**Algorithm 1** Attribution Tree Construction

**Input:** $[a_{i,j}^l]_{n \times n}$: Attribution scores

$\{E^l\}_{l=1}^{|l|}$: Retained attribution edges

**Output:** $\mathcal{V}, \mathcal{E}$: Node set and edge set of Attr tree

1: ▷ *Initialize the state of all tokens, each token has three states:* NotAppear, Appear, Fixed
2: **for** $i \leftarrow n, \cdots, 1$ **do**
3: $\quad State_i = $ NotAppear
4: ▷ *Choose the top node of the attribution tree*
5: $[AttrAll_i]_n = \sum_{l=1}^{|l|} \sum_{j=1, j \neq i}^{n} a_{i,j}^l$
6: $TopNode = argmax([AttrAll_i]_n)$
7: $\mathcal{V} \leftarrow \{TopNode\}; State_{TopNode} = $ Appear
8: ▷ *Build the attribution tree downward*
9: **for** $l \leftarrow |l| - 1, \cdots, 1$ **do**
10: $\quad$ **for** $(i,j)_{i \neq j}^l \in E^l$ **do**
11: $\quad\quad$ **if** $State_i$ is Appear and $State_j$ is NotAppear **then**
12: $\quad\quad\quad \mathcal{E} \leftarrow \mathcal{E} \bigcup \{(i,j)\}$
13: $\quad\quad\quad \mathcal{V} \leftarrow \mathcal{V} \bigcup \{j\}$
14: $\quad\quad\quad State_i = $ Fixed
15: $\quad\quad\quad State_j = $ Appear
16: $\quad\quad$ **if** $State_i$ is Fixed and $State_j$ is NotAppear **then**
17: $\quad\quad\quad \mathcal{E} \leftarrow \mathcal{E} \bigcup \{(i,j)\}$
18: $\quad\quad\quad \mathcal{V} \leftarrow \mathcal{V} \bigcup \{j\}$
19: $\quad\quad\quad State_j = $ Appear
20: ▷ *Add the terminal of the information flow*
21: $\mathcal{V} \leftarrow \{[\text{CLS}]\}$
22: **for** $j \leftarrow n, \cdots, 1$ **do**
23: $\quad$ **if** $State_j \in \{$Appear, Fixed$\}$ **then**
24: $\quad\quad \mathcal{E} \leftarrow \mathcal{E} \bigcup \{([\text{CLS}], j)\}$
25: **return** $Tree = \{\mathcal{V}, \mathcal{E}\}$

Sentence: Seldom has a movie so closely matched the spirit of a man and his work

Model input is the attribution scores between each token in the sentence and the edges which tokens are related

## Algorithm 1 Attribution Tree Construction

**Input:** $[a_{i,j}^l]_{n \times n}$: Attribution scores

$\{E^l\}_{l=1}^{|l|}$: Retained attribution edges

**Output:** $\mathcal{V}, \mathcal{E}$: Node set and edge set of Attr tree

1: ▷ *Initialize the state of all tokens, each token has three states:* NotAppear, Appear, Fixed
2: **for** $i \leftarrow n, \cdots, 1$ **do**
3:      $State_i = $ NotAppear
4: ▷ *Choose the top node of the attribution tree*
5: $[AttrAll_i]_n = \sum_{l=1}^{|l|} \sum_{j=1, j \neq i}^{n} a_{i,j}^l$
6: $TopNode = argmax([AttrAll_i]_n)$
7: $\mathcal{V} \leftarrow \{TopNode\}; State_{TopNode} = $ Appear
8: ▷ *Build the attribution tree downward*
9: **for** $l \leftarrow |l| - 1, \cdots, 1$ **do**
10:      **for** $(i,j)_{i \neq j}^l \in E^l$ **do**
11:          **if** $State_i$ is Appear and $State_j$ is NotAppear **then**
12:              $\mathcal{E} \leftarrow \mathcal{E} \bigcup \{(i,j)\}$
13:              $\mathcal{V} \leftarrow \mathcal{V} \bigcup \{j\}$
14:              $State_i = $ Fixed
15:              $State_j = $ Appear
16:          **if** $State_i$ is Fixed and $State_j$ is NotAppear **then**
17:              $\mathcal{E} \leftarrow \mathcal{E} \bigcup \{(i,j)\}$
18:              $\mathcal{V} \leftarrow \mathcal{V} \bigcup \{j\}$
19:              $State_j = $ Appear
20: ▷ *Add the terminal of the information flow*
21: $\mathcal{V} \leftarrow \{[\text{CLS}]\}$
22: **for** $j \leftarrow n, \cdots, 1$ **do**
23:      **if** $State_j \in \{$Appear, Fixed$\}$ **then**
24:          $\mathcal{E} \leftarrow \mathcal{E} \bigcup \{([\text{CLS}], j)\}$
25: **return** $Tree = \{\mathcal{V}, \mathcal{E}\}$

## Status NotAppear Tokens

| | | |
|---|---|---|
| seldom | closely | a |
| has | matched | man |
| a | the | and |
| movie | spirit | his |
| so | of | work |

**Algorithm 1** Attribution Tree Construction

**Input:** $[a_{i,j}^l]_{n \times n}$: Attribution scores

$\{E^l\}_{l=1}^{|l|}$: Retained attribution edges

**Output:** $\mathcal{V}, \mathcal{E}$: Node set and edge set of Attr tree

1: ▷ *Initialize the state of all tokens, each token has three states:* NotAppear, Appear, Fixed
2: **for** $i \leftarrow n, \cdots, 1$ **do**
3: $\quad State_i = $ NotAppear
4: ▷ *Choose the top node of the attribution tree*
5: $[AttrAll_i]_n = \sum_{l=1}^{|l|} \sum_{j=1, j \neq i}^{n} a_{i,j}^l$
6: $TopNode = argmax([AttrAll_i]_n)$
7: $\mathcal{V} \leftarrow \{TopNode\}; State_{TopNode} = $ Appear
8: ▷ *Build the attribution tree downward*
9: **for** $l \leftarrow |l| - 1, \cdots, 1$ **do**
10: $\quad$ **for** $(i,j)_{i \neq j}^l \in E^l$ **do**
11: $\quad\quad$ **if** $State_i$ is Appear and $State_j$ is NotAppear **then**
12: $\quad\quad\quad \mathcal{E} \leftarrow \mathcal{E} \bigcup \{(i,j)\}$
13: $\quad\quad\quad \mathcal{V} \leftarrow \mathcal{V} \bigcup \{j\}$
14: $\quad\quad\quad State_i = $ Fixed
15: $\quad\quad\quad State_j = $ Appear
16: $\quad\quad$ **if** $State_i$ is Fixed and $State_j$ is NotAppear **then**
17: $\quad\quad\quad \mathcal{E} \leftarrow \mathcal{E} \bigcup \{(i,j)\}$
18: $\quad\quad\quad \mathcal{V} \leftarrow \mathcal{V} \bigcup \{j\}$
19: $\quad\quad\quad State_j = $ Appear
20: ▷ *Add the terminal of the information flow*
21: $\mathcal{V} \leftarrow \{[\text{CLS}]\}$
22: **for** $j \leftarrow n, \cdots, 1$ **do**
23: $\quad$ **if** $State_j \in \{$Appear, Fixed$\}$ **then**
24: $\quad\quad \mathcal{E} \leftarrow \mathcal{E} \bigcup \{([\text{CLS}], j)\}$
25: **return** $Tree = \{\mathcal{V}, \mathcal{E}\}$

**Attribution Tree**

matched

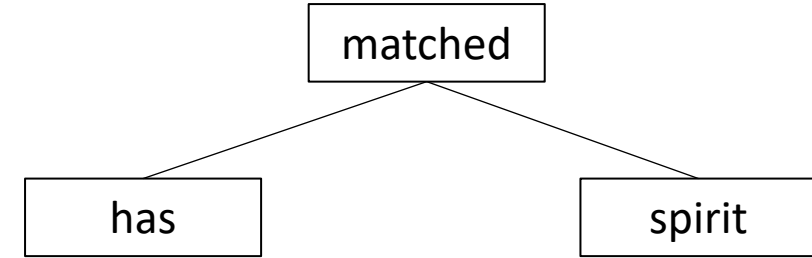**Algorithm 1** Attribution Tree Construction

**Input:**  $[a_{i,j}^l]_{n \times n}$: Attribution scores

$\{E^l\}_{l=1}^{|l|}$: Retained attribution edges

**Output:**  $\mathcal{V}, \mathcal{E}$: Node set and edge set of Attr tree

 1: ▷ *Initialize the state of all tokens, each token has three states:* `NotAppear, Appear, Fixed`
 2: **for** $i \leftarrow n, \cdots, 1$ **do**
 3:     $State_i = $ `NotAppear`
 4: ▷ *Choose the top node of the attribution tree*
 5: $[AttrAll_i]_n = \sum_{l=1}^{|l|} \sum_{j=1, j \neq i}^{n} a_{i,j}^l$
 6: $TopNode = argmax([AttrAll_i]_n)$
 7: $\mathcal{V} \leftarrow \{TopNode\}; State_{TopNode} = $ `Appear`
 8: ▷ *Build the attribution tree downward*
 9: **for** $l \leftarrow |l| - 1, \cdots, 1$ **do**
10:     **for** $(i,j)_{i \neq j}^l \in E^l$ **do**
11:        **if** $State_i$ is `Appear` and $State_j$ is `NotAppear` **then**
12:           $\mathcal{E} \leftarrow \mathcal{E} \bigcup \{(i,j)\}$
13:           $\mathcal{V} \leftarrow \mathcal{V} \bigcup \{j\}$
14:           $State_i = $ `Fixed`
15:           $State_j = $ `Appear`
16:        **if** $State_i$ is `Fixed` and $State_j$ is `NotAppear` **then**
17:           $\mathcal{E} \leftarrow \mathcal{E} \bigcup \{(i,j)\}$
18:           $\mathcal{V} \leftarrow \mathcal{V} \bigcup \{j\}$
19:           $State_j = $ `Appear`
20: ▷ *Add the terminal of the information flow*
21: $\mathcal{V} \leftarrow \{$`[CLS]`$\}$
22: **for** $j \leftarrow n, \cdots, 1$ **do**
23:     **if** $State_j \in \{$`Appear, Fixed`$\}$ **then**
24:        $\mathcal{E} \leftarrow \mathcal{E} \bigcup \{($`[CLS]`$, j)\}$
25: **return** $Tree = \{\mathcal{V}, \mathcal{E}\}$



**Attribution Tree**
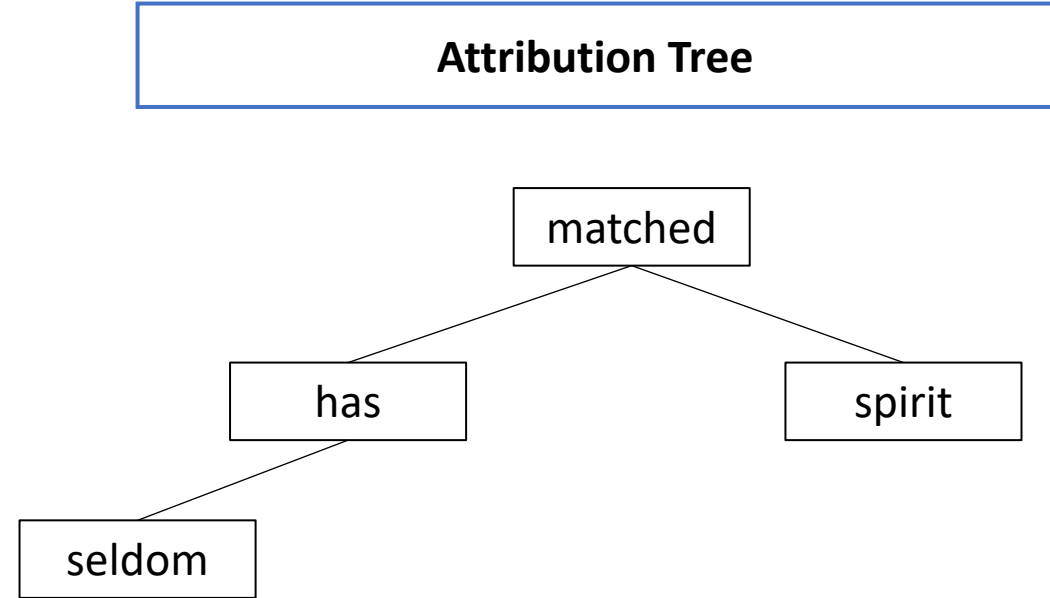
matched

has      spirit

**Algorithm 1** Attribution Tree Construction

**Input:**    $[a_{i,j}^l]_{n \times n}$: Attribution scores

            $\{E^l\}_{l=1}^{|l|}$: Retained attribution edges

**Output:**  $\mathcal{V}, \mathcal{E}$: Node set and edge set of Attr tree

  1: ▷ *Initialize the state of all tokens, each token has three states:* `NotAppear`, `Appear`, `Fixed`

  2: **for** $i \leftarrow n, \cdots, 1$ **do**

  3:     $State_i = $ `NotAppear`

  4: ▷ *Choose the top node of the attribution tree*

  5: $[AttrAll_i]_n = \sum_{l=1}^{|l|} \sum_{j=1, j \neq i}^{n} a_{i,j}^l$

  6: $TopNode = argmax([AttrAll_i]_n)$

  7: $\mathcal{V} \leftarrow \{TopNode\}; State_{TopNode} = $ `Appear`

  8: ▷ *Build the attribution tree downward*

  9: **for** $l \leftarrow |l| - 1, \cdots, 1$ **do**

10:     **for** $(i, j)_{i \neq j}^l \in E^l$ **do**

11:         **if** $State_i$ is `Appear` and $State_j$ is `NotAppear` **then**

12:             $\mathcal{E} \leftarrow \mathcal{E} \bigcup \{(i, j)\}$

13:             $\mathcal{V} \leftarrow \mathcal{V} \bigcup \{j\}$

14:             $State_i = $ `Fixed`

15:             $State_j = $ `Appear`

16:         **if** $State_i$ is `Fixed` and $State_j$ is `NotAppear` **then**

17:             $\mathcal{E} \leftarrow \mathcal{E} \bigcup \{(i, j)\}$

18:             $\mathcal{V} \leftarrow \mathcal{V} \bigcup \{j\}$

19:             $State_j = $ `Appear`

20: ▷ *Add the terminal of the information flow*

21: $\mathcal{V} \leftarrow \{$`[CLS]`$\}$

22: **for** $j \leftarrow n, \cdots, 1$ **do**

23:     **if** $State_j \in \{$`Appear`, `Fixed`$\}$ **then**

24:         $\mathcal{E} \leftarrow \mathcal{E} \bigcup \{($`[CLS]`$, j)\}$

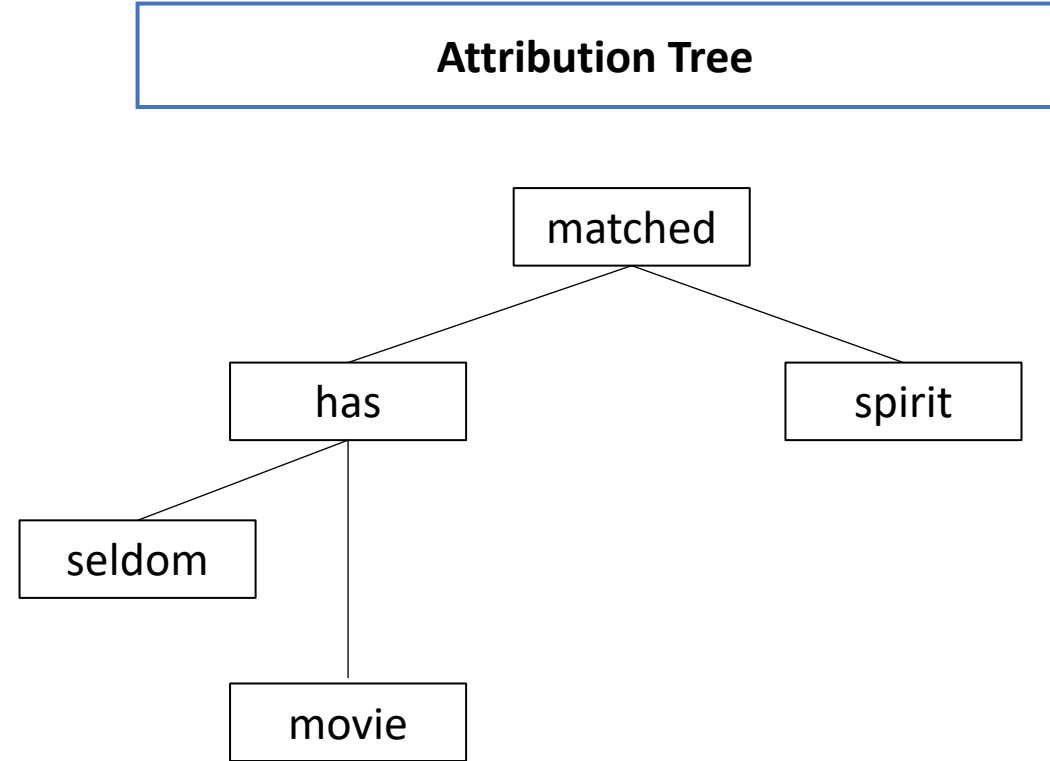25: **return** $Tree = \{\mathcal{V}, \mathcal{E}\}$



Attribution Tree

**Algorithm 1** Attribution Tree Construction

**Input:** $[a_{i,j}^l]_{n \times n}$: Attribution scores

$\{E^l\}_{l=1}^{|l|}$: Retained attribution edges

**Output:** $\mathcal{V}, \mathcal{E}$: Node set and edge set of Attr tree

1:  ▷ *Initialize the state of all tokens, each token has three states:* NotAppear, Appear, Fixed
2:  **for** $i \leftarrow n, \cdots, 1$ **do**
3:      $State_i = $ NotAppear
4:  ▷ *Choose the top node of the attribution tree*
5:  $[AttrAll_i]_n = \sum_{l=1}^{|l|} \sum_{j=1, j \neq i}^{n} a_{i,j}^l$
6:  $TopNode = argmax([AttrAll_i]_n)$
7:  $\mathcal{V} \leftarrow \{TopNode\}; State_{TopNode} = $ Appear
8:  ▷ *Build the attribution tree downward*
9:  **for** $l \leftarrow |l| - 1, \cdots, 1$ **do**
10:     **for** $(i, j)_{i \neq j}^l \in E^l$ **do**
11:         **if** $State_i$ is Appear and $State_j$ is NotAppear **then**
12:             $\mathcal{E} \leftarrow \mathcal{E} \bigcup \{(i, j)\}$
13:             $\mathcal{V} \leftarrow \mathcal{V} \bigcup \{j\}$
14:             $State_i = $ Fixed
15:             $State_j = $ Appear
16:         **if** $State_i$ is Fixed and $State_j$ is NotAppear **then**
17:             $\mathcal{E} \leftarrow \mathcal{E} \bigcup \{(i, j)\}$
18:             $\mathcal{V} \leftarrow \mathcal{V} \bigcup \{j\}$
19:             $State_j = $ Appear
20: ▷ *Add the terminal of the information flow*
21: $\mathcal{V} \leftarrow \{[CLS]\}$
22: **for** $j \leftarrow n, \cdots, 1$ **do**
23:     **if** $State_j \in \{$Appear, Fixed$\}$ **then**
24:         $\mathcal{E} \leftarrow \mathcal{E} \bigcup \{([CLS], j)\}$
25: **return** $Tree = \{\mathcal{V}, \mathcal{E}\}$



**Attribution Tree**

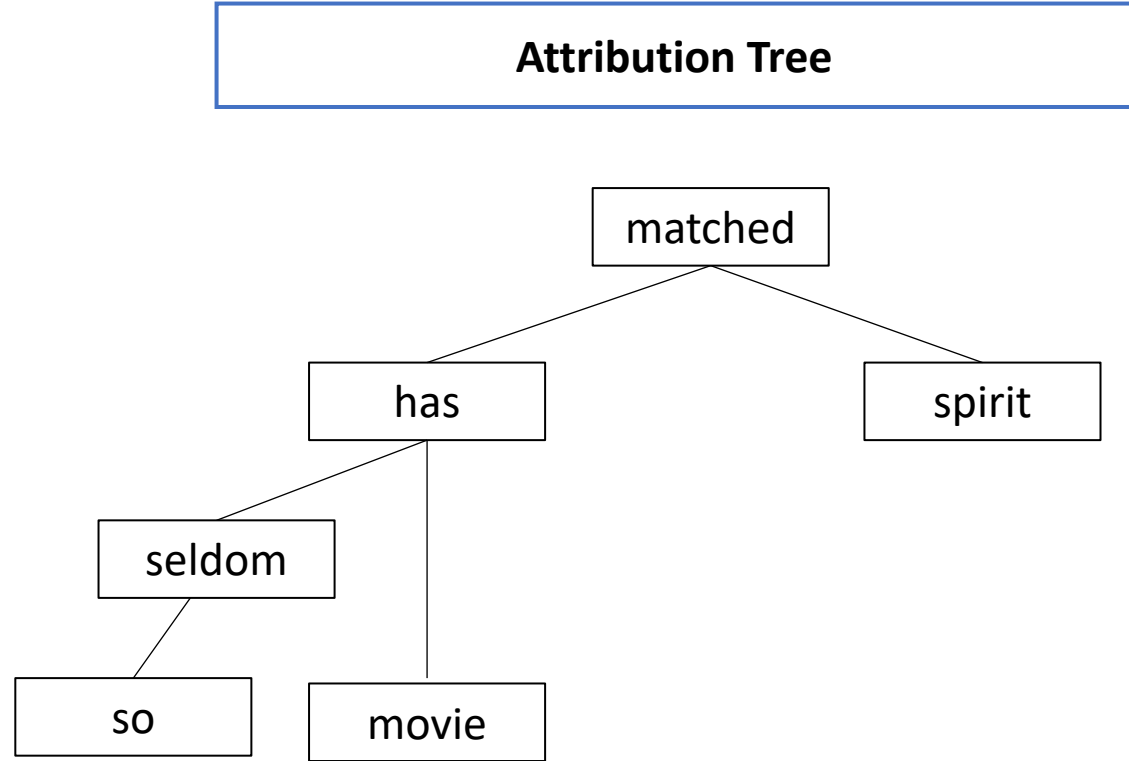matched
├── has
│   ├── seldom
│   └── movie
└── spirit

**Algorithm 1** Attribution Tree Construction

**Input:**  $[a_{i,j}^l]_{n \times n}$ : Attribution scores

$\{E^l\}_{l=1}^{|l|}$ : Retained attribution edges

**Output:**  $\mathcal{V}, \mathcal{E}$ : Node set and edge set of Attr tree

1: ▷ *Initialize the state of all tokens, each token has three states:* `NotAppear`, `Appear`, `Fixed`

2: **for** $i \leftarrow n, \cdots, 1$ **do**

3:      $State_i = $ `NotAppear`

4: ▷ *Choose the top node of the attribution tree*

5: $[AttrAll_i]_n = \sum_{l=1}^{|l|} \sum_{j=1, j \neq i}^{n} a_{i,j}^l$

6: $TopNode = argmax([AttrAll_i]_n)$

7: $\mathcal{V} \leftarrow \{TopNode\}; State_{TopNode} = $ `Appear`

8: ▷ *Build the attribution tree downward*

9: **for** $l \leftarrow |l| - 1, \cdots, 1$ **do**

10:      **for** $(i,j)_{i \neq j}^l \in E^l$ **do**

11:          **if** $State_i$ is `Appear` and $State_j$ is `NotAppear` **then**

12:              $\mathcal{E} \leftarrow \mathcal{E} \bigcup \{(i,j)\}$

13:              $\mathcal{V} \leftarrow \mathcal{V} \bigcup \{j\}$

14:              $State_i = $ `Fixed`

15:              $State_j = $ `Appear`

16:          **if** $State_i$ is `Fixed` and $State_j$ is `NotAppear` **then**

17:              $\mathcal{E} \leftarrow \mathcal{E} \bigcup \{(i,j)\}$

18:              $\mathcal{V} \leftarrow \mathcal{V} \bigcup \{j\}$

19:              $State_j = $ `Appear`

20: ▷ *Add the terminal of the information flow*

21: $\mathcal{V} \leftarrow \{$ `[CLS]` $\}$

22: **for** $j \leftarrow n, \cdots, 1$ **do**

23:      **if** $State_j \in \{$ `Appear`, `Fixed` $\}$ **then**

24:          $\mathcal{E} \leftarrow \mathcal{E} \bigcup \{($ `[CLS]` $, j)\}$

25: **return** $Tree = \{\mathcal{V}, \mathcal{E}\}$
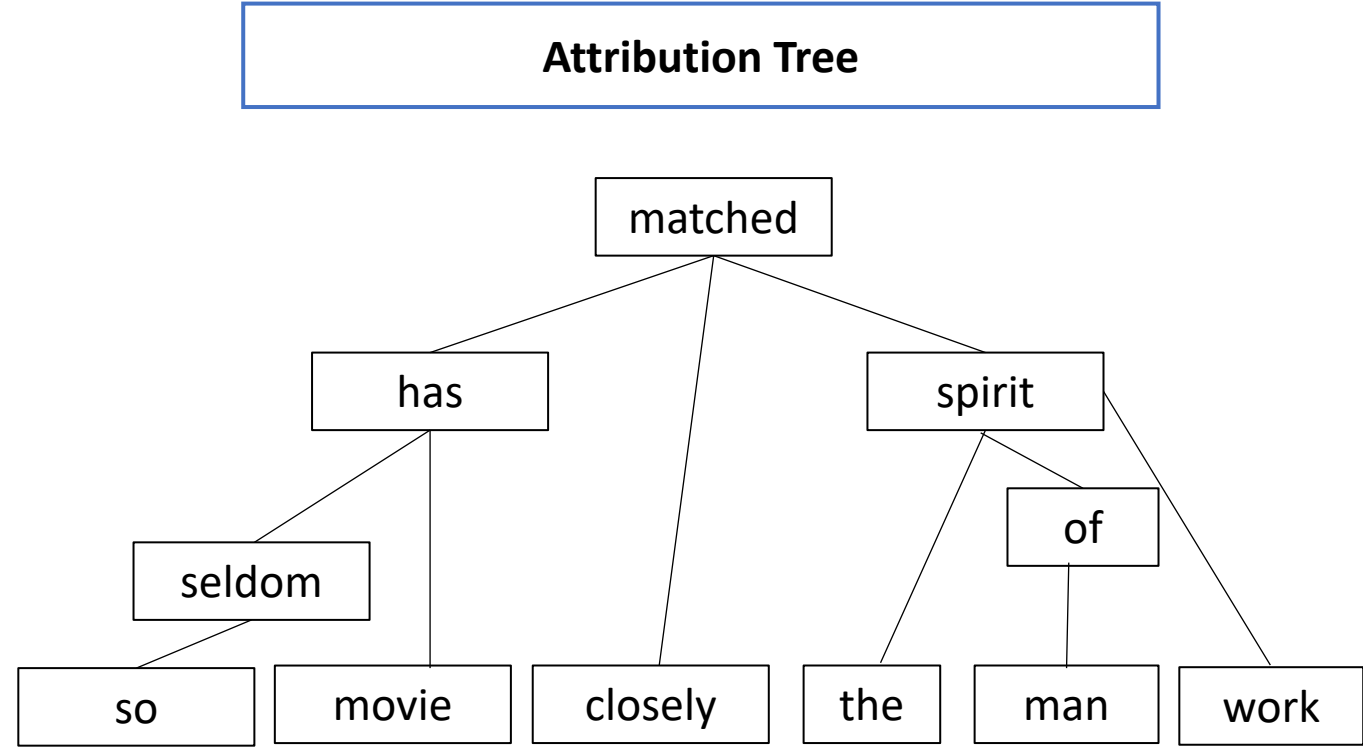


Attribution Tree

**Algorithm 1** Attribution Tree Construction

**Input:**    $[a_{i,j}^l]_{n \times n}$: Attribution scores

           $\{E^l\}_{l=1}^{|l|}$: Retained attribution edges

**Output:**  $\mathcal{V}, \mathcal{E}$: Node set and edge set of Attr tree

  1: ▷ *Initialize the state of all tokens, each token has three states:* `NotAppear`, `Appear`, `Fixed`

  2: **for** $i \leftarrow n, \cdots, 1$ **do**

  3:     $State_i = $ `NotAppear`

  4: ▷ *Choose the top node of the attribution tree*

  5: $[AttrAll_i]_n = \sum_{l=1}^{|l|} \sum_{j=1, j \neq i}^{n} a_{i,j}^l$

  6: $TopNode = argmax([AttrAll_i]_n)$

  7: $\mathcal{V} \leftarrow \{TopNode\}; State_{TopNode} = $ `Appear`

  8: ▷ *Build the attribution tree downward*

  9: **for** $l \leftarrow |l| - 1, \cdots, 1$ **do**

10:     **for** $(i,j)_{i \neq j}^l \in E^l$ **do**

11:         **if** $State_i$ is `Appear` and $State_j$ is `NotAppear` **then**

12:             $\mathcal{E} \leftarrow \mathcal{E} \bigcup \{(i,j)\}$

13:             $\mathcal{V} \leftarrow \mathcal{V} \bigcup \{j\}$

14:             $State_i = $ `Fixed`

15:             $State_j = $ `Appear`

16:         **if** $State_i$ is `Fixed` and $State_j$ is `NotAppear` **then**

17:             $\mathcal{E} \leftarrow \mathcal{E} \bigcup \{(i,j)\}$

18:             $\mathcal{V} \leftarrow \mathcal{V} \bigcup \{j\}$

19:             $State_j = $ `Appear`

20: ▷ *Add the terminal of the information flow*

21: $\mathcal{V} \leftarrow \{[\texttt{CLS}]\}$

22: **for** $j \leftarrow n, \cdots, 1$ **do**

23:     **if** $State_j \in \{$`Appear`, `Fixed`$\}$ **then**

24:         $\mathcal{E} \leftarrow \mathcal{E} \bigcup \{([\texttt{CLS}], j)\}$

25: **return** $Tree = \{\mathcal{V}, \mathcal{E}\}$



**Attribution Tree**

**Algorithm 1** Attribution Tree Construction

**Input:**    $[a_{i,j}^l]_{n \times n}$: Attribution scores

            $\{E^l\}_{l=1}^{|l|}$: Retained attribution edges

**Output:**   $\mathcal{V}, \mathcal{E}$: Node set and edge set of Attr tree

  1: ▷ *Initialize the state of all tokens, each token has three states:* NotAppear, Appear, Fixed

  2: **for** $i \leftarrow n, \cdots, 1$ **do**

  3:      $State_i = $ NotAppear

  4: ▷ *Choose the top node of the attribution tree*

  5: $[AttrAll_i]_n = \sum_{l=1}^{|l|} \sum_{j=1, j \neq i}^n a_{i,j}^l$

  6: $TopNode = argmax([AttrAll_i]_n)$

  7: $\mathcal{V} \leftarrow \{TopNode\}; State_{TopNode} = $ Appear

  8: ▷ *Build the attribution tree downward*

  9: **for** $l \leftarrow |l| - 1, \cdots, 1$ **do**

10:     **for** $(i,j)_{i \neq j}^l \in E^l$ **do**

11:        **if** $State_i$ is Appear and $State_j$ is NotAppear **then**

12:           $\mathcal{E} \leftarrow \mathcal{E} \bigcup \{(i,j)\}$

13:           $\mathcal{V} \leftarrow \mathcal{V} \bigcup \{j\}$

14:           $State_i = $ Fixed

15:           $State_j = $ Appear

16:        **if** $State_i$ is Fixed and $State_j$ is NotAppear **then**

17:           $\mathcal{E} \leftarrow \mathcal{E} \bigcup \{(i,j)\}$

18:           $\mathcal{V} \leftarrow \mathcal{V} \bigcup \{j\}$

19:           $State_j = $ Appear

20: ▷ *Add the terminal of the information flow*

21: $\mathcal{V} \leftarrow \{[\text{CLS}]\}$

22: **for** $j \leftarrow n, \cdots, 1$ **do**

23:    **if** $State_j \in \{$Appear, Fixed$\}$ **then**

24:        $\mathcal{E} \leftarrow \mathcal{E} \bigcup \{([\text{CLS}], j)\}$

25: **return** $Tree = \{\mathcal{V}, \mathcal{E}\}$



Attribution Tree