# PARALLEL ZIP

**Ajinkya Pawar    Jitender Kumar**
**Palak Purohit     Sagar Bisen**

**Instructor - Prof. Nipun Batra**
**Mentor - Rishiraj Adhikari**

# INDEX

# MOTIVATION

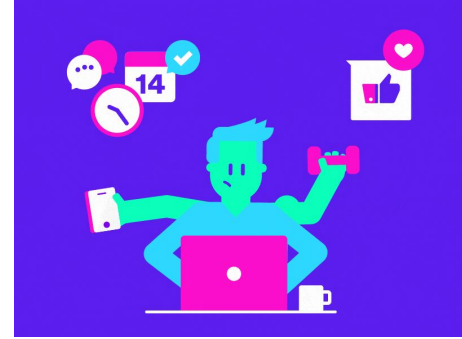❖ **Why zip files parallely?**

- **Multi-threaded programs have several advantages**

  - **Makes a single process faster**
  - **Threads don't require new address spaces**

- **Consequence**
  - **Improved CPU performance**
    - **Reduced Turnaround Time**

**Already stressed out CPU**



**We try to give some relief**



3

# APPROACH

- **Compression Algorithms**

- **Sequential zip**

- **Parallel Zip**

- **BETTER PARALLEL ZIP**

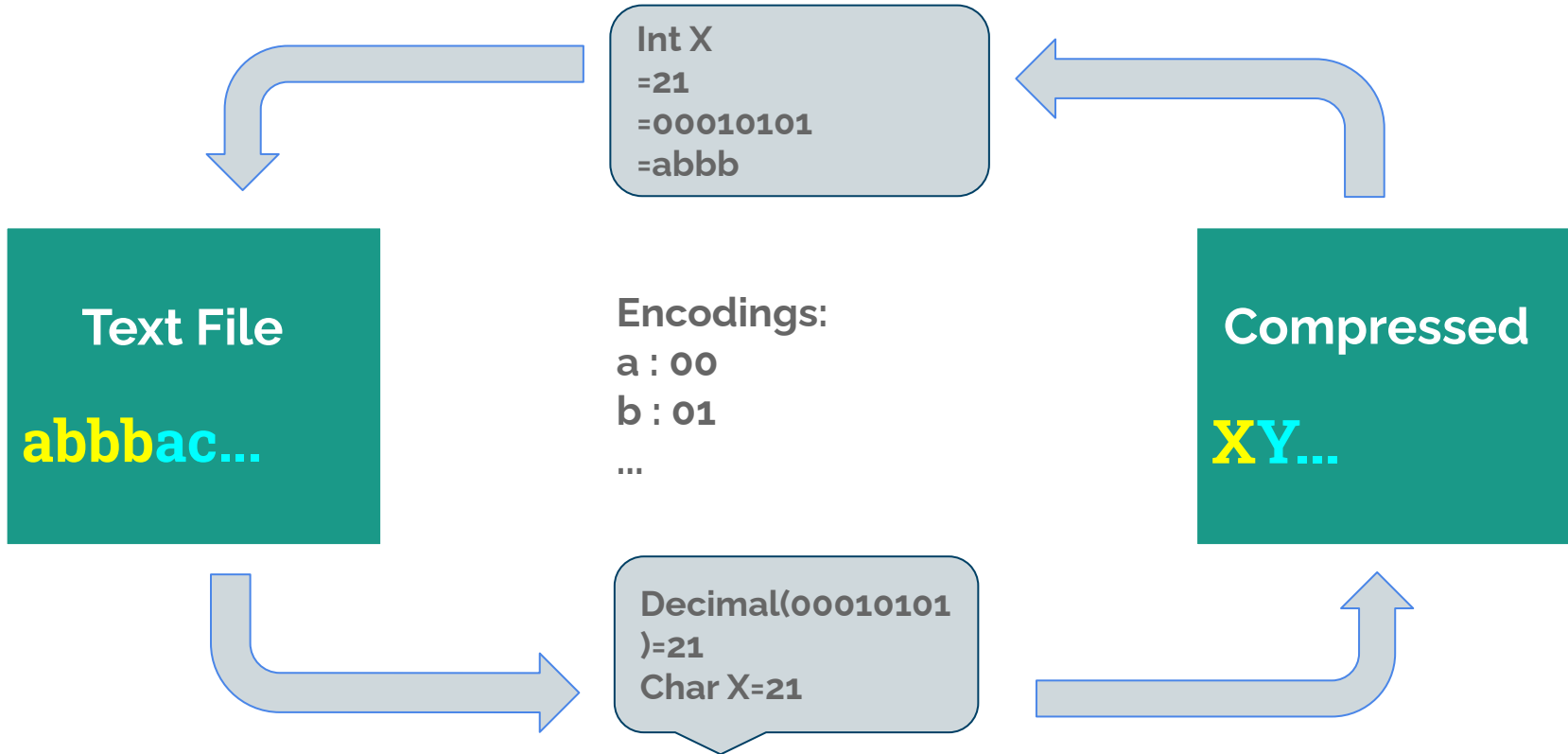**Well, keep an eye on what's attracting your eyes here**

# Two encoding schemes

**Huffman Encoding**

**Run length encoding(RLE)**

# COMPRESSION - HUFFMAN CODING



Int X
=21
=00010101
=abbb

**Text File**

**abbbac...**

Encodings:
a : 00
b : 01
...

**Compressed**

**XY...**

Decimal(00010101
)=21
Char X=21

# RUN LENGTH ENCODING (RLE)

hhllo → hhllo

hhllo → 2h2l1o

❖ **For run length >2, encoding makes sense**

wwwwbbb → wwwwbbbb

wwwwbbb → 4w4b
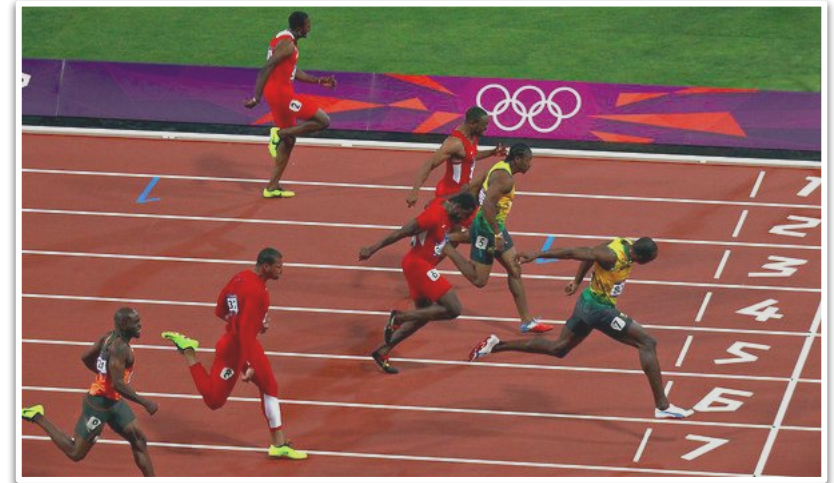
# INITIAL MODELS

- **SEQUENTIAL ZIP**
  - **Pseudo code**
    - **For i in range (no. of files)**
      **zip (ith file)**


- **PARALLEL ZIP**
  - **Pseudo code**
    - **Pthread_t threads[no. of files];**
    - **For i in range (no. of files)**
      **Pthread_create**
      **(ith thread, zip , ith file)**



sequential vs parallel zip

# STRAGGLERS PROBLEM

- **What are Stragglers ?**
- **Time**
  - **Sequential**
    - **File1 + File2 + ...**
  - **Parallel**
    - **Max ( core 1, core 2, ...)+ thread overhead.**
- **Can sequential be faster than parallel !?**

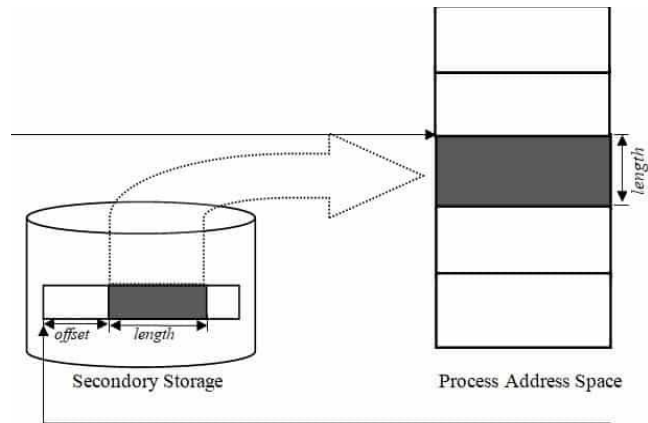# SOLUTION - PRODUCER & CONSUMERS


CAT-GIFs.com

- **Single producer**
- **Multiple consumers**

- **Producer:**
  - **Adds to the buffer**
  - **Each buffer object has**
    - **File number**
    - **Page number**
    - **Memory mapping**
- **Consumers:**
  - **Eat from the buffer**
  - **Each zipped object has**
    - **Compressed data**
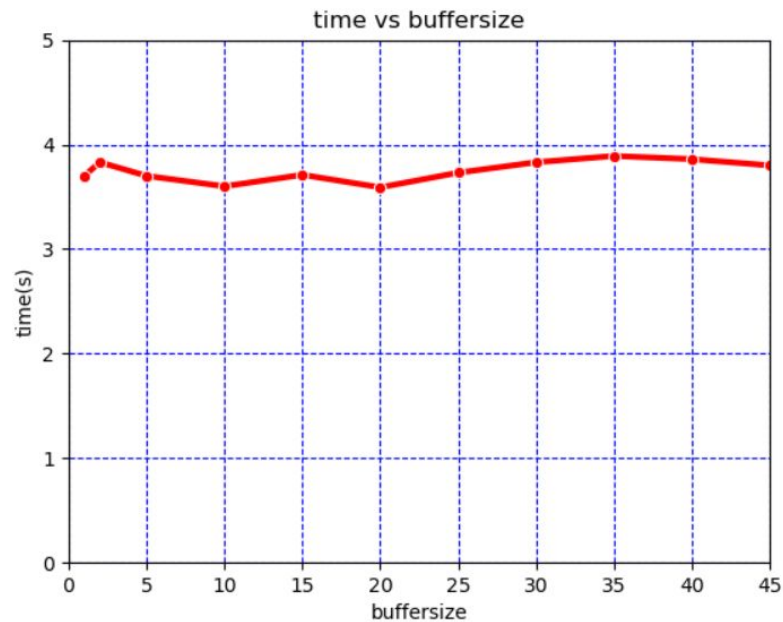    - **Data size**

# Read/Write System Calls vs MMAP

- **Traditional Read**
  - **Process can't access disc**
  - **Entire data read and copied from target file into I/O buffer followed by process buffers**

- **MMAP**
  - **Loading relevant file sections to the RAM**
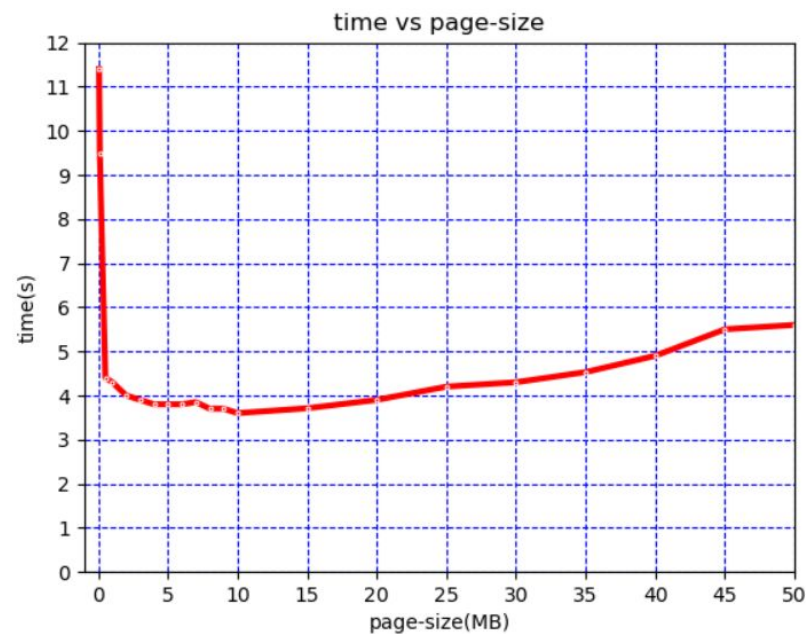  - **File is loaded as pages into RAM**



offset  length

Secondary Storage

length

Process Address Space

# RESULTS

❖ **Varying** BUFFER SIZE

❖ **Varying** PAGE SIZE

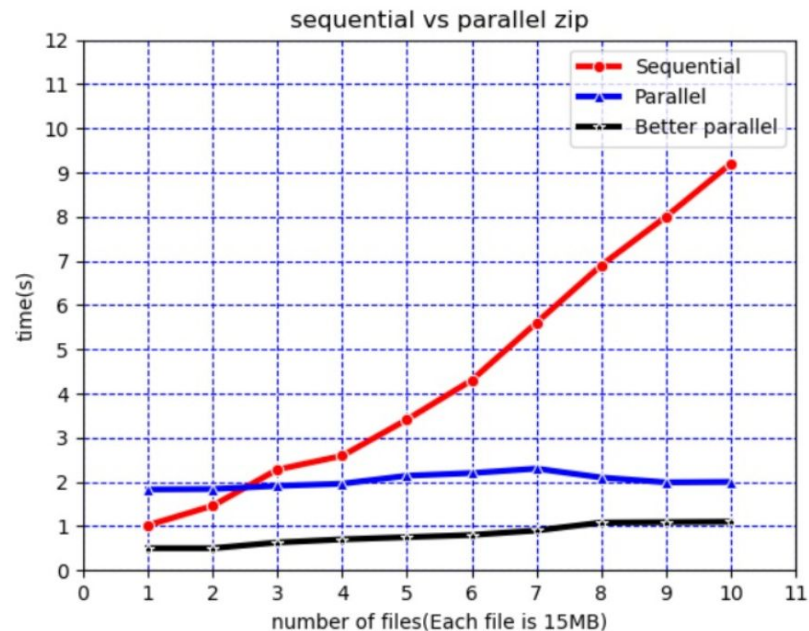● **Workload - 3 Files (100 MB, 300 MB, 400 MB)**



time vs buffersize



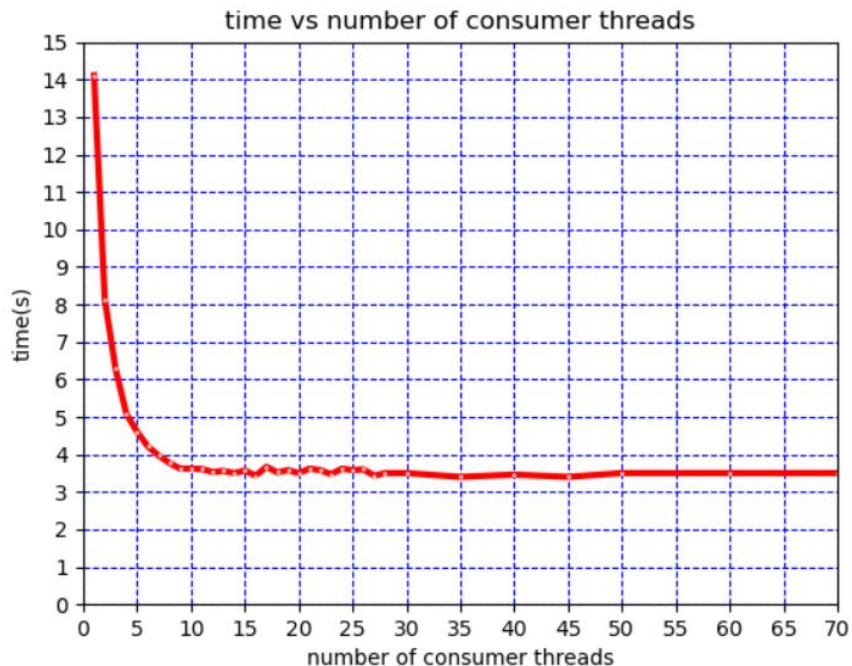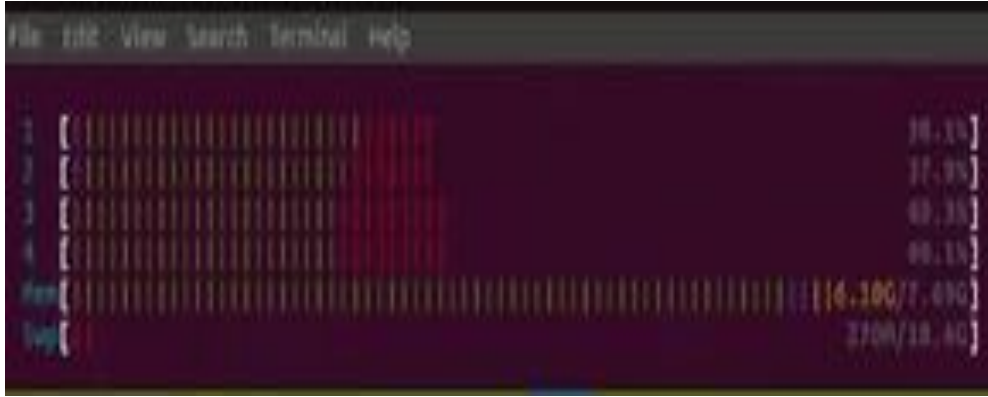time vs page-size

# RESULTS

❖ **Varying** NUMBER OF CONSUMER THREADS

# RESULTS



**CPU Utilization of all 4 cores on running code**

❖ **Performance ∝ No. of cores**



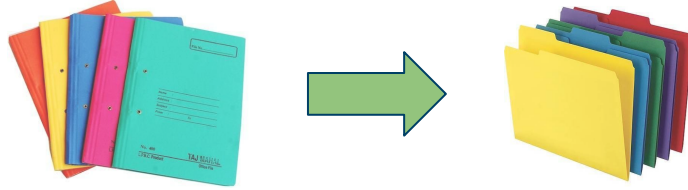Zip 5 files(16MB, 30MB, 50MB, 110MB, 116MB)

❖ **Performance:**
**Sequential<Parallel<<Better Parallel**

# LIMITATIONS AND SCOPE

- **Scope**



- **Limitations**
    - **MMAP wastes space, wastage=( int*page_size - file_size)**
    - **Parallel Unzip**

# THANK YOU