

## COMPRESSING DOCUMENTS

Use algorithms to reduce the bit used to represent a document, in order to minimize the resources required to store and transmit data.



SPACE-TIME COMPLEXITY TRADEOFF = You have to consume resources in the compression / decompression phase

Types of compression

**LOSSLESS** = DATA DOESN'T LOSE INFORMATION. IT EXPLOITS STATISTICAL REDUNDANCY (EXAMPLE: IF I HAVE 200 CONTAINS RED PIXEL, I DON'T HAVE TO CODE THEM ALL AS "REO, REO, REO..." BUT I CAN DO INSTEAD "200 REO")

LZ

**LOSSY** = DROPS NON-ESSENTIAL DETAILS OF THE DATA TO SAVE MORE SPACE (EXAMPLE: ROUND UP THE PIXEL COLORS THAT ARE NOT PERCEIVED WELL BY HUMANS)

|                    |   |            |
|--------------------|---|------------|
| R:7.2<br>(r:3.222) | → | G:3<br>B:4 |
| B:7                |   |            |

SPEC

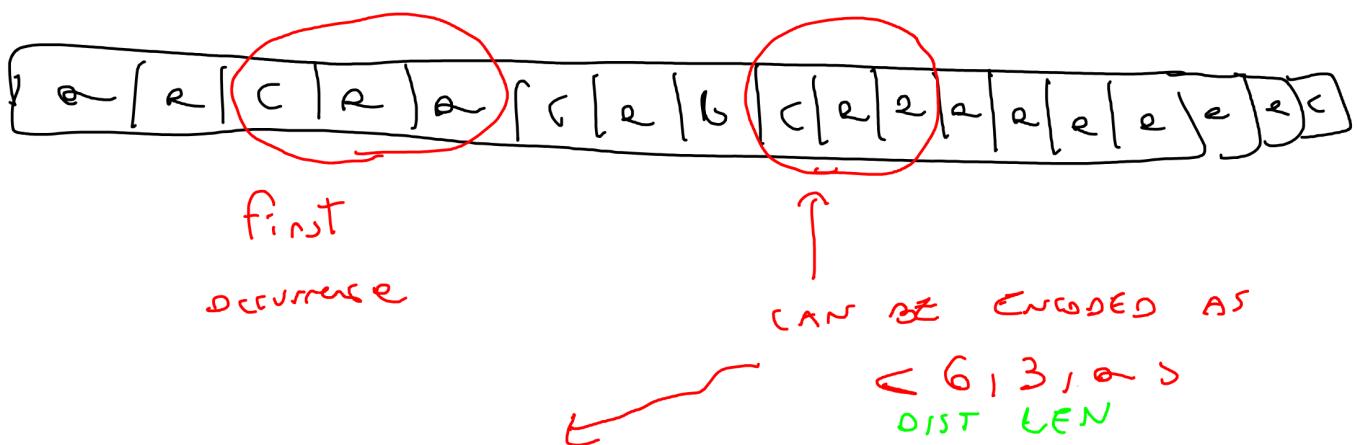
## LZ77

WE ALSO REFER TO LZ AS DEFL

LOSSLESS DATA COMPRESSION ALGR.

IT REPLACES REPEATED OCCURRENCES OF DATA WITH A REFERENCE TO A SINGLE COPY OF THAT DATA. IN ORDER TO FIND MATCHES, THE ALGR. KEEPS TRACK OF SOME AMOUNT OF THE RECENT DATA (CURRENT SLIDING WINDOW). A MATCH IS ENCODED AS A LENGTH-INSTANCE PAIR (IN THE EXAMPLE WE ALSO KEEP THE VALUE OF THE NEXT CHAR, BECAUSE IT MAKES DECODING EASIER)

EXAMPLE : LET'S IMAGINE THAT THE WINDOW IS THE WHOLE STRING



THE NEXT 3 CHARS ARE EQUAL TO THE CHARS YOU SAW "6 CHARS AGO" THEN THERE IS AN OR AFTER

## EXAMPLE

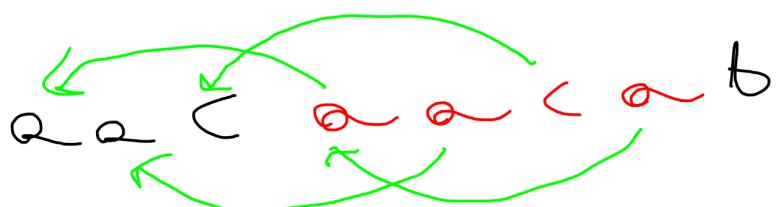
IT IS ALSO POSSIBLE TO HAVE PAIRS WITH LENGTH DIST WHICH IS USED TO SAY TO THE DECODER THAT THE DATA WANT TO DECODE NEEDS TO BE PASTED UNTIL NO SPACE REMAINS



THE DECODER WILL FIND THIS:

$a a c \langle 3, 4, 6 \rangle$

AND WILL DECODE AS.



Example: with window = 6

< dist, len, next char >

THE STRINGS:



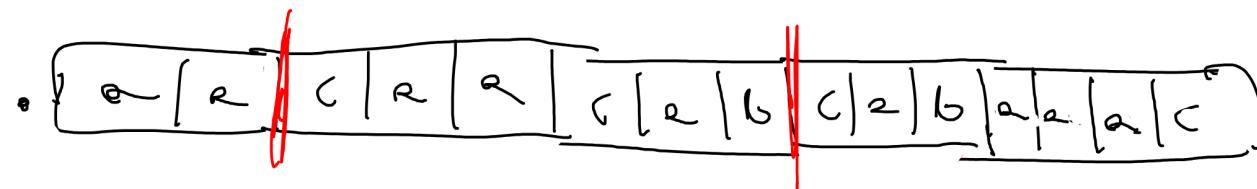
OUTPUT = < 0, 0, a >



OUTPUT = < 1, 1, c >      a matches with a



OUTPUT = < 3, 4, b >      a & c | a matches with each other



OUTPUT = < 3, 3, a >



OUTPUT = < 1, 2, c >

FINISH!

FINAL RESULT : < 0, 0, a > < 1, 1, c > < 3, 4, b > < 3, 3, a > < 1, 2, c >

## EXAMPLE

DECODE THE RESULT ABOVE

→ TRACE THE TRIPLES, ONE BY ONE

→ DO WHAT THIS ORDER

$\langle 0, 0, 0 \rangle \langle 1, 1, c \rangle \langle 3, 4, b \rangle \langle 3, 3, a \rangle \langle 1, 2, c \rangle$

$\alpha \langle 1, 1, c \rangle \langle 3, 4, b \rangle \langle 3, 3, a \rangle \langle 1, 2, c \rangle$

$\alpha \alpha \langle 3, 4, b \rangle \langle 3, 3, a \rangle \langle 1, 2, c \rangle$

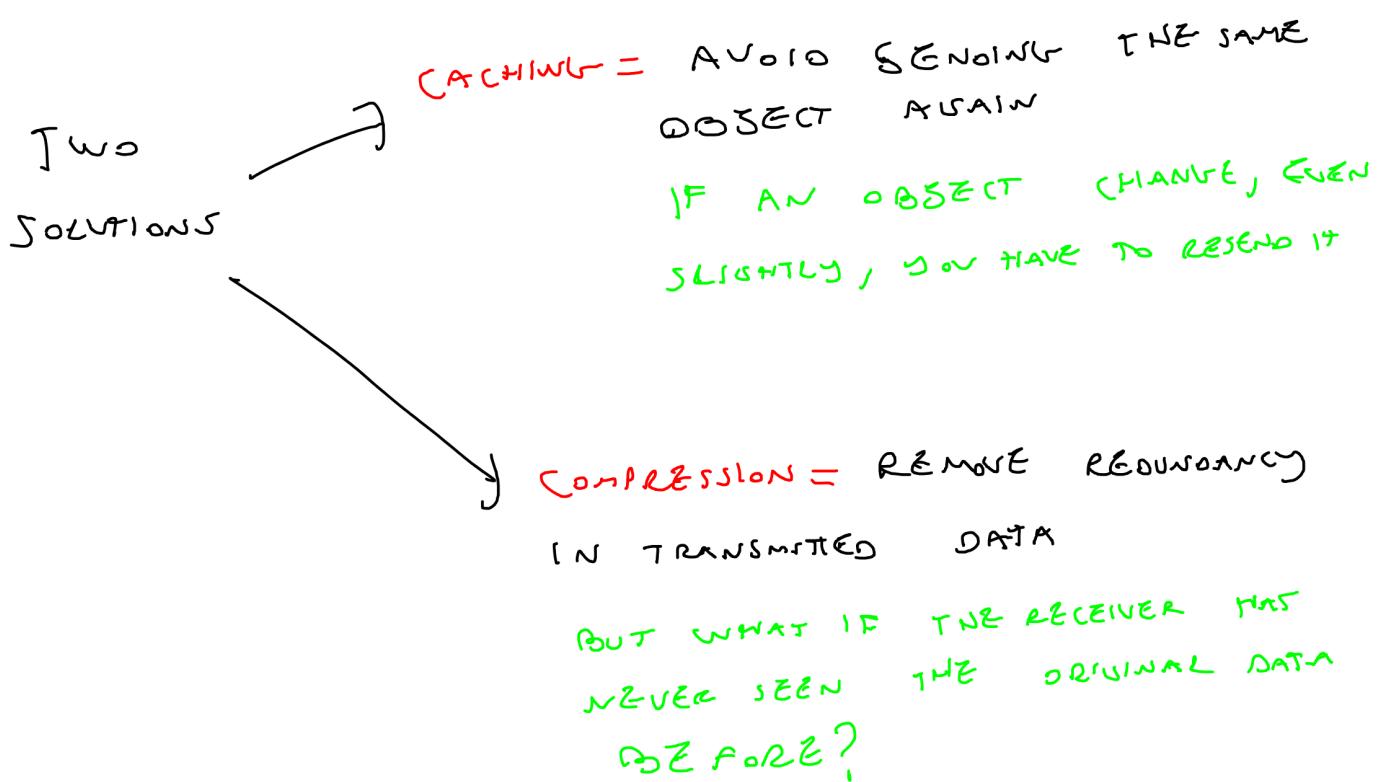
$\alpha \alpha \langle \alpha \alpha \langle \alpha b \langle 3, 3, a \rangle \langle 1, 2, c \rangle$

$\alpha \alpha \langle \alpha \alpha \langle \alpha b \langle \alpha b \alpha \langle 1, 2, c \rangle$

$\alpha \alpha \langle \alpha \alpha \langle \alpha b \langle \alpha b \alpha \alpha \langle$

## COMPRESSION & NETWORKING

WE WANT TO MAKE DATA TRANSFER AS TRANSPARENT AS POSSIBLE TO THE USERS.  
EVEN IF INTERNET IS NOW FAST, SOMEONE MAY BE STILL USING A SLOW CLIENT (MOBILE) OR BATTERY BASED.



→ IF BOTH SENDER & RECEIVER KNOWS THE FILE (UNSTRUCTURED)

→ **DELTA COMPRESSION (diff)**

SENDS ONLY THE DIFFERENCES BETWEEN THE REQUESTED PART AND THE ONES IN CACHE

### FORMATS

Z-DELTA COMPRESSION

Given two files  $f_{known}$ , known to both parties, and  $f_{new}$ , known only to sender, we want to compute  $f_d$  of minimum size, such that  $f_{new}$  can be derived by the receiver using  $f_{known}$  and  $f_d$ .  
FIND AN OPTIMAL COVERING OF  $f_{new}$  BASED ON  $f_{known}$

### OUR APPROACH

COMPRESS THE TEXT AS  $f_{known} f_{new}$  USING THE ZCZT SCHEME, STARTING THE ENCODE FROM  $f_{new}$ , while  $f_{known}$  acts as the "previously encoded text"

**EFFICIENT AND  
OPTIMAL SOLUTION<sup>90%</sup>**

USES ZDELTAS USED ALSO IN BACKUPS AND VERSION CONTROL

## CLUSTER-BASED

## DELTA COMPRESSION

2 DELTA CAN BE ALSO USED TO COMPRESS  
A GROUP OF FILES (USEFUL ON  
A DYNAMIC COLLECTION OF WEB PAGES OR  
BACKUPS...).

THE IDEA IS TO FIND A GOOD  
REFERENCE FOR EACH FILE FILE AND  
APPLY A PAIRWISE DELTA



THE PROBLEM REDUCES TO THE MIN. BRANCHING  
PROBLEM ON DAGS (DIRECTED ACYCLIC GRAPHS)

→ BUILD A (COMPLETE?) WEIGHTED GRAPH  
GRAPH (NODES = FILES, WEIGHTS =  $\frac{\text{delta size}}{\# \text{triples created during encoding}}$ )

→ INSERT A DUMMY NODE CONNECTED TO ALL  
AND THE WEIGHTS ARE THE ZIP SIZE  
OF THE RESPECTIVE NODES

→ COMPUTE THE DIRECTED SPANNING TREE OF  
MIN TOTAL COST THAT COVERS ALL  
5'S NODES (A PATH THAT CONNECTS ALL THE NODES)

SEE EXAMPLE

**PROBLEM** = CONSTRUCTION  $\hookrightarrow$  IS COSTLY,  
 BECAUSE IT REQUIRES  $n^2$  edge calculation  
 (the number of delta executions)

SOLUTION

PRUNING?

- COLLECTION ANALYSIS = CLUSTER THE FILES THAT APPEAR SIMILAR, HENCE GOOD FOR ZDELTA, THEN BUILD A SPARSE WEIGHTED GRAPH  $G_F$  CONTAINING ONLY EDGES BETWEEN PAIRS OF FILES IN THE SAME CLUSTER

- ASSIGN WEIGHTS = ESTIMATE APPROPRIATE EDGE WEIGHTS FOR  $G_F$ , HENCE SPARING ZDELTA EXEC (HOWEVER, STILL  $n^2$  TIMES)

### SOME PERFORMANCE...

NO PRUNING

|                | SPACE | TIME     |
|----------------|-------|----------|
| UNCOMPRESS.    | 30MB  | ...      |
| Tgz            | 20%   | $O(n)$   |
| CLUSTER ZDELTA | 8%    | $O(n^2)$ |

PRUNING

|                | SPACE | TIME  |
|----------------|-------|-------|
| UNCOMPRESS.    | 20MB  | ...   |
| Tgz            | 12%   | 2min  |
| CLUSTER ZDELTA | 8%    | 16min |

## EXAMPLE

GIVEN 3 STRINGS  $s_1 = "abac"$ ,  $s_2 = "dad"$   
 and  $s_3 = "coaba"$ , SHOW THE COMPRESSION ORDER  
 $\chi$ -DELTa APPLIED TO THE GROUP OF THESE  
 STRINGS USING ZIP ( $\zeta\chi\chi$ ) AND  
 GIVES THEIR CONCATENATION ORDER

① Apply zip to the strings

$$\text{zip}(s_1) = \langle \langle 0, 0, a \rangle \langle 0, 0, b \rangle \langle 2, 1, c \rangle \langle 0, 0, 0 \rangle \rangle \quad \text{size 4}$$

$$\text{zip}(s_2) = \langle \langle 0, 0, d \rangle \langle 0, 0, a \rangle \langle 2, 2, \$ \rangle \rangle \quad \text{size 3}$$

$$\text{zip}(s_3) = \langle \langle 0, 0, c \rangle \langle 0, 0, 0 \rangle \langle 0, 0, 0 \rangle \langle 0, 0, b \rangle \langle 2, 1, \$ \rangle \rangle \quad \text{size 5}$$

② Apply  $\chi$ -delta pair-wise ( ENCODE ONLY SECOND TERM )

$$\chi\text{-delta}(s_1 | s_2) = \langle \langle 0, 0, d \rangle \langle 7, 1, d \rangle \langle 2, 1, \$ \rangle \rangle \quad \text{size 3}$$

$$\chi\text{-delta}(s_1 | s_3) = \langle \langle 2, 2, \$ \rangle \langle 7, 2, \$ \rangle \rangle \quad \text{size 2}$$

$$\chi\text{-delta}(s_2 | s_1) = \langle \langle 1, 1, b \rangle \langle 2, 1, c \rangle \langle 0, 0, 0 \rangle \rangle \quad \text{size 3}$$

$$\chi\text{-delta}(s_2 | s_3) = \langle \langle 0, 0, c \rangle \langle 0, 0, 0 \rangle \langle 3, 1, b \rangle \langle 2, 1, \$ \rangle \rangle \quad \text{size 4}$$

$$\chi\text{-delta}(s_3 | s_1) = \langle \langle 3, 3, c \rangle \langle 8, 1, \$ \rangle \rangle \quad \text{size 2}$$

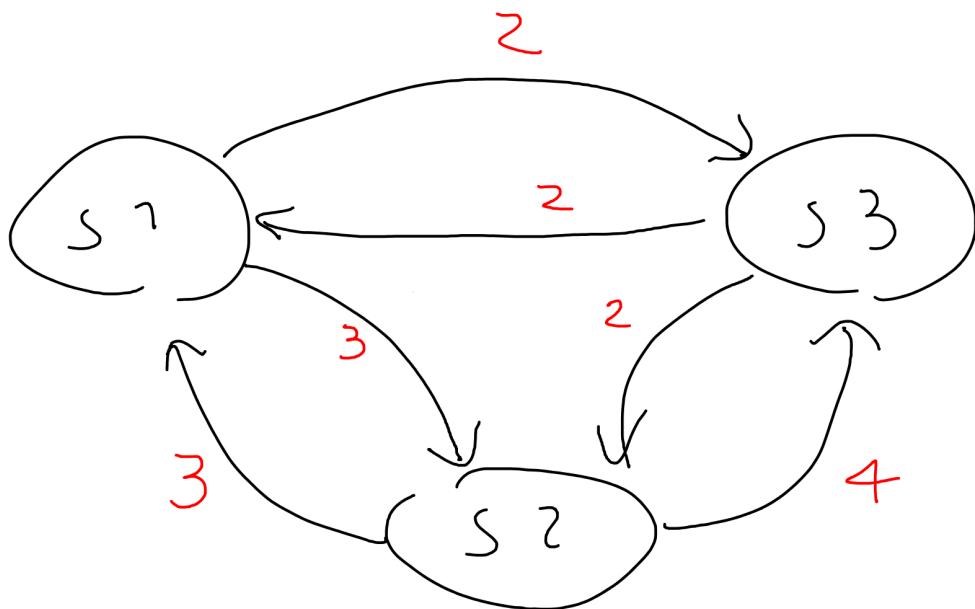
$$\chi\text{-delta}(s_3 | s_2) = \langle \langle 0, 0, d \rangle \langle 2, 3, \$ \rangle \rangle \quad \text{size 2}$$

NICE OVERLAP

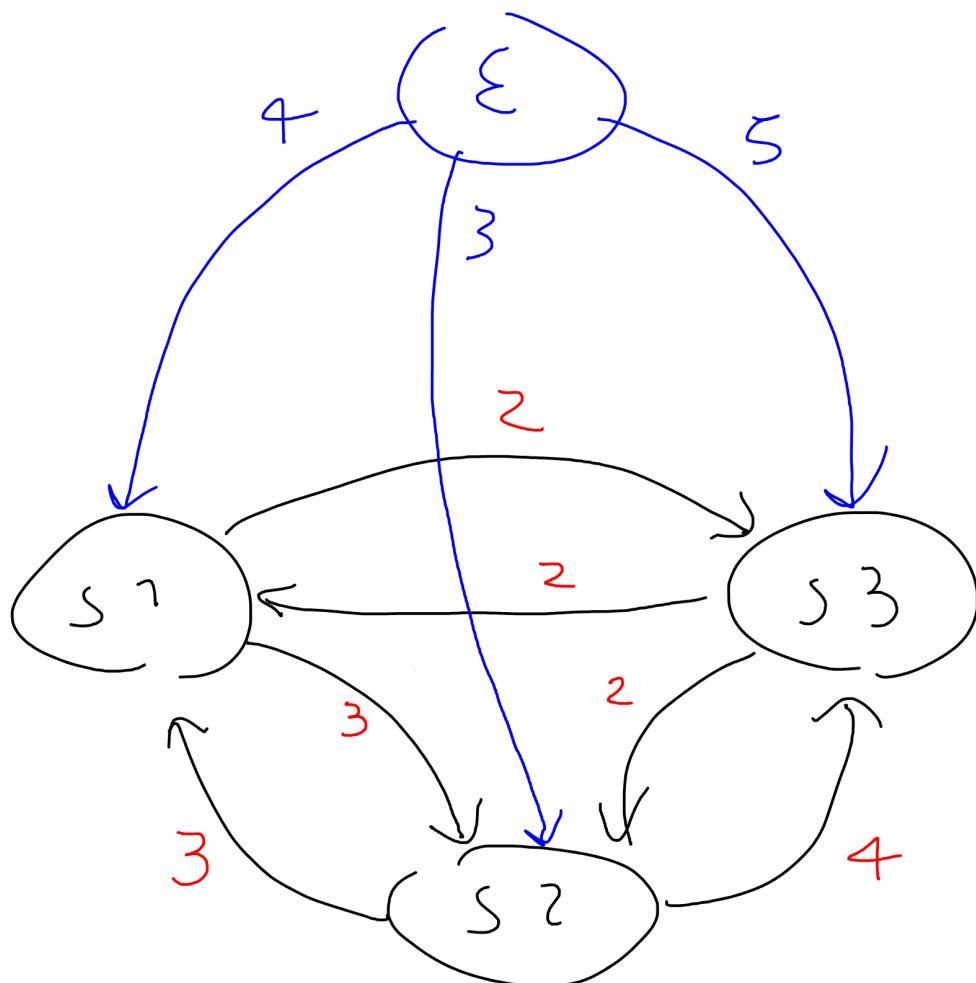
COABA DA DA

(3) BUILD THE WEIGHTED GRAPH

(nodes = files, weights = delta size)



(4) ADD DUMMY NODE  $\epsilon$



(5)

FIND THE DIRECTED SPANNING TREE OF  
MIN TOTAL COST

Note = You HAVE TO START FROM {

TWO POSSIBILITIES

- { → s<sub>2</sub> → s<sub>1</sub> → s<sub>3</sub>      COST      8
- { → s<sub>1</sub> → s<sub>3</sub> → s<sub>2</sub>      COST      8

Now taking the second path

zdelta compression will compress tree

gzip(s<sub>1</sub>), zdelta(s<sub>1</sub>/s<sub>3</sub>), zdelta(s<sub>3</sub>/s<sub>2</sub>)

→ IF THERE IS ONLY A PARTIAL KNOWLEDGE  
BETWEEN SERVER & RECEIVER

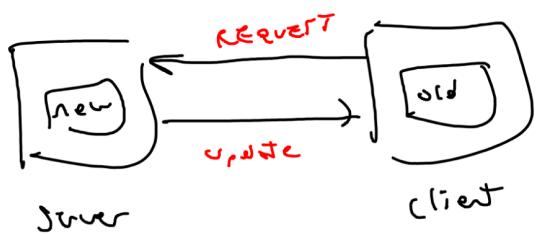


UNSTRUCTURED  
FILES

RECORD-BASED  
DATA

### FILE SYNCHRONIZATION (rsync, zsync)

Client request for a new version  
of a file. Server sends only  
the updated parts, without knowing  
the old file



### SET RECONCILIATION

CLIENT UPDATES

STRUCTURED FILE,  
WHICH IS AVAILABLE  
ON A SERVER

### FILE SYNCHRONIZATION WITH RSYNC

- ① Client divides old file in block of size  $B$  and hashes each of them
- ② Client sends  $\frac{\text{file size}}{B}$  hashes to the server
- ③ The server identifies the mismatching pairs of the files by hash rolling  $B$  chars: from left to right, substitute the hashes where possible
- ④ Server sends back a concatenation of literals and hashes known by the client, which forms the updated file

- Single round-trip
- Not good in theory (granularity) problem

## EXAMPLE

Server has  $f_{new} = CABADDACCCC$  and client  
 $f_{old} = ABACCCDAC$ ,  $B = 3$

(1)  $f_{old} = \underbrace{ABA}_{h_1} \mid \underbrace{CCC}_{h_2} \mid \underbrace{DAC}_{h_3}$

(2) Client  $\xrightarrow{h_1, h_2, h_3}$  Server

(3) Hash rolling  $f_{new}$  to compute  $f_d$

$$f_d = \cancel{CABAD} \cancel{DAC} C C C$$

(4) Server  $\xrightarrow{ch_1, ch_3, C C}$  Client

Client is able to build  $f_{new}$

# FILE SYNCHRONIZATION WITH ZSYNC

The idea is to reduce the workload on the server (sort of symmetrical approach to rsync)

(1) Server divides new file in  $B$  blocks and hashes each of them

(2) Server sends  $\frac{new\ size}{B}$  hashes to client

(3) Client identifies the missing parts by hash rolling  $B$  chars and then sends to the server & bitmask to request

then ( $size = \# \text{hashes received} = \frac{new\ size}{B}$ )  
( $B[i] = 0$  if client doesn't own hash  $i$ , else  $1$ )

(4) Using the bitmask, the server is able to understand what the client only and what not, therefore it delta compresses the missing blocks according to the owned ones and sends the encoded file.

## Example

Server has  $f_{new} = \text{ABRACACALABRIA}$  and client  
 $f_{old} = \text{ABRACADABRA}$   $B = 3$

(1)  $f_{new} = \text{ABR} | \text{A C A} | \text{LA B R I A}$   
 $h_1 \quad h_2 \quad h_3 \quad h_4$

(2) Server  $\xrightarrow{h_1, h_2, h_3, h_4}$  Client

(3) HASH rolling  $f_{old}$

A B R A C A D A B R A  
 $h_1 \quad h_2 \quad h_3$

bitmask = 1100

Client  $\xrightarrow{1100}$  Server

(4) Server does  $\rightarrow$  delta compression of  
 $h_3 h_4$ , taking  $h_1 h_2$  as known text

$\text{gzip}(h_1 h_2 | h_3 h_4) = \text{gzip}(\text{ABRACA} | \text{LABRIA}) =$   
 $= \langle 0, 0, L \rangle \langle 7, 3, 1 \rangle \langle 4, 1, \$ \rangle$

Server  $\xrightarrow{\langle 0, 0, L \rangle \langle 7, 3, 1 \rangle \langle 4, 1, \$ \rangle}$  Client

CLIENT IS NOW ABLE TO BUILD  $f_{new}$