

程序代写代做 CS编程辅导

---



# COMP3349 Web Data Compression and Search

Assignment Project Exam Help

JSX & XBW;  
Email: [tutorcs@163.com](mailto:tutorcs@163.com)  
(Web) Graph Compression  
QQ: 749389476

<https://tutorcs.com>

# ISX Requirements



1. Space does matter for many applications
2. Generally reducing space improves cache locality  
WeChat: cstutorcs
3. Indirection is expensive  
Assignment Project Exam Help
4. Support fast navigations  
Email: tutorcs@163.com
5. Support fast insertion and deletion  
QQ: 749389476
6. Support efficient joins  
<https://tutorcs.com>
7. Separate topology, text and schema

# ISX Goal

程序代写代做 CS编程辅导



- To find a space-efficient storage scheme for XML data without compromising both query and update performances

WeChat: cstutorcs  
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

# Proposed Storage Structure



**Topology Layer**



WeChat: cstutorcs

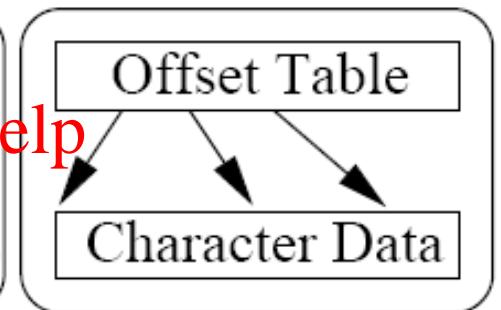
Assignment Project Exam Help

Email: tutorcs@163.com

**Internal Node Layer**  
(Tags)

Symbol Table,  
Topology Labels +  
Text Data Signatures

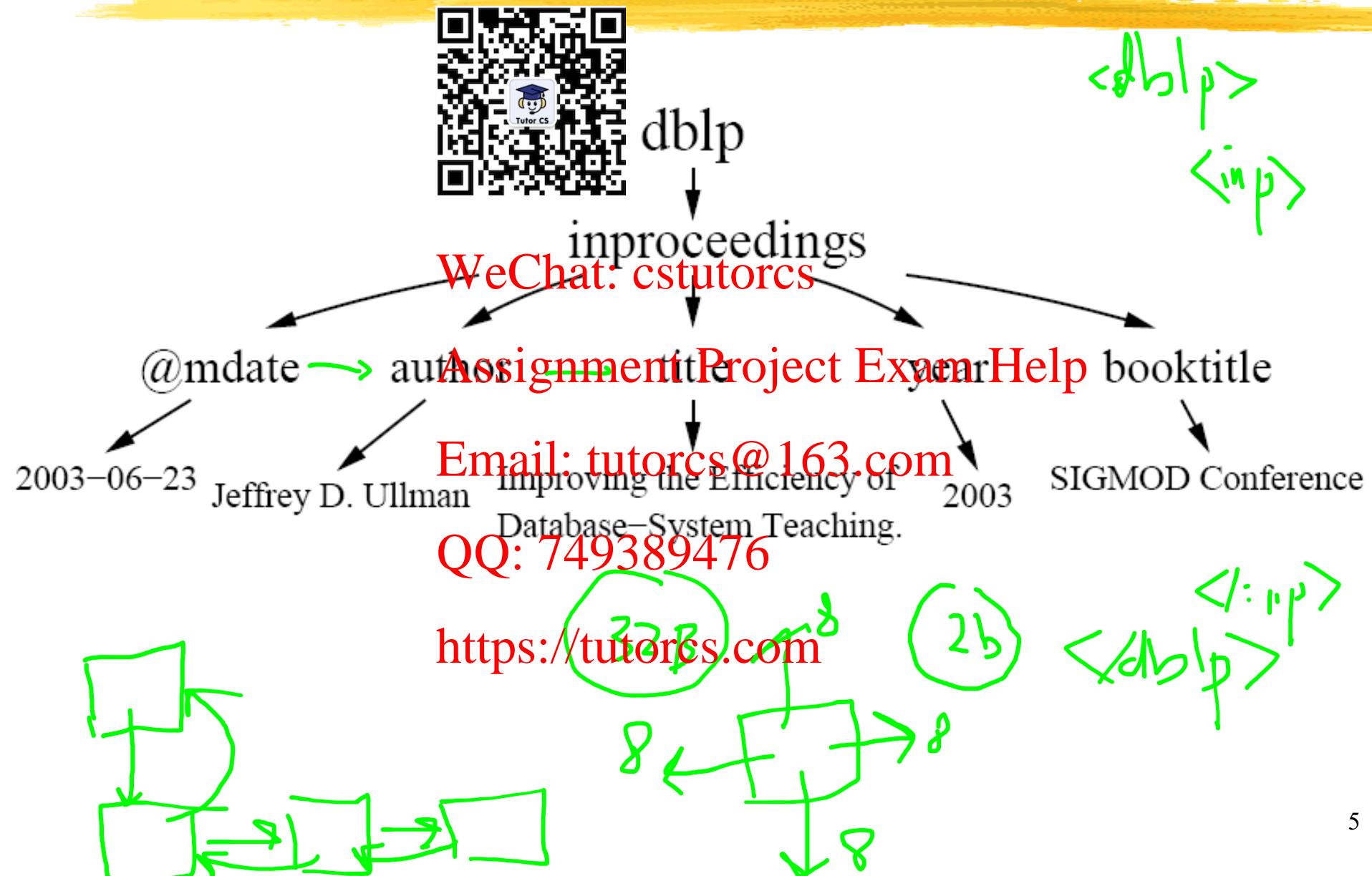
**Leaf Node Layer**  
(Text Data)



QQ: 749389476

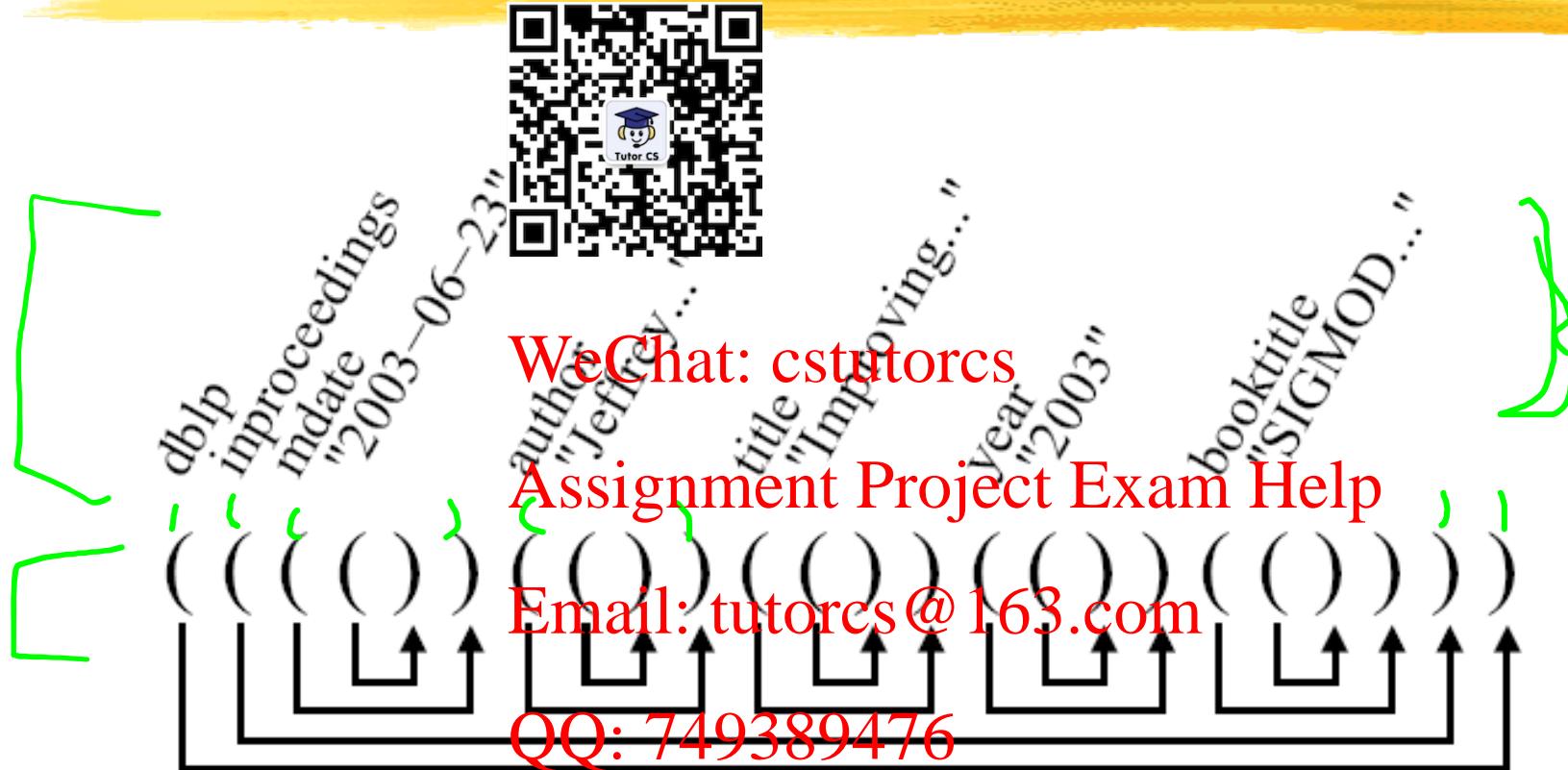
The ISX Structure  
<https://tutorcs.com>

# Sample DBLP XML Fragment



# Balanced Parenthesis Encoding

程序代写代做 CS编程辅导

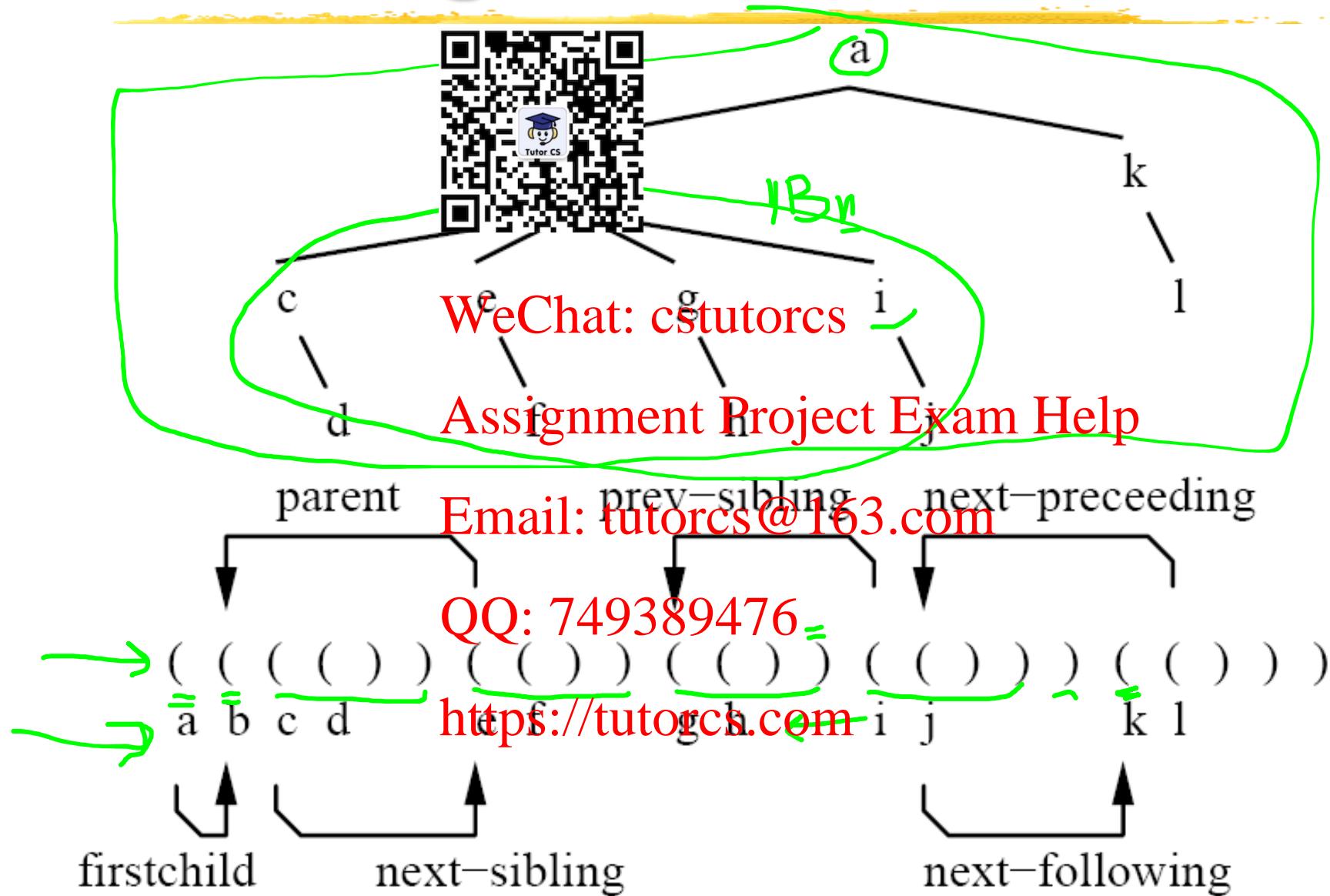


<https://tutorcs.com>

→ 0 0 0 0 1 1 0 0 1 1 0 0 1 1 0 0 1 1 1 1

# Node Navigations

程序代写代做 CS 编程辅导



# Primitive operators



## Algorithm

### FORWARDEXCESS(*start*, *end*)

```
1 for each current from start to end do
2   if (tier0[current] is an open parenthesis) then
3     k ← k + 1
4   else
5     k ← k + 1
6     if (k = 0) then
7       return current
8 return NOT-FOUND
```

### BACKWARDEXCESS(*start*, *end*, *k*)

```
1 for each current from start to end step -1 do
2   if (tier0[current] is an open parenthesis) then
3     k ← k - 1
4   else
5     k ← k + 1
6   if (k = 0) then
7     return current
8 return NOT-FOUND
```

### PREV(*node*)

```
1 if (node > 0) then return node - 1 else return NOT-FOUND
```

### NEXT(*node*)

```
1 if (node < |tier0|) then return node + 1 else return NOT-FOUND
```

### FINDCLOSE(*node*)

```
1 return FORWARDEXCESS(node, |tier0|, 0)
```

### FINDOPEN(*node*)

```
1 return BACKWARDEXCESS(node, |tier0|, 0)
```

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

# Tier 2 excess 程序代写代做 CS编程辅导



## Algorithm 3 Calculate Local Excess

## Excess in a Tier 2 Block

TIER2LOCALEXCESS( $t_2$ )

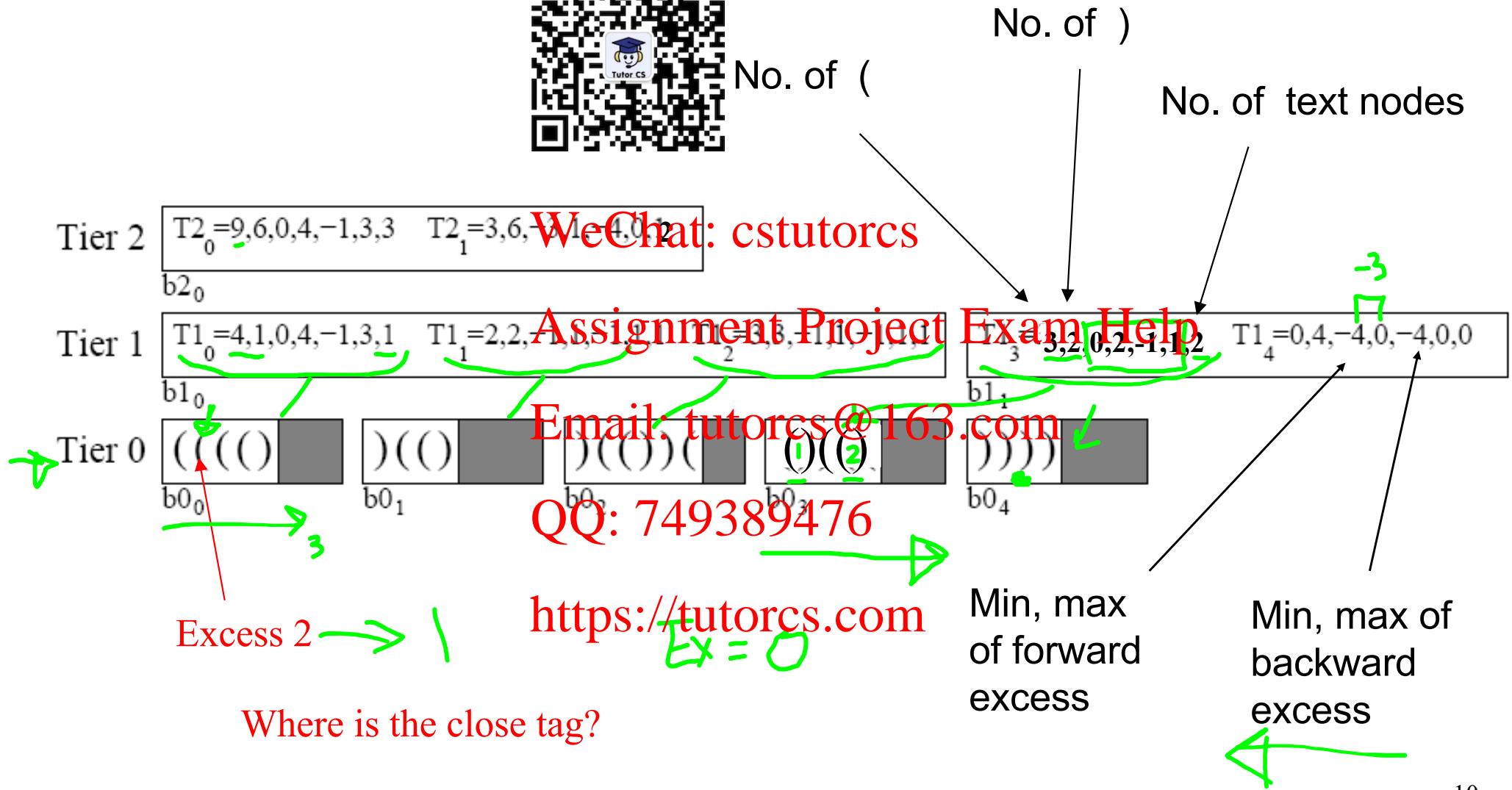
```
1    $\{t_{1\_start}, t_{1\_end}\} \leftarrow \{\frac{t_2 * |T^2|}{|T|}, \frac{(t_2 + 1) * |T^2|}{|T|} - 1\}$ 
2    $\{tier2[t_2].m, tier2[t_2].M\} \leftarrow \{tier1[t_{1\_start}].m, tier1[t_{1\_start}].M\}$ 
3    $excess \leftarrow tier1[t_{1\_start}] = tier1[t_{1\_start}] \cdot R$ 
4   for each  $t_1$  from  $t_{1\_start} + 1$  to  $t_{1\_end}$  do
5       if ( $excess + tier1[t_1].m < tier2[t_2].M$ ) then
6            $tier1[t_1].m \leftarrow excess + tier1[t_1].m$ 
7       if ( $excess + tier1[t_1].M > tier2[t_2].M$ ) then
8            $tier1[t_1].M \leftarrow excess + tier1[t_1].M$ 
9    $excess \leftarrow excess + tier1[t_1].L - tier1[t_1].R$ 
```

WeChat: cstutorcs  
Assignment Project Exam Help

Email: tutorcs@163.com  
QQ: 749389476

<https://tutorcs.com>

# Topology Tiers



# Efficient Updates



Density Threshold      Depth

[0.50, 0.75]      0

$$d_9 = 37.5\%$$

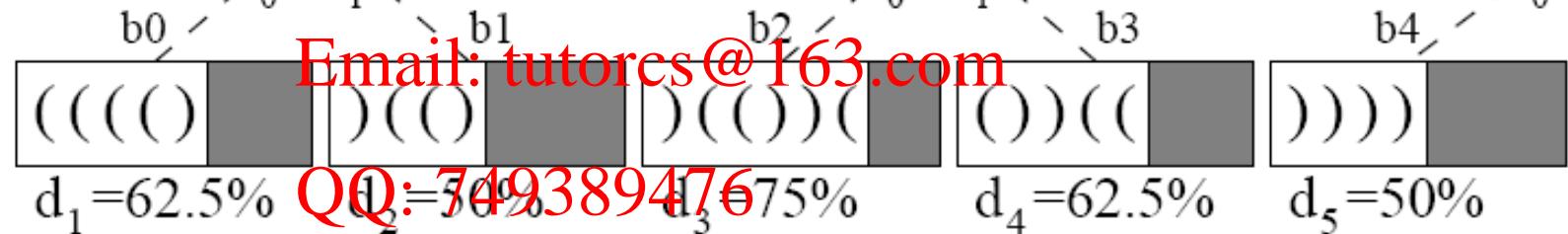
[0.42, 0.83]      1

WeChat: cstutorcs

[0.33, 0.92]      2

[0.25, 1.00]      3

Assignment Project Exam Help



d: density within a range of blocks      height of virtual binary trie: 3  
<https://tutores.com>

# Example

程序代写代做 CS编程辅导



- 100 MB DBLP document
  - 5 million XML nodes
  - ISX: 1MB topology
- WeChat: cstutorcs  
Assignment Project Exam Help  
Email: tutorcs@163.com  
QQ: 749389476  
<https://tutorcs.com>

# Another example



- Core Duo 1
- 1GB RAM
- 5400 RPM Harddrive

• MS Vista

WeChat: cstutorcs

5M DBLP	MSXML	ISX
Runtime (loading)	15MB	4MB
Loading time	0.54s	QQ: 749389476
Runtime (//www)	21MB	4MB
//www	0.096s	0.004s

100M DBLP	MSXML	ISX
Runtime (loading)	329MB	67MB
Loading time	17.8s	0.67s
Runtime (//www)	333MB	67MB
//www	1.814s	0.143s

# ISX Features

程序代写代做 CS编程辅导



Features	XMill	XGrind	NoK	TIMBER	ISX
Compression	WeChat: cstutorcs				✓
Document Traversal	✓	✓	✓	✓	✓
Node Navigation of All Axes					✓
Update Operation				✓	✓
Support XPath Query	Email: tutorcs@163.com	✓	✓	✓	✓

QQ: 749389476

COMPARISON OF SUPPORTED FEATURES

<https://tutorcs.com>

# Experiments

程序代写代做 CS编程辅导

## Setup



- Fixed at **64MB memory buffer** WeChat: cstutorcs
- Up to 16 GB XML document Assignment Project Exam Help
- E.g. 16 GB DBLP contains Email: tutorcs@163.com **>770 million nodes**
- **NO** index or query optimization has been employed for ISX (*except for ISX Stream where TurboXPath algorithm has been employed*) QQ: 749389476 <https://tutorcs.com>

# Storage Size (ISX vs NoK)



Document Size (MB)			PSD		TreeBank	
	NoK	ISX	NoK	ISX	NoK	ISX
5	18	3.64	17.91	3.36	19	3.21
10	35	7.25	36.12	6.82	38	6.36
50	181	36.1	182.42	34.14	196	31.78
100	367	72.1	377.52	68.74	389	63.43
250	918	180.2	950	171.9	974	159

WeChat: cstutorcs  
Assignment Project Exam Help  
Email: tutorcs@163.com

TABLE II  
QQ: 749389476

STORAGE SIZE OF ISX vs. NoK

<https://tutorcs.com>

# Storage Size (ISX, XMill, XGrind): DBLP

程序代写代做 CS编程辅导



Source Data (MB)	ISX (MB)	ISX Compressed (MB)	ISX Uncompressed (MB)	XMill (MB)	Source Data (MB)	ISX (MB)	ISX Compressed (MB)	XMill (MB)	XGrind 1 (MB)
1	1	0.4	0.1	0.3	256	182	82.7	31.5	75.0
2	1	0.7	0.3	0.6	500	363	163.7	62.6	Failed
5	3	1.5	0.6	1.3	750	549	249.7	94.0	Failed
8	5	2.5	0.9	2.1	1000	726	327.5	125.3	Failed
16	10	5	1.8	4.3	2000	1452	654.9	250.5	Failed
32	21	10	3.7	8.6	4000	2903	1309.8	501.0	Failed
64	42	20	7.2	17.4	8000	5807	2619.6	978.48	Failed
128	87	40.2	14.9	35.8	16000	9411	4629.9	1952.81	Failed

WeChat: cstutorcs  
Assignment Project Exam Help

Email: tutorcs@163.com

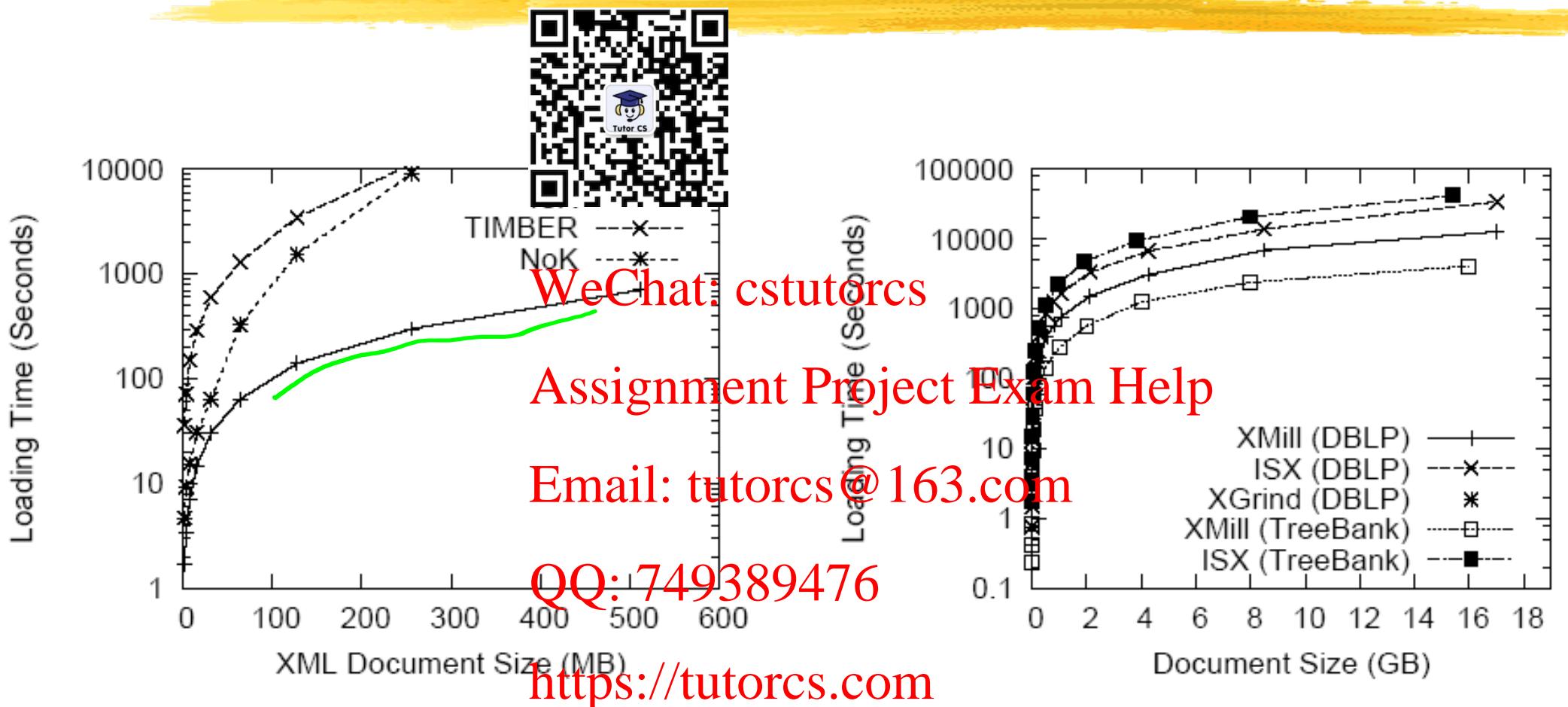
TABLE III

STORAGE SIZE OF ISX (WITH AND WITHOUT TEXT COMPRESSION), XMILL AND XGRIND ON DBLP

QQ: 749389476

<https://tutorcs.com>

# Bulk Loading Performance



# Queries

程序代写代做 CS编程辅导



Query #	XPath Expression	1 GB	2 GB	4 GB	8 GB	16 GB
		Final	Final	Final	Final	Final
Q1	//inproceedings	402667	981484	2012761	4160339	8453066
Q2	//mastersthesis	74	156	315	627	1251
Q3	/dblp/article	442184	717449	1379945	2630711	5135130
Q4	//inproceedings/title	402667	981484	2012761	4160339	8453066
Q5	//article[./month/text() = "July"]//title	812209	172419	3454708	6920136	13848372
Q6	//inproceedings[./ee]//pages	796742	1607116	3210628	6430194	12868471

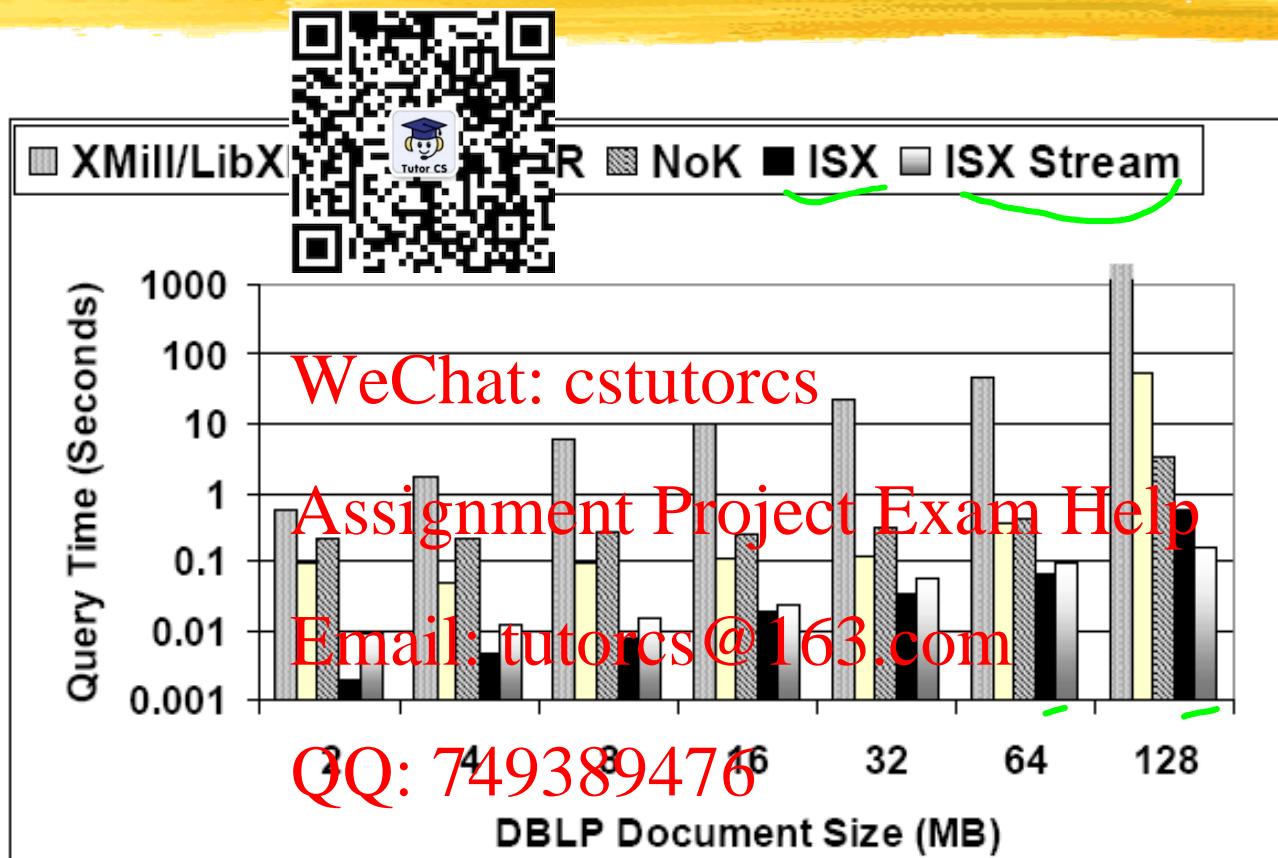
WeChat: cstutorcs  
Assignment Project Exam Help

Table 4: Test Queries and Final Result Sizes

QQ: 749389476

<https://tutorcs.com>

# Q1: //inproceedings

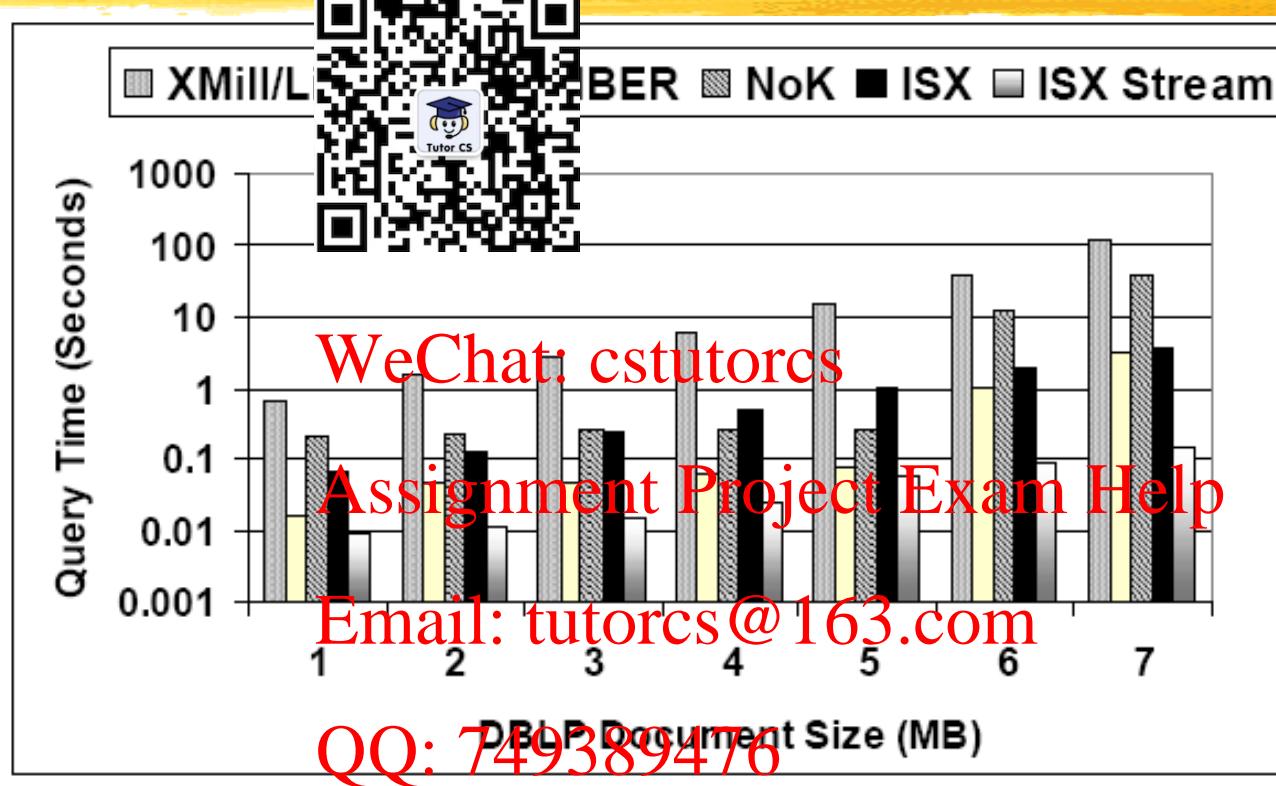


<https://tutorcs.com>

(a) DBLP Q1

# Q5: //article[.//month/text() = “July”]//title

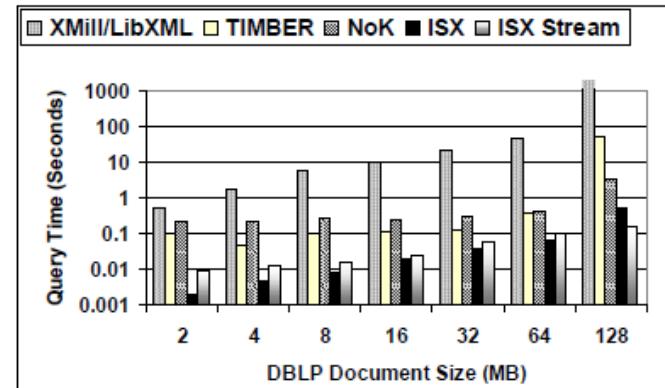
程序代写代做 CS编程辅导



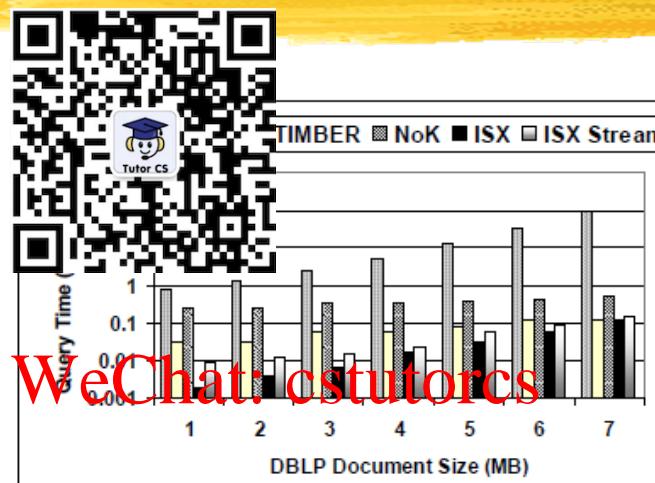
<https://tutorcs.com>

(e) DBLP Q5

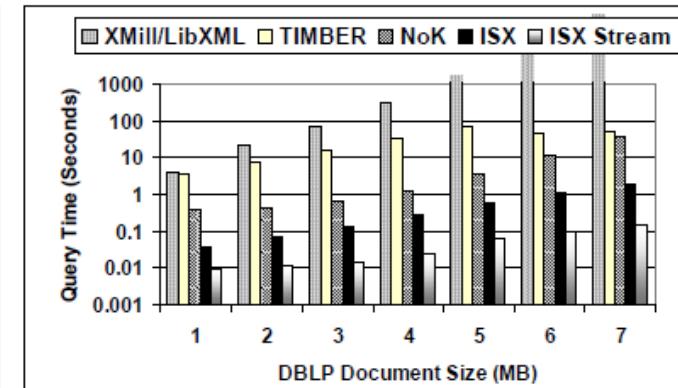
# Other queries



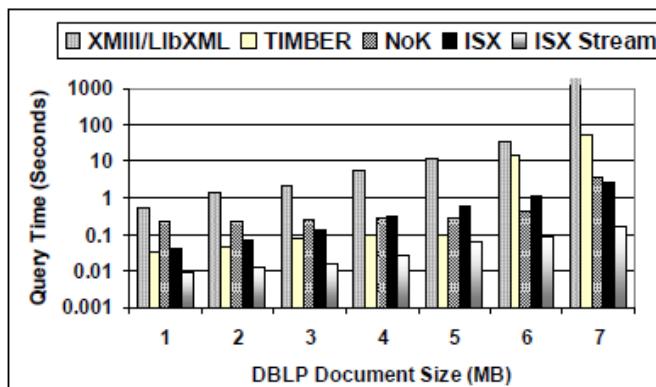
9(a) DBLP Q1



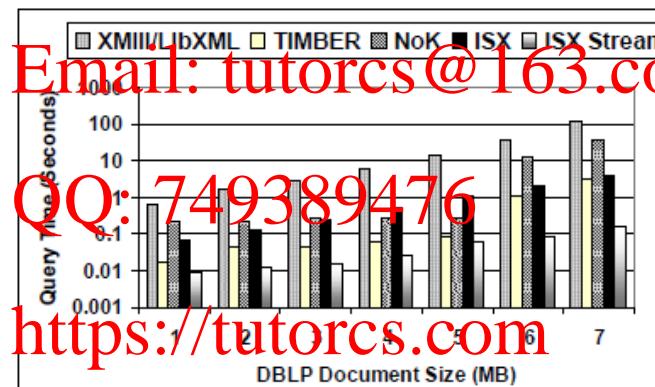
9(b) DBLP Q2



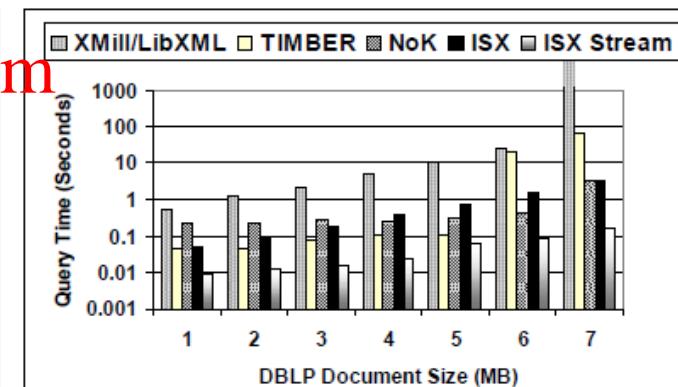
9(c) DBLP Q3



9(d) DBLP Q4



9(e) DBLP Q5



9(f) DBLP Q6

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

# XPath 13 axes

We can navigate along 13 axes:

ancestor

ancestor-or-self

attribute

child

descendant

descendant-or-self

following

following-sibling

namespace

parent

preceding

preceding-sibling

self



WeChat: cstutorcs

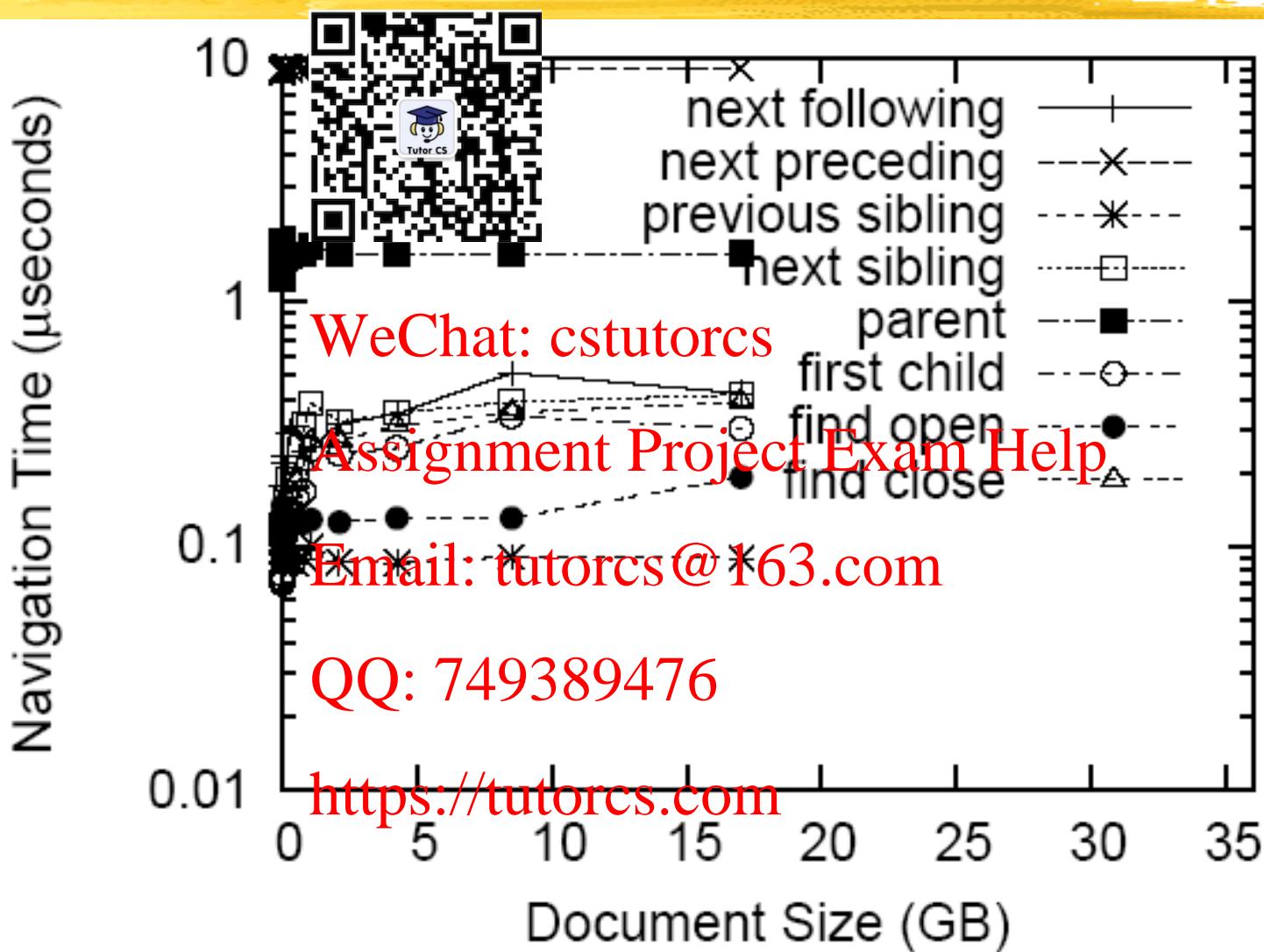
Assignment Project Exam Help

Email: tutorcs@163.com

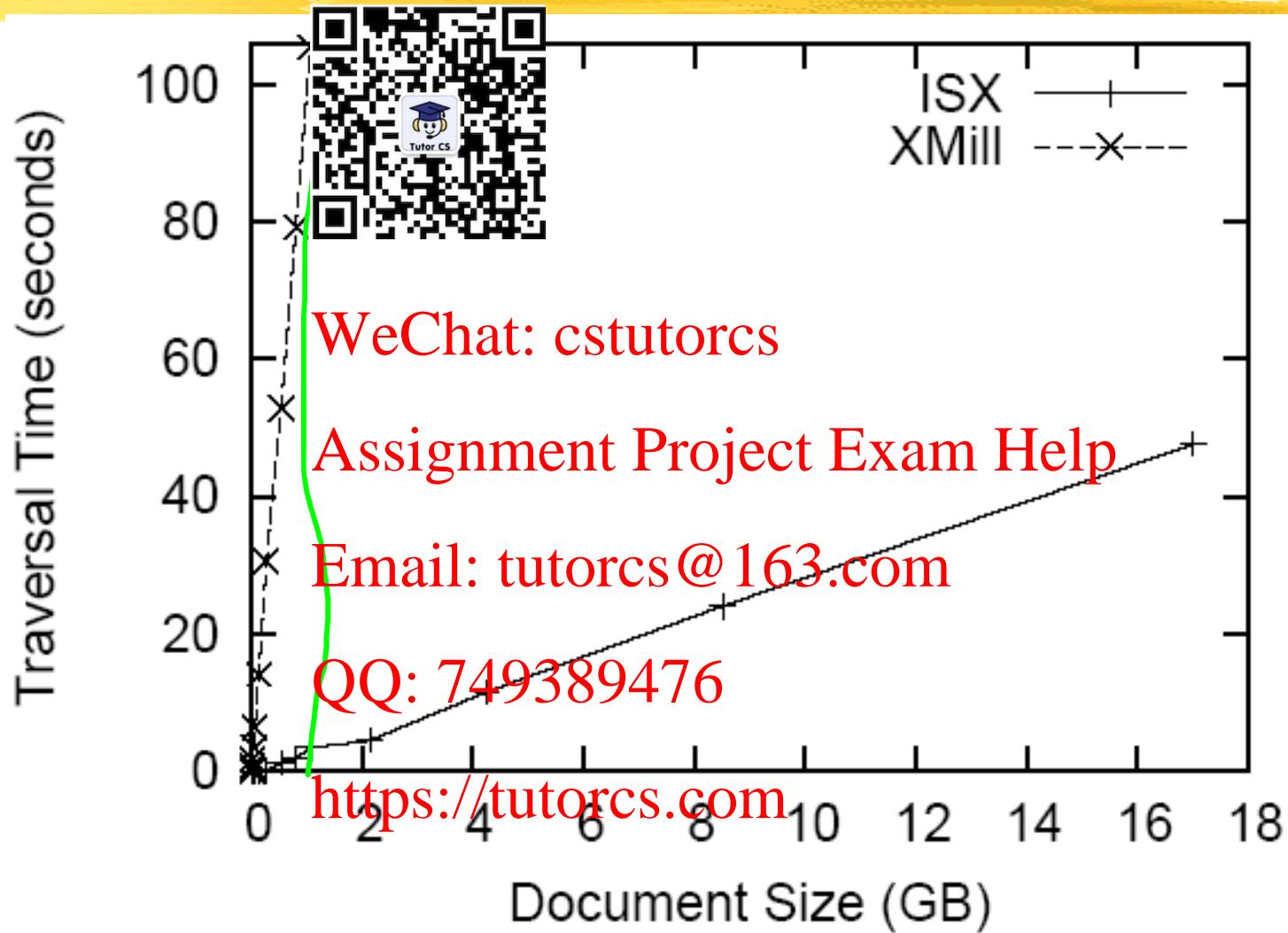
QQ: 749389476

<https://tutorcs.com>

# Node Navigation

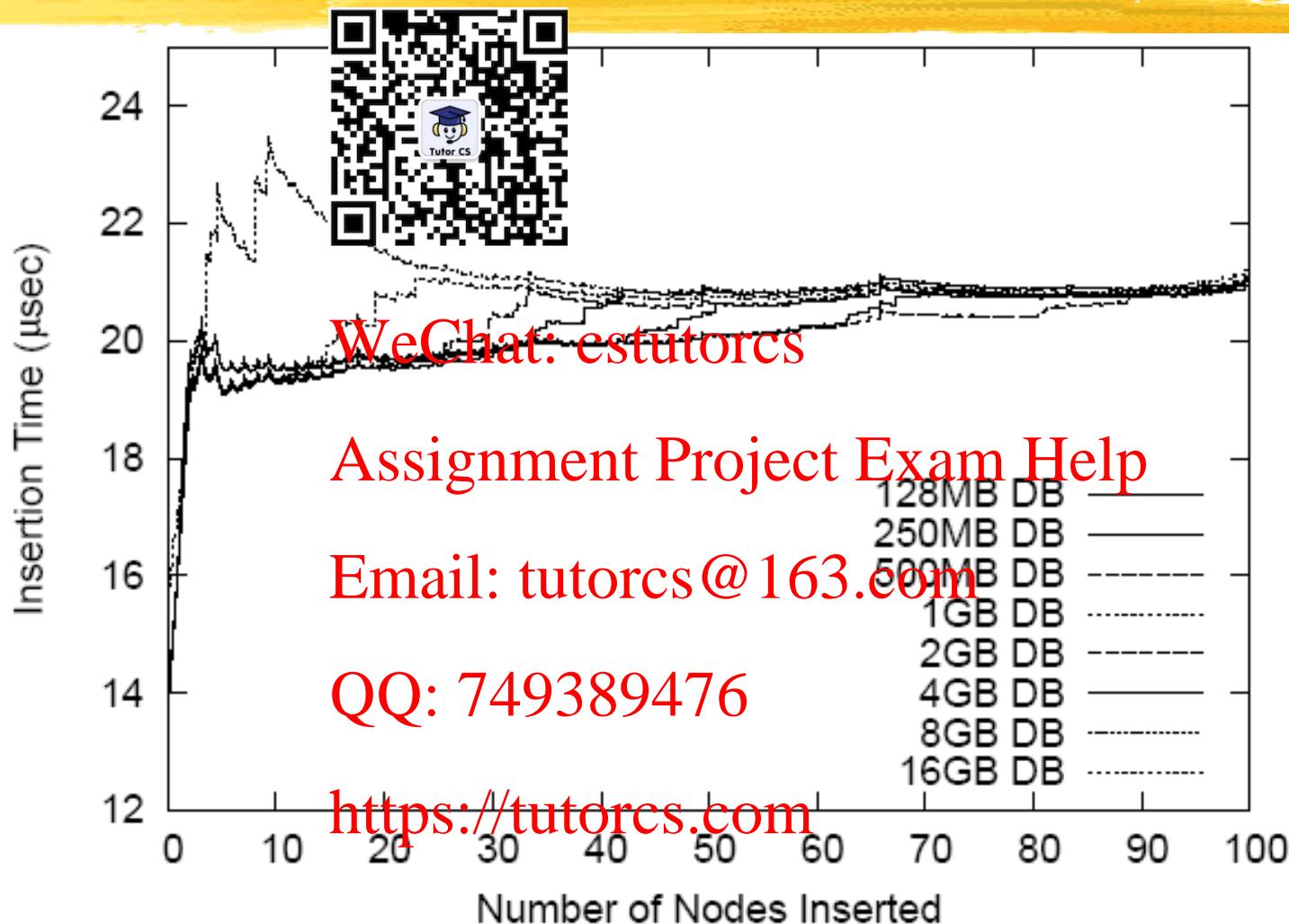


# Full document traversal



# Update (Insertion) Performance

程序代写代做 CS 编程辅导



# ISX Summary



- Small storage footprint
- Small runtime footprint
- Fast and consistent performance on navigational access
- Superior query performance (further indexing / query optimization can be added)
- Superior update performance

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

# Compressing and Searching XML Data Via Two Zips



WeChat: cstutorcs

Assignment Project Exam Help

Paolo Ferragina et al.  
Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

# A transform for “labeled trees”

[Ferragina et al, IEEE Focs '05]

程序代写代做 CS编程辅导



- ✓ We proposed the **XBW Transform** that mimics on trees the nice structural properties of the Burrows-and-Wheeler Transform on strings

WeChat: cstutorcs

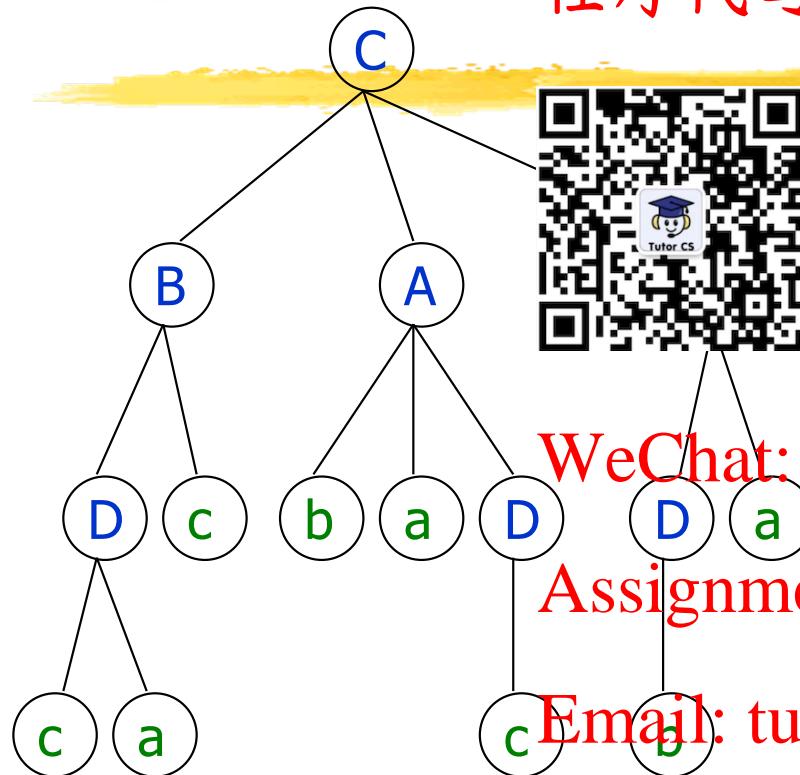
- ✓ The XBW **linearizes** the tree  $T$  in  $\mathcal{O}(n \log n)$  s.t.

- ✓ the **compression** of  $T$  reduces to use any compressor (*gzip*, *bzip*, ... ) over these two arrays  
QQ: 749389476

- ✓ the **indexing** of  $T$  reduces to implement simple **rank/select** query operations over these two arrays  
<https://tutorcs.com>

# The XBW-Transform

程序代写代做 CS 编程辅导



WeChat: cstutorcs  
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

Step 1.

Visit the tree in pre-order.

For each node, write down its label  
and the labels on its upward path

<https://tutorcs.com>

Permutation  
of tree nodes

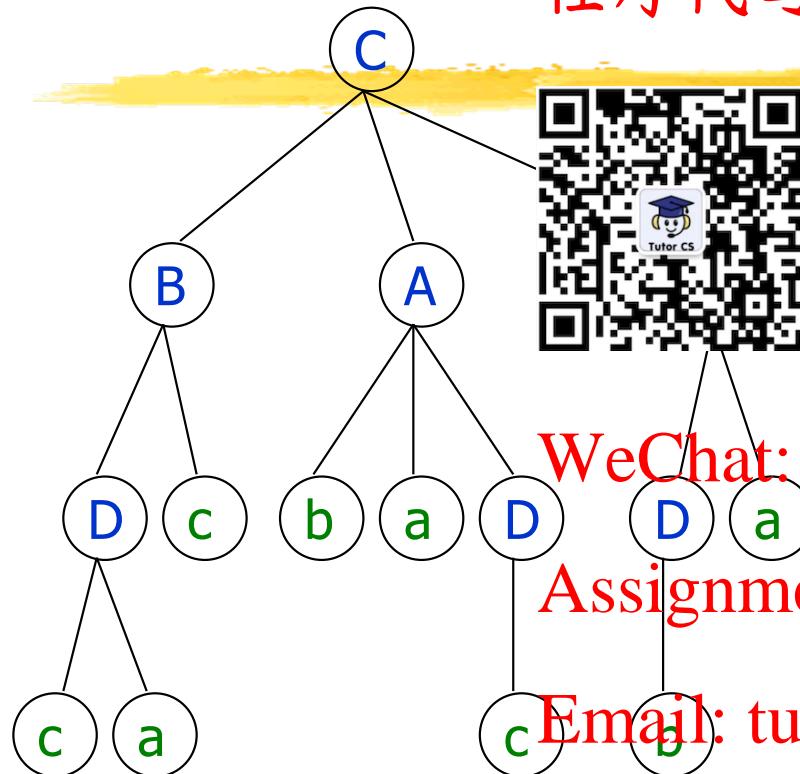
$S_\alpha$        $S_\pi$

C	$\epsilon$
B	C
D	BC
c	DBC
a	DBC
c	BC
A	C
b	A
a	AC
D	AC
c	DAC
B	C
D	BC
b	BC
a	BC

upward labeled paths

# The XBW-Transform

程序代写代做 CS编程辅导



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Step 2.

Stably sort according to  $S_\pi$

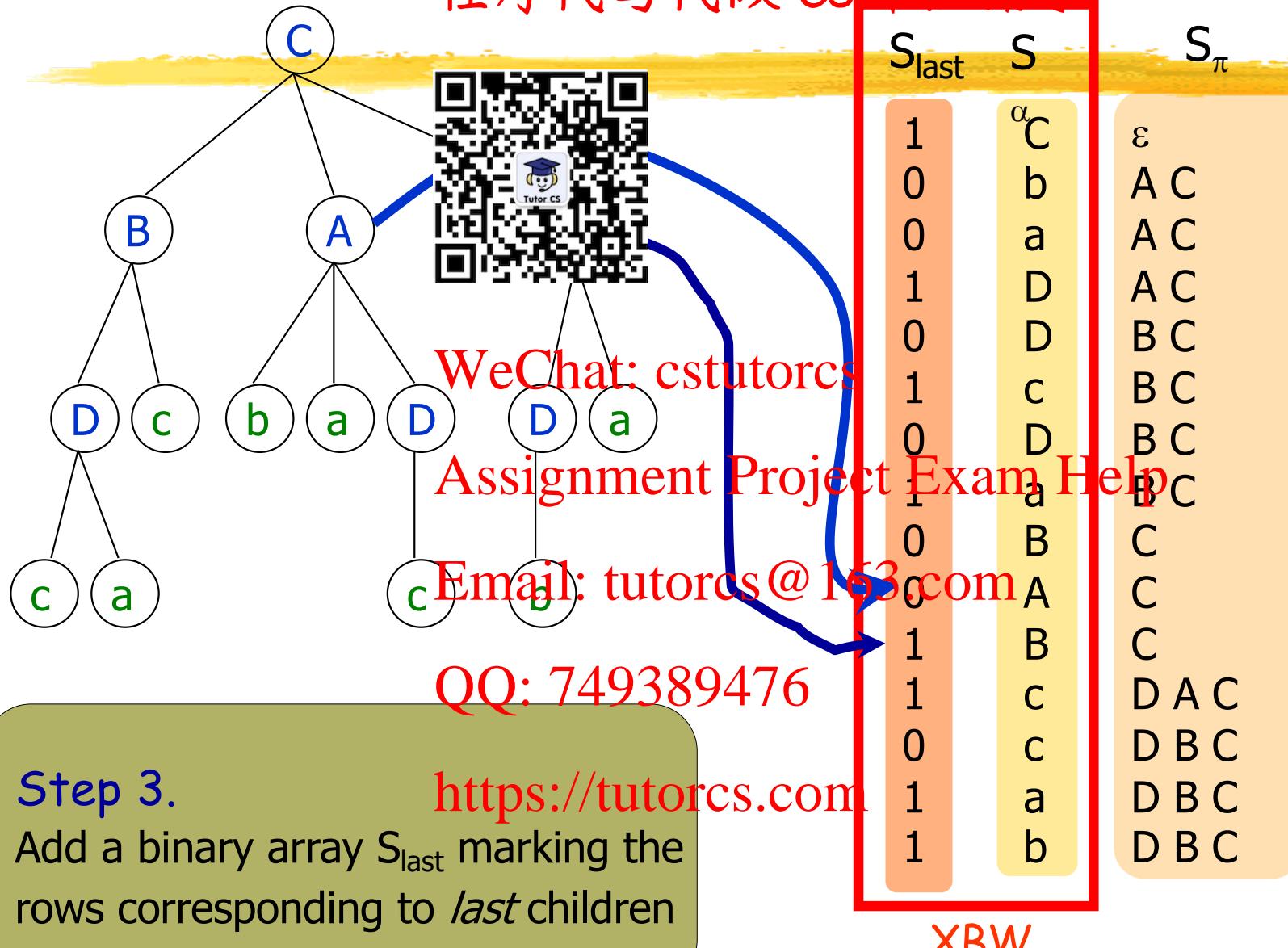
$S_\alpha$        $S_\pi$

C	$\epsilon$
b	A C
a	A C
D	A C
D	B C
c	B C
D	B C
a	B C
c	B C
D	C
a	C
B	C
A	C
B	C
c	D A C
c	D B C
a	D B C
b	D B C

↑  
upward labeled paths

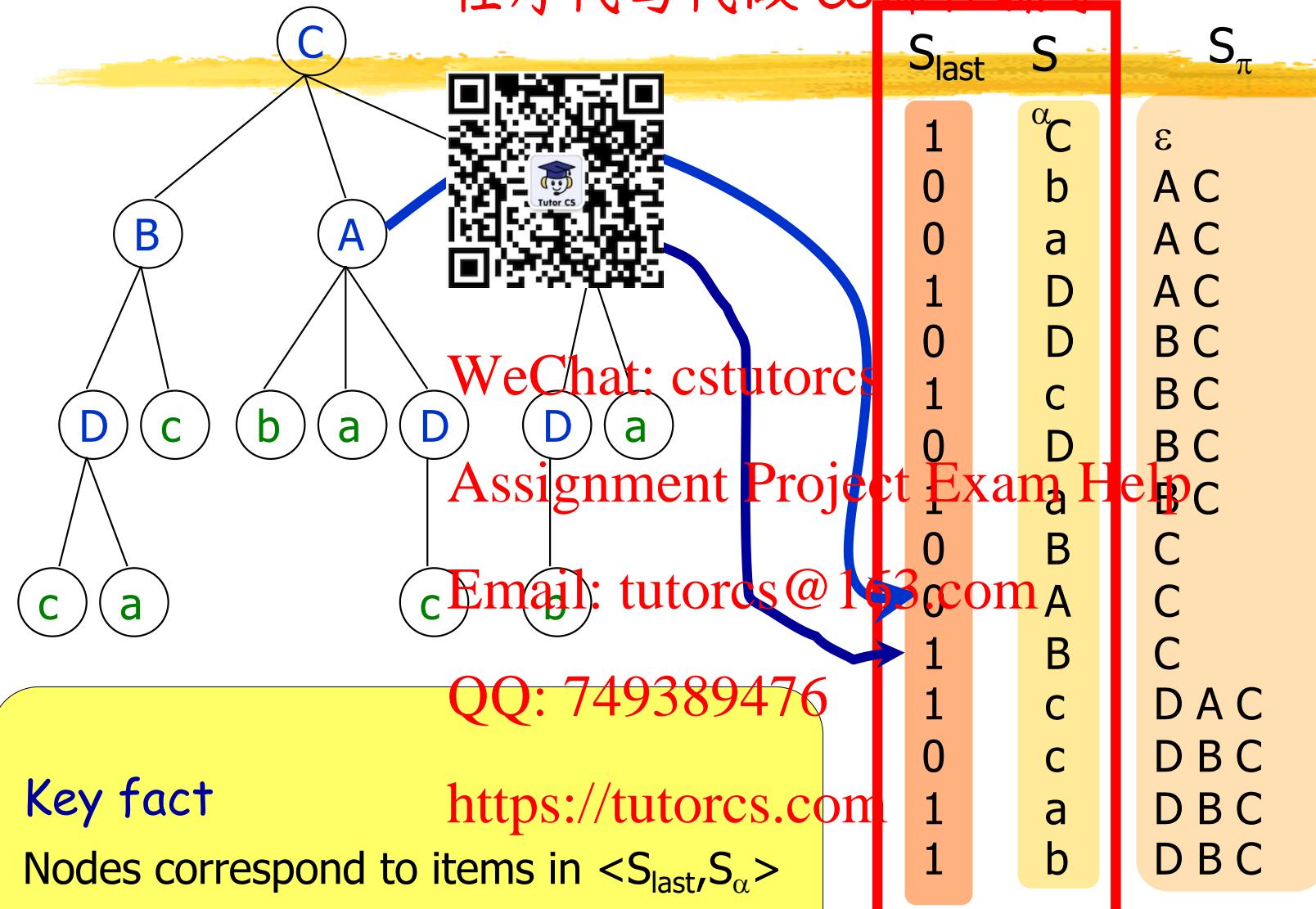
# The XBW-Transform

程序代写代做 CS 编程辅导



# The XBW-Transform

程序代写代做 CS 编程辅导



# XBzip – a simple XML compressor

程序代写代做 CS编程辅导



WeChat: cstutorcs

Rk	$S_{last}$	$S_\alpha$	$S_\pi$
1	1	<biblio	empty string
2	1	-	<author><book><biblio
3	1	-	<author><book><biblio
4	0	<book	<biblio
5	1	<book	<biblio
6	0	@id	<book><biblio
7	0	<author	<book><biblio
8	1	<title	<book><biblio
9	0	@id	<book><biblio
10	0	<author	<book><biblio
11	1	<title	<book><biblio
12	1	-	<title><book><biblio
13	1	-	<title><book><biblio
14	1	-	@id<book><biblio
15	1	-	@id<book><biblio
16	ØJ. Austin	-<author><book><biblio	
17	ØC. Bronte	-<author><book><biblio	
18	ØEmma	-<title><book><biblio	
19	ØJane Eyre	-<title><book><biblio	
20	Ø1	-@id<book><biblio	
21	Ø2	-@id<book><biblio	

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476 Tools: Attributes and symbol =

<https://tutorcs.com>

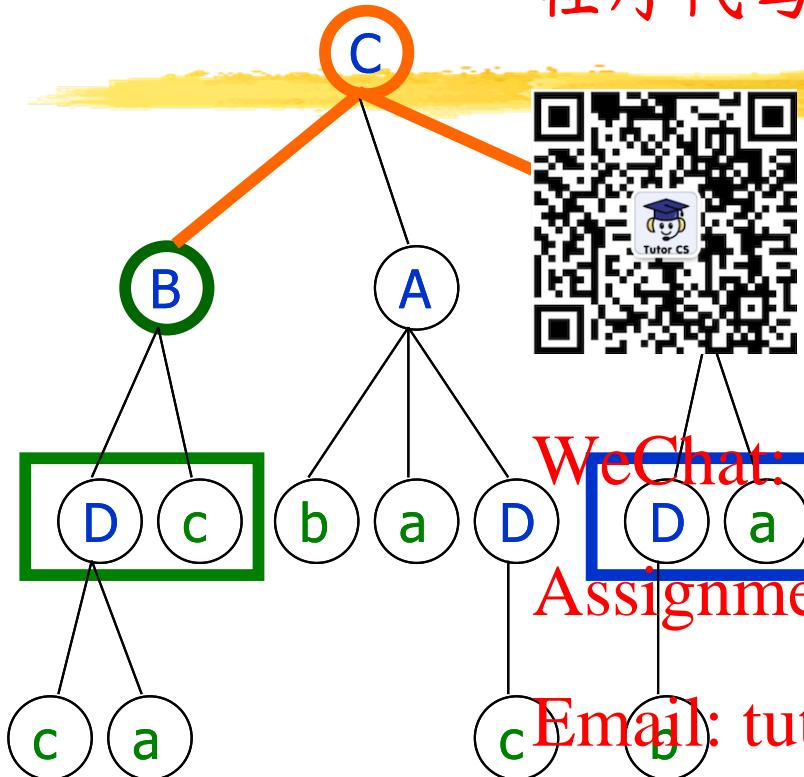
PCDATA

**XBW is compressible:**

- ①  $S_\alpha$  and  $S_{pcdata}$  are locally homogeneous
- ②  $S_{last}$  has some structure

# Some structural properties

程序代写代做 CS 编程辅导



WeChat: cstutorcs  
Assignment Project Exam Help

Email: tutorcs@163.com  
QQ: 749389476

Two useful properties:

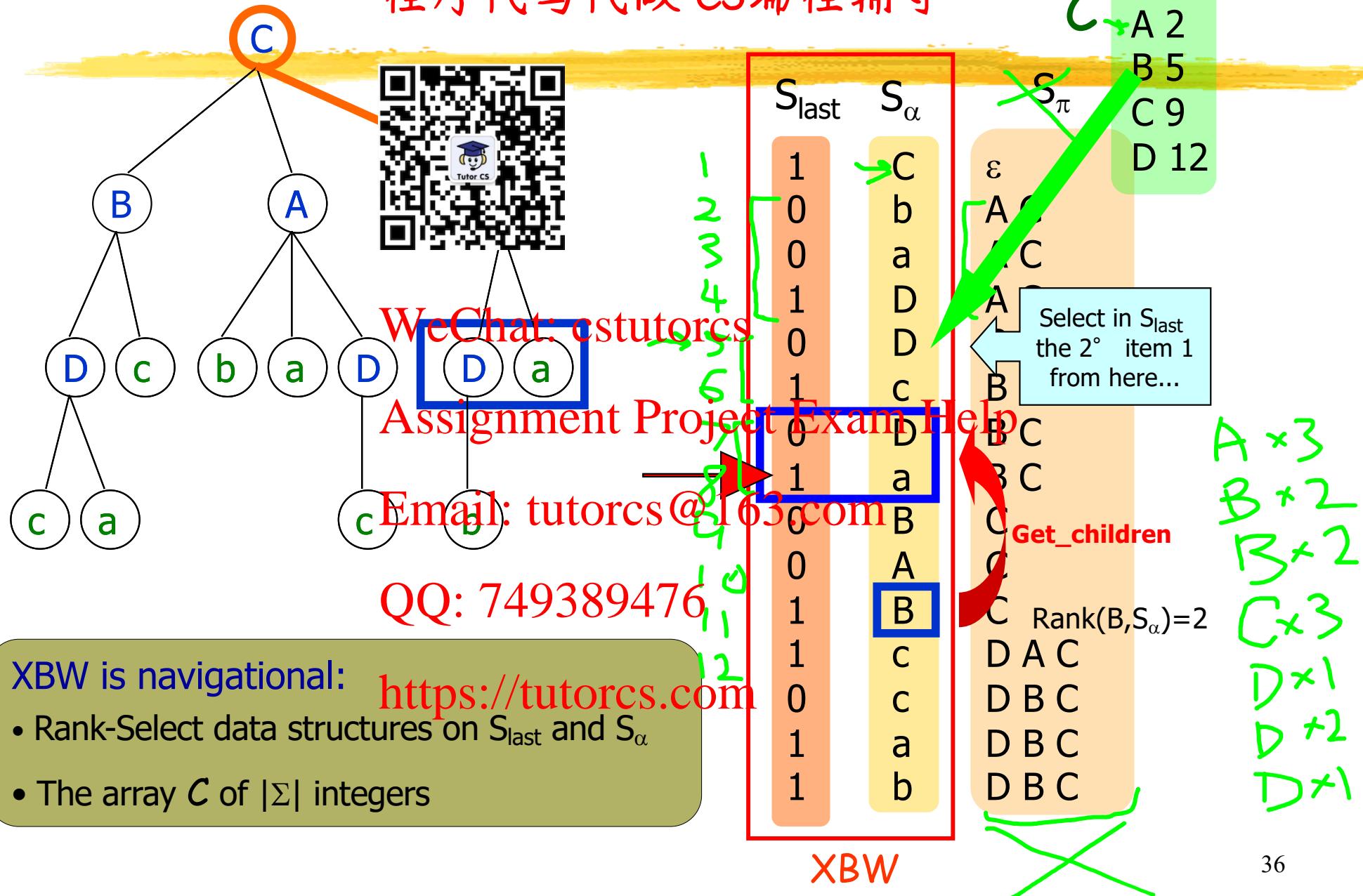
- Children are contiguous and delimited by 1s
- Children reflect the order of their parents

<https://tutorcs.com>

$S_{\text{last}}$	$S_{\alpha}$	$S_{\pi}$
1	C	A
0	b	C
0	a	A
1	D	C
0	D	B
1	c	C
0	D	B
1	a	C
0	B	B
0	A	C
1	B	C
1	A	A
0	B	B
0	C	C
1	C	C
1	D	D
0	A	A
1	B	B
1	C	D
0	C	B
1	A	C
1	B	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1	B	B
1	C	B
0	B	D
0	C	B
1	C	C
1	D	D
0	A	A
1		

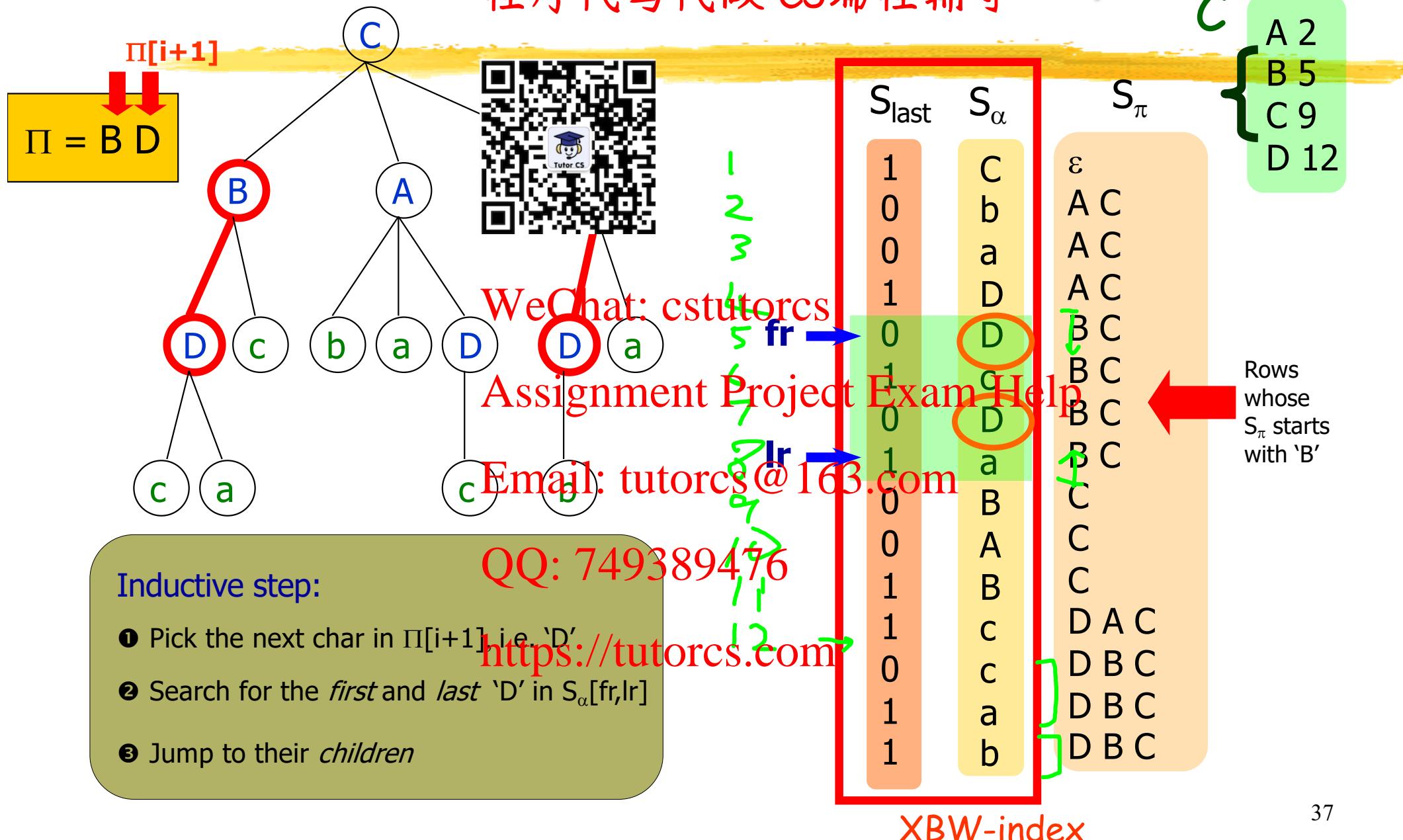
# XBW is navigational

程序代写代做 CS 编程辅导



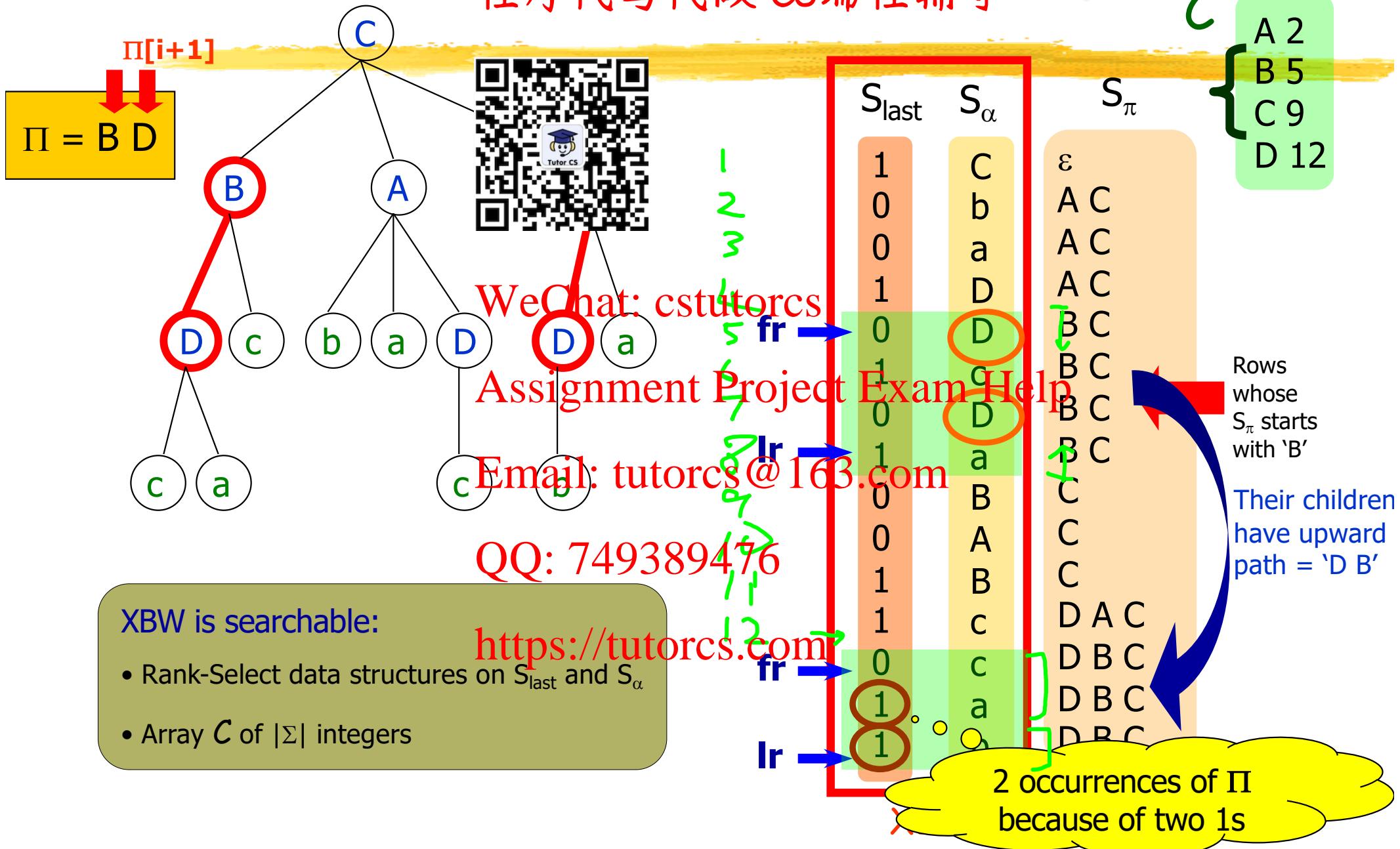
# XBW is searchable (count subpaths)

# 程序代写代做 CS端程辅导



# XBW is searchable (count subpaths)

程序代写代做CS编程辅导



# Graph compression



- Useful for many graph based applications such as:

WeChat: cstutorcs

- Web graph

Assignment Project Exam Help

- Social network

Email: tutors@163.com

QQ: 749389476

- Chemical & biological applications

<https://tutorcs.com>

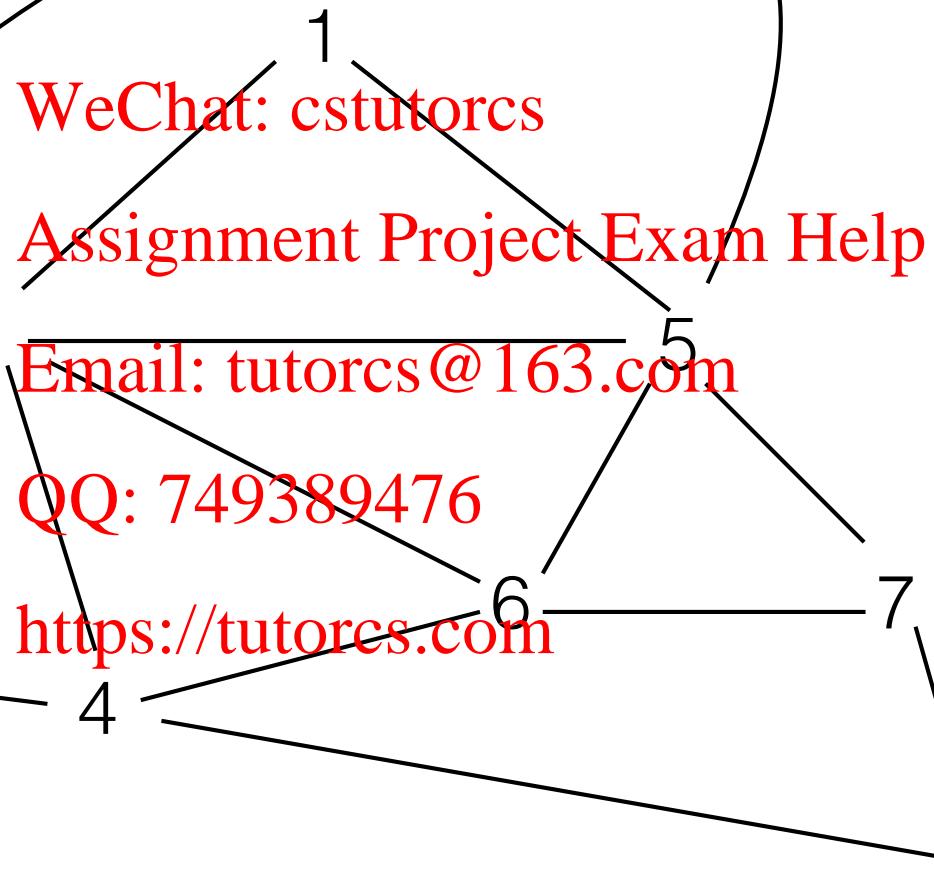
- Graph visualisation and analysis

# Graph compression

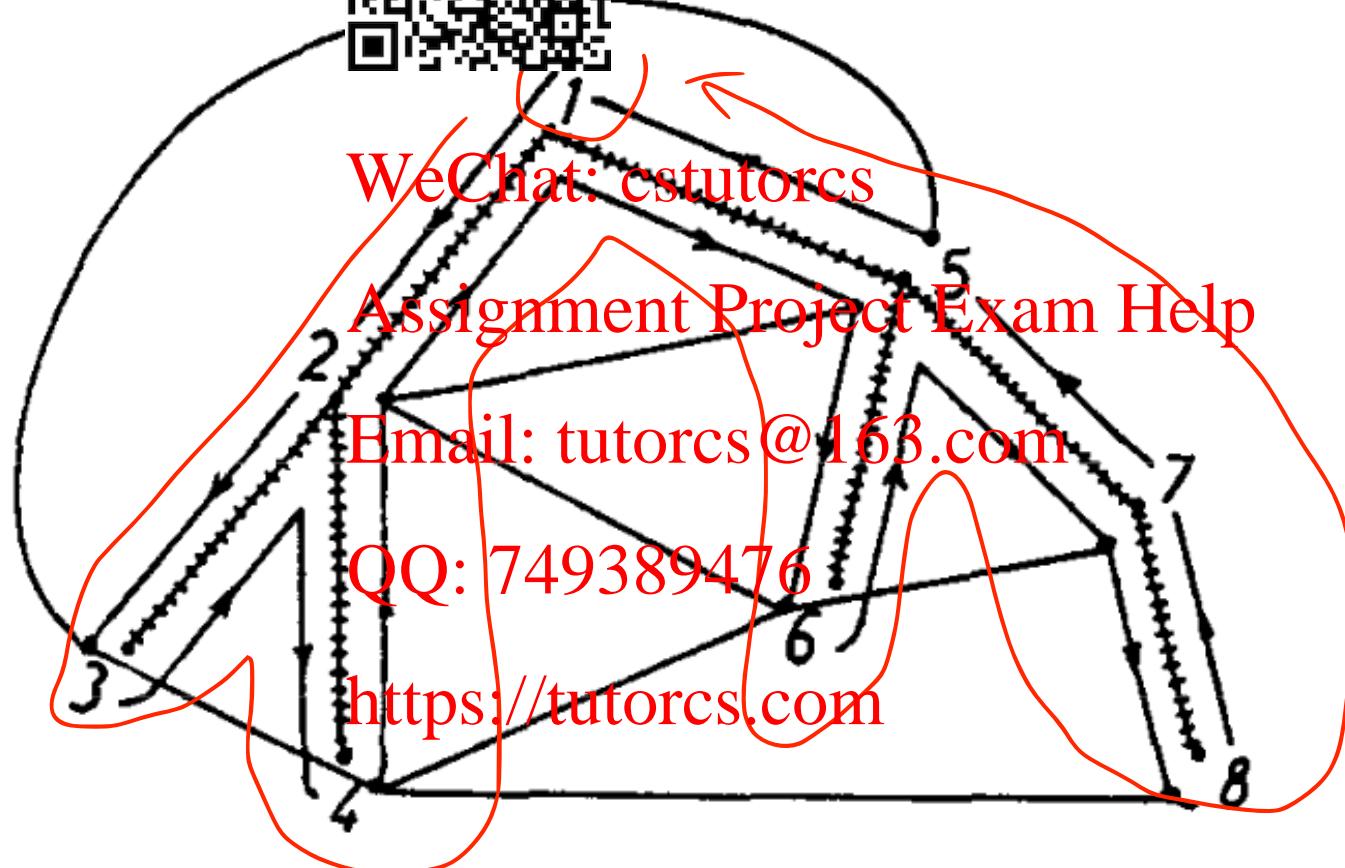


- Many techniques, e.g.,  
**WeChat: cstutorcs**
  - Succinct graph representation  
**Assignment Project Exam Help**
  - Adjacency matrix  
**Email: tutorcs@163.com**  
**QQ: 749389476**
  - Adjacency list  
**<https://tutorcs.com>**

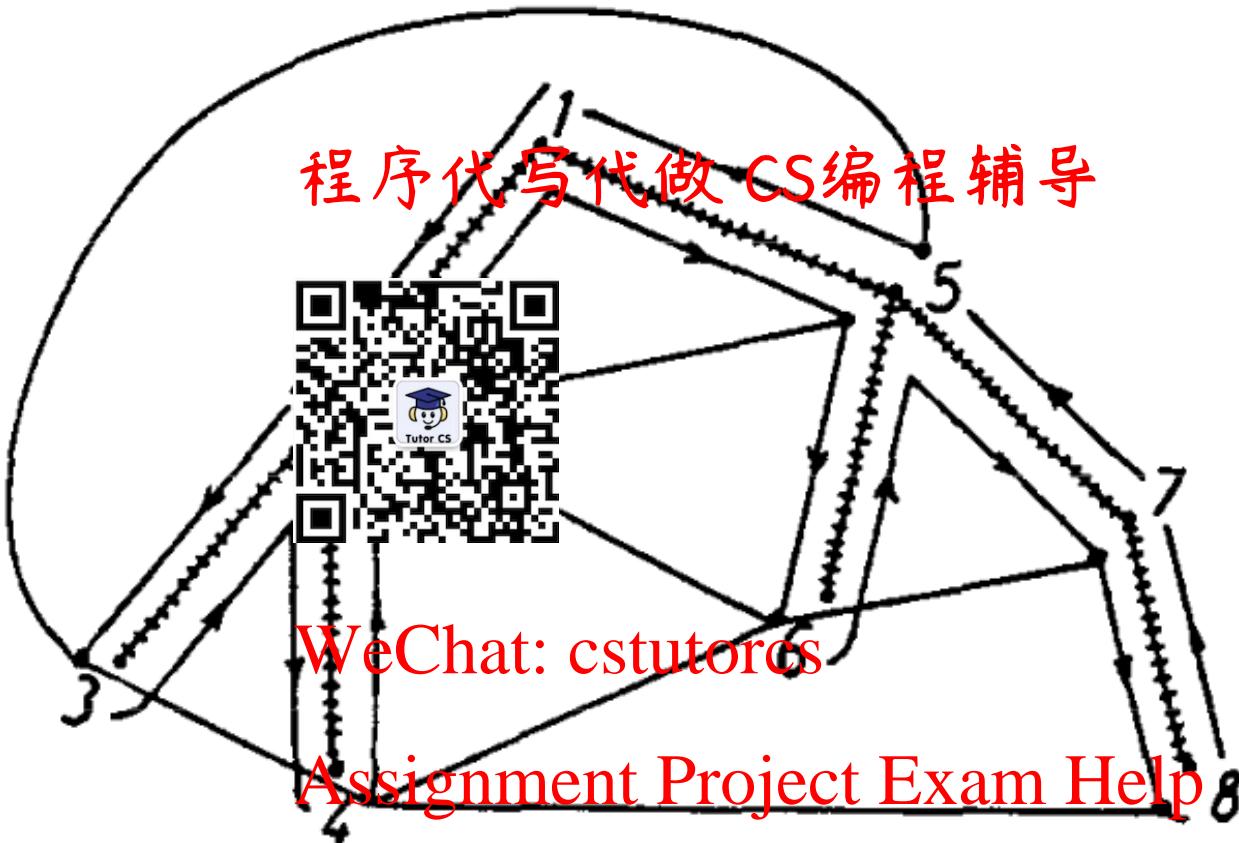
# A planar graph G



# A spanning tree of G



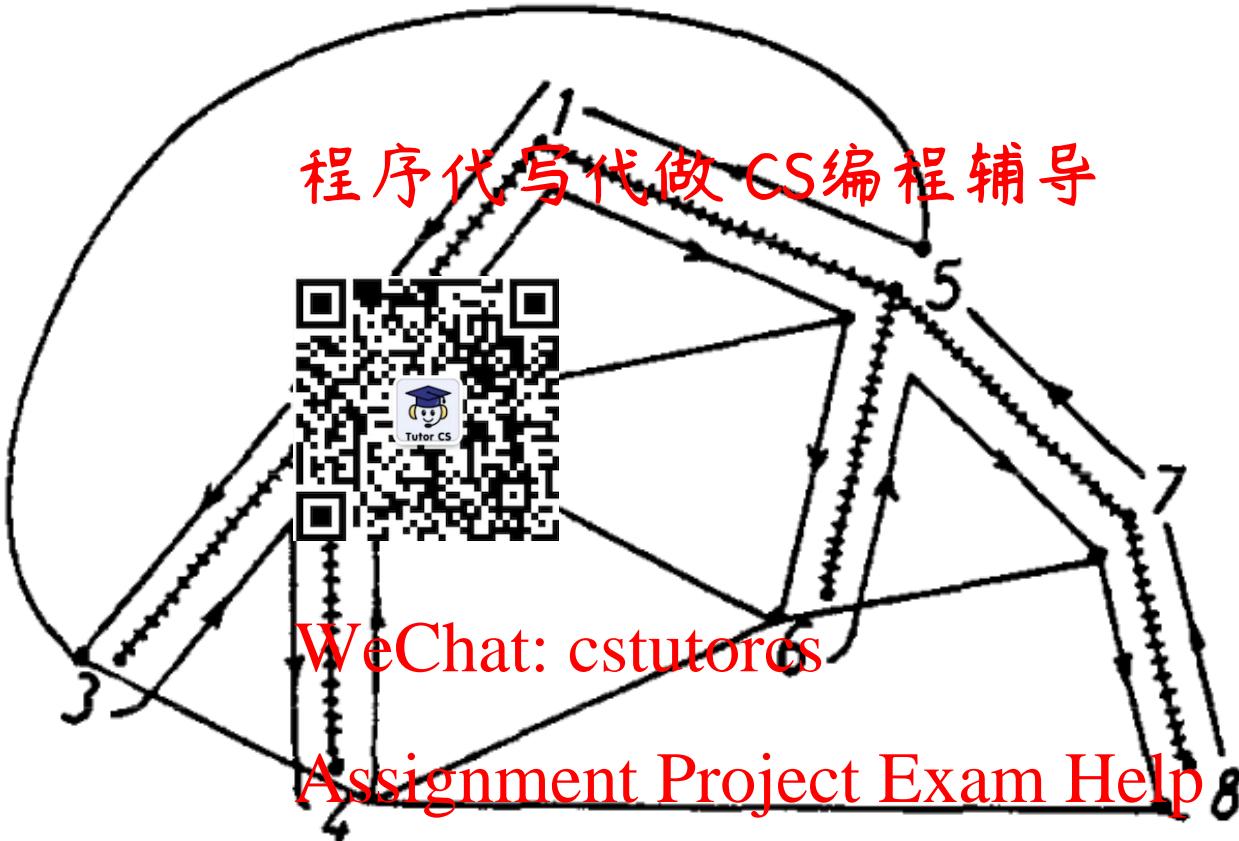
程序代写代做 CS编程辅导



Email: tutorcs@163.com

1 2 3 2 4 2 1 5 6 5 7 8 7 5 1

QQ: 749389476  
<https://tutorcs.com>



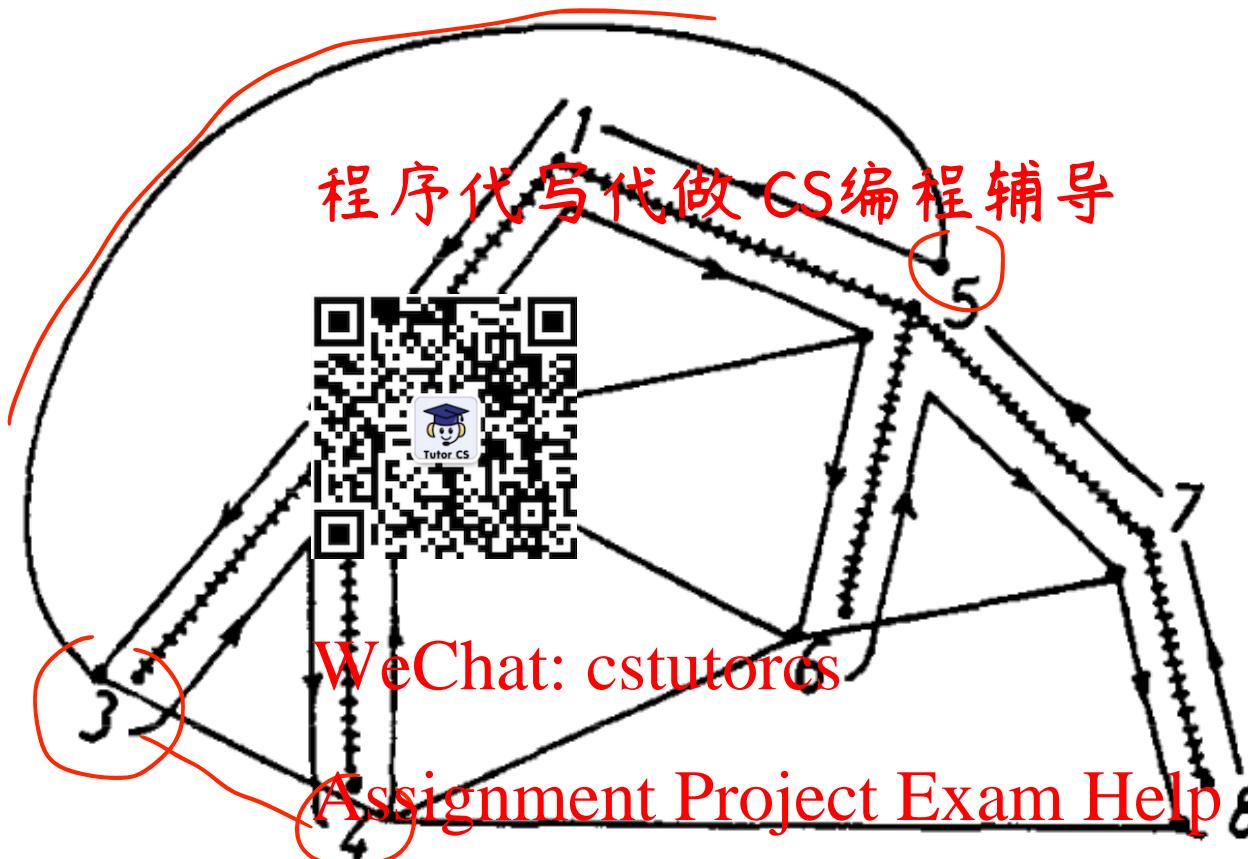
Email: tutorcs@163.com

QQ: 749389476

1 2 3 2 4 2 1 5 6 5 7 8 7 5 1

<https://tutorcs.com>

- - + - + + - - + - - + + +



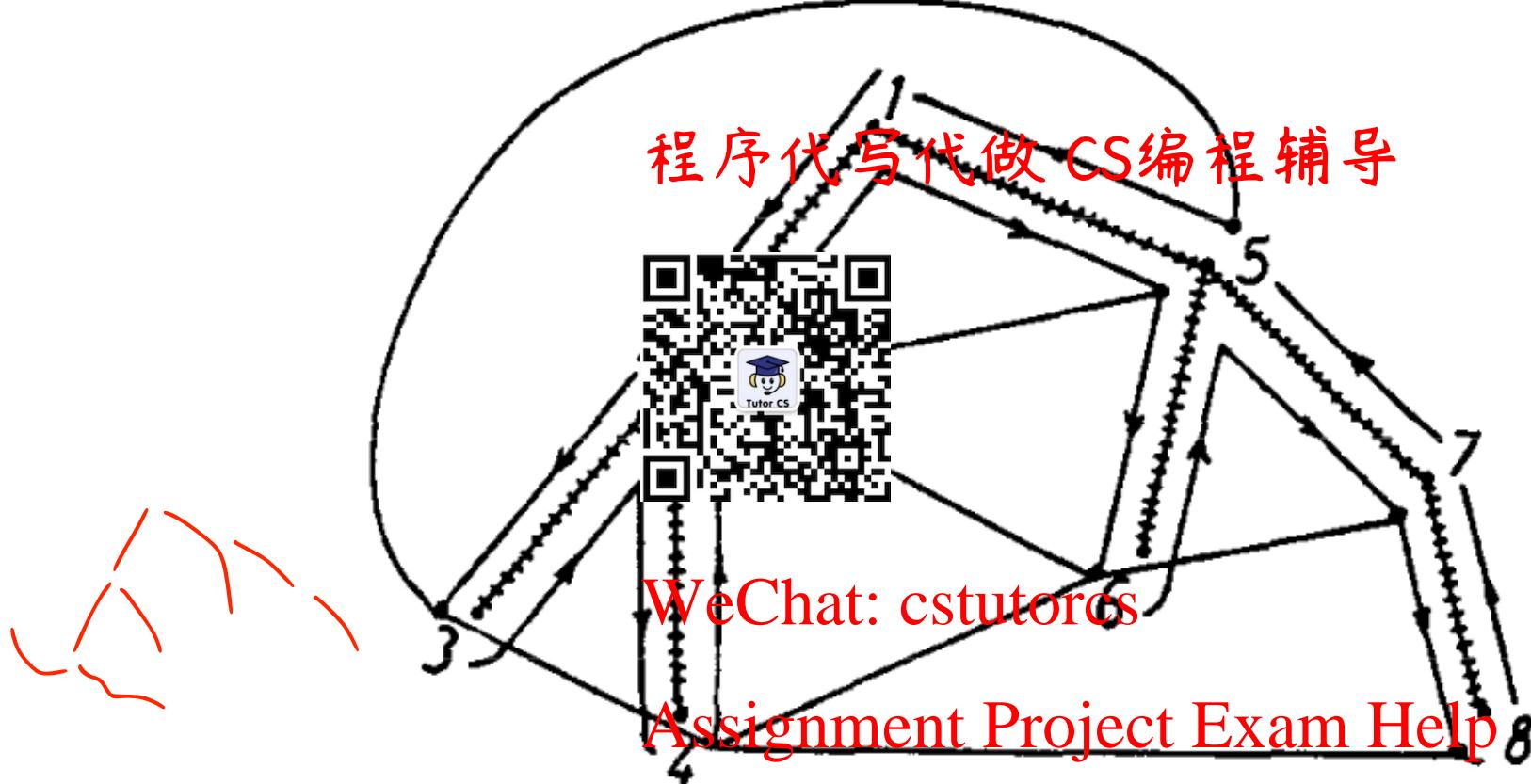
Email: tutorcs@163.com

1 2 3 2 4 2 1 5 6 5 7 8 7 5 1

QQ: 749389476  
<https://tutorcs.com>

- - + - + + - - + - - + + +

(( ))(( (( ) ) ))( ( ) ) { }



Email: tutorcs@163.com

- - + - + + - - + - - + + +  
(( ))(( (( ))))(( )) ) ) )

QQ: 749389476  
<https://tutorcs.com>

- - (( + - ))(( + (( + - ) - ))(( + - ) - ) + + ) +

程序代写代做 CS编程辅导



1

2

5

3

WeChat: cstutorcs

Assignment Project Exam Help

6

Email: tutorcs@163.com

QQ: 749389476

8

<https://tutorcs.com>

- - (( + - )(( + (( + - ) - ))( + - ) - ) + + ) +

00: -    01: +    10: (    11: )

# Succinct representation of G



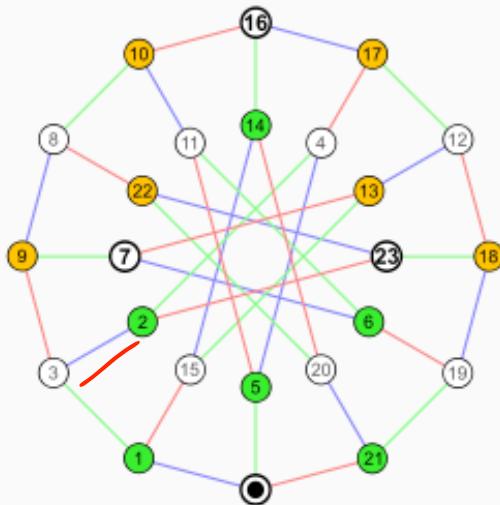
- - (( + - )(( + (( + - ) - ))( + - ) - ) + + ) +  
**WeChat: cstutorcs**

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

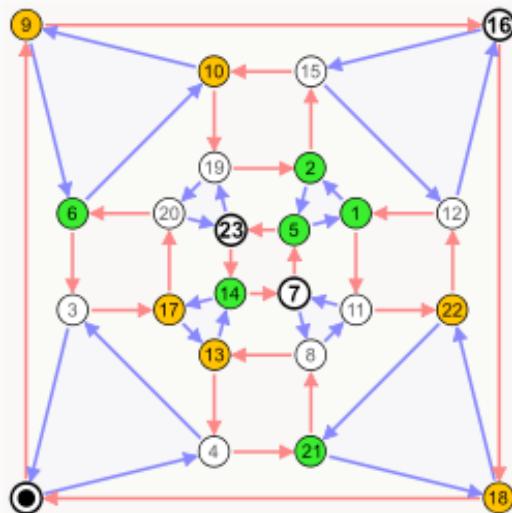
<https://tutorcs.com>



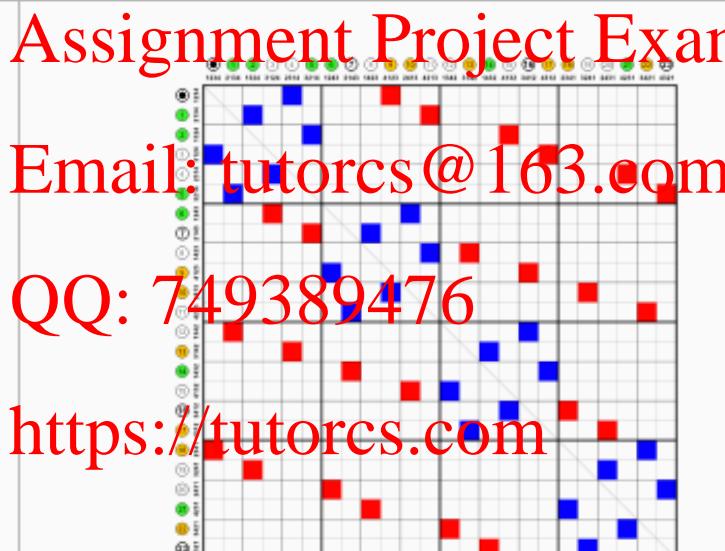
Nauru graph



程序代写代做 CS 编程辅导



Directed Cayley graph of  $S_4$



Coordinates are 0–23.

As the graph is directed, the matrix is not symmetric.

WeChat: estutores  
Coordinates are 0–23.  
White fields are zeros, colored fields are ones.

Assignment Project Exam Help  
Adjacency matrix

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

# Adjacency List



WeChat: cstutorcs

- Each vertex associated with an (sorted / unsorted) array of adjacent vertices

**Assignment Project Exam Help**  
Email: tutorcs@163.com

- More space efficient for sparse graph  
**QQ: 749389476**

<https://tutorcs.com>

# 程序代写代做 CS编程辅导



# Web Graphics representation and compression

# Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Slides modified from L.S. Buriol and D. Donato's

程序代写代做 CS编程辅导

# Internet/Web as Graphs

- Graph of the Internet: physical layer with routers, computers as nodes and physical connections as edges
  - It is limited
  - Does not capture the graphical connections associated with the information on the Internet
- Web Graph where nodes represent web pages and edges are associated with hyperlinks



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS编程辅导

# Web Graph



WeChat: cstutorcs  
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

© 2003 TouchGraph LLC

<http://www.touchgraph.com/TGoogleBrowser.html>

# Web Graph Considerations



- Graph is highly dynamic
  - Nodes and edges are added/deleted often
  - Content of existing nodes is also subject to change
  - Pages and hyperlinks created on the fly
- Apart from primary connected component there are also smaller disconnected components

<https://tutorcs.com>

程序代写代做 CS编程辅导

# Why the Web Graph?

- Example of large, complex, dynamic and distributed
- Possibly similar to other complex graphs in social, biological and other systems
- Reflects how humans organize information (relevance, ranking) and their societies
- Efficient navigation algorithms
- Study behavior of users as they traverse the web graph (e-commerce)



WeChat: cstutorcs

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS编程辅导

## Statistics of Interest

- Size and connectivity of the graph
- Number of connected components
- Distribution of pages per site  
WeChat: cstutorcs
- Distribution of incoming and outgoing connections per site  
Assignment Project Exam Help
- Average and maximal length of the shortest path between any two vertices (diameter)  
Email: tutorcs@163.com  
QQ: 749389476  
<https://tutorcs.com>

# 程序代写代做CS编程辅导 Web Graph



A web graph  $G = (V, E)$  to a set of URLs is a directed graph having those URLs as the set of nodes. An arc  $u \rightarrow v$  is identified for each hyperlink from a URL  $u$  towards a URL  $v$ .

WeChat: cstutorcs  
Assignment Project Exam Help

Email: tutorcs@163.com

URLs that do not appear either as sources or in more than  $T$  (4) pages are ignored;

QQ: 749389476

The URLs are normalized by converting hostnames to lower case, canonicalizes port number, re-introducing them where they need, and adding a trailing slash to all URLs that do not have it.

# Main features of Web Graphs



**Locality:** usually most of the hyperlinks are local, i.e, they point to other URLs on the same host. WeChat: cstutorcs  
The literature reports that on average 80% of the hyperlinks are Assignment Project Exam Help

**Consecutivity:** links within same page are likely to be consecutive respecting to the lexicographic order.

QQ: 749389476

<https://tutorcs.com>

# Main features of WebGraphs

**Similarity:** Pages on the same host tend to have many links pointing to the same pages.



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

# 程序代写代做CS编程辅导



Connectivity Center (1998) – *Digital Systems Research Center* – Stanford University – K. Bharat, A. Broder, M. Henzinger, P. Kumar, S. Venkatasubramanian;  
WeChat: csutorcs

Assignment Project Exam Help  
Link Database (2001) - Compaq Systems Research Center – K. Randall, R. Stata, R. Wickremesinghe, J. Wiener;  
Email: tutorcs@163.com  
QQ: 749389476

WebGraph Framework (2002) – Universita degli Studi di Milano – P. Boldi, S. Vigna.

# 程序代写 游戏代做 CS编程辅导

# Connectivity Server

- Tool for graphs visualisation, analysis (connectivity, ranking pages) and URLs compression
- Used by Alta Vista;
- Links represented by an outgoing and an incoming adjacency lists; Assignment Project Exam Help

Email: tutorcs@163.com

- Composed of:

QQ: 749389476

URL Database: URL, fingerprint, URL-id;

Host Database: <https://tutorcs.com> group of URLs based on the hostname portion;

Link Database: URL, outlinks, inlinks.

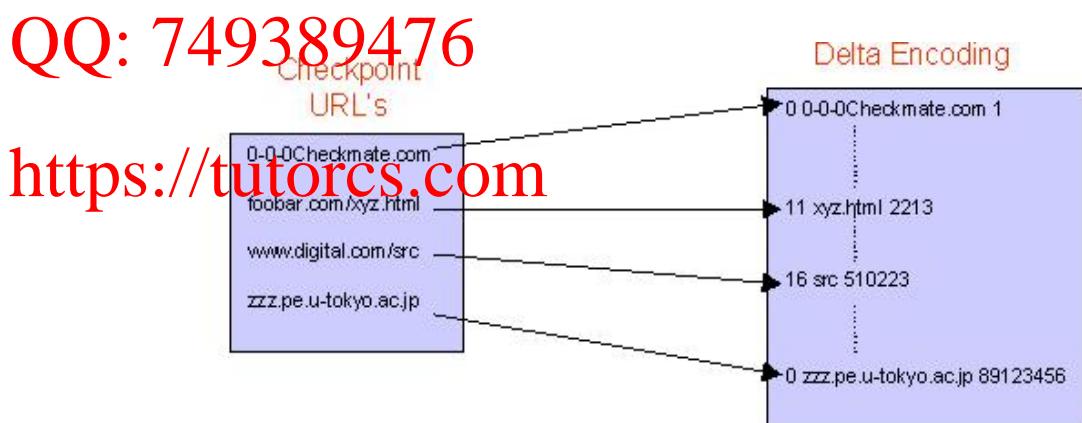
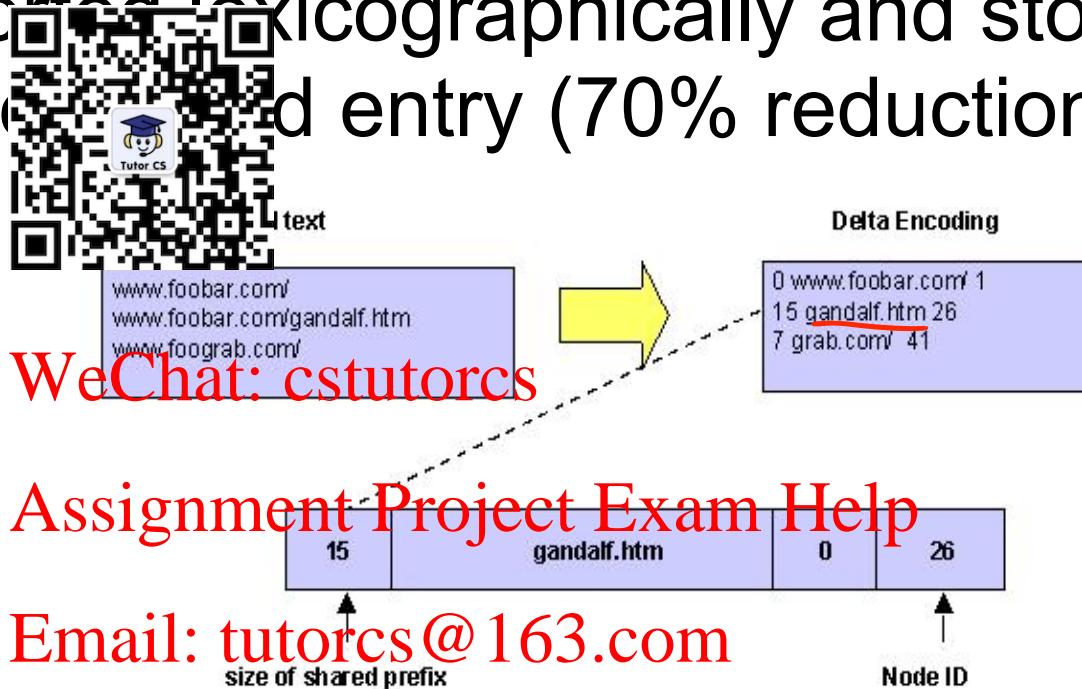


# Connectivity Server 做 URL compression

URLs are sorted lexicographically and stored as a delta encoded entry (70% reduction).

URLs delta encoding

Indexing the delta encoding



# Link1: first version of Link Database



No compression. simple representation of outgoing and incoming adjacency lists of links.  
Assignment Project Exam Help

Email: tutorcs@163.com

Avg. inlink size: 34 bits

QQ: 749389476

Avg. outlink size: 24 bits

<https://tutorcs.com>

# Link2: second version of Link

程序代写 代码代做 CS 编程辅导

## Database



Single list compression and starts  
WeChat: cstutorcs  
compression

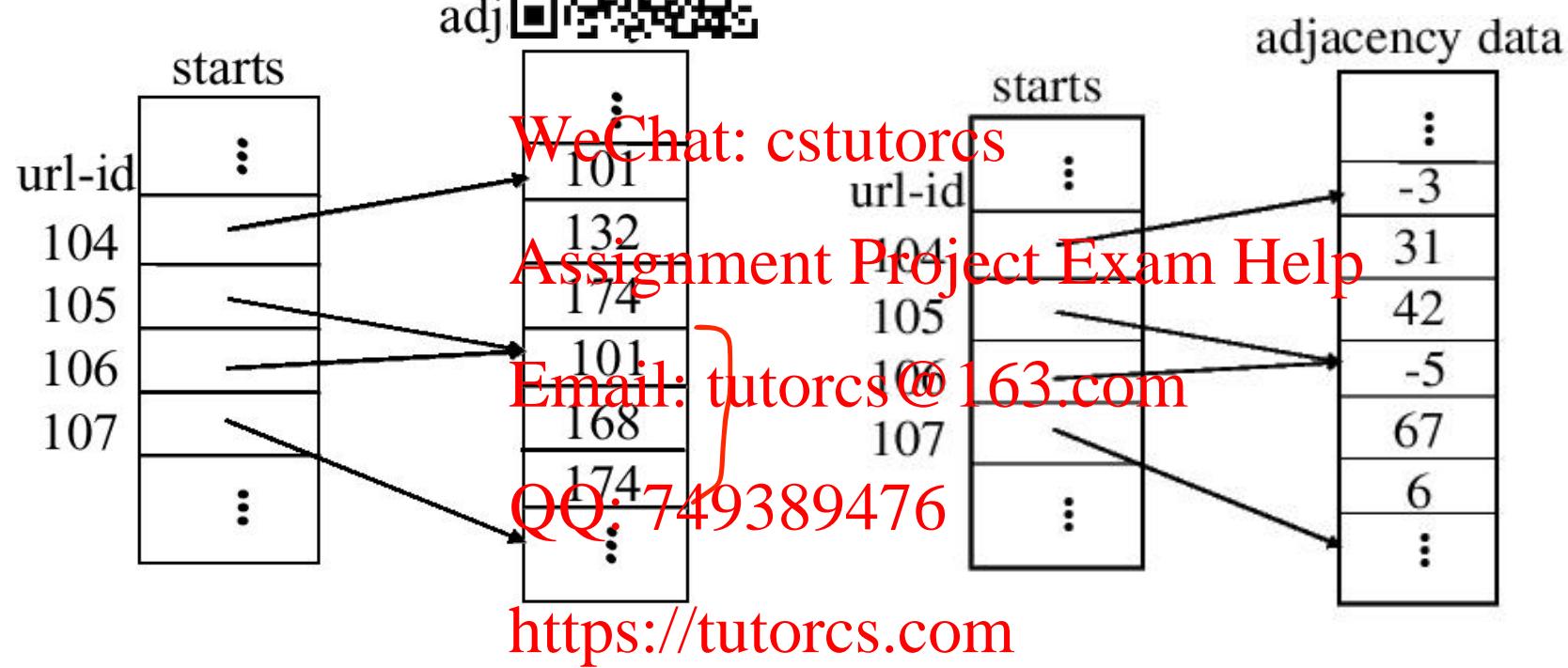
Assignment Project Exam Help

Avg. inlink size: 8.9 bits

Avg. outlink size: 11.03 bits

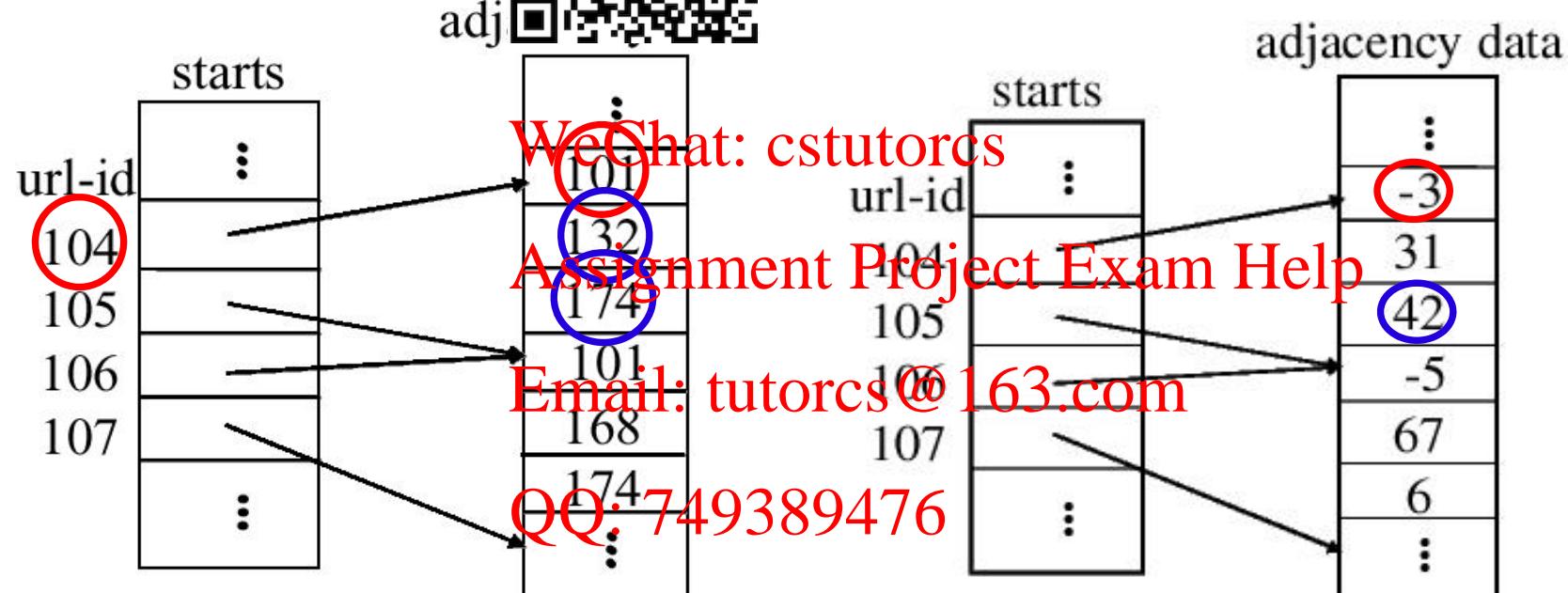
<https://tutorcs.com>

# Delta Encoding of the Adjacency Lists



Each array element is 32 bits long.

# Delta Encoding of the Adjacency Lists



$-3 = 101 - 104$  (first item)

$42 = 174 - 132$  (other items)

# Starts 程序代写 代码代做 CS 编程辅导 array compression



- The URLs are divided into three partitions based on their degree:

Assignment Project Exam Help

Email: tutorcs@163.com

- The literature reports that 74% of the entries are in the low-degree partition.

~~QQ: 749389476~~  
<https://tutorcs.com>

# Link3: third version of Link Database



Interlist compression with representative  
WeChat: cstutorcs  
list

Assignment Project Exam Help

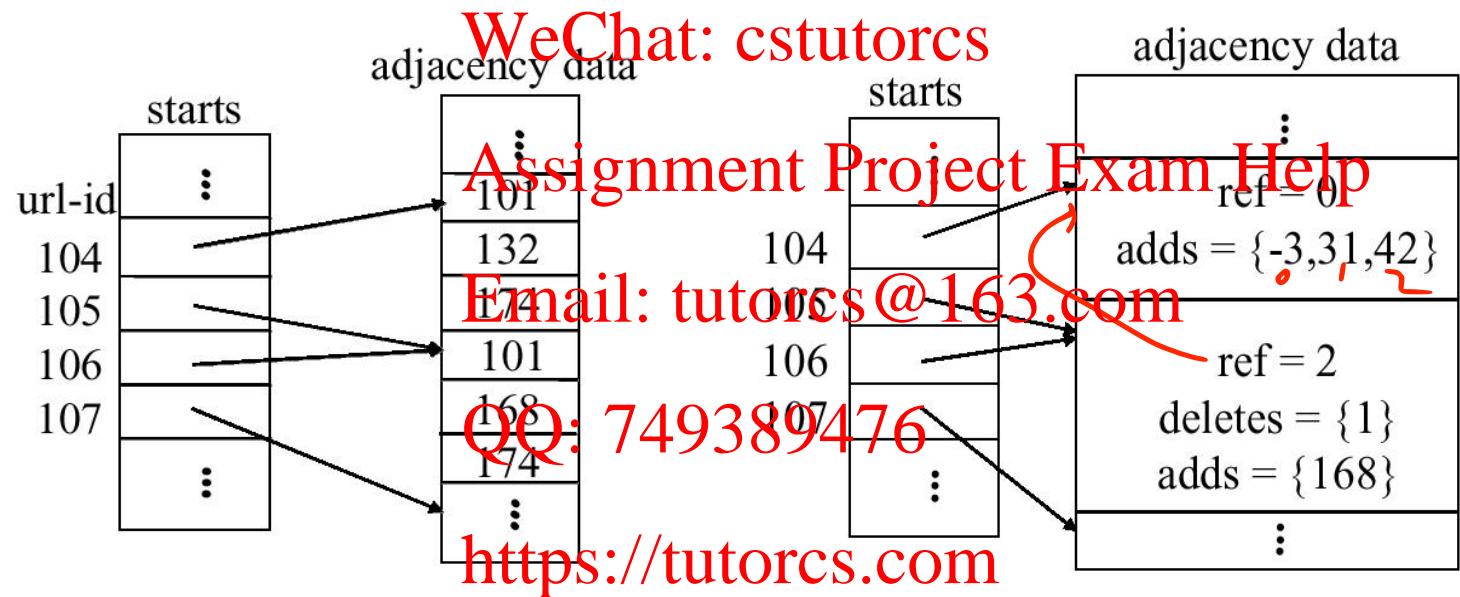
Avg. inlink size: 5.60 bits

Avg.outlink size: 5.61 bits

<https://tutorcs.com>

# Intelligent Compression

*ref* : relative index of the URL in the representative adjacency list;  
*deletes*: set of URL-ids to be removed from the representative list;  
*adds*: set of URL-ids to be added to the representative list.



**LimitSelect-K-L**: chooses the best representative adjacency list from among the previous  $K$  (8) URL-ids' adjacency lists and only allows chains of fewer than  $L$  (4) hops.

# $\zeta$ -codes (WebGraph 程序代写代做 CS 编程辅导 Framework)



Interlist compression with representative  
WeChat: cstutorcs  
list

Assignment Project Exam Help

Avg. inlink size: 3.08 bits

Avg.outlink size: 2.89 bits

<https://tutorcs.com>

# 程序代写代做CS编程辅导 Using copy lists

Uncompressed  
adjacency list



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

WeChat: cstutorcs

Adjacency list with  
copy lists.

| 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                             |
| 18   | 5     | 3    | 111100000000 | 50   |
|      |       |      |              | ...  |

QQ: 749389476

Each bit on the copy list informs whether the corresponding successor of  $y$  is also a successor of  $x$ ;

The reference list index  $ref.$  is chosen as the value between 0 and  $W$  (window size) that gives the best compression.

<https://tutorcs.com>

# 程序代寫代做 CS 編程輔導 Using copy blocks

Adjacency list with copy lists.

| Ref. | Copy list    | Extra nodes                                    |
|------|--------------|--|
| 17   | ...          | 13, 15, 16, 17, 18, 19, 23, 24, 203, 315, 1034 |
| 18   | 011100011010 | 22, 316, 317, 3041                             |
| 5    | 111100000000 | 50   |
| 3    | ...          | ...  |

WeChat: cstutorcs

Adjacency list with copy blocks.

| Node | Outd. | Ref. | # blocks | Copy blocks         | Extra nodes                                    |
|------|-------|------|----------|---------------------|--|
| 15   | 11    | 0    | 1        | 0                   | 13, 15, 16, 17, 18, 19, 23, 24, 203, 315, 1034 |
| 16   | 10    | 1    | 7        | 0, 0, 2, 1, 1, 0, 0 | 22, 316, 317, 3041                             |
| 17   | 0     | 0    | 1        | 0                   | ...  |
| 18   | 5     | 3    | 1        | 4                   | 50   |
| ...  | ...   | ...  | ...      | ...                 | ...  |

QQ: 749389476

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

The last block is omitted;

The first copy block is 0 if the copy list starts with 0;

The length is decremented by one for all blocks except the first one.

# 程序代写代做CS编程辅导

# Conclusions

The compression techniques are specialized for Web Graphs.



The average link size decreases with the increase of the graph.

WeChat: cstutorcs  
Assignment Project Exam Help

The average link access time increases with the increase of the graph.

Email: tutorcs@163.com  
QQ: 749389476

The  $\zeta$ -codes seems to have the best trade-off between avg. bit size and access time.

<https://tutorcs.com>