

# Implementation of Hopcroft's Algorithm

Hang Zhou

ENS

7 January 2009

1 Introduction

2 Data Structures

3 Analysis of Time Complexity

## Definition

Let  $\mathcal{P}$  be a partition of  $Q$ . The number of classes in  $\mathcal{P}$  is bounded by  $N$ . Hence every class of  $\mathcal{P}$  can be entitled with a unique **name** between 1 and  $N$ .

## Definition

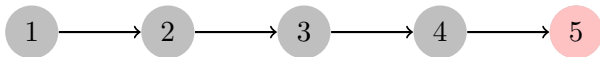
Let  $\mathcal{A}$  be a deterministic finite automaton  $(Q, A, E, I, F)$  and let  $a \in A$  and  $B, C \subseteq Q$ . Then  $B$  is **split** by  $(C, a)$ , if  $B \cdot a \not\subseteq C$  and  $B \cdot a \cap C \neq \emptyset$  with  $B \cdot a = \{x \cdot a | x \in B\}$ .

# Hopcroft's Algorithm - First Version

ALGORITHM1()

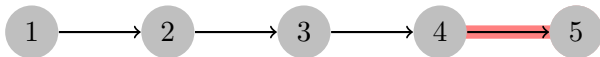
```
1   $\mathcal{P} \leftarrow \{F, F^c\}$ 
2   $S \leftarrow \emptyset$ 
3  for  $a \in A$ 
4  do INSERT ( $\min(F, F^c), a$ ) to  $S$ 
5  while  $S \neq \emptyset$ 
6  do DELETE ( $C, a$ ) from  $S$ 
7      $\mathcal{T} \leftarrow \{B \mid B \text{ is split by } (C, a)\}$ 
8     for each  $B \in \mathcal{T}$ 
9     do  $B', B'' \leftarrow \text{SPLIT}(B, C, a)$ 
10     REPLACE  $B$  by  $B'$  and  $B''$  in  $\mathcal{P}$ 
11     for  $b \in A$ 
12     do if  $(B, b) \in S$ 
13         then REPLACE  $(B, b)$  by  $(B', b)$  and  $(B'', b)$  in  $S$ 
14         else INSERT ( $\min(B', B''), b$ ) to  $S$ 
```

# Extreme Condition



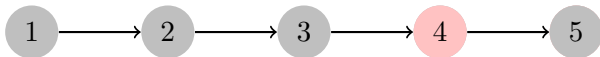
Partition in this moment : 1234 - 5

# Extreme Condition



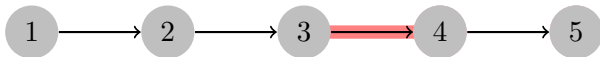
Partition in this moment : 1234 - 5

# Extreme Condition



Partition in this moment : 123- 4 - 5

# Extreme Condition



Partition in this moment : 123- 4 - 5



# Extreme Condition



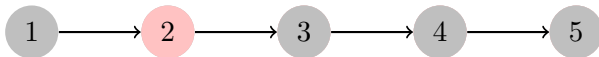
Partition in this moment : 12 - 3 - 4 - 5

# Extreme Condition



Partition in this moment : 12 - 3 - 4 - 5

# Extreme Condition



Partition in this moment : 1 - 2 - 3 - 4 - 5

# Extreme Condition



Partition in this moment : 1 - 2 - 3 - 4 - 5

# Extreme Condition



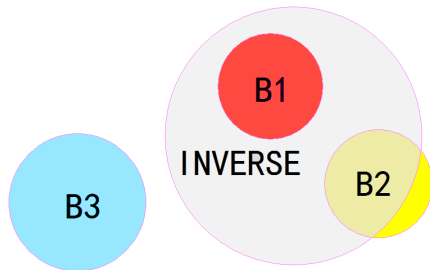
Partition in this moment : 1 - 2 - 3 - 4 - 5

# The Question Throughout the Writing

What's the execution time of the *while* loop if it is carried out efficiently?

Idea:

- define **INVERSE** =  $a^{-1}C$
- wish to prove  $O(|\text{INVERSE}|)$  complexity of *while* loop

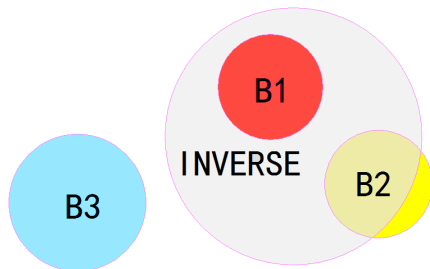


# Difficulty

We cannot directly obtain  $\mathcal{T}$  in  $O(|\text{INVERSE}|)$ , what to do?

Solution:

- Instead, we calculate its **approximate substitute**  $\mathcal{T}'$  :  
$$\mathcal{T}' = \{B \mid B \cap \text{INVERSE} \neq \emptyset\}$$
- $O(1)$  time to check whether an element in  $\mathcal{T}'$  is also in  $\mathcal{T}$



# Hopcroft's Algorithm - Second Version

ALGORITHM2()

```
1   $\mathcal{P} \leftarrow \{F, F^c\}$ 
2   $S \leftarrow \emptyset$ 
3  for  $a \in A$ 
4  do INSERT ( $\min(F, F^c), a$ ) to  $S$ 
5  while  $S \neq \emptyset$ 
6  do DELETE ( $C, a$ ) from  $S$ 
7     INVERSE  $\leftarrow a^{-1}C$ 
8      $\mathcal{T}' \leftarrow \{B \mid B \cap \text{INVERSE} \neq \emptyset\}$ 
9     for each  $B \in \mathcal{T}'$ 
10    do if  $B \cdot a \not\subseteq C$ 
11        then  $B', B'' \leftarrow \text{SPLIT}(B, C, a)$ 
12            REPLACE  $B$  by  $B'$  and  $B''$  in  $\mathcal{P}$ 
13            for  $b \in A$ 
14                do if  $(B, b) \in S$ 
15                    then REPLACE  $(B, b)$  by  $(B', b)$  and  $(B'', b)$  in  $S$ 
16                    else INSERT ( $\min(B', B''), b$ ) to  $S$ 
```



# Calculate INVERSE efficiently

By definition,  $\text{INVERSE} = a^{-1}C$ .

We use a two-dimensional array  $\text{Inv}$  :

- $\text{Inv}[x,a]$  :  
a pointer to the linked list of all states  $y$ , such that  $y \cdot a = x$

- So we have :

$$\begin{array}{c} \text{for all } x \in \text{INVERSE} \\ \updownarrow \\ \text{for all } y \in C \text{ and all } x \in \text{Inv}[y, a] \end{array}$$

# Support list operations efficiently

List operations includes:

- **search** an element
- **insert** an element
- **delete** an *arbitrary* element
- test whether the list is **empty**

We use a **linked list S** and an **array InList** together, s.t.

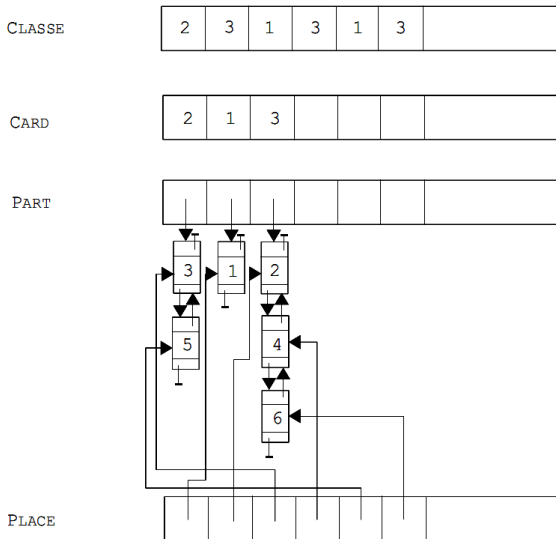
$$\text{InList}[C, a] = \begin{cases} \text{a pointer to the position of } (C, a) \text{ in } S & \text{if } (C, a) \in S \\ \text{Null} & \text{if } (C, a) \notin S \end{cases}$$

**$O(1)$  time** of all operations above

# Characterize a Partition

- **Class** [ $x$ ] : the name of the class which contains state  $x$ ;
- **Part** [ $i$ ] : a pointer to the doubly linked list with all states in class  $i$ ;
- **Card** [ $i$ ] : the size of class  $i$ ;
- **Place** [ $x$ ] : a pointer to the position of state  $x$  in the linked list **Part** [ $p$ ], where  $p = \text{Class}[x]$ .

# Example: $Q = \{1, \dots, 6\}$ and $\mathcal{P}: 35-1-246$



# Split the Classes

- **Involved** : a linked list of the names of classes in  $\mathcal{T}'$ ;
- **Size[i]** : the size of  $B \cap a^{-1}C$ , where  $B$  is the class  $i$ ;
- **Twin[i]** : the name of the newly created class while splitting the class  $i$ .

*Noting:* They should be cleared every time the *while* loop is executed.

# Time Complexity - Part 1

ALGORITHM2()

1  $\mathcal{P} \leftarrow \{F, F^c\}$

2  $S \leftarrow \emptyset$

3 **for**  $a \in A$

4 **do** INSERT ( $\min(F, F^c), a$ ) to  $S$

$O(M)$

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

# Time Complexity - Part 2

ALGORITHM2()

1  $\mathcal{P} \leftarrow \{F, F^c\}$

2  $S \leftarrow \emptyset$

3 **for**  $a \in A$

4 **do** INSERT ( $\min(F, F^c), a$ ) to  $S$

$O(M)$

5

6 **while**  $S \neq \emptyset$

7 **do** DELETE ( $C, a$ ) from  $S$

$O(MN)?$

8

9

10

11

12

13

14

15

16

17

18

19

# Time Complexity - Part 3

ALGORITHM2()

1  $\mathcal{P} \leftarrow \{F, F^c\}$

2  $S \leftarrow \emptyset$

3 **for**  $a \in A$

4 **do** INSERT ( $\min(F, F^c), a$ ) to  $S$

$O(M)$

5

6 **while**  $S \neq \emptyset$

7 **do** DELETE ( $C, a$ ) from  $S$

$O(MN)?$

8

9     INVERSE  $\leftarrow a^{-1}C$

10      $\mathcal{T}' \leftarrow \{B \mid B \cap \text{INVERSE} \neq \emptyset\}$

11     **for** each  $B \in \mathcal{T}'$

12     **do if**  $B \cdot a \notin C$

13         **then**  $B', B'' \leftarrow \text{SPLIT}(B, C, a)$

14         REPLACE  $B$  by  $B'$  and  $B''$  in  $\mathcal{P}$

$O(MN \log N)?$

15

16

17

18

19



# Time Complexity - Part 4

ALGORITHM2()

```
1   $\mathcal{P} \leftarrow \{F, F^c\}$ 
2   $S \leftarrow \emptyset$ 
3  for  $a \in A$ 
4  do INSERT ( $\min(F, F^c), a$ ) to  $S$   $O(M)$ 
5
6  while  $S \neq \emptyset$ 
7  do DELETE ( $C, a$ ) from  $S$   $O(MN)?$ 
8
9      INVERSE  $\leftarrow a^{-1}C$ 
10      $T' \leftarrow \{B \mid B \cap \text{INVERSE} \neq \emptyset\}$ 
11     for each  $B \in T'$ 
12     do if  $B \cdot a \notin C$ 
13         then  $B', B'' \leftarrow \text{SPLIT}(B, C, a)$ 
14             REPLACE  $B$  by  $B'$  and  $B''$  in  $\mathcal{P}$   $O(MN \log N)?$ 
15
16         for  $b \in A$ 
17         do if  $(B, b) \in S$ 
18             then REPLACE  $(B, b)$  by  $(B', b)$  and  $(B'', b)$  in  $S$ 
19             else INSERT ( $\min(B', B''), b$ ) to  $S$   $O(MN)?$ 
```

# Analysis of Time Complexity

## Proposition

*The number of pairs  $(C, a)$  inserted into  $S$  is at most  $2MN$ .*

## Lemma

*The number of classes created during the execution is at most  $2N - 1$ .*

## Corollary

*The number of iterations of the while loop is at most  $2MN$ .*

# Analysis of Time Complexity

Why does **Part 3** (i.e. bottleneck) have global execution time  $O(MN \log_2 N)$ ?

We take the following two steps:

- Evaluate  $\sum |\text{INVERSE}|$   
(i.e. the sum of  $|\text{INVERSE}|$  in all executions of the *while* loop)
- Ensure the global execution time to be  $O(\sum |\text{INVERSE}|)$

## Proposition

*When  $a \in A$  and  $p \in Q$  are fixed, the number of  $(C, a)$  being **deleted** from list  $S$ , such that  $p \in C$ , is bounded by  $\log_2 N$ .*

## Corollary

$\sum |\text{INVERSE}|$  is bounded by  $MN \log_2 N$ .

# Bottleneck Analysis - Implementation - Part 3.1

ALGORITHM3()

```
1  create an empty list Involved
2  for all  $y \in C$  and all  $x \in \text{Inv}[y, a]$ 
3  do  $i \leftarrow \text{Class}[x]$ 
4     if  $\text{Size}[i] = 0$ 
5         then  $\text{Size}[i] \leftarrow 1$ 
6             insert  $i$  to Involved
7         else  $\text{Size}[i] \leftarrow \text{Size}[i] + 1$ 
8
9
10
11
12
13
14
15
16
17
18
19
```

# Bottleneck Analysis - Implementation - Part 3.2

ALGORITHM3()

```
1  create an empty list Involved
2  for all  $y \in C$  and all  $x \in \text{Inv}[y, a]$ 
3  do  $i \leftarrow \text{Class}[x]$ 
4     if  $\text{Size}[i] = 0$ 
5         then  $\text{Size}[i] \leftarrow 1$ 
6             insert  $i$  to Involved
7         else  $\text{Size}[i] \leftarrow \text{Size}[i] + 1$ 
8
9  for all  $y \in C$  and all  $x \in \text{Inv}[y, a]$ 
10 do  $i \leftarrow \text{Class}[x]$ 
11    if  $\text{Size}[i] < \text{Card}[i]$ 
12        then if  $\text{Twin}[i] = 0$ 
13            then  $\text{Counter} \leftarrow \text{Counter} + 1$ 
14                 $\text{Twin}[i] \leftarrow \text{Counter}$ 
15                delete  $x$  from class  $i$  and insert  $x$  into class  $\text{Twin}[i]$ 
16
17
18
19
```

# Bottleneck Analysis - Implementation - Part 3.3

ALGORITHM3()

```
1  create an empty list Involved
2  for all  $y \in C$  and all  $x \in \text{Inv}[y, a]$ 
3  do  $i \leftarrow \text{Class}[x]$ 
4     if  $\text{Size}[i] = 0$ 
5         then  $\text{Size}[i] \leftarrow 1$ 
6             insert  $i$  to Involved
7         else  $\text{Size}[i] \leftarrow \text{Size}[i] + 1$ 
8
9  for all  $y \in C$  and all  $x \in \text{Inv}[y, a]$ 
10 do  $i \leftarrow \text{Class}[x]$ 
11     if  $\text{Size}[i] < \text{Card}[i]$ 
12         then if  $\text{Twin}[i] = 0$ 
13             then  $\text{Counter} \leftarrow \text{Counter} + 1$ 
14                  $\text{Twin}[i] \leftarrow \text{Counter}$ 
15                 delete  $x$  from class  $i$  and insert  $x$  into class  $\text{Twin}[i]$ 
16
17 for all  $j \in \text{Involved}$ 
18 do  $\text{Size}[j] \leftarrow 0$ 
19      $\text{Twin}[j] \leftarrow 0$ 
```

The global execution time of **Part 3** is  $O(MN\log_2N)$ .

So we have:

## Theorem

*Let  $A$  be an alphabet of  $M$  letters and let  $\mathcal{A}$  be a DFA on this alphabet with  $N$  states, then Hopcroft's algorithm using previous data structures is in time  $O(MN\log_2N)$  in the worst case.*

*Remark:* This bound can also be proved to be tight.