

# Compression techniques

## Graph

Paolo Boldi, DSI, Università degli studi di Milano

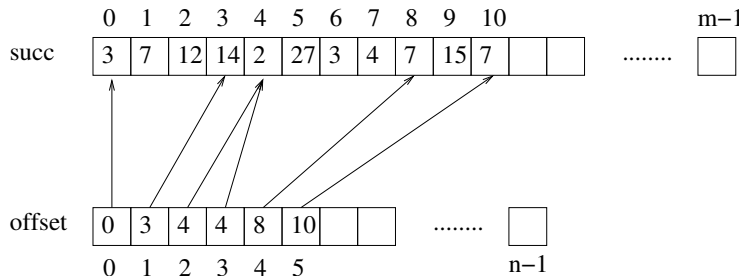
## What do we mean by. . .

... “storing the web graph”?

- ▶ Having a data structure that allows you, for a given node, to know its successors.
- ▶ Possibly: having a way to know which URL corresponds to a given node and vice versa.

We already studied good data structures for the URL  $\mapsto$  node map. The other one (less useful) is not covered here.

## Naive representation



The offset vector tells, for each given node  $x$ , where the successor list of  $x$  starts from. Implicitly, it also gives the degree of each node.

## Naive representation

How much space does this representation take?

- ▶ Successor array:  $m$  elements (arcs), each containing a node ( $\log n$  bits); with 32 bits, we can store up to 4 billion nodes (half of it, if we don't have unsigned types)
- ▶ Offset array:  $n$  elements (nodes), each containing an index in the successor array ( $\log m$  bits); with 32 bits, we can store up to 4 billion arcs.

All in all,  $32(n + m)$  bits. If we assume  $m = 8n$  (a very modest assumption on the outdegree), we need  $288n$  bits, i.e., 288 bits/node, 36 bits/arc.

We show how to reduce this of an order of magnitude.

# Idea

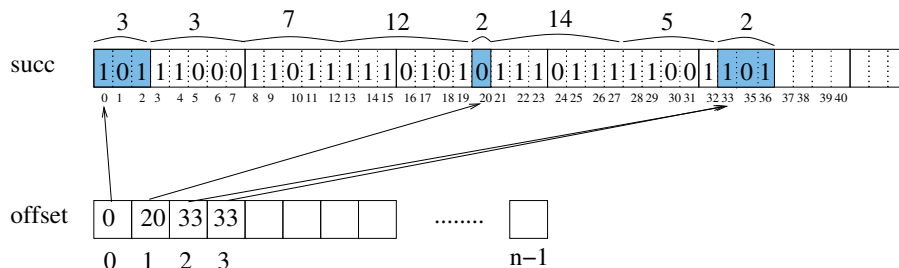
Use a variable-length representation for successors. Such a representation should obviously...

- ▶ be instantaneously decodable
- ▶ minimize the expected bitlength.

What about the offset array?

- ▶ bit displacement vs. byte displacement (with alignment)
- ▶ we have to keep an explicit representation of the node degrees (e.g., in the successor array, before every successor list).

## Variable-length representation



Node degrees (blue background), followed by successors. Each number is represented using an instantaneous code (possibly, different for degree and successors).

## Instantaneous code

- ▶ An *instantaneous (binary) code* for the set  $S$  is a function  $c : S \rightarrow \{0, 1\}^*$  such that, for all  $x, y \in S$ , if  $c(x)$  is a prefix of  $c(y)$ , then  $x = y$ .
- ▶ Let  $l_x$  be the length (in bits) of  $c(x)$ .
- ▶ Kraft-McMillan: there exists an instantaneous code with lengths  $l_x$  ( $x \in S$ ) if and only if

$$\sum_{x \in S} 2^{-l_x} \leq 1.$$

## Intended distribution

- ▶ If  $p : S \rightarrow [0, 1]$  is the probability distribution of the source, than the ideal code for  $S$  is such that (Shannon)

$$l_x = -\log p(x)$$

(all logarithms are in base 2).

- ▶ So a code with lengths  $l_x$  has *intended distribution*

$$p(x) = 2^{-l_x}.$$

- ▶ The choice of the code to use will be based on the expected distribution of the data.



## Fixed-length coding

- ▶ If  $S = \{1, 2, \dots, N\}$ , to represent an element of  $S$  it is sufficient to use  $\lceil \log N \rceil$  bits.
- ▶ The fixed-length representation for  $S$  uses exactly that number of bits for every element (and represents  $x$  using the standard binary coding of  $x - 1$  on  $\lceil \log N \rceil$  bits).
- ▶ Intended distribution:

$$p(x) = 2^{-\lceil \log N \rceil} \quad \text{uniform distribution.}$$

# Unary coding

- ▶ If  $S = \mathbf{N}$ , one can represent  $x \in S$  writing  $x$  zeroes followed by a one.
- ▶ So  $l_x = x + 1$ , and the intended distribution is

$$p(x) = 2^{-x-1} \quad \text{geometric distribution of ratio } 1/2.$$

0	1
1	01
2	001
3	0001
4	00001

## A more general viewpoint

Unary coding can be seen as a special case of a more general kind of coding for  $\mathbf{N}$ . Suppose you group  $\mathbf{N}$  into *slots*: every slot is made by consecutive integers; let

$$V = \langle s_1, s_2, s_3, \dots \rangle$$

be the slot sizes (in the unary case  $s_1 = s_2 = \dots = 1$ ).

Then, to represent  $x \in \mathbf{N}$  one can

- ▶ encode *in unary* the index  $i$  of the slot containing  $x$ ;
- ▶ encode *in binary* the offset of  $x$  within its slot (using  $\lceil \log s_i \rceil$  bits).

# Golomb coding

*Golomb coding with modulus  $b$*  is obtained choosing

$$V = \langle b, b, b, \dots \rangle.$$

To represent  $x \in \mathbf{N}$  you need to specify the slot where  $x$  falls (that is,  $\lfloor x/b \rfloor$ ) in unary, and then represent the offset using  $\lceil \log b \rceil$  bits (or  $\lfloor \log b \rfloor$  bits).

So

$$l_x = \left\lfloor \frac{x}{b} \right\rfloor + \lceil \log b \rceil.$$

The intended distribution is

$$p(x) = 2^{-l_x} \propto (2^{1/b})^{-x} \quad \text{geometric distribution of ratio } 1/\sqrt[b]{2}.$$

## More precisely. . .

A finer analysis shows that Golomb coding is optimal (=Huffman) for a geometric distribution of ratio  $p$ , provided that  $b$  is chosen as

$$b = \left\lceil \frac{\log(2 - p)}{-\log(1 - p)} \right\rceil.$$

0	<b>10</b>
1	<b>110</b>
2	<b>111</b>
3	<b>010</b>
4	<b>0110</b>
5	<b>0111</b>
6	<b>0010</b>

Elias'  $\gamma$ 

Elias'  $\gamma$  coding of  $x \in \mathbf{N}^+$  is obtained by representing  $x$  in binary preceded by a unary representation of its length (minus one).  
More precisely, to represent  $x$  we write in unary  $\lfloor \log x \rfloor$  and then in binary  $x - 2^{\lfloor \log x \rfloor}$  (on  $\lfloor \log x \rfloor$  bits). So

$$l_x = 1 + 2\lfloor \log x \rfloor \quad \implies \quad p(x) \propto \frac{1}{2x^2} \text{ (Zipf of exponent 2)}$$

1	<b>1</b>
2	<b>010</b>
3	<b>011</b>
4	<b>00100</b>
5	<b>00101</b>

Elias'  $\delta$ 

Elias'  $\delta$  coding of  $x \in \mathbf{N}^+$  is obtained by representing  $x$  in binary preceded by a representation of its length in  $\gamma$ .

So

$$l_x = 1 + 2\lfloor \log \log x \rfloor + \lfloor \log x \rfloor \quad \implies \quad p(x) \propto \frac{1}{2x(\log x)^2}$$

1	<b>1</b>
2	<b>0100</b>
3	<b>0101</b>
4	<b>01100</b>
5	<b>01101</b>
6	<b>00100000</b>
7	<b>00100001</b>

## An alternative way. . .

. . . to think of  $\gamma$  coding is that  $x$  is represented using its usual binary representation (except for the initial “1”, which is omitted), with every bit “coming with” a continuation bit, that tells whether the representation continues or whether it stops there.

For example (up to bit permutation)  $\gamma$  coding of 724 (in binary: 1011010100) is

0	1	1	1	1	0	1	1	1	0	1	1	1	0	1	0	0
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---



## $k$ -bit-variable coding

What happens if we group digits  $k$  by  $k$ ?

0 1 1 1 1 1 0 1 1 1 0 1 1 1 0 1 0 0

0 0 1 1 1 1 0 1 1 0 1 1 0 0 0

0 1 1 1 0 1 0 1 1 0 0 0

1 1 0 1 1 0 1 0 0 0

## $k$ -bit-variable coding (cont'd)

For  $x$ , we use  $\lceil \log(x)/k \rceil$  bits for the unary part, and the same number of bits multiplied by  $k$  for the binary part.

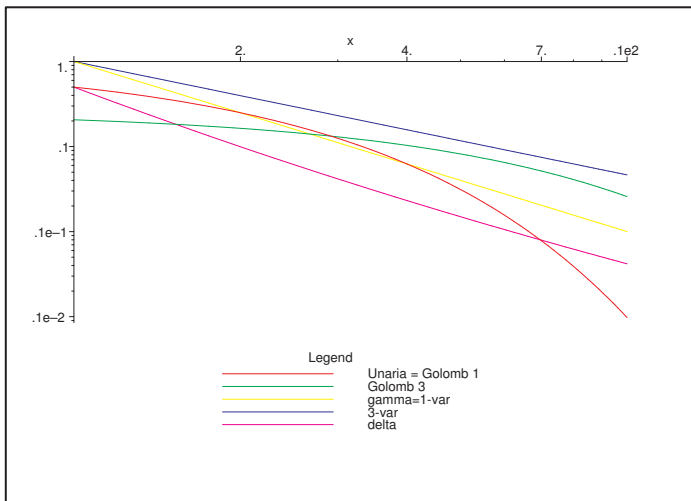
So

$$l_x = (k+1)(\lceil \log(x)/k \rceil) \implies p(x) \propto x^{-(k+1)/k} \text{ (Zipf } (k+1)/k)$$

A more efficient variant: the  $\zeta_k$  codes (for Zipf  $1 \rightarrow 2$ ).

	$\gamma = \zeta_1$	$\zeta_2$	$\zeta_3$	$\zeta_4$
1	1	10	100	1000
2	010	110	1010	10010
3	011	111	1011	10011
4	00100	01000	1100	10100
5	00101	01001	1101	10101
6	00110	01010	1110	10110
7	00111	01011	1111	10111
8	0001000	011000	0100000	11000

# Comparing codings



## Coding techniques. . .

. . . alone do not improve on compression: we have first to guarantee that the data we represent have a distribution close to the intended one (depending on the coding we are going to use). In particular, they have to enjoy a monotonic distribution (smaller values are more probable than larger ones).

- ▶ BTW: some codings (e.g., Elias  $\gamma$  and  $\delta$ ) are universal: for whatever monotonic distribution, they guarantee an expected length that is only within a constant factor of the optimal one.
- ▶ Degrees are distributed like a Zipf of exponent  $\approx 2.7$ : they can be safely encoded using  $\gamma$ .
- ▶ What about successors? Let us assume that successors of  $x$  are  $y_1, \dots, y_k$ : how should we encode  $y_1, \dots, y_k$ ?

## Locality

In general, we cannot say much about their distribution, unless we make some assumption on the way in which nodes are ordered.

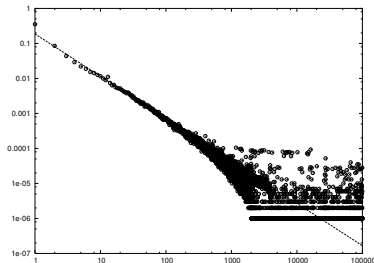
- ▶ Many hypertextual links contained in a web page are *navigational* (“home”, “next”, “up”...). If we compare the URL they refer to with that of the page containing them, they share a long common prefix. This property is known as *locality* and it was first observed by the authors of the Connectivity Server.
- ▶ To exploit this property, assume that URLs are ordered lexicographically (that is, node 0 is the first URL in lexicographic order, etc.). Then, if  $x \rightarrow y$  is an arc, most of the times  $|x - y|$  will be “small”.

# Exploiting locality

If  $x$  has successors  $y_1 < y_2 < \dots < y_k$ , we represent its successor list through the gaps (*differentiation*):

$$y_1 - x, y_2 - y_1 - 1, \dots, y_k - y_{k-1} - 1$$

(only the first value can be negative: we make it into a natural number...). How are such differences distributed?



Zipf with exponent 1.2  $\implies$  we use  $\zeta_3$ .

# Similarity

URLs close to each other (in lexicographic order) have similar successor sets: this fact (known as *similarity*) was exploited for the first time in the Link database.

We may encode the successor list of  $x$  as follows:

- ▶ we write the differences with respect to the successor list of some previous node  $x - r$  (called the *reference node*)
- ▶ we explicitly encode (as before) only the successors of  $x$  that were not successors of  $x - r$ .

## Similarity (cont'd)

More explicitly, the successor list of  $x$  is encoded as (*referencing*):

- ▶ an integer  $r$  (reference): if  $r > 0$ , the list is described by difference with respect to the successor list of  $x - r$ ; in this case, we write a bitvector (of length equal to  $d^+(x - r)$ ) discriminating the elements in  $N^+(x - r) \cap N^+(x)$  from the ones in  $N^+(x - r) \setminus N^+(x)$
- ▶ an explicit list of *extra nodes*, containing the elements of  $N^+(x) \setminus N^+(x - r)$  (or the whole  $N^+(x)$ , if  $r = 0$ ), encoded as explained before.



## Referencing example

Node	Outdegree	Successors
...	...	...
15	11	13, 15, 16, 17, 18, 19, 23, 24, 203, 315, 1034
16	10	15, 16, 17, 22, 23, 24, 315, 316, 317, 3041
17	0	
18	5	13, 15, 16, 17, 50
...	...	...

Node	Outd.	Ref.	Copy list	Extra nodes
...	...	...	...	...
15	11	0		13, 15, 16, 17, 18, 19, 23, 24, 203, 315, 1034
16	10	1	01110011010	22, 316, 317, 3041
17	0			
18	5	3	11110000000	50
...	...	...	...	...

Blocks (*differential compression*)

Instead of using a bitvector, we use run-length encoding, telling the length of successive runs (blocks) of “0” and “1”:

Node	Outdegree	Successors			
...	...	...			
15	11	13, 15, 16, 17, 18, 19, 23, 24, 203, 315, 1034			
16	10	15, 16, 17, 22, 23, 24, 315, 316, 317, 3041			
17	0				
18	5	13, 15, 16, 17, 50			
...	...	...			

Node	Outd.	Ref.	# blocks	Copy blocks	Extra nodes
...	...	...	...	...	...
15	11	0	...	...	13, 15, 16, 17, 18, 19, 23, ...
16	10	1	7	0, 0, 2, 1, 1, 0, 0	22, 316, ...
17	0				
18	5	3	1	4	50
...	...	...	...	...	...

# Consecutivity

Among the extra nodes, many happen to sport the *consecutivity* property: they appear in clusters of consecutive integers. This phenomenon, observed empirically, have some possible explanations:

- ▶ most pages contain groups of navigational links that correspond to a certain hierarchical level of the website, and are often consecutive to one another;
- ▶ in the transpose graph, moreover, consecutivity is the dual of similarity with reference 1: when there is a cluster of consecutive pages with many similar links, in the transpose there are intervals of consecutive outgoing links.

## Consecutivity (cont'd)

To exploit consecutivity, we use a special representation for the extra node list called *intervalization*, that is:

- ▶ sufficiently long (say  $\geq T$ ) intervals of consecutive integers are represented by their left extreme and their length minus  $T$ ;
- ▶ other extra nodes, if any, are called *residual nodes* and are represented alone.

# Intervalization example

Node	Outdegree	Successors
...	...	...
15	11	13, 15, 16, 17, 18, 19, 23, 24, 203, 315, 1034
16	10	15, 16, 17, 22, 23, 24, 315, 316, 317, 3041
17	0	
18	5	13, 15, 16, 17, 50
...	...	...

Node	Outd.	Ref.	# bl.	Copy bl.s	# int.	Lft extr.	Lth	Residuals
...	...	...	...	...	...	...	...	...
15	11	0			2	0, 2	3, 0	5, 189, 111, 718
16	10	1	7	0, 0, ...	1	600	0	12, 3018
17	0							
18	5	3	1	4	0			50
...	...	...	...	...	...	...	...	...

## Reference window

When the reference node is chosen, how far back in the “past” are we allowed to go? We need to keep track of a window of the last  $W$  successor lists. The choice of  $W$  is critical:

- ▶ a large  $W$  guarantees better compression, but increases compression time and space
- ▶ after  $W = 7$  there is no significant improvement in compression.

The choice of  $W$  does not impact on decompression time.

## Reference chain length

Referencing involves recursion: to decode the successor list of  $x$ , we need first to decompress the successor list of  $x - r$ , etc. This chain is called the *reference chain* of  $x$ : decompression speed depends on the length of such chains.

During compression, it is possible to limit their length keeping into account of how long is the reference chain for every node in the window and avoiding to use nodes whose reference chain is already of a given maximum length  $R$ .

The choice of  $R$  influences the compression ratio (with  $R = \infty$  giving the best possible compression) but also on decompression speed ( $R = \infty$  may produce access time that can be two orders of magnitude larger than  $R = 1$  — it may even produce stack overflows).