

WEB GRAPH

- WE CAN SEE THE STATIC WEB

(WEB CONTAINING STATIC PAGES, THOSE WHOSE CONTENT DOES NOT VARY FROM ONE REQUEST TO ANOTHER) AS A DIRECTED GRAPH IN WHICH EACH WEB PAGE IS A NODE AND EACH HYPERLINK IS A DIRECTED EDGE, CALLED THE WEB GRAPH.

- OF COURSE, THE WEB GRAPH IS NOT STRONGLY CONNECTED; YOU CAN'T REACH ANY PAGE YOU WANT FROM ONE PAGE, USING ONE HYPERLINK.

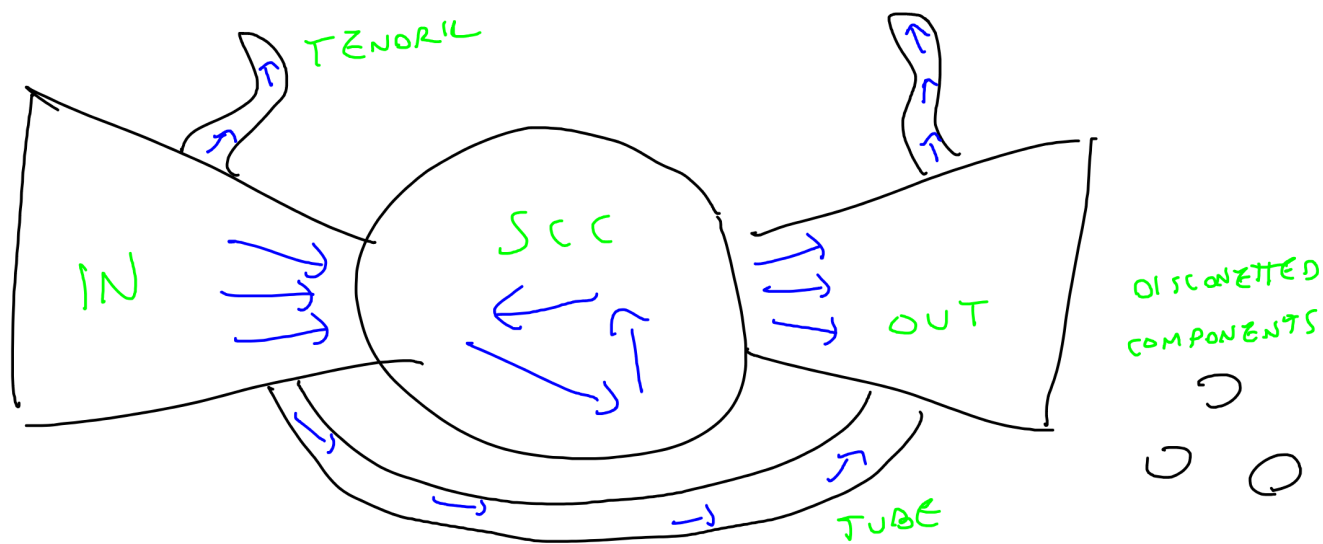
- WE CALL THE HYPERLINKS INTO A PAGE in-links, AND WE CALL out-links THOSE OUTSIDE

- THE NUMBER OF IN IN-LINKS TO A PAGE IS CALLED IN-DEGREE, AND THE NUMBER OF LINKS OUT OF THE PAGE IS CALLED OUT-DEGREE

→ THE NUMBER OF LINKS IS NOT RANDOMLY DISTRIBUTED (DOESN'T FOLLOW THE POISSON DISTRIBUTION, THAT WOULD HAVE BEEN THE CASE IF EACH PAGE WAS ABLE TO PICK A LINK'S DESTINATION UNIFORMLY AT RANDOM)

RATHER, THE DISTRIBUTION FOLLOWS THE POWER LAW: THE TOTAL NUMBER OF WEB-PAGES WITH IN-DEGREE i IS PROPORTIONAL TO $\frac{1}{i^\alpha}$ (FEW NODES WITH LOTS OF LINKS, LOTS OF NODES WITH FEW LINKS)

→ THE WEB GRAPH IS BOW-TIE SHAPED



3 TYPES OF WEB PAGES =

IN = YOU CAN PASS FROM ANY PAGE IN IN TO ANY PAGE IN SCC (BUT NOT VICEVERSA)

SCC = YOU CAN PASS FROM ANY PAGE IN SCC TO ANY PAGE IN SCC OR OUT (BUT, IN THE LATTER CASE, YOU CAN'T GO BACK TO SCC PAGES)

OUT = NO HYPERLINKS HERE

THE WEB GRAPH, MORE FORMALLY

→ Directed graph $G = (V, E)$, where
 $V = \text{URLS}$, $E = (u, v)$ where u has an
hyperlink to v

→ SKEWED DISTRIBUTION

$$P(\text{in-degree}(u) = k) \propto \frac{1}{k^\alpha} \quad \text{where } \alpha = 2.1$$

POWER LAW DISTR.

NOTE: THIS IS ALSO TRUE FOR OUT-DEGREE
AND THE SIZE OF THE 3 CATEGORIES

→ LOCALITY

MOST OF THE LINKS FROM URL u POINTS TO
OTHER URLS IN THE SAME HOST OF u (80%)

→ SIMILARITY

IF URLS u AND v ARE CLOSE IN
LEXICOGRAPHIC ORDER, THEN THEY TEND
TO SHARE MANY HYPERLINKS

Other characteristics

- 1 trillion pages available, 50 billions of which are crawled
- 5-40Kb per page
- Size grows everyday
- It's a dynamic graph :
 - Life time of a link is 10 days
 - 8% new pages, 25% new links daily
- Average 10 links per page

EXPLOIT LOCALITY AND SIMILARITY TO COMPRESS THE ADJACENT LISTS OF THE GRAPH

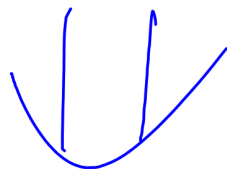
SAME DOMAIN \Rightarrow CLOSE TO EACH OTHER IN
LEXICOGRAPHIC ORDER

\Downarrow
CLOSE NODES

\Downarrow
COMPRESS THE NODES WITH
GO P-ENCODING AND VARIABLE-LENGTH
REPRESENTATIONS

UNCOMPRESSED ADJACENCY LIST

NODE	OUTdegree	Successors
\vdots	\vdots	\vdots
15	11	13, 15, 16, 17, 18, 19, 23, 24, 203, 315, 1034
16	10	15, 16, 17, 22, 23, 24, 315, 316, 317, 3047
17	0	
18	5	13, 15, 16, 17, 50
\vdots	\vdots	\vdots



COMPRESSED ADJACENCY LIST

NODE	OUTDEG	REF	COPYLIST	EXTRA NODES
⋮	⋮	⋮	⋮	
15	11	0		SAME AS BEFORE
16	10	1	01110011010	24, 316, 317, 30+1
17	0			
18	5	3	11110000000	50
⋮	⋮	⋮		

→ The reference index ($REF=0$) is the one that gives the best compression

→ EACH BIT OF THE COPYLIST INFORMS WHETHER THE CORRESPONDING SUCCESSOR OF

$(REF \neq 0)$ IS ALSO A SUCCESSOR OF X ($REF=0$)

WE CAN ALSO USE THE RUN-LENGTH

ENCODING TO THE BIT SEQUENCES OF THE

COPY LIST IN ORDER TO COMPUTE THE

ADJACENT LIST WITH COPY BLOCKS

COPY BLOCKS \equiv RLE (COPY LIST)

- ① → TAKE THE COPY LIST: EACH GROUP OF CONSECUTIVE EQUAL BIT IS A BLOCK
- ② → IF THE COPYLIST STARTS WITH A (group of) ZERO, THE FIRST BIT OF THE COPY-BLOCK IS A ZERO, OTHERWISE IT'S A ONE
- ③ → FOR EACH OTHER BIT, COUNTS THE GROUP OF EQUAL BITS (A BLOCK), AND PUT IN THE COPY-BLOCK THEIR LENGTH MINUS ONE
- ④ → THE LAST BLOCK IS OMITTED, SINCE WE KNOW THE OUT-DEGREE OF THE NODE

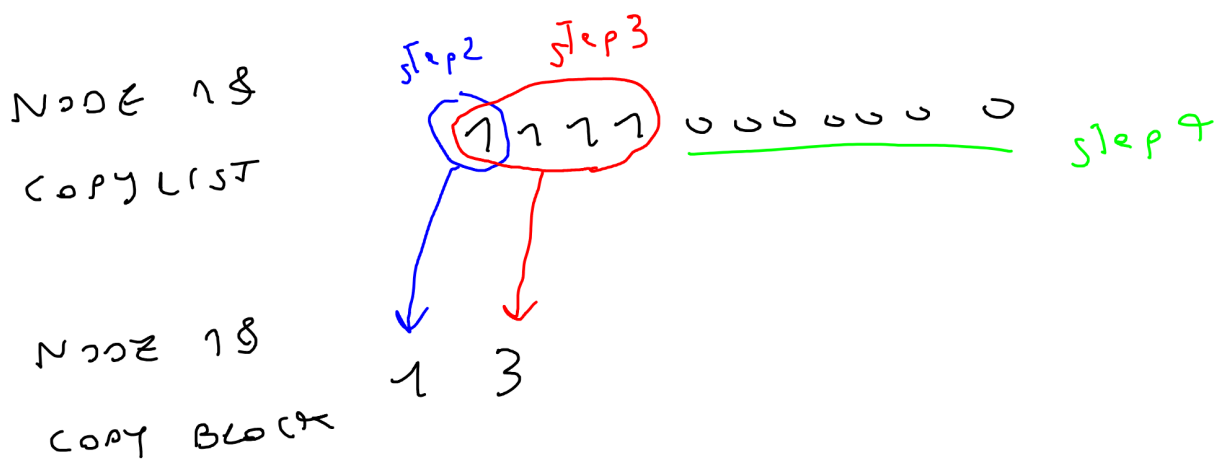
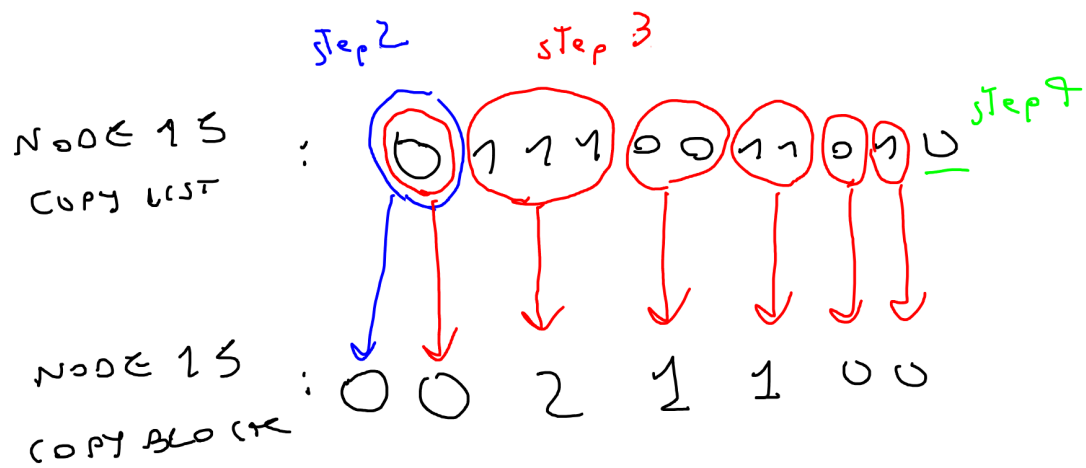
COPY-LIST

NODE	OUTDEG.	REF	COPYLIST	EXTRA NODES
⋮	⋮	⋮	⋮	
15	17	0		SAME AS BEFORE
16	10	1	01110011010	24, 316, 317, 30+1
17	0			
18	5	3	11110000000	50
⋮	⋮	⋮		



COPY BLOCK

NODE	OUTDEG.	REF	#BLOCKS	COPY BLOCKS	EXTRA NODES
⋮	⋮	⋮	⋮	⋮	
15	17	0			SAME AS BEFORE
16	10	1	7	0, 0, 2, 1, 1, 0, 0	24, 316, 317, 30+1
17	0				
18	5	3	2	1, 3	50
⋮	⋮	⋮			



EXAMPLE

GIVEN TWO ADJACENT LISTS FOR TWO WEBAVES (NODES) :

14 \rightarrow 3, 10, 11, 13, 14, 17, 18, 21, 25

15 \rightarrow 5, 10, 11, 12, 14, 17, 19, 20, 21, 24, 33

SHOW HOW TO COMPRESS LIST 15 USING

COPY LISTS (1) AND THEN COPY BLOCKS (2)

(1)

Node	Ref	copy list	EXTRA NODES
14	0		AS BEFORE
15	1	0 1 1 0 1 1 1 1 0	5, 12, 20, 24, 33

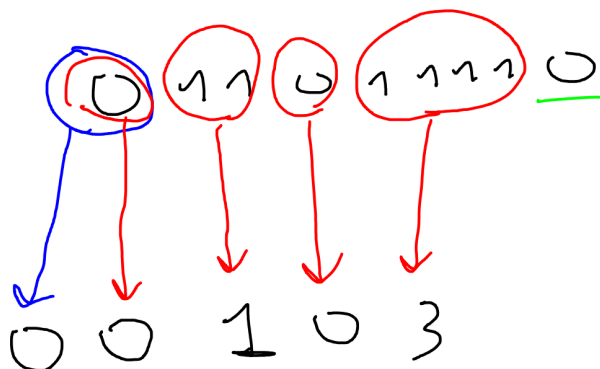
(2)

NODE 15

COPY-LIST

NODE 15

COPY BLOCK



blocks = 5