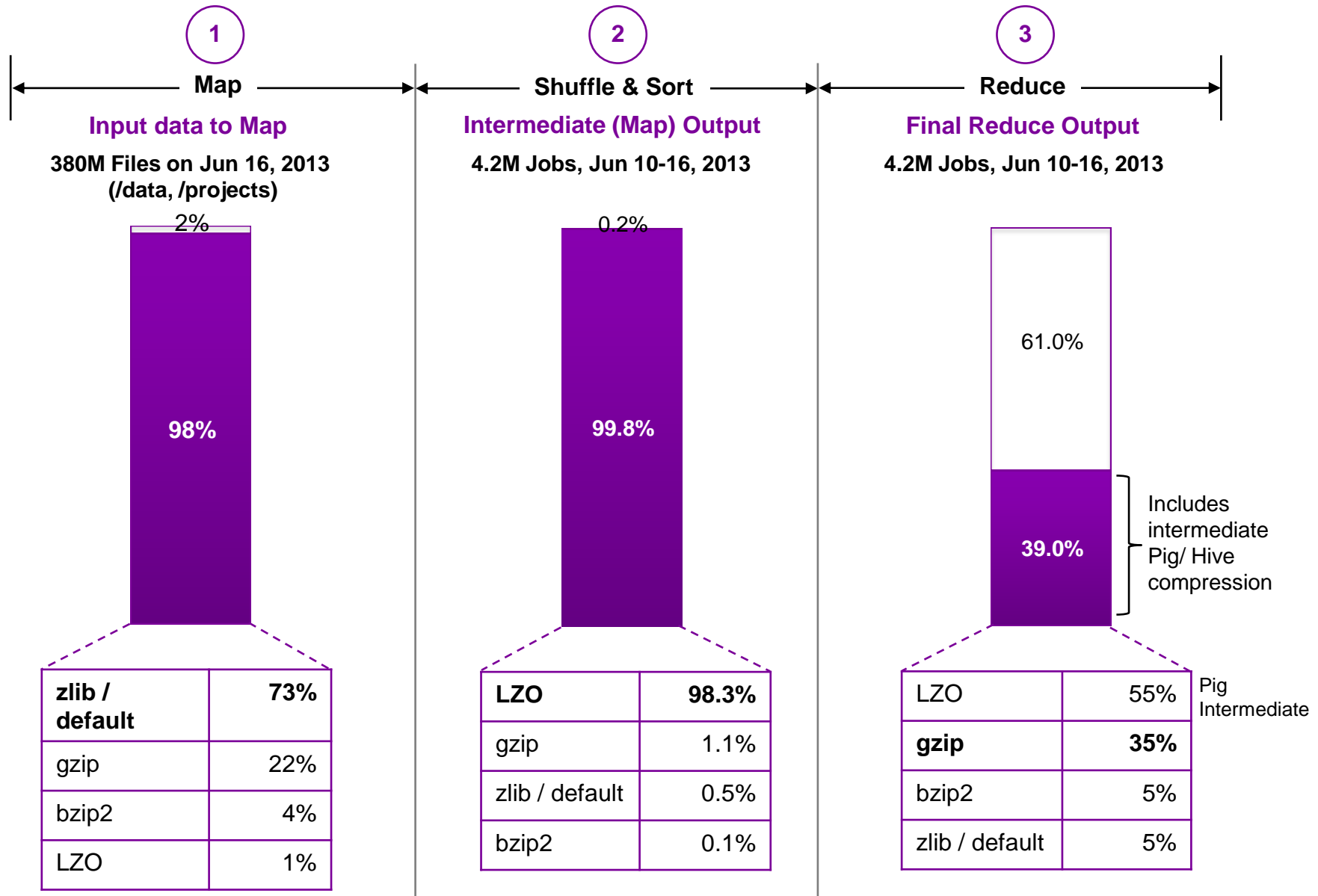


 Compress  Decompress

Compression in the Hadoop Ecosystem

| Component | When to Use | What to Use |
|---|---|--|
| Pig  | <ul style="list-style-type: none">▪ Compressing data between MR job▪ Typical in Pig scripts that include joins or other operators that expand your data size | Enable compression and select the codec: <code>pig.tmpfilecompression = true</code> <code>pig.tmpfilecompression.codec = gzip, lzo</code> |
| Hive  | <ul style="list-style-type: none">▪ Intermediate files produced by Hive between multiple map-reduce jobs▪ Hive writes output to a table | Enable intermediate or output compression: <code>hive.exec.compress.intermediate = true</code> <code>hive.exec.compress.output = true</code> |
| HBase  | <ul style="list-style-type: none">▪ Compress data at the CF level (support for LZO, gzip, Snappy, and LZ4) | List required JNI libraries: <code>hbase.regionserver.codecs</code> Enabling compression: <code>create 'table', { NAME => 'colfam', COMPRESSION => 'LZO' }</code> <code>alter 'table', { NAME => 'colfam', COMPRESSION => 'LZO' }</code> |

Compression in Hadoop at Yahoo!



Compression for Data Storage Efficiency

- DSE considerations at Yahoo!
- RCFile instead of SequenceFile
- Faster implementation of bzip2
- Native-code bzip2 codec
- HADOOP-8462¹, available in 0.23.7
- Substituting the IPP library

¹ Native-code bzip2 implementation done in collaboration with Jason Lowe, Hadoop Core PMC member

IPP Libraries

- Integrated Performance Primitives from Intel
- Algorithmic and architectural optimizations
- Processor-specific variants of each function
- Applications remain processor-neutral
- Compression: LZ, RLE, BWT, LZO
- High level formats include: zlib, gzip, bzip2 and LZO

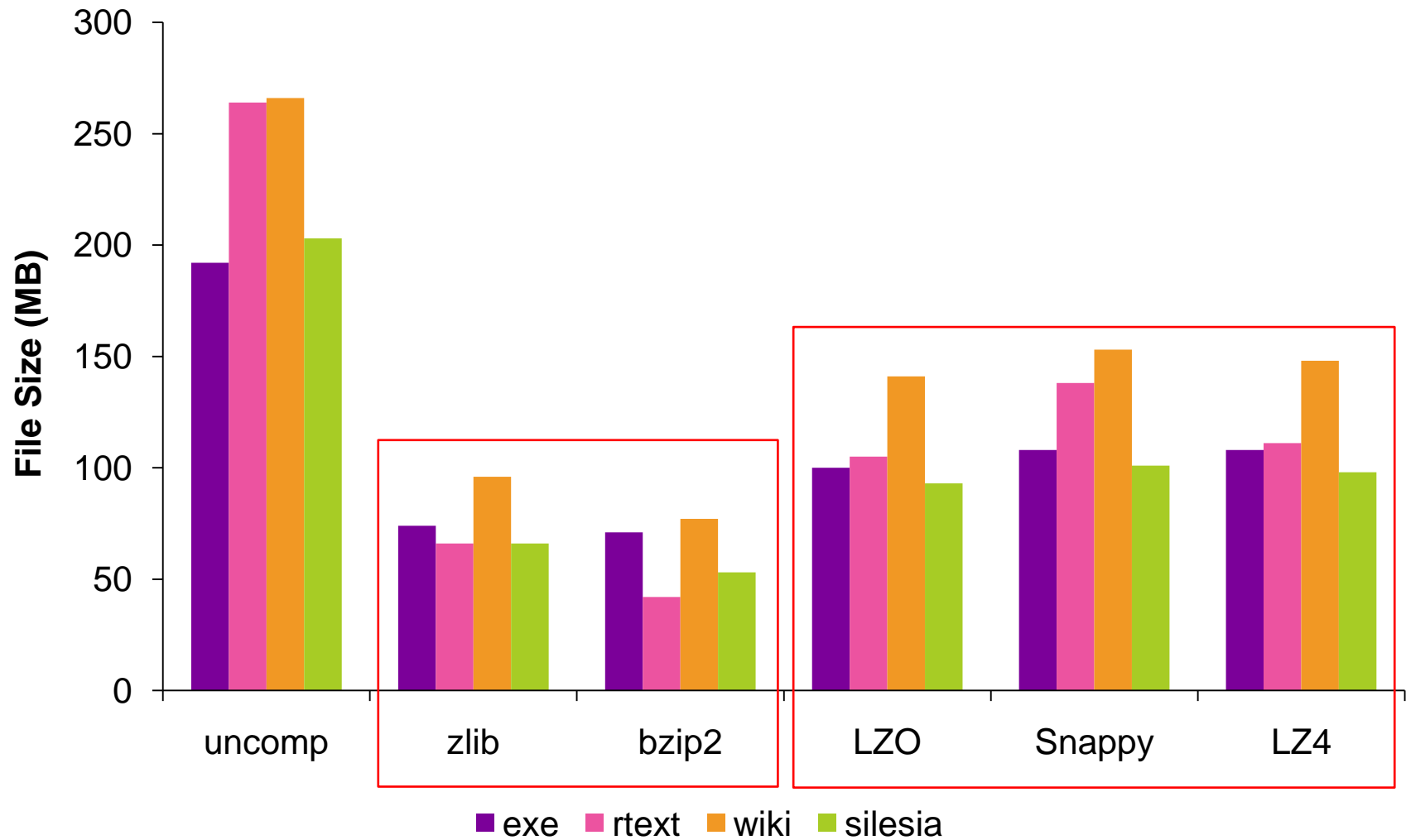
Measuring Standalone Performance

- Standard programs (gzip, bzip2) used
- Driver program written for other cases
- 32-bit mode
- Single-threaded
- JVM load overhead discounted
- Default compression level
- Quad-core Xeon machine

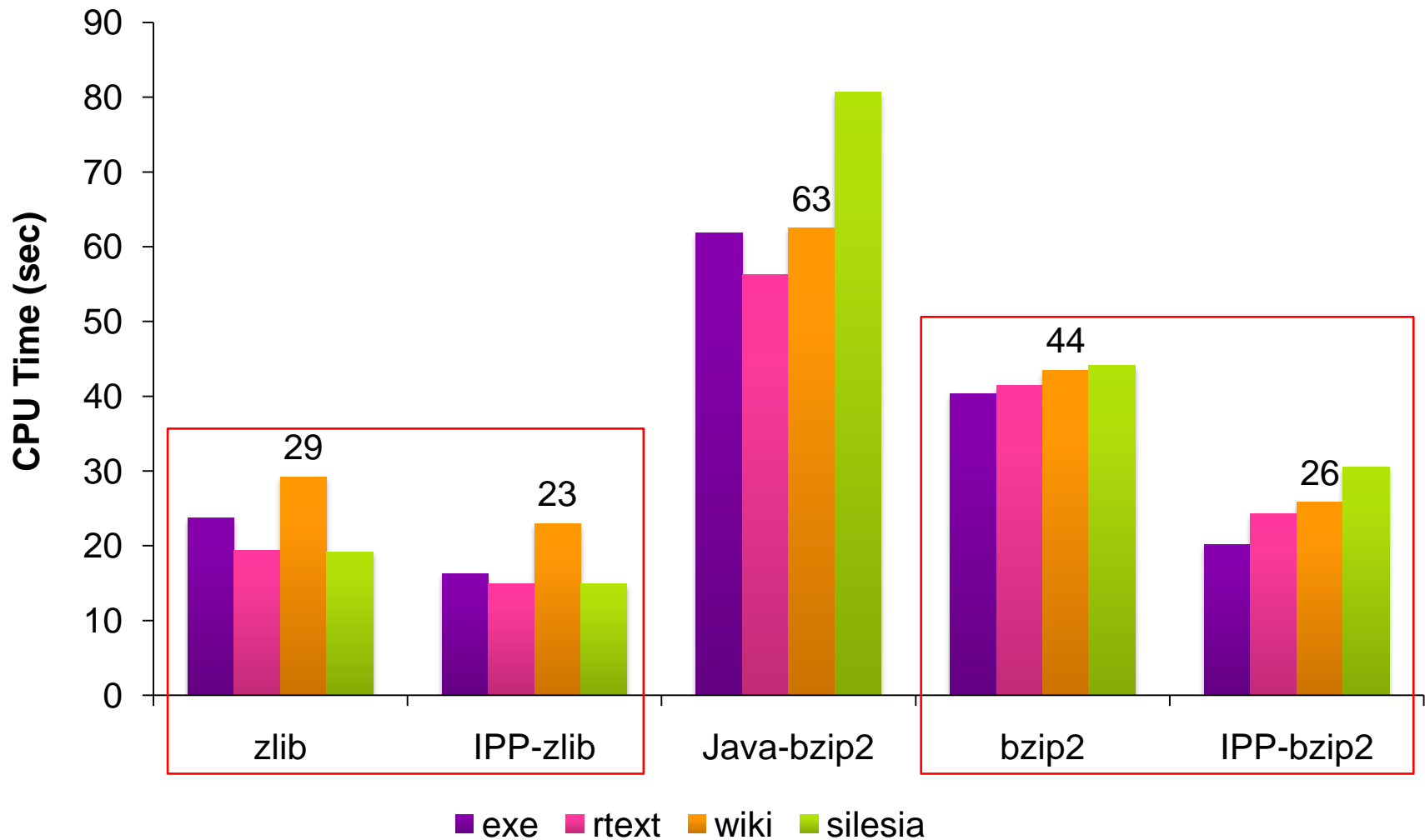
Data Corpora Used

- Binary files
- Generated text from randomtextwriter
- Wikipedia corpus
- Silesia corpus

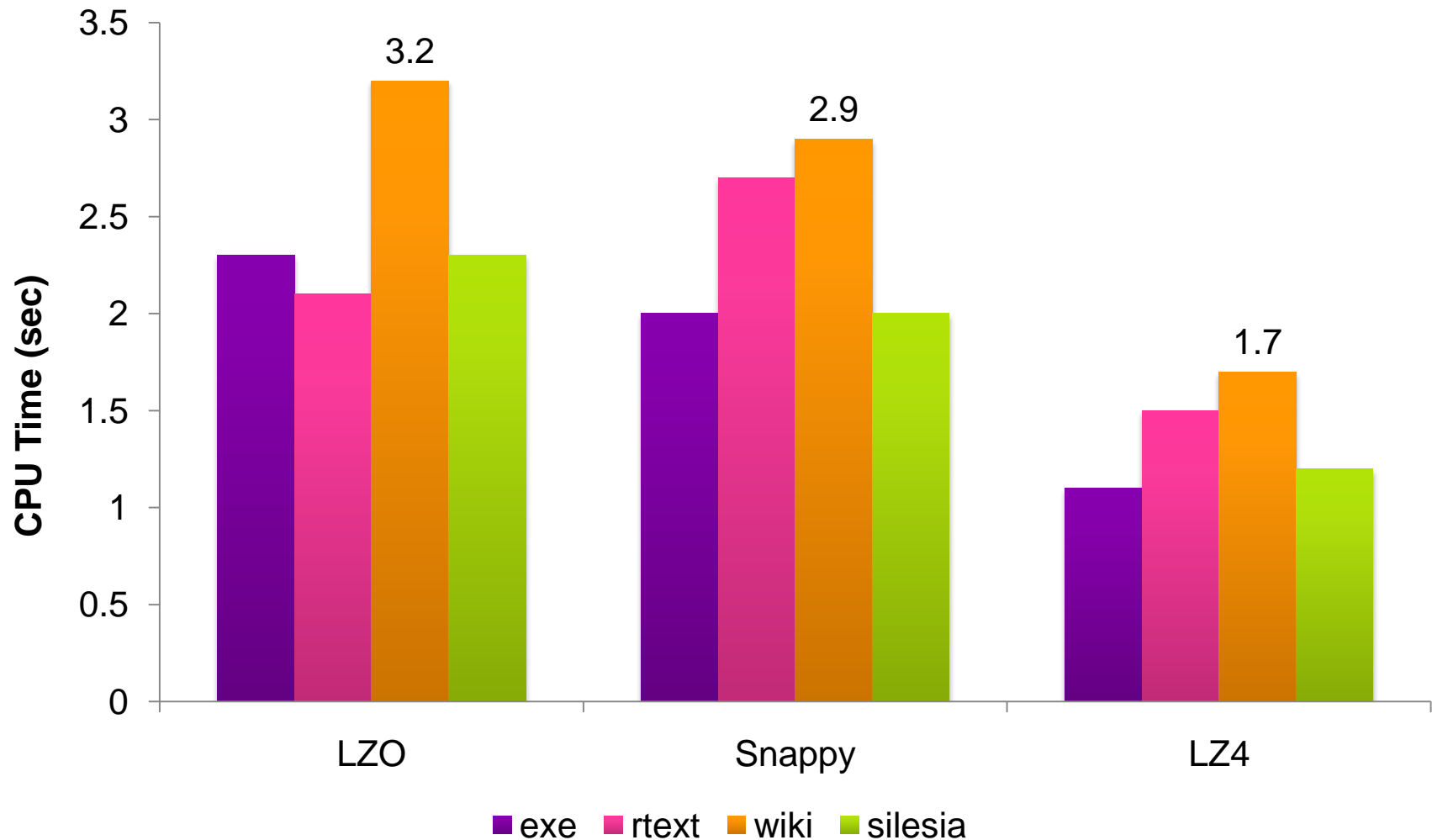
Compression Ratio



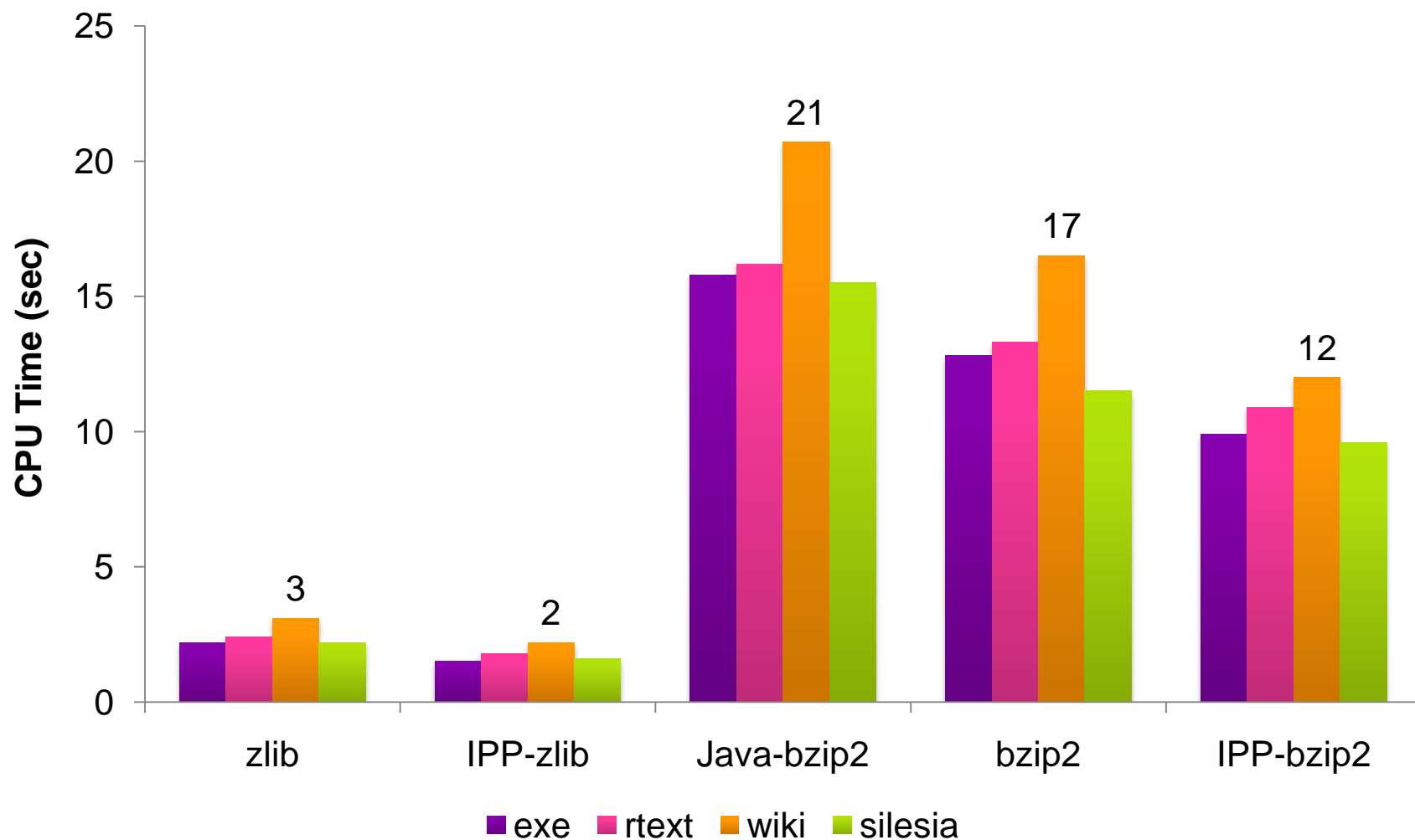
Compression Performance



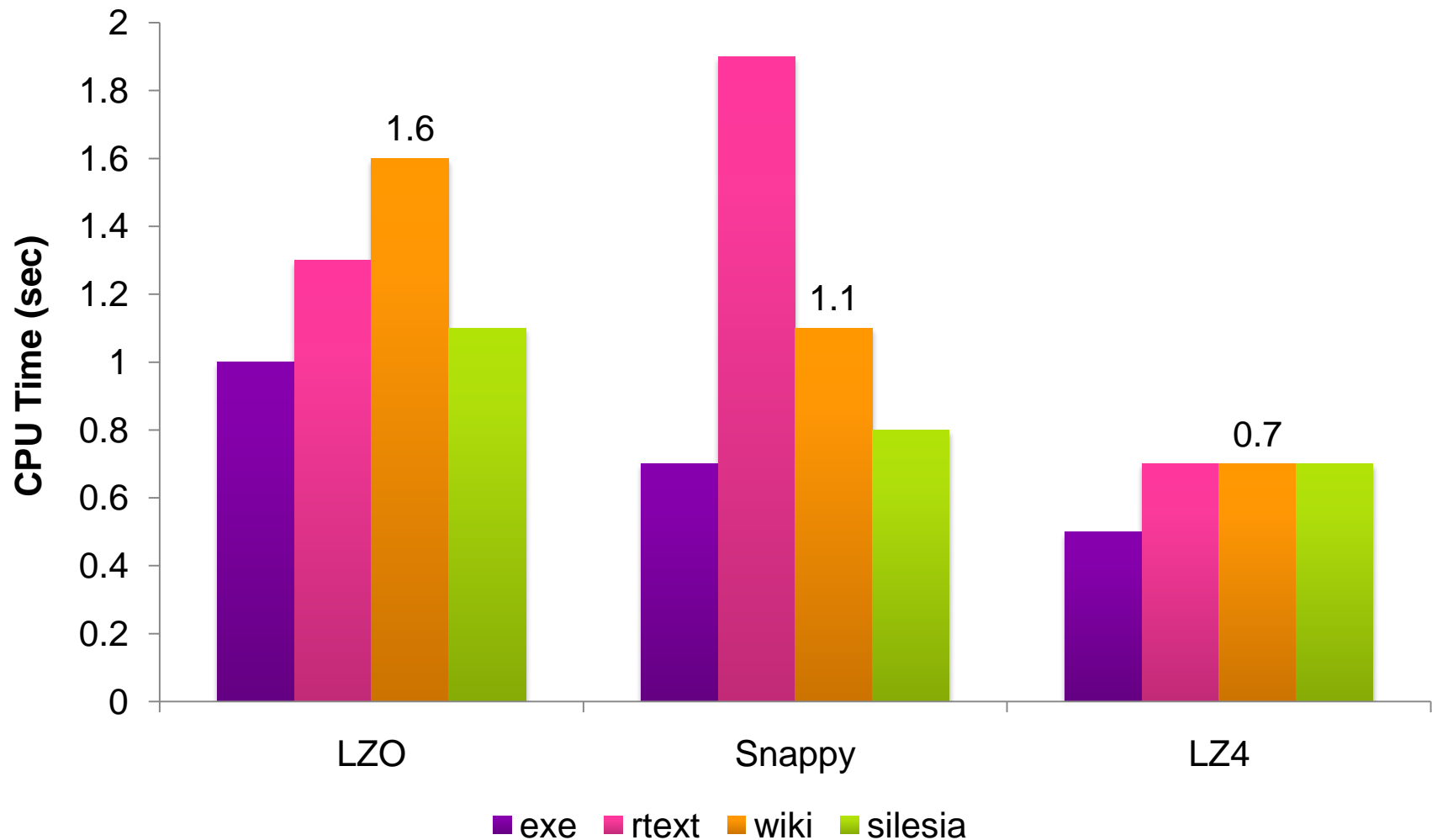
Compression Performance (Fast Algorithms)



Decompression Performance



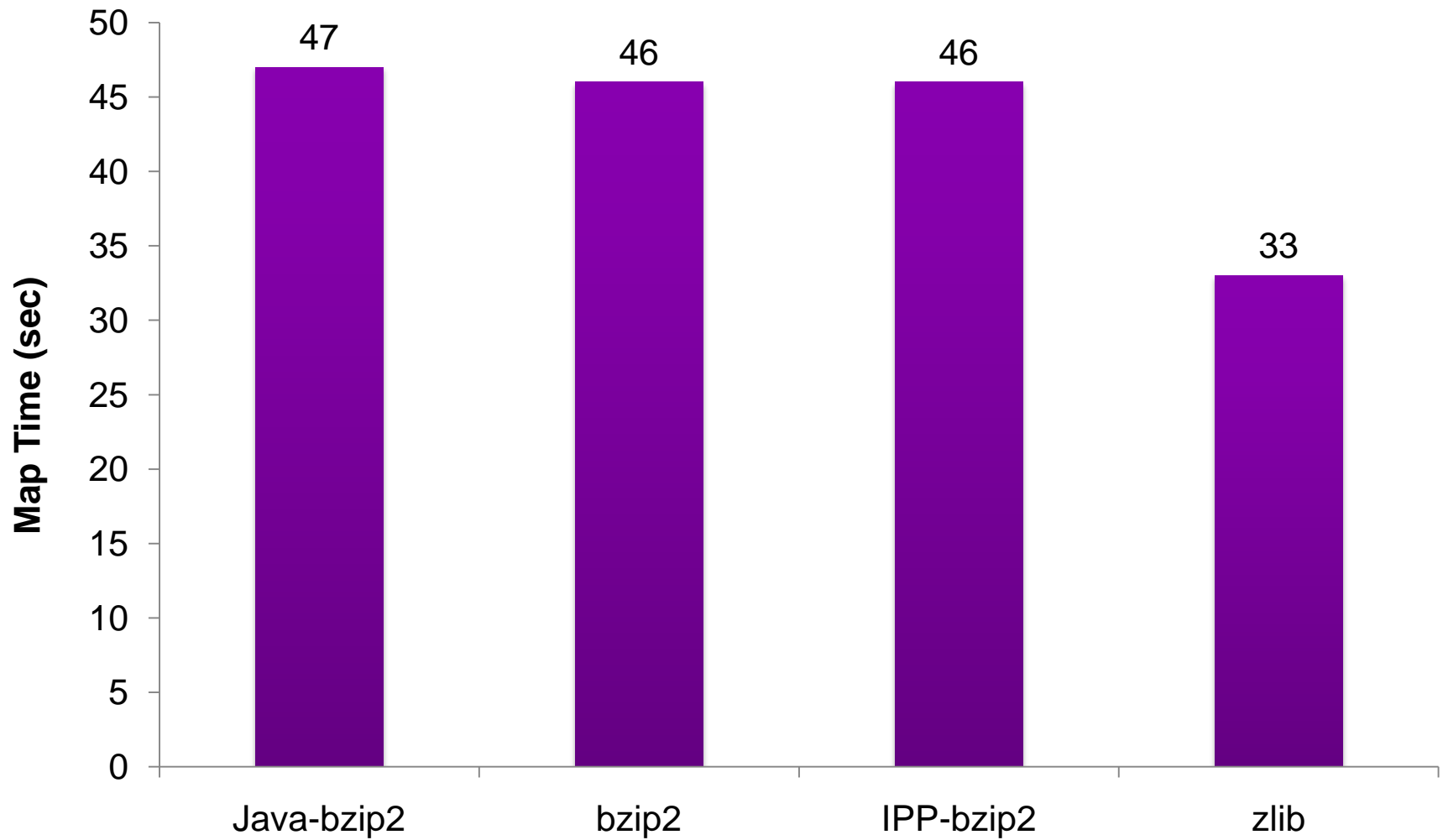
Decompression Performance (Fast Algorithms)



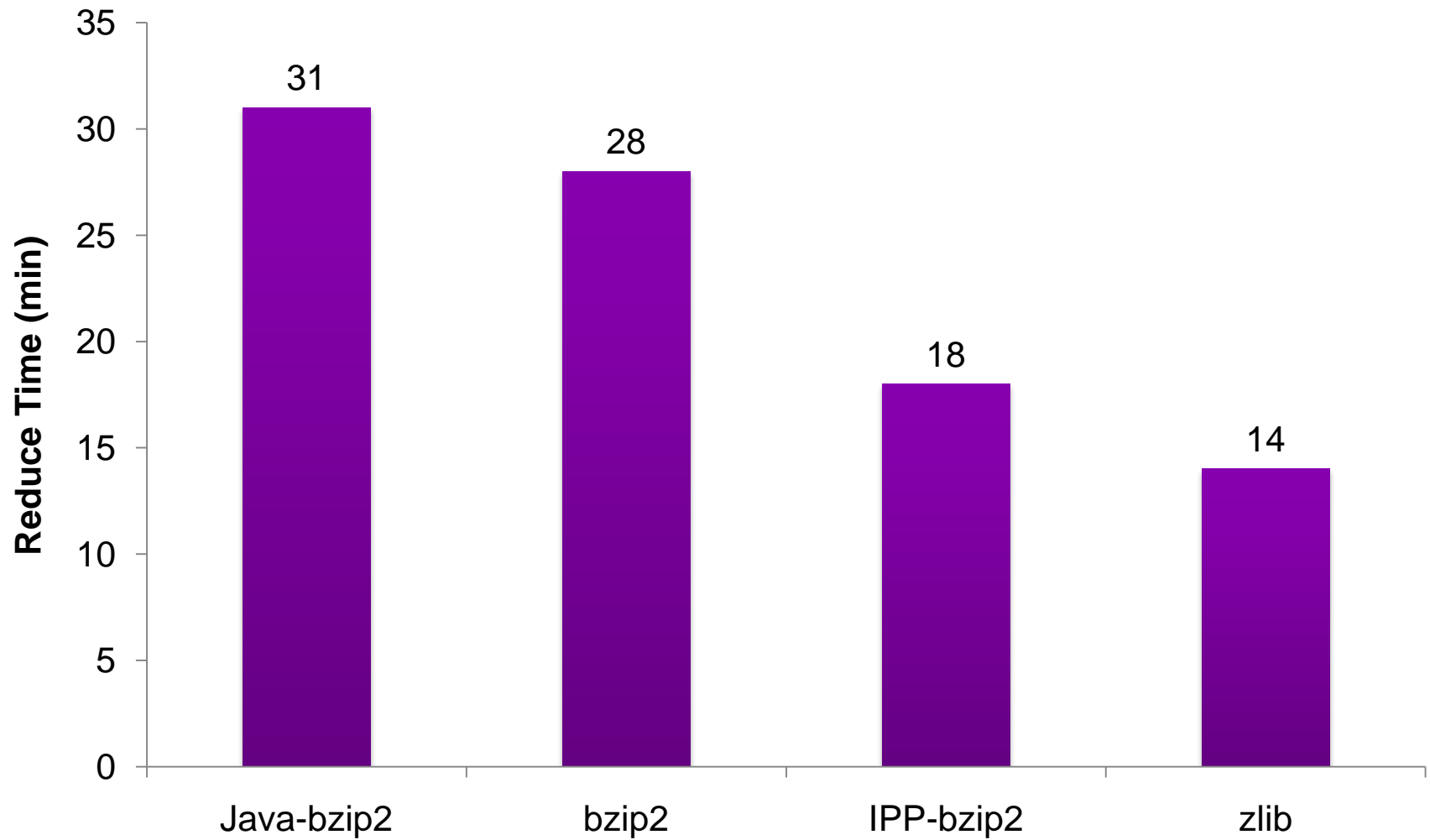
Compression Performance within Hadoop

- Daytona performance framework
- GridMix v1
- Loadgen and sort jobs
- Input data compressed with zlib / bzip2
- LZO used for intermediate compression
- 35 datanodes, dual-quad-core machines

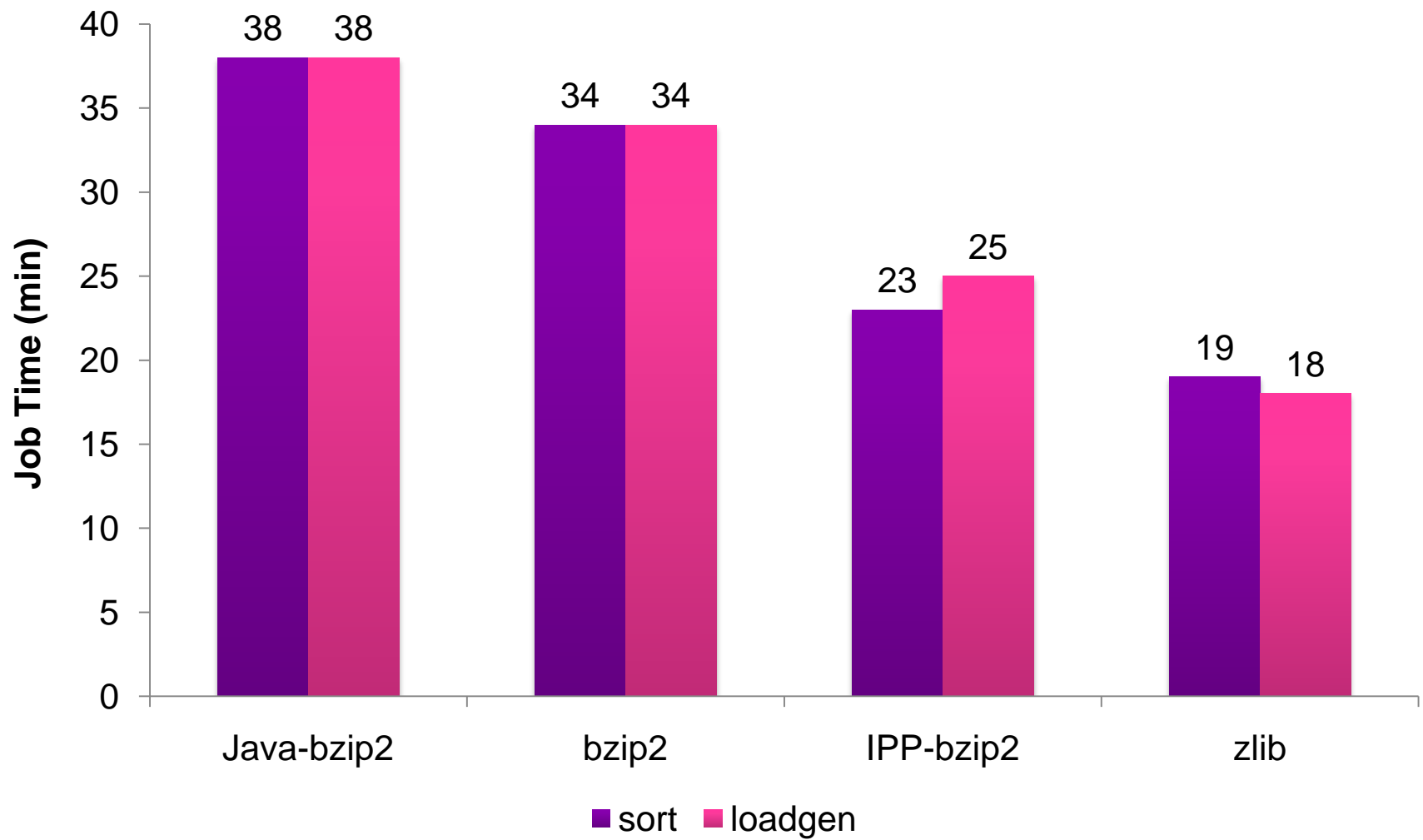
Map Performance



Reduce Performance



Job Performance



Future Work

- Splittability support for native-code bzip2 codec
- Enhancing Pig to use common bzip2 codec
- Optimizing the JNI interface and buffer copies
- Varying the compression effort parameter
- Performance evaluation for 64-bit mode
- Updating the zlib codec to specify alternative libraries
- Other codec combinations, such as zlib for transient data
- Other compression algorithms

Considerations in Selecting Compression Type

- Nature of the data set
- Chained jobs
- Data-storage efficiency requirements
- Frequency of compression vs. decompression
- Requirement for compatibility with a standard data format
- Splittability requirements
- Size of the intermediate and final data
- Alternative implementations of compression libraries

YAHOO!