

---

Universidad Nacional Autónoma de México  
Facultad de Ciencias  
Algoritmos Paralelos

## Proyecto final: Implementación de Run-length encoding (RLE)

Autor: Medina Peralta Joaquín

---

## Problema

Run-Length encoding es una forma de compresión de datos sin pérdida, donde la secuencia de símbolos de la entrada se almacenan como una simple ocurrencia con el número simbolos consecutivos.

Definimos, para una cadena  $w \in \Sigma^*$ , una serie de pares  $(r, l, s)$  de datos consecutivos dadas por:

- El carácter de control  $r \notin \Sigma$ .
- La longitud de la secuencia  $l \in \mathbb{N}$
- El símbolo del alfabeto  $s \in \Sigma$

Por ejemplo, tenemos una cadena perteneciente al alfabeto en inglés, ABABBBC y la lista de tuplas con el caracter de control  $\alpha$ , se tiene:

$$\begin{aligned} & (\alpha, 1, A), (\alpha, 1, B), (\alpha, 1, A) \\ & (\alpha, 3, B), (\alpha, 1, C) \end{aligned}$$

Para poder representarlo en un archivo, los símbolos con un solo elemento consecutivo, es decir que no se repiten, se representarán como  $s$ . Con el ejemplo anterior, la cadena comprimida sería  $ABA(\alpha, 3, B)C$ .

Por lo tanto, el algoritmo es eficiente para cadenas o contenido donde hay muchas secuencias consecutivas, es decir, archivos de texto, binarios representación de una imagen por medio de pixeles o componentes de bloques largos de archivos de sonido. (Pu, 2005).

Para este trabajo, se consideran únicamente archivos de texto para comprimir que estén en **utf-8**, debido a que, se usará una combinación para representar las tuplas. Para los símbolos sin repetición o secuencia muy cortas, , como se menciona en (Pu, 2005), se implementará un modo literal donde los bytes se escriben directamente sin tupla. Esto evita la expansión del archivo, permitiendo que la compresión se reserve únicamente para las secuencias que si generen un ahorro significativo.



## Algoritmo secuencial

Como se menciona en (Pu, 2005), se define caracteres de control de repetición  $r_3, r_4, \dots$ , los cuales para  $r_i$  representa la cantidad  $i$  de símbolos consecutivos. De esta manera, se usará el carácter de control únicamente cuando la longitud sea mayor a 3.

Así, definimos el algoritmo de encriptación como:

- Repetir hasta llegar al final del archivo:
  - Leer el símbolo  $S$  de la secuencia:
    - Contar el número de símbolos consecutivos iguales  $i$  hasta encontrar uno diferente.
    - Guardar en el archivo de salida el par  $r_iS$  solo si  $i > 2$ , en otro caso guardar el símbolo  $i$

Para la implementación, se debe de considerar tres apuntadores  $i, j$  donde  $j = i + 1$  donde los usaremos para recorrer el arreglo de símbolos  $S$ , para lo cual, mientras que  $S[i] == S[j]$  entonces podemos incrementar el contador, de esta manera verificamos si es mayor al umbral, de esta manera podemos representarlo en forma de tupla y en otro caso, se escribe el carácter  $S[q]$  y  $S[p]$  en la cadena comprimida. (El pseudocódigo se encuentra en el anexo 1)

De esta manera, definimos el algoritmo de desencriptación como:

- Repetir hasta llegar al final del archivo:
  - Leer el símbolo  $S$  de la secuencia:
    - Si el símbolo  $S = r_i$ , se debe de leer el símbolo consecutivo y escribir  $i$  veces en la salida.
    - En otro caso, escribir el símbolo actual.

De la misma manera, es necesario un apuntador  $p$  para recorrer el archivo comprimido, donde si  $S[p]$  es el carácter de control definido  $r$ , se toma el valor de  $S[p + 1] = n$  como un número y  $S[p + 2]$  como el símbolo a repetir  $n$  veces, en otro caso se guarda  $S[p]$  en la salida. (El pseudocódigo se encuentra en el anexo 2)

A su vez, se considera el uso de **FLAG\_RLE** para iniciar la tupla definida en la sección anterior y en dado caso que el símbolo  $s$  sea igual a **FLAG\_RLE**, se considera a **FLAG\_LITERAL** para escapar y escribir el valor de **FLAG\_RLE**.

## Algoritmo en paralelo

Para esto, se toma en inspiración del artículo (Manchev, 2006), donde se propone el dividir la secuencia de entrada en bloques procesados en paralelo donde cada proceso identifica las corridas locales dentro de su bloque, para posteriormente realizar una fase de combinación donde se revisa si los límites entre bloques contiguos se cruzan con corridas de símbolos, asegurando que las secuencias iguales se representen correctamente.



De esta manera, se debe de utilizar el uso de memoria compartida y distribuida para poder ejecutar el algoritmo y por lo tanto, nos esta dando la pista de usar MPI.

Antes de pasar a la implementación, debemos mencionar que esta estrategia de parallelizar el algoritmo secuencial RLE sigue el principio de **Fork Join** donde se divide el trabajo entre los distintos procesadores y al final, juntar los resultados de cada procesador en uno global.

Se puede observar la técnica de **Fork Join** de manera explicita al dividir la secuencia de entrada en bloques para que cada procesador realice RLE de manera local y al final se junta cada bloque revisando si no hay cruces entre cada uno.

De esta manera, proponemos el pseudocódigo (3) donde se puede observar la compresión en RLE en paralelo y en (4) la manera de descomprimir en paralelo.

## Experimentación

Para esto, trabajaremos con 3 archivos de 100 MB con una alta, media y baja repetitividad los cuales se calculará el tiempo de comprimir de manera secuencial y paralela con 2, 4, 6, 20, 50 procesadores.

De este modo, se realizará la prueba en el siguiente equipo:

```
-----
Sistema operativo: Darwin MacBook-Pro-****.local 22.6.0 Darwin
Kernel Version 22.6.0: Tue Jul 15 08:22:28 PDT 2025; root:xnu
-8796.141.3.713.2~2/
RELEASE_X86_64 x86_64
CPU: Intel(R) Core(TM) i5-7360U CPU @ 2.30GHz
Nucleos fisicos: 2
Nucleos logicos: 4
Memoria RAM total: 8,00 GB
-----
```

## Resultados

## Conclusión

## Referencias

- Manchev, N. (2006). Parallel Algorithm for Run Length Encoding. *Proceedings of the Third International Conference on Information Technology: New Generations (ITNG 2006)*. <https://doi.org/10.1109/ITNG.2006.106>
- Pu, I. M. (2005). *Fundamental Data Compression* [pp. 56–64]. Elsevier Butterworth-Heinemann.



## Anexo

---

**Algoritmo 1** Algoritmo para comprimir en RLE con modo literal (Umbral 3)

---

```
1: procedure COMPRIMIR( $S$ )                                 $\triangleright$  Cadena de símbolos  $S$ 
2:    $N \leftarrow longitud(S)$ 
3:    $O \leftarrow []$                                           $\triangleright$  Cadena de salida comprimida
4:    $p \leftarrow 0$ 
5:   while not  $EOF$  do
6:      $i \leftarrow 1$ 
7:      $q \leftarrow p + 1$                                       $\triangleright$  Fase 1: Contar la secuencia
8:     while  $q < N$  y  $S[q] = S[p]$  y  $i < 255$  do           $\triangleright$  Máximo conteo en 1 byte
9:        $i \leftarrow i + 1$ 
10:       $q \leftarrow q + 1$ 
11:    end while                                          $\triangleright$  Fase 2: Aplicar el umbral y codificar
12:    if  $i \geq 3$  then                                  $\triangleright$  Guardar como tupla RLE:  $\langle \text{FLAG\_RLE}, i, S_p \rangle$ 
13:       $O \leftarrow O \parallel \text{FLAG\_RLE} \parallel i \parallel S[p]$ 
14:    else                                               $\triangleright$  Guardar como secuencia literal:  $S[p]$  se repite  $i$  veces
15:      for  $k \leftarrow 0$  hasta  $i - 1$  do
16:         $O \leftarrow O \parallel S[p]$ 
17:      end for
18:    end if
19:     $p \leftarrow p + i$                                       $\triangleright$  Mover el apuntador principal al inicio de la siguiente secuencia
20:  end while
21:  Salida:  $O$ 
22: end procedure
```

---

Figura 1: Algoritmo secuencial para comprimir



---

**Algoritmo 2** Algoritmo para Descomprimir en RLE (Modo Extendido)

---

```
1: procedure DESCOMPRIMIR( $C$ )                                 $\triangleright$  Cadena de símbolos comprimidos  $C$ 
2:    $N \leftarrow longitud(C)$ 
3:    $O \leftarrow []$                                                   $\triangleright$  Cadena de salida descomprimida
4:    $p \leftarrow 0$ 
5:   while  $p < N_C$  do
6:      $S \leftarrow C[p]$                                           $\triangleright$  Leer el símbolo actual (byte de control o literal)
7:     if  $S = \text{FLAG\_RLE}$  then
8:       if  $p + 2 \geq N_C$  then
9:         break                                               $\triangleright$  Error: Buffer incompleto
10:        end if
11:         $i \leftarrow C[p + 1]$                                       $\triangleright$  Leer el conteo (longitud)
12:         $V \leftarrow C[p + 2]$                                       $\triangleright$  Leer el valor del símbolo
13:        for  $k \leftarrow 1$  hasta  $i$  do
14:           $O \leftarrow O \parallel V$                                   $\triangleright$  Escribir el valor  $V$  el número de veces  $i$ 
15:        end for
16:         $p \leftarrow p + 3$                                           $\triangleright$  Avanzar el apuntador 3 posiciones (FLAG, Conteo, Valor)
17:      else if  $S = \text{FLAG\_LITERAL}$  then                       $\triangleright$  Símbolo FLAG escapado
18:        if  $p + 1 \geq N_C$  then
19:          break                                               $\triangleright$  Error: Buffer incompleto
20:        end if
21:         $V \leftarrow C[p + 1]$                                       $\triangleright$  Leer el siguiente byte (que es un FLAG real)
22:         $O \leftarrow O \parallel V$                                   $\triangleright$  Escribir el byte escapado
23:         $p \leftarrow p + 2$                                           $\triangleright$  Avanzar 2 posiciones (FLAG_LITERAL, Valor)
24:      else                                                  $\triangleright$  Símbolo sin repetición
25:         $O \leftarrow O \parallel S$                                   $\triangleright$  Escribir el símbolo literal actual  $S$ 
26:         $p \leftarrow p + 1$                                           $\triangleright$  Avanzar 1 posición
27:      end if
28:    end while
29:    Salida:  $O$ 
30:  end procedure
```

---

Figura 2: Algoritmo secuencial para descomprimir



---

**Algoritmo 3** Algoritmo para comprimir en RLE en paralelo

---

```
1: procedure COMPRIMIR_PARALELO( $S, N$ )            $\triangleright S$  entrada y  $N$  número de procesos
2:                                          $\triangleright$  Lógica del Coordinador (Pre-procesamiento)
3:      $Global\_Size \leftarrow longitud(S)$ 
4:     for  $id \leftarrow 0 \dots N - 1$  do
5:         OffsetStart
6:          $SizeChunk \leftarrow calculate\_chunk(id, Global\_Size, S)$ 
7:          $Chunkid \leftarrow S[OffsetStart \dots OffsetStart + SizeChunk + 1]$             $\triangleright$  Incluye 1 byte de frontera
8:         Enviar( $Chunkid$ ) a  $N\_id$ 
9:     end for
10:    Sincronizar()                                      $\triangleright$  Lógica del Trabajador
11:    InputT  $\leftarrow$  Recibir( $Chunkid$ )
12:    InputMain  $\leftarrow$  InputT[datos principales]
13:    ByteNext  $\leftarrow$  InputT[byte de frontera]
14:    Outputlocal  $\leftarrow$  Comprimir_Secuencial(InputMain)            $\triangleright$  Usa RLE Extendido
15:                                          $\triangleright$  Corrección de Fronteras (Comunicación/Sincronización)
16:    if  $id < N - 1$  then
17:        Enviar_Front(Outputlocal.último_bloque) a Trabajadorid+1
18:        Recibir_Front(SeñalFus) de Trabajadorid+1
19:        if SeñalFus = FUSIONADO then
20:            Outputlocal.Eliminar_último_bloque()
21:        end if
22:    end if
23:    if  $id > 0$  then
24:        Recibir_Front(BloquePrev) de Trabajadorid-1
25:        Intentar_Fusion(Outputlocal.primero, BloquePrev)
26:        Enviar_Front(SeñalFus) a Trabajadorid-1
27:    end if
28:    Enviar(Outputlocal) al Coordinador
29:    Sincronizar()                                      $\triangleright$  Lógica del Coordinador (Post-procesamiento)
30:    Recibir_Buffers(Todos los trabajadores)
31:    LenAcumulada  $\leftarrow$  Calcular_Suma_Prefijo(Longitudes_locales)
32:     $O_{comp} \leftarrow$  Concatenar(Buffers recibidos, LenAcumulada)
33:    Salida:  $O_{comp}$ 
34: end procedure
```

---

Figura 3: Algoritmo paralelo para comprimir



---

**Algoritmo 4** Algoritmo para descomprimir en RLE en paralelo

---

```
1: procedure DESCOMPRIMIR_PARALELO( $C, N$ )  $\triangleright C$  entrada comprimida y  $N$  número de procesos
2:  $Par_{count} \leftarrow read\_Metadato(C)$   $\triangleright$  Lógica del Coordinador (Pre-procesamiento)
3: for  $id \leftarrow 0 \dots N - 1$  do
4:    $Offset_{par}, Size_{par} \leftarrow Calcular\_Particion\_Pares(id, Par_{count})$   $\triangleright$  Número total de pares/bloques
5:    $Chunk_{id} \leftarrow C[Offset_{par} \dots Offset_{par} + Size_{par}]$   $\triangleright$  Lectura de porción de pares
6:   Enviar( $Chunk_{id}$ ) a  $N\_id$ 
7: end for
8: Sincronizar()
9:  $Syncronizar()$   $\triangleright$  Lógica del Trabajador
10:  $Input_T \leftarrow Recibir(Chunk_{id})$ 
11:  $Output_{local} \leftarrow Descomprimir\_Secuencial(Input_T)$   $\triangleright$  Genera bytes sin comprimir
12:  $Len_{local} \leftarrow longitud(Output_{local})$ 
13:  $Offset_{Final} \leftarrow Calcular\_Suma\_Prefijo(Todas\_Len_{local})[id]$   $\triangleright$  Cálculo del Offset Global (Suma Prefijo)
14: Enviar( $Len_{local}$ ) a Concentrador
15: Sincronizar()
16: Recibir(Todas_Lenlocal)
17:  $Offset_{Final} \leftarrow Calcular\_Suma\_Prefijo(Todas\_Len_{local})[id]$   $\triangleright$  Cada trabajador conoce dónde escribir
18: Escribir_en_Posicion_Global( $O_{descomp}, Offset_{Final}, Output_{local}$ )  $\triangleright$  Escritura dispersa
19: Sincronizar()
20:  $Syncronizar()$   $\triangleright$  Lógica del Coordinador (Post-procesamiento)
21: Finalizar_Archivo( $O_{descomp}$ )
22: Salida:  $O_{descomp}$ 
23: end procedure
```

---

Figura 4: Algoritmo paralelo para descomprimir