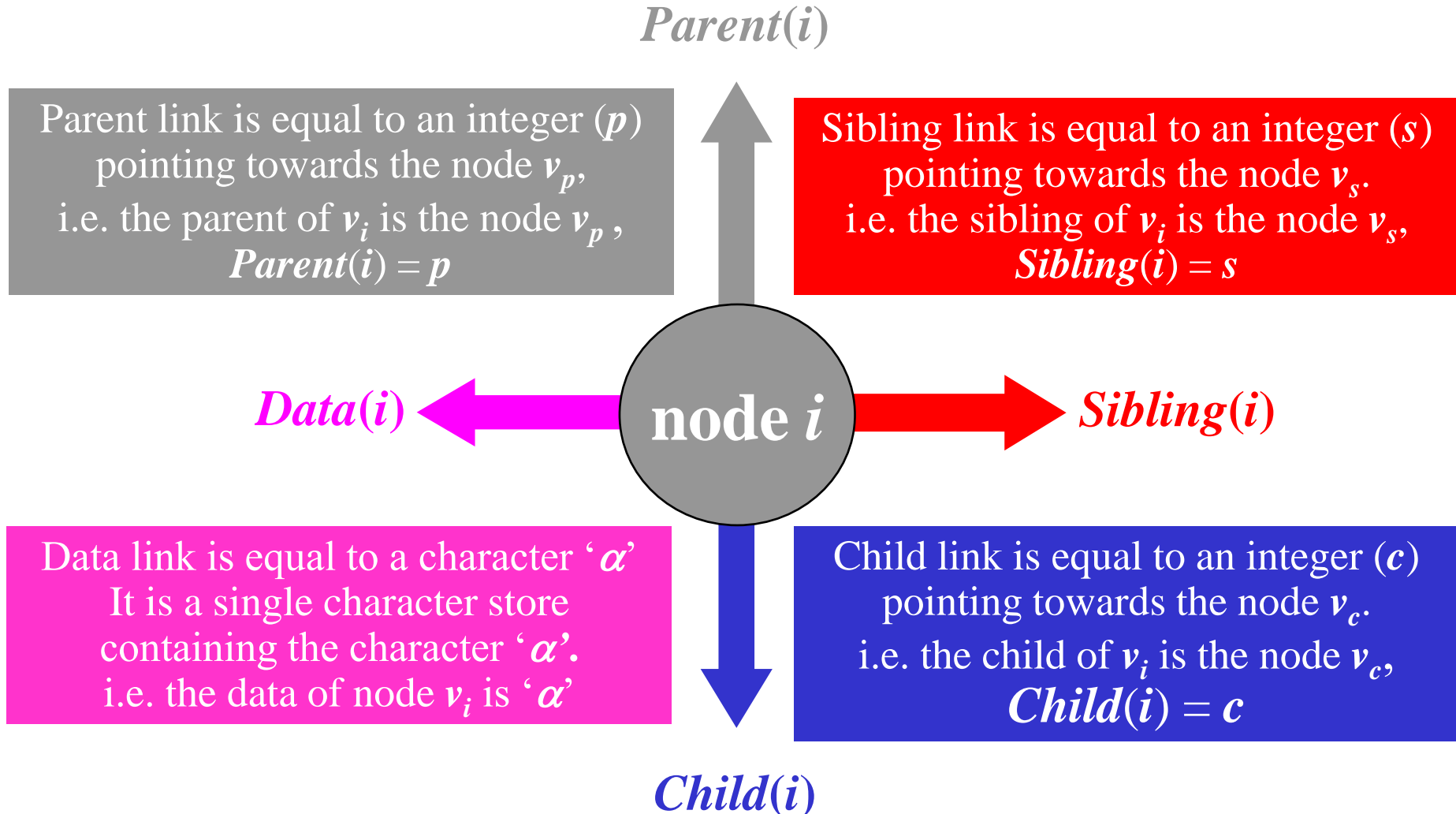*S&M*

# Split and Merge Compression Algorithm

# By

# Abdullah Hashim

## *Compression: The Tree Structure*

# *The Tree Structure*

## **Data Tree (DT) links**:
If a link has a value equal to -1, the link is a null link, leading to no node; the node does not exist.

*Parent(i)*

Parent link is equal to an integer $(p)$ pointing towards the node $v_p$, i.e. the parent of $v_i$ is the node $v_p$, $Parent(i) = p$

Sibling link is equal to an integer $(s)$ pointing towards the node $v_s$. i.e. the sibling of $v_i$ is the node $v_s$, $Sibling(i) = s$

*Data(i)*

**node i**

*Sibling(i)*

Data link is equal to a character '$\alpha$' It is a single character store containing the character '$\alpha$'. i.e. the data of node $v_i$ is '$\alpha$'

Child link is equal to an integer $(c)$ pointing towards the node $v_c$. i.e. the child of $v_i$ is the node $v_c$, $Child(i) = c$
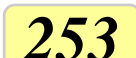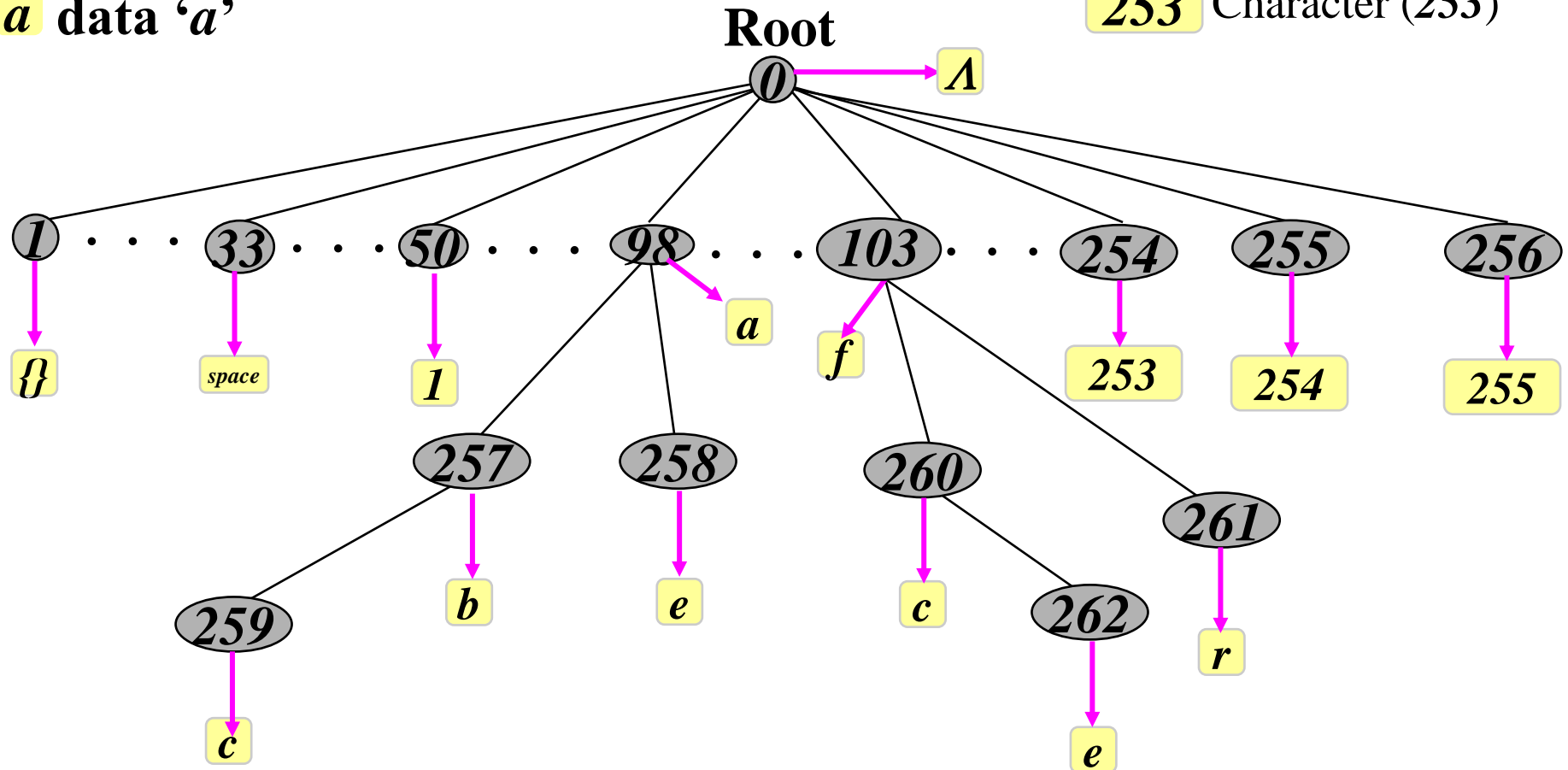
*Child(i)*

# The Tree Structure

Let $D = \{0, 1, 2, \ldots, 253, 254, 255, ab, ae, abc, fc, fr, fce\}$

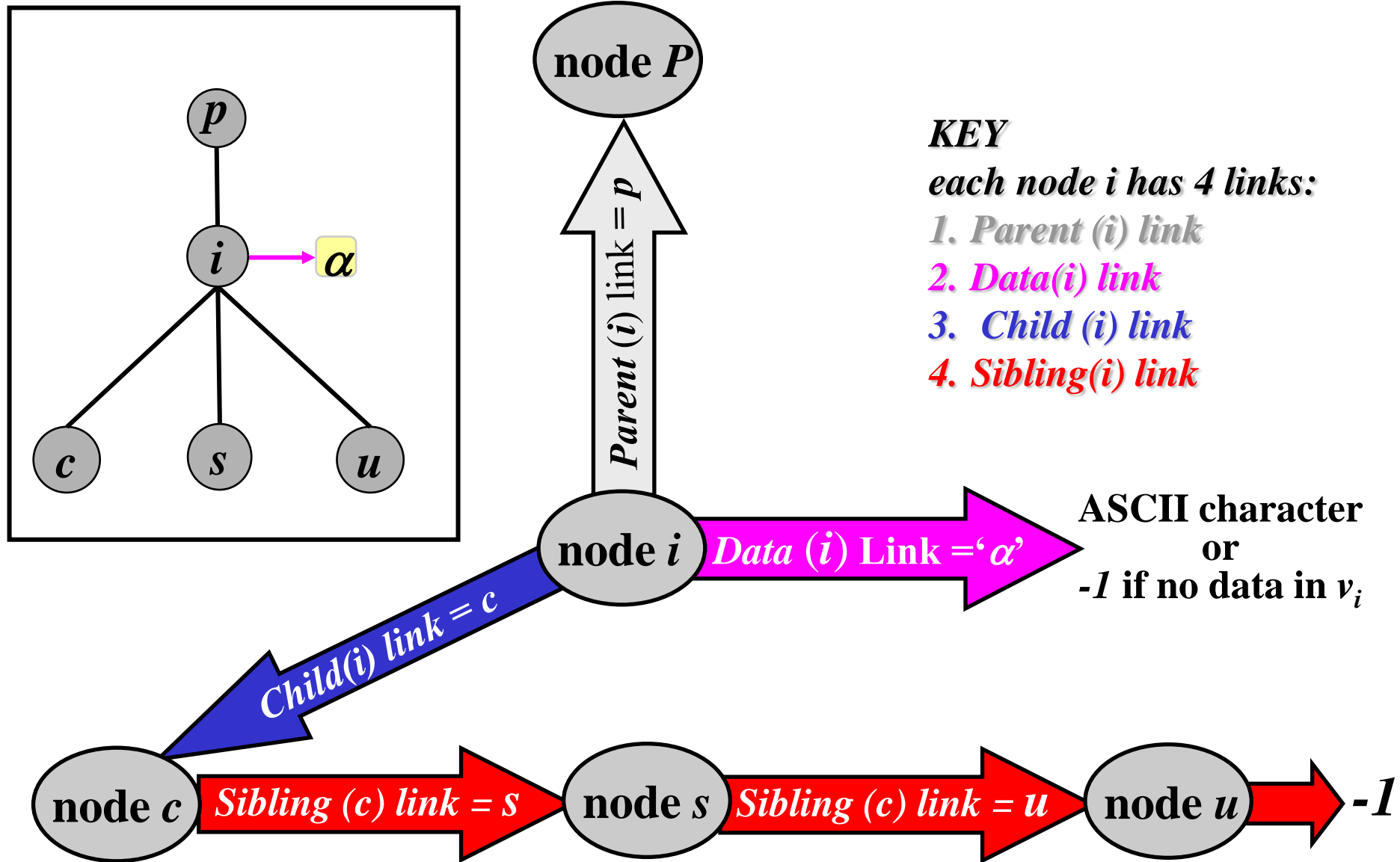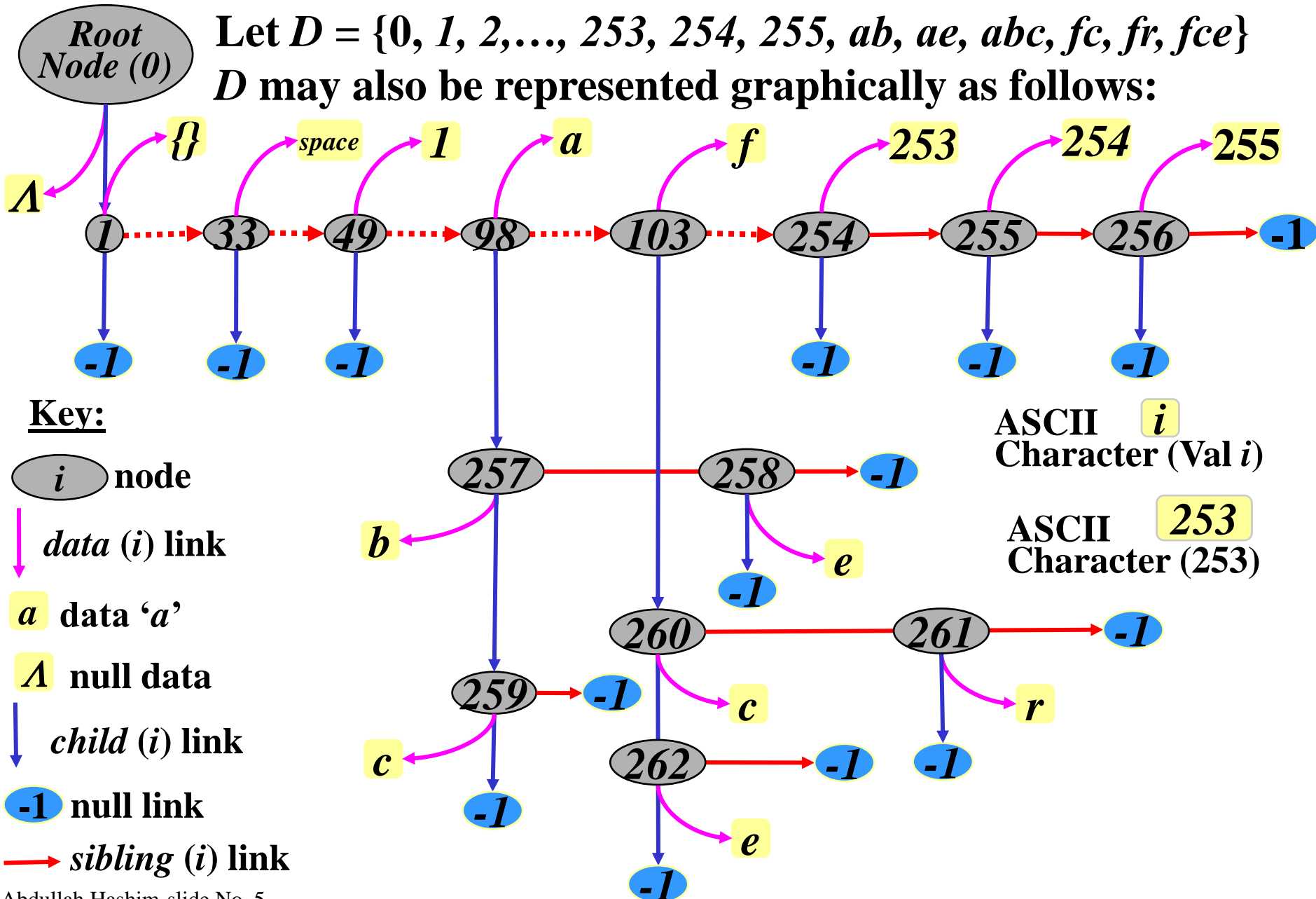The dictionary $D$ may be represented graphically as follows:

# *The Tree Structure*



**KEY**
*each node i has 4 links:*
1. *Parent (i) link*
2. *Data(i) link*
3. *Child (i) link*
4. *Sibling(i) link*

node $P$

$Parent$ $(i)$ link = $p$

node $i$   *Data* $(i)$ **Link =** '$\alpha$'

**ASCII character**
**or**
*-1* **if no data in** $v_i$

*Child(i) link = c*

node $c$   *Sibling (c) link = S*   node $s$   *Sibling (c) link = u*   node $u$   *-1*

# *Tree Structure*

**Let *D* = {0, 1, 2,…, 253, 254, 255, ab, ae, abc, fc, fr, fce}**
***D* may also be represented graphically as follows:**



**Key:**

- *i* node
- *data* (*i*) link
- *a* data '*a*'
- *Λ* null data
- *child* (*i*) link
- -1 null link
- *sibling* (*i*) link

ASCII *i* Character (Val *i*)

ASCII 253 Character (253)

# *The Tree Structure*

**ED**: *Extended Dictionary*:  **ED** is a dictionary containing:

1. All single ASCII characters, $A = \{\alpha_1, \alpha_2, \ldots, \alpha_N\}$.

2. All multiple character words, $W = \{w_1, w_2, \ldots w_M\}$.

3. All command "control" words, $C = \{\delta_1, \delta_2, \ldots, \delta_C\}$, usually $|C|$ is in the range of two to three characters.

4. All the null words, $\{\} = \{\Lambda_1, \Lambda_2, \ldots, \Lambda_i, \ldots, \Lambda_B\}$.

Therefore: total words in **ED** is $|ED|_{max} = N + M + C + B$

# *The Tree Structure*

**For example let the extended dictionary be given by the set *ED*:**

$\{0, 1, 2, \ldots, 253, 254, 255, \delta_1, \delta_2, ab, ae, abc, fc, fr, fce, \Lambda_1, \Lambda_2\}$
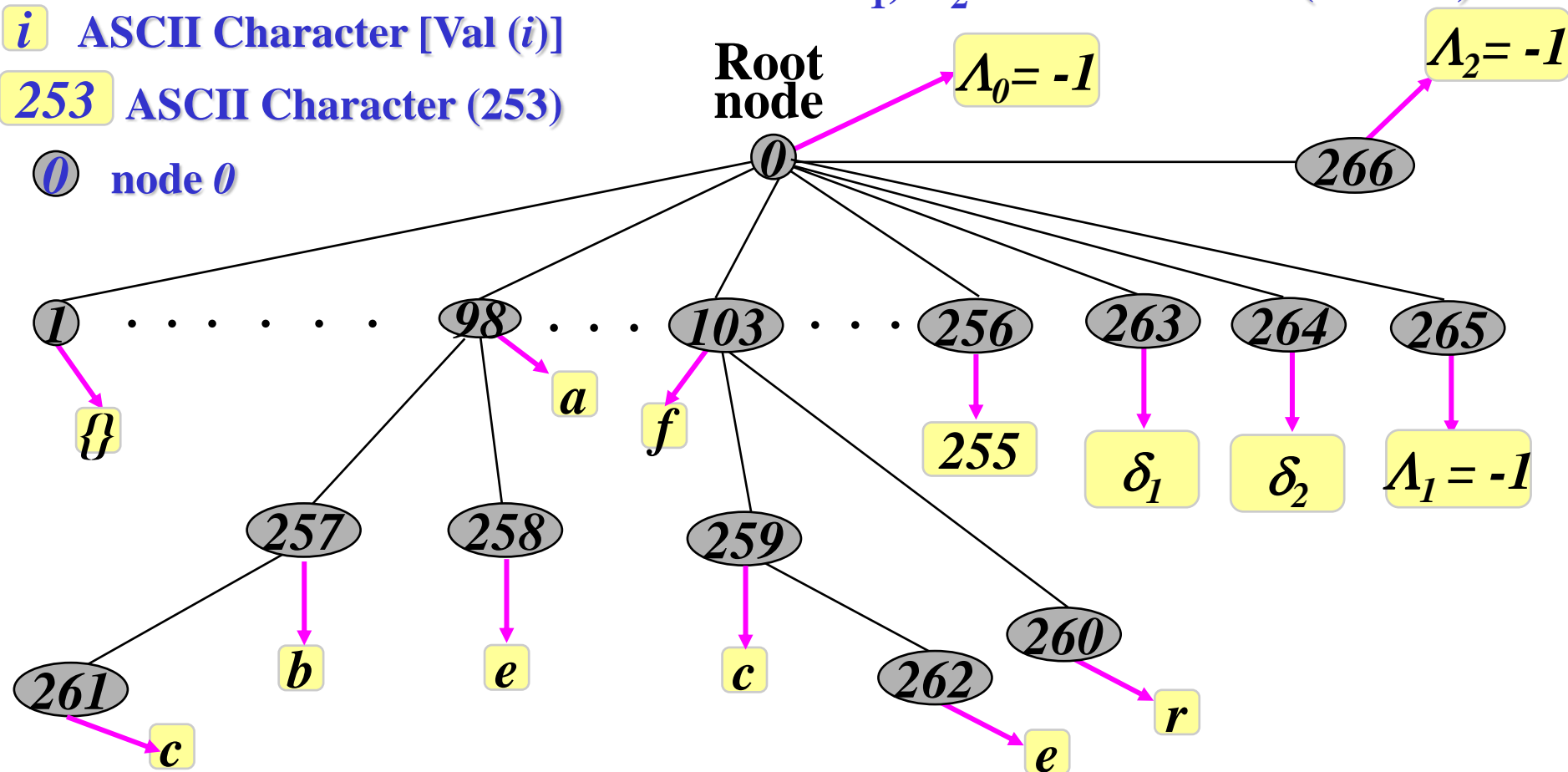
$\delta_1, \delta_2$ : control characters.

$\Lambda_1, \Lambda_2$ : null words = -1 (no data).

*Data*(*i*) **link**

| *i* | **ASCII Character [Val (*i*)]** |

| *253* | **ASCII Character (253)** |

| *0* | **node *0*** |



**Root node**

$\Lambda_0 = -1$

$\Lambda_2 = -1$

*0* ——— *266*

*1* · · · · · · *98* · · · *103* · · · *256*     *263*     *264*     *265*

*{}*

*a*

*f*

*255*

$\delta_1$

$\delta_2$

$\Lambda_1 = -1$

*257*     *258*     *259*

*260*

*261*

*b*     *e*     *c*     *262*

*r*

*c*     *e*

# *The Tree Structure*

**OD**     *Ordered Dictionary*. The ordered dictionary is a list of words. The words are ordered according to some parameter or function; the leftmost word ($w_1$) has the highest rank while the rightmost word ($w_M$) has the lowest rank.

$$OD = <w_1, w_2, \ldots, w_j, \ldots, w_{(M-1)}, w_M>$$

**ODT**     *Ordered Dictionary Tree*. A graphical presentation of **OD**. Let $ODT = <v_1, v_2, \ldots, v_j, \ldots, v_{(M-1)}, v_M>$, where $v_j$ is a node in **ODT** corresponding to the word $w_j$ in the **OD**.

To represent **ODT** graphically, two extra links will be needed. The first is left node link of node *i*, the second is right node link of node *i*. The links of node *i* leading to named nodes are:

1. ***Parent***(*i*) link;     2. ***Child***(*i*) link;     3. ***Sibling***(*i*) link;

4. ***Data***(*i*) link;     5. ***Left*** (*i*) link and    6. ***Right***(*i*) link.
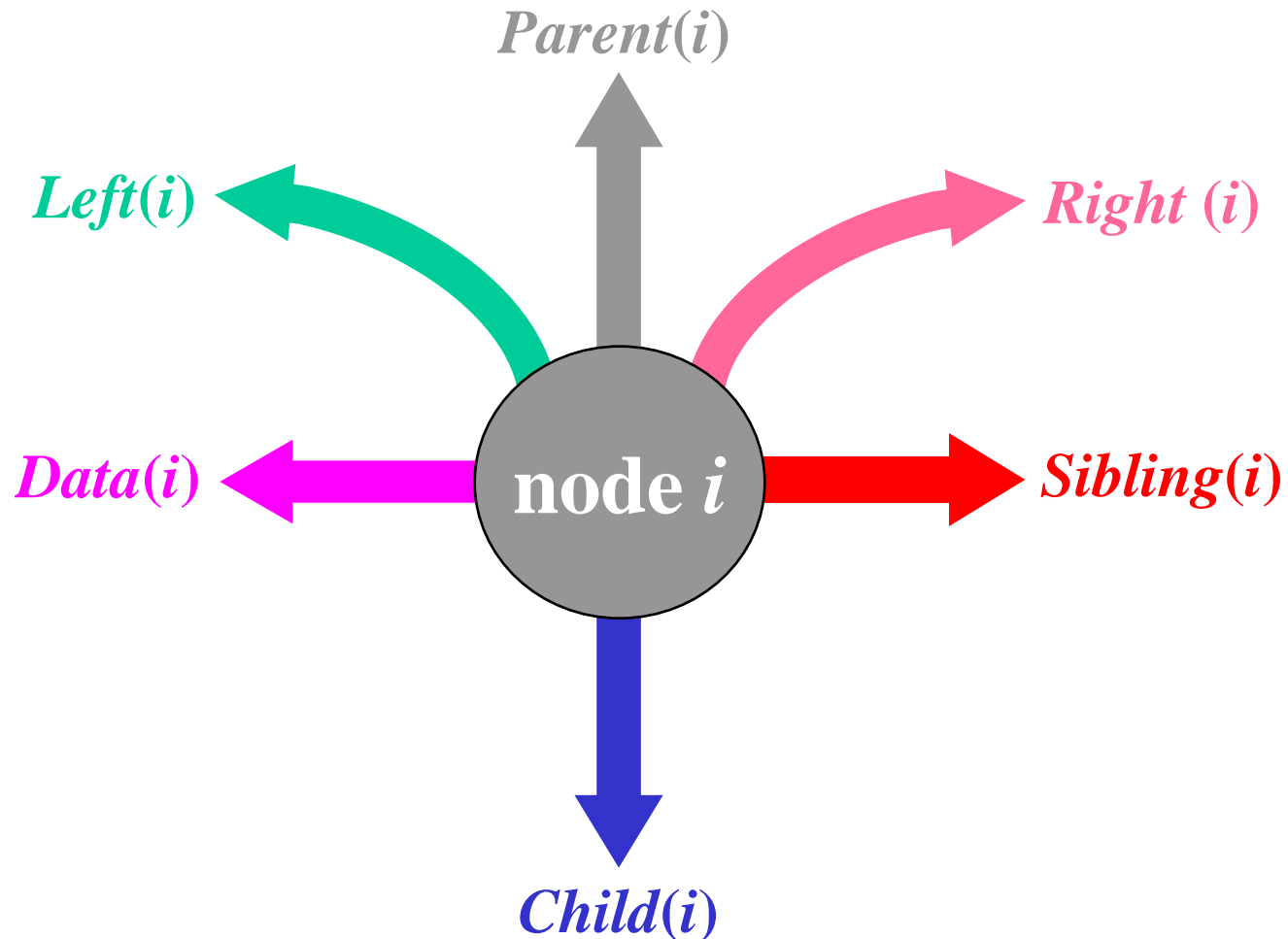
# *The Tree Structure*

## *ODT* **links**:
### With left and right links
If a link has a value equal to -1, the link is a null link.



*Parent*(*i*)

*Left*(*i*)　　　　　*Right* (*i*)

*Data*(*i*)　　**node *i***　　*Sibling*(*i*)

*Child*(*i*)

# *The Tree Structure*

## Graphical Representation of *ODT*

Let *OD* = <*3 , 257,  260, 256, 261, 1, 258, 259, ……*>

**$Data$(i) link**      **$Right$(i) link**      **$Left$(i) link**

**Left Most Node = 3**      **Right Most Node = •**

3    257    260    256    261    1    258    259    • •
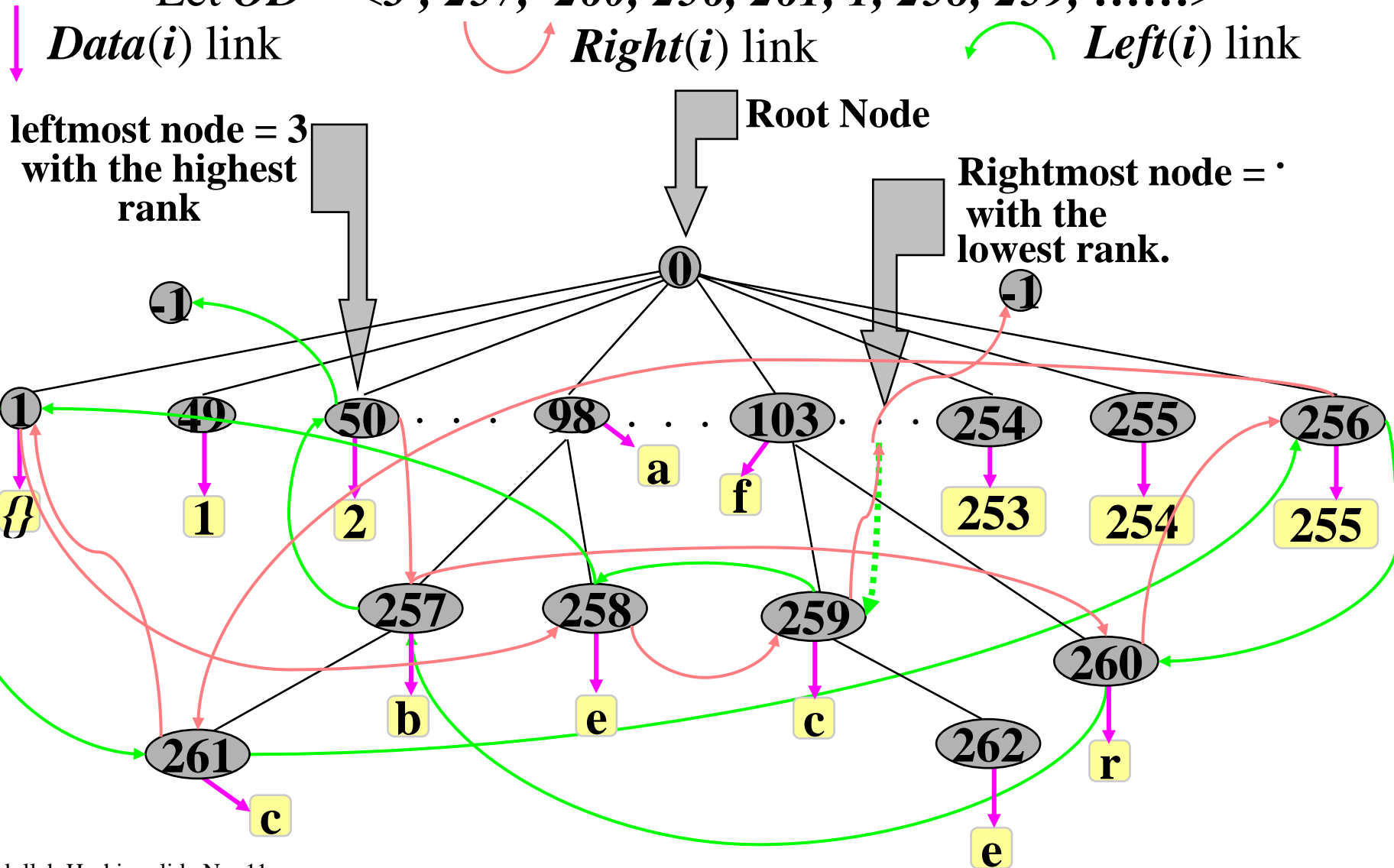
2    b    r    255    c    0    e    c

Node **3** has a higher rank than node **257**, node **257** has a higher rank than node **260, 260** has higher rank than **256**, ……etc.

# *The Tree Structure*

## Graphical Representation of *ODT*

Let *OD* = <*3 , 257,  260, 256, 261, 1, 258, 259, ……*>

*Data*(*i*) link        *Right*(*i*) link        *Left*(*i*) link

**Root Node**

**leftmost node = 3 with the highest rank**

**Rightmost node = ·  with the lowest rank.**

# *Lempel & Ziv algorithm*

Let string $s_n$ be an input string from a file of characters in alphabet $A$ over language $\Theta$, at interval $n$. Match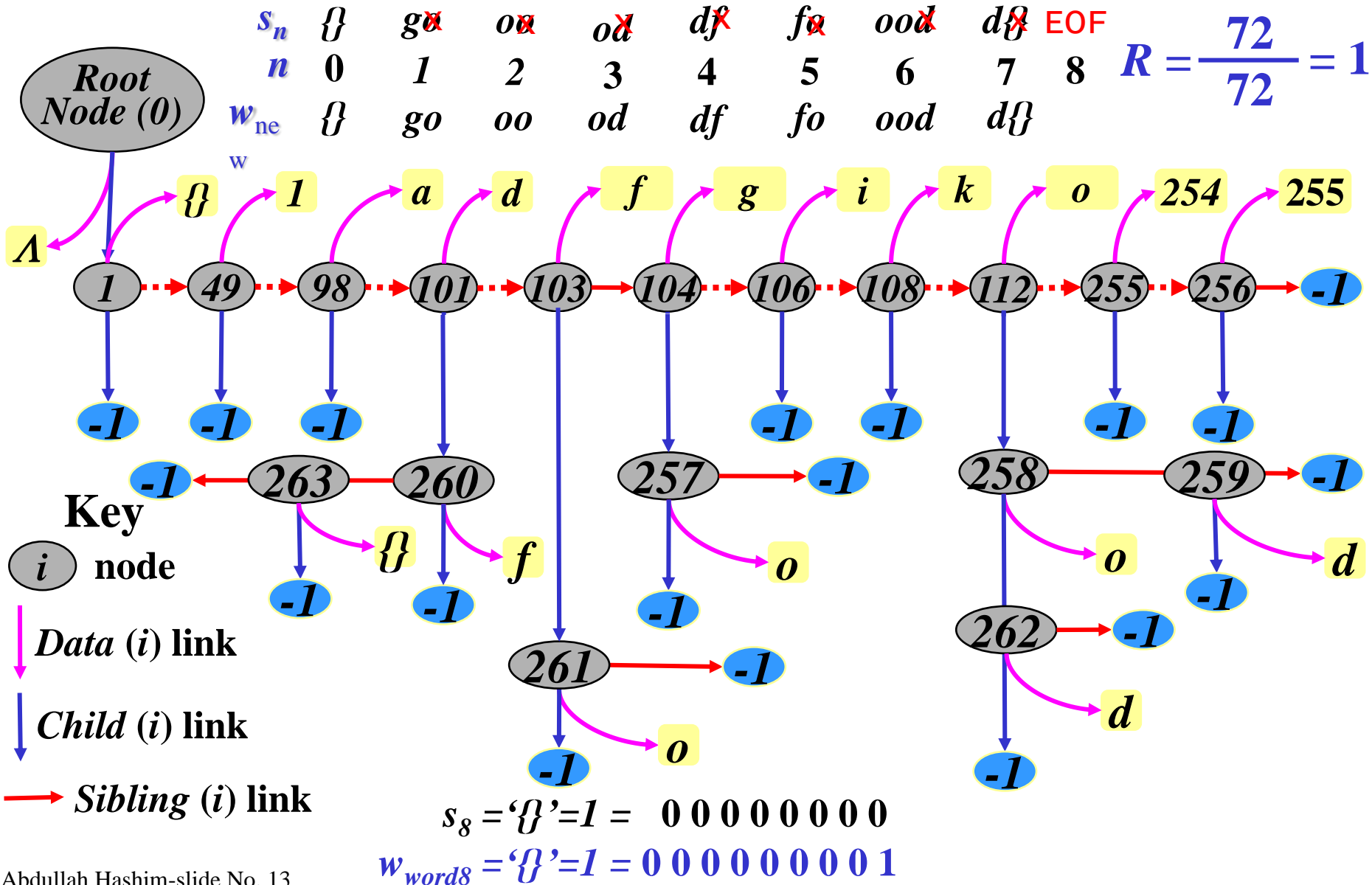 $s_n$ to the longest word $w_n$ in the dictionary given by: $s_n = w_n = <e_1 e_2 \ldots e_i \ldots e_p>$. If $e_{p+1}$ is the unmatched character resulting from the matching process, then the new word added to the dictionary at interval $n$ is $w_{new} = <e_1 e_2 \ldots e_p e_{p+1}>$.

1. At interval $n$, an input string of source symbols is matched with the longest string in the dictionary ($w_n$).

2. The matched word $w_n$ is coded by a binary codeword to form the compressed word [$w_{word}(w_n)$].

3. The input string $s_n = w_n$ is appended to the unmatched character resulting from the matching process in 2 above to form the new word ($w_{new}$).

4. The new word $w_{new}$ is added to the dictionary.

5. The process is repeated starting from step one until EOF.

# *Lempel & Ziv algorithm*

## A practical simulation for input string <goodfood{}>



| $s_n$ | {} | g̶o̶ | o̶o̶ | o̶d̶ | d̶f̶ | f̶o̶ | o̶o̶d̶ | d̶{̶}̶ | EOF |
|-------|----|----|----|----|----|----|-----|-----|-----|
| $n$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| $w_{ne}$ | {} | go | oo | od | df | fo | ood | d{} | |

$$R = \frac{72}{72} = 1$$

**Key**

- (i) node
- **Data** (*i*) link
- **Child** (*i*) link
- **Sibling** (*i*) link

$$s_8 = \text{'}\{\}\text{'} = 1 = \ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0$$

$$w_{word8} = \text{'}\{\}\text{'} = 1 = 0\ 0\ 0\ 0\ 0\ 0\ 0\ 1$$

# *Lempel & Ziv algorithm*

Lempel and Ziv in their original paper (On the Complexity of Finite sequences. IEEE Trans. IT-22.1 Jan 1976, pp75-81) have shown that if an initial dictionary of single character words of an alphabet is allowed to grow indefinitely by adding a new word at every successive interval in time and the dictionary words are coded by equal length codewords, asymptotically the compression ratio of such a coding scheme will tend to infinity as the dictionary size tends to infinity. Many practical implementations of this basic theory have been reported in the literature of different approaches and complexity. However, all suffer from the same basic inherent feature that once a practical limit is imposed on the dictionary size and its wordlength, the compression ratio will be degraded rapidly. In fact all practical implementations of the dictionary result in an average codeword length even higher than second order entropy. The low compression ratios are due to the practical implementations of the scheme which fill the dictionary mainly with two-character words of low frequency of occurrence instead of long words with high frequency.