

Data Compression Techniques for Big Data

¹Ms.Poonam Bonde, ²Mr. Sachin Barahate

¹P.G Student, ²Assistent Professor in I.T. Department

¹Student of YTGIOFOE, Mumbai, India

²Padmabhushan Vasantdada Patil Pratishthan's College Of Engineering, Sion, Mumbai, India

Abstract—Due to rapidly increasing size of data which comes from different heterogeneous sources, storage, maintenance, processing & real time transfer of big data is becoming a challenging task. Many times data files contain irrelevant and redundant data which can be removed to reduce size of file manageable. Various techniques to achieve this task are called compression techniques. These techniques can be utilized to text, image, audio, video etc data. Compressed data not only increases transfer speed but also improves knowledge discovery as well as processing time.

Index Terms—Compression, text compression, image compression, audio-video compression, big data compression, lossy-lossless compression

I. INTRODUCTION

What is Big Data:

Big data means huge amount of data which maybe in structured or in unstructured format .This data is very complex and in high volume so it becomes very difficult for traditional data processing applications.

Challenges or issues due to lage size of data:

1. Storage and transport issues
2. Data Management Issues
3. Processing issues
4. Data access and information sharing
5. Relevant data searching
6. Data analytics

Compression is the process of reduction in size of data in order to save space or transmission time.

Reasons to Compress

- To Reduce File Size
- To Save storage space
- To Increase transfer speed at a given data rate
- To Allow real-time transfer at a given data rate

lossless algorithms, which can reconstruct the original message exactly from the compressed message as they don't involve loss of information, and lossy algorithms, which can only reconstruct an approximation of the original message as they involves certain amount of loss of information.

II. LITERATURE REVIEW

Compression algorithms for Text data

1. Arithmetic coding
2. Huffman Algorithm
3. LZW
4. Run length encoding

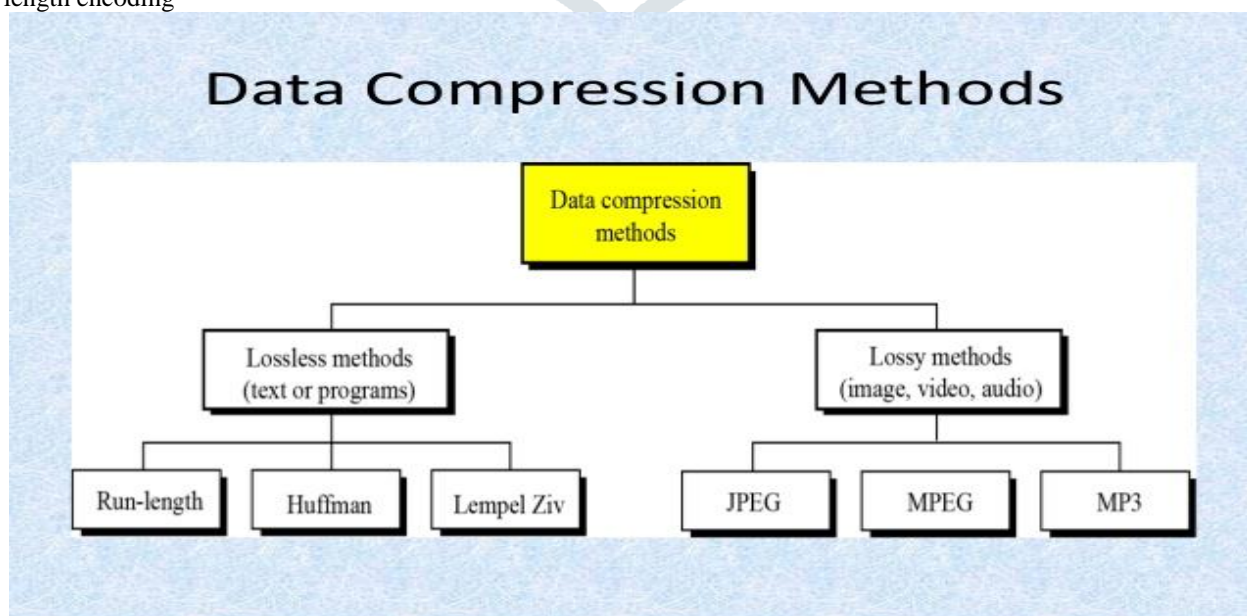


Fig. Data Compression Methods

The general purpose of data compression algorithms on text files is to convert a string of characters into a new string which contains the same information but with new length as small as possible.

Compression algorithms for image data

Methods for lossless image compression

1. Run length encoding:
2. Entropy encoding
3. Adaptive dictionary algorithms such as LZW – used in GIF and TIFF

Methods for lossy image compression

1. Chroma subsampling
2. Transform coding
3. Fractal compression

Compression algorithms for audio data

1. MPEG

Compression algorithms for video data

1. MPEG
2. H.264

Basic types of compression techniques

1. Lossy
2. Lossless[1]

Combining Compression techniques in big data

We enlisted above some issues regarding big data. Main issues are storage space & transformation cost. Using compression techniques in big data to reduce size of files can provide solution in some extent. We are focussing here on text & image compression of big data.

Text Compression

Burrows Wheeler transform

The Burrows-Wheeler Algorithm was published in the year 1994 by Michael Burrows and David Wheeler in the research report “A Block-sorting Lossless Data Compression Algorithm”. This research report is based on an unpublished work by David Wheeler from the year 1983.^[2]

Structure of the Burrows-Wheeler Algorithm

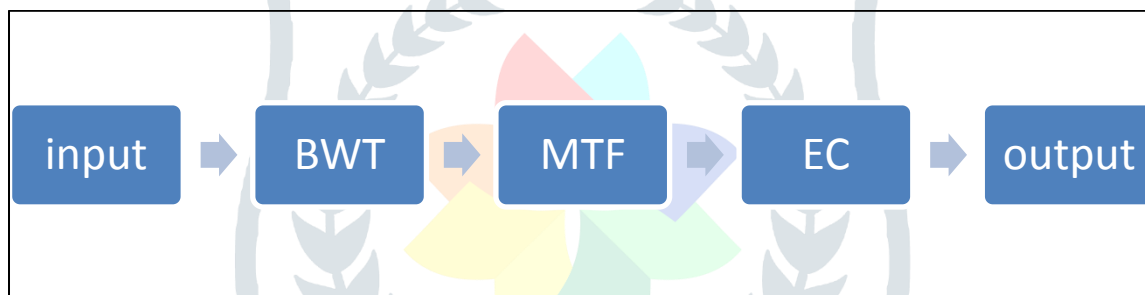


Fig. Structure of Burrows Wheeler Algorithm

The original algorithm consists of three stages.

- 1) Transform Coding
- 2) Move-To-Front Transform.
- 3) Zero Run Length Coding.
- 4) Backward Search Algorithm

Image Compression

Image compression is the process of transforming an image file in such a way that it consumes less space than the original file. It is a compression technique that reduces the size of an image file without degrading its quality to a greater extent.

Anamorphic Stretch Transform:

An anamorphic stretch transform (AST) is a mathematical transformation in which analog or digital data is stretched and warped in a specific fashion such that after down sampling, the volume of data is reduced without loss of pertinent information.^[3]

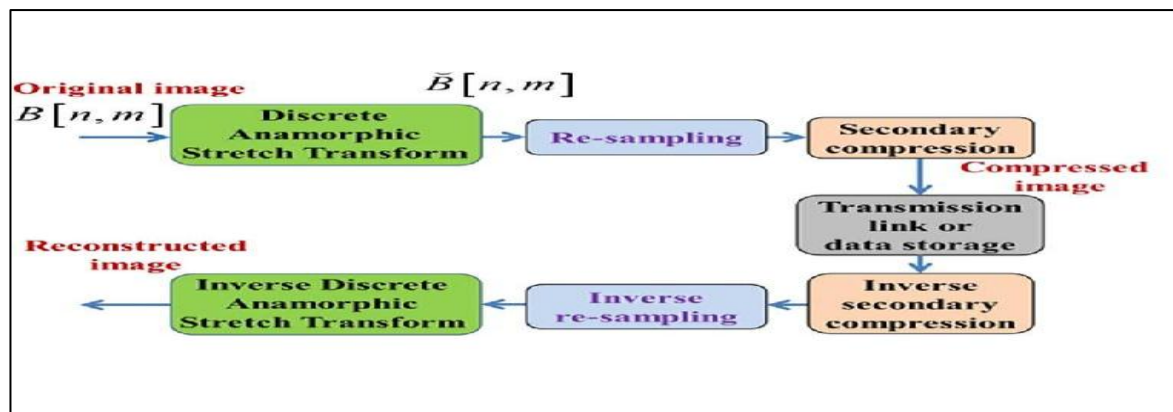


Fig. Different steps for implementation of DAST for application to image compression

The discrete implementation of this technique, dubbed Discrete Anamorphic Stretch Transform (DAST) represents a new solution to image compression [4]. The transform reshapes the image before uniform re-sampling in such a way that sharp features are naturally stretched more than coarse features. This causes sharp features to experience higher sampling density than coarse features.

III. FRAMEWORK OF OUR APPROACH

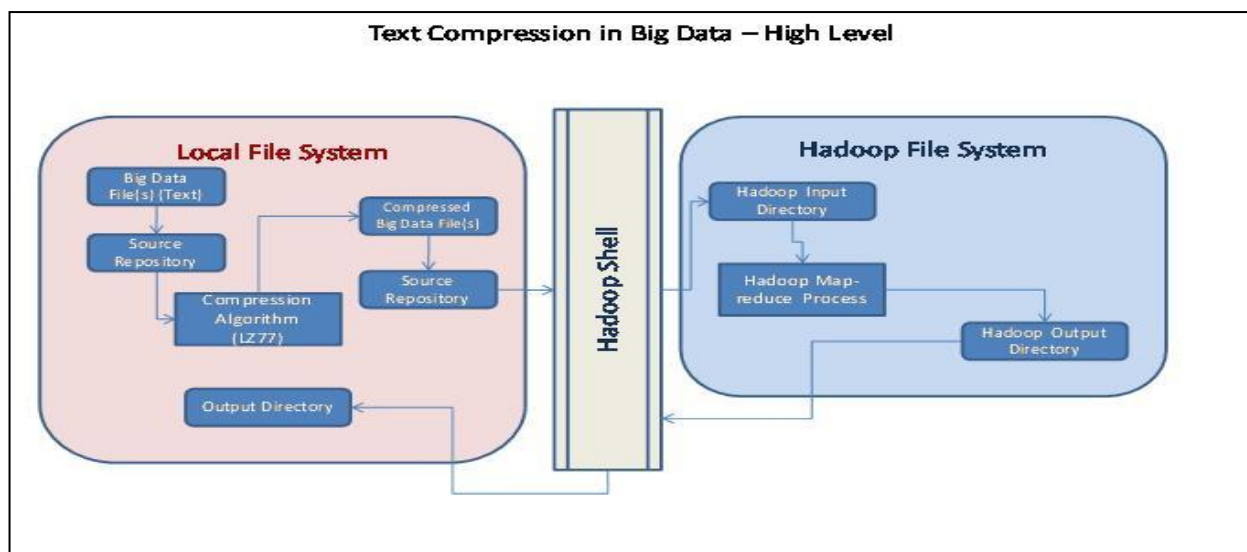


Fig . Text Compression in Big Data

Steps:

1. Get big text source file(s) (public classified information only) (input from various online social networking sites, ecommerce transaction files, supply chain transaction files, geospatial records, healthcare information etc).
2. Place the file(s) in source directory of local file system (In my case, I have used Ubuntu operating system).
3. Use a compression program (Working based on LZ77, Burrows wheeler algorithm) to compress the input file(s).
4. Using Hadoop Shell commands; transfer the compressed file to Hadoop file system.
5. Run the analytical program to analyze the files.
6. The analytics program works just fine as in case of uncompressed file (This is an observation).
7. Get the output files from the Hadoop file system to Local file system to read conveniently.

Image Compression:

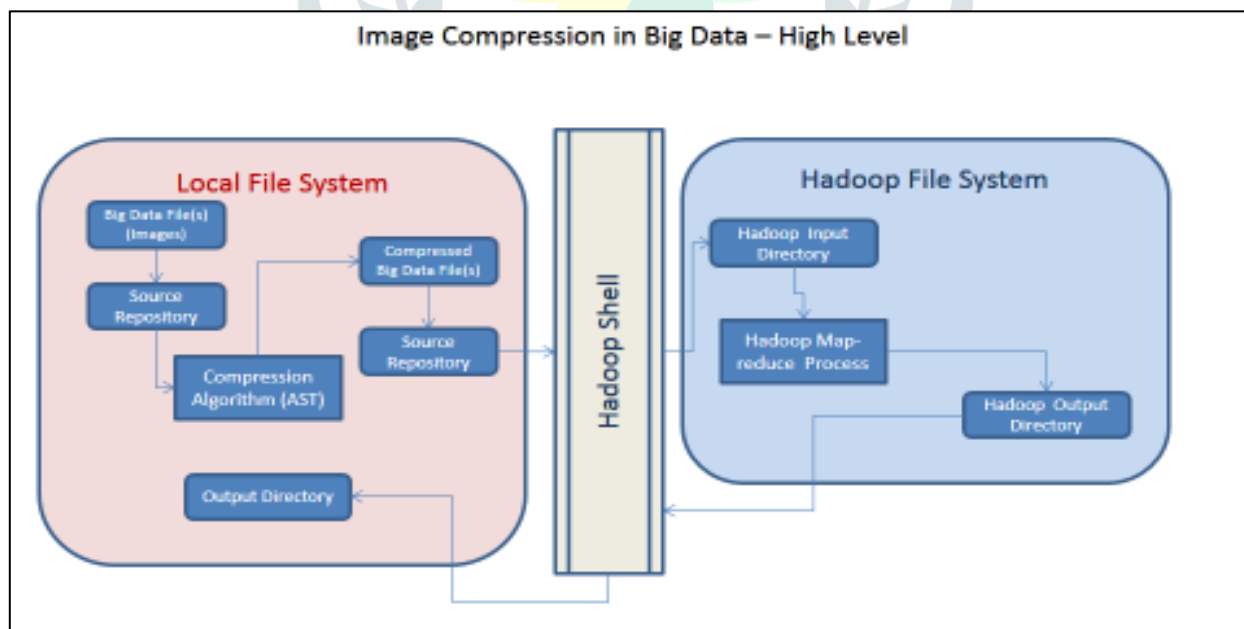


Fig. Image Compression in Big Data

Steps:

1. Get big source image file(s) (public classified information only) (input from various online social networking sites, ecommerce transaction files, supply chain transaction files, geospatial records, healthcare information etc).
2. Place the file(s) in source directory of local file system (In my case, I have used Ubuntu operating system).
3. Use a compression program which works on Anamorphic Stretch Transform to compress the input file(s).
4. Using Hadoop Shell commands; transfer the compressed file to Hadoop file system.
5. Run the analytical program to analyze the files.

Observations of using text compression techniques in Big Data using Hadoop

In my experiment, I took a large text file, using Hadoop file system, I transferred the file from local to Hadoop file system, ran a simple wordcount analytical program to check the results.

The observations are as follows

The compression technique used for following observations is modified LZ77 which is used by GZIP utility.

1. Environment Preparation and File Transfer

```

poonam@poonam-Compaq-Presario-C700-Notebook-PC: /hdc
poonam@poonam-Compaq-Presario-C700-Notebook-PC:/hdc$ . prepenv.sh > prepenv.result
14/11/30 13:35:48 INFO fs.TrashPolicyDefault: Namenode trash configuration: Deletion interval = 0 minutes, Emptier interval = 0 minutes.
rm: Cannot delete /user/poonam/sen_text_compressed. Name node is in safe mode.
14/11/30 13:35:48 INFO fs.TrashPolicyDefault: Namenode trash configuration: Deletion interval = 0 minutes, Emptier interval = 0 minutes.
rm: Cannot delete /user/poonam/sen_text_compressed-out. Name node is in safe mode.
14/11/30 13:35:48 INFO fs.TrashPolicyDefault: Namenode trash configuration: Deletion interval = 0 minutes, Emptier interval = 0 minutes.
rm: Cannot delete /user/poonam/sen_text_uncompressed. Name node is in safe mode.
14/11/30 13:35:48 INFO fs.TrashPolicyDefault: Namenode trash configuration: Deletion interval = 0 minutes, Emptier interval = 0 minutes.
rm: Cannot delete /user/poonam/sen_text_uncompressed-out. Name node is in safe mode.
mkdir: /user/poonam/sen_text_uncompressed: File exists
mkdir: /user/poonam/sen_text_compressed: File exists
rm: cannot remove '/hdc/sen_text_compressed/*': No such file or directory
rm: cannot remove '/hdc/sen_text_uncompressed-out/*': No such file or directory
rm: cannot remove '/hdc/sen_text_compressed-out/*': No such file or directory
poonam@poonam-Compaq-Presario-C700-Notebook-PC:/hdc$ . filetransfer.sh > filetransfer.result
rm: cannot remove '/hdc/sen_text_compressed/*.*': No such file or directory
14/11/30 13:36:24 INFO fs.TrashPolicyDefault: Namenode trash configuration: Deletion interval = 0 minutes, Emptier interval = 0 minutes.
14/11/30 13:36:27 INFO fs.TrashPolicyDefault: Namenode trash configuration: Deletion interval = 0 minutes, Emptier interval = 0 minutes.
tar: Removing leading '/' from member names
poonam@poonam-Compaq-Presario-C700-Notebook-PC:/hdc$
    
```

Fig. Environment Preparation & File Transfer

2. Results of File Transfer

```

poonam@poonam-Compaq-Presario-C700-Notebook-PC: /hdc
-----
Removing all files from Hadoop File System and compressed files from Local File System
-----
Deleted /user/poonam/sen_text_compressed
Deleted /user/poonam/sen_text_uncompressed
File system has been cleaned
-----
Moving Uncompressed File to Hadoop File System
-----
Start Time:Sun Nov 30 13:36:28 IST 2014
Uncompressed File has been transferred to Hadoop File System
End Time:Sun Nov 30 13:37:11 IST 2014
-----
Compressing File on Local File System using tar.gz format
-----
Start Time:Sun Nov 30 13:37:11 IST 2014
/hdc/sen_text_uncompressed/Uly.txt
End Time:Sun Nov 30 13:37:27 IST 2014
Compressed the File on Local File System using tar.gz format
-----
Transferring gz File to Hadoop File System
-----
Start Time:Sun Nov 30 13:37:27 IST 2014
End Time:Sun Nov 30 13:37:33 IST 2014
Compressed File has been transferred to Hadoop File System
-----
~
~
    
```

Fig. Results of File Transfer

3. Results of running the analytics (wordcount program in this case)

```

poonam@poonam-Compaq-Presario-C700-Notebook-PC: /hdc
-----
Running wordcount analytics on uncompressed File
Start Time:Sun Nov 30 13:38:59 IST 2014
End Time:Sun Nov 30 13:41:25 IST 2014
-----
Running wordcount analytics on compressed File
Start Time:Sun Nov 30 13:41:25 IST 2014
End Time:Sun Nov 30 13:43:48 IST 2014
-----
~
~
~
~
    
```

Fig. Results of running the analytics

Time

File Transfer			
File Type	Start Time	End Time	Time Taken by Process
Uncompressed	13:36:28	13:37:11	0:00:43
Compressed	13:37:27	13:37:33	0:00:06

Fig. File Transfer Time

File Compression	
Start Time	13:37:11
End Time	13:37:27
Time taken by Compression Process	0:00:16

Fig. File Compression

Analytics Processing Time			
	Start Time	End Time	Process Time
Uncompressed File	13:38:59	13:41:25	0:02:26
Compressed File	13:41:25	13:43:48	0:02:23

Fig. Processing Time

Total Time	
Uncompressed File	0:03:09
Compressed File	0:02:45

Time saved due to compression	0:00:24
-------------------------------	---------

Fig. Total Time

Size

File Size	Bytes	MBs
Before Compression	552311760	526.7255
After Compression	4743196	4.523464

Fig. Size of files

Outputs

Uncompressed File

```

poonam@poonam-Compaq-Presario-C700-Notebook-PC: /hdc
type="text/javascript" 25410
type="text/javascript"> 25410
type="text/plain" 25410
typeof="pgterms:agent" 25410
typeof="pgterms:ebook" 25410
typeof="pgterms:file"&gt; 101640
types 25410
url(/www.gutenberg.org/cache/epub/4300/pg4300.qrcode.png) 25410
url(/pics/sprite.png?1416701893) 25410
us 50820
use." 25410
value="" 25410
value="XKAL6BZL3VPSN" 25410
value="s_xclick" 25410
var 228690
version="XHTML+RDFa" 25410
we 25410
which 25410
why 25410
without 50820
xml:lang="en" 50820
xml:lang="en"&gt;14&lt;/th&gt; 25410
xml:lang="en"&gt;#74&lt;/td&gt; 25410
xml:lang="en"&gt;#&lt;/th&gt; 25410
xml:lang="en"&gt;( 25410
xml:lang="en"&gt;&lt;/th&gt; 25410
xml:lang="en"&gt;&lt;/th&gt; 25410
xml:lang="en"&gt;a, shakespeare&lt;/td&gt; 25410
xml:lang="en"&gt;cat.&lt;/th&gt; 25410
xml:lang="en"&gt;jane 25410
xml:lang="en"&gt;juvenile 25410
xml:lang="en"&gt;l.&lt;/th&gt; 25410
xml:lang="en"&gt;love 50820
xml:lang="en"&gt;n.&lt;/th&gt; 25410
xml:lang="en"&gt;qui.&lt;/td&gt; 25410
xml:lang="en"&gt;s.&lt;/th&gt; 25410
xml:lang="en"&gt;s, shakespeare&lt;/td&gt; 25410
xml:lang="en"&gt;shakespeare 25410
xml:lang="en"&gt;t.&lt;/th&gt; 25410
xml:lang="en"&gt;verne 25410
xml:lang="en"&gt;|&lt;/th&gt; 25410
xml:lang="en_US" 25410
xmlns:dcterms="http://purl.org/dc/terms/" 25410
xmlns:ebook="http://www.gutenberg.org/ebooks/" 25410
xmlns:marcrel="https://www.loc.gov/loc/terms/relators/" 25410
xmlns:og="http://opengraphprotocol.org/schema/" 25410
xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#" 25410
xmlns:xsd="http://www.w3.org/2001/XMLSchema#" 25410
xmlns="http://www.w3.org/1999/xhtml" 25410
xx_XX 25410
your 50820
[ 50820
] 25410
) 50820
poonam@poonam-Compaq-Presario-C700-Notebook-PC: /hdc$
    </pre>
</div>
<div data-bbox="412 937 580 953" data-label="Caption">
<p>Fig. Uncompressed File</p>
</div>
<div data-bbox="52 967 170 983" data-label="Page-Footer">
<p>JETIR1703041</p>
</div>
<div data-bbox="183 967 839 983" data-label="Page-Footer">
<p>Journal of Emerging Technologies and Innovative Research (JETIR) <a href="http://www.jetir.org">www.jetir.org</a></p>
</div>
<div data-bbox="883 967 923 982" data-label="Page-Footer">
<p>197</p>
</div>
```


Compressed File

```

poonam@poonam-Compaq-Presario-C700-Notebook-PC: /hdc
type="text/javascript" 25410
type="text/javascript"> 25410
type="text/plain" 25410
typeof="pgterms:agent" 25410
typeof="pgterms:ebook" 25410
typeof="pgterms:file"&gt; 101640
types 25410
url(/www.gutenberg.org/cache/epub/4300/pg4300_qrcode.png) 25410
url(/pics/sprite.png:1416701893) 25410
us 50820
use." 25410
value=" 25410
value="XKAL6BZL3VPSN" 25410
value="_s_xclick" 25410
var 228690
version="HTML-RDFa" 25410
we 25410
which 25410
why 25410
without 50820
xml:lang="en" 50820
xml:lang="en"&gt;|&lt;/th&gt; 25410
xml:lang="en"&gt;#74&lt;/td&gt; 25410
xml:lang="en"&gt;#&lt;/th&gt; 25410
xml:lang="en"&gt;(&lt; 25410
xml:lang="en"&gt;&lt;/th&gt; 25410
xml:lang="en"&gt;a.&lt;/th&gt; 25410
xml:lang="en"&gt;cat.&lt;/th&gt; 25410
xml:lang="en"&gt;jane 25410
xml:lang="en"&gt;juvenile 25410
xml:lang="en"&gt;l.&lt;/th&gt; 25410
xml:lang="en"&gt;love 50820
xml:lang="en"&gt;on.&lt;/th&gt; 25410
xml:lang="en"&gt;qui.&lt;/td&gt; 25410
xml:lang="en"&gt;s.&lt;/th&gt; 25410
xml:lang="en"&gt;s.shakespeare&lt;/td&gt; 25410
xml:lang="en"&gt;t.&lt;/th&gt; 25410
xml:lang="en"&gt;verne 25410
xml:lang="en"&gt;|&lt;/th&gt; 25410
xml:lang="en US 25410
xmlns:dcterms="http://purl.org/dc/terms/" 25410
xmlns:ebook="http://www.gutenberg.org/ebooks/" 25410
xmlns:marc-rel="http://www.loc.gov/loc/terms/relators/" 25410
xmlns:og="http://opengraphprotocol.org/schema/" 25410
xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#" 25410
xmlns:xsd="http://www.w3.org/2001/XMLSchema#" 25410
xmlns="http://www.w3.org/1999/xhtml" 25410
xx_xx 25410
your 50820
} 50820
| 25410
) 50820
} 50820
poonam@poonam-Compaq-Presario-C700-Notebook-PC: /hdc$
</pre>
</div>
<div data-bbox="422 353 573 369" data-label="Caption">Fig. Compressed File</div>
<div data-bbox="19 381 193 395" data-label="Section-Header"><b>Results and Conclusion</b></div>
<div data-bbox="79 395 680 492" data-label="List-Group">
<ol>
<li>1. Compressed files take less time to transfer from local system to Hadoop file system.</li>
<li>2. Both Compressed files and uncompressed files yield the similar results of analytics.</li>
<li>3. Hence use of compression techniques is feasible in Big Data by using Hadoop.</li>
<li>4. Following Algorithms can be used in Hadoop for text compression
<ol type="a">
<li>a. LZ77 (Used by Gzip in unix)</li>
<li>b. Burrows Wheeler (Used by Bzip2 in unix)</li>
<li>c. Along with these according to Hadoop documents, it can use Snappy.</li>
</ol>
</li>
</ol>
</div>
<div data-bbox="19 491 310 506" data-label="Section-Header"><b>Benefits of File Compression in Hadoop</b></div>
<div data-bbox="79 505 693 533" data-label="List-Group">
<ol>
<li>1. Requires less storage space.</li>
<li>2. File transfer and ultimately the analytics are faster compared with uncompressed files.</li>
</ol>
</div>
<div data-bbox="19 532 117 546" data-label="Section-Header"><b>Work Ahead</b></div>
<div data-bbox="79 546 700 560" data-label="List-Group">
<ol>
<li>1. Using AST (Anamorphic Stretch Transform) to compress images and use it in Hadoop.</li>
</ol>
</div>
<div data-bbox="19 572 138 587" data-label="Section-Header"><b>REFERENCES</b></div>
<div data-bbox="49 586 977 683" data-label="List-Group">
<ol>
<li>[1] International Journal of Engineering Inventions e-ISSN: 2278-7461, p-ISSN: 2319-6491 Volume 2, Issue 4 (February 2013) www.ijejournal.com<br/>Image Compression Techniques: A Survey Athira B. Kaimal, S. Manimurugan, C.S.C .Devadass</li>
<li>[2] The Burrows-Wheeler Algorithm Daniel Schiller August 5, 2012</li>
<li>[3] <a href="http://en.wikipedia.org/wiki/Anamorphic_stretch_transform">http://en.wikipedia.org/wiki/Anamorphic_stretch_transform</a></li>
<li>[4] B. Jalali and M. H. Asghari, "The anamorphic stretch transform: putting the squeeze on Big Data", Optics and Photonics News Magazine, February issue cover story, Optical Society of America, Feb. 2014</li>
</ol>
</div>
<div data-bbox="51 966 169 983" data-label="Page-Footer">JETIR1703041</div>
<div data-bbox="185 966 837 983" data-label="Page-Footer">Journal of Emerging Technologies and Innovative Research (JETIR) <a href="http://www.jetir.org">www.jetir.org</a></div>
<div data-bbox="876 966 916 982" data-label="Page-Footer">198</div>
```