

## Warm Up Problem

- What is the definition of a NFA? (Try it without looking!)
- Write an NFA over  $\Sigma = \{a, b, c\}$  that accepts  $L = \{abc\} \cup \{w : w \text{ ends with } cc\}$ .

## CS 241 Lecture 8

Non-Deterministic Finite Automata with  $\epsilon$ -transitions

With thanks to Brad Lushman, Troy Vasiga and Kevin Lanctot

# Non-Deterministic Finite Automata

The above idea can be mathematically described as follows:

## Definition

An **NFA** is a 5-tuple  $(\Sigma, Q, q_0, A, \delta)$ :

- $\Sigma$  is a finite non-empty set (alphabet).
- $Q$  is a finite non-empty set of states.
- $q_0 \in Q$  is a start state
- $A \subseteq Q$  is a set of accepting states
- $\delta : (Q \times \Sigma) \rightarrow 2^Q$  is our [total] transition function. Note that  $2^Q$  denotes the *power set* of  $Q$ , that is, the set of all subsets of  $Q$ . This allows us to go to multiple states at once!

## Extending $\delta$ For an NFA

Again we can extend the definition of  $\delta : (Q \times \Sigma) \rightarrow 2^Q$  to a function  $\delta^* : (2^Q \times \Sigma^*) \rightarrow 2^Q$  via:

$$\begin{aligned}\delta^* : (2^Q \times \Sigma^*) &\rightarrow 2^Q \\ (S, \epsilon) &\mapsto S \\ (S, aw) &\mapsto \delta^* \left( \bigcup_{q \in S} \delta(q, a), w \right)\end{aligned}$$

where  $a \in \Sigma$ . Analogously, we also have:

### Definition

An NFA given by  $M = (\Sigma, Q, q_0, A, \delta)$  **accepts a string**  $w$  if and only if  $\delta^*({q_0}, w) \cap A \neq \emptyset$ .

# Simulating an NFA

---

**Algorithm 1** Algorithm to Simulate an NFA

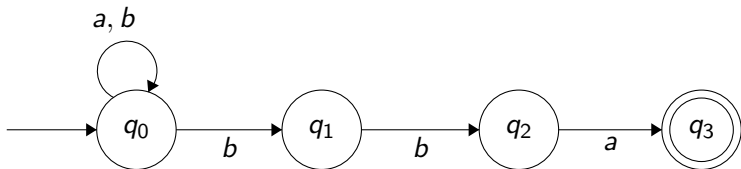
---

```
1:  $S = \{q_0\}$ 
2: while not EOF do
3:    $c = \text{read\_char}()$ 
4:    $S = \bigcup_{q \in S} \delta(q, c)$ 
5: end while
6: if  $S \cap A \neq \emptyset$  then
7:   Accept
8: else
9:   Reject
10: end if
```

---

## Practice Simulating $w = abbba$

Example:  $\Sigma = \{a, b\}$ ,  $L = \{w : w \text{ ends with } bba\}$



Processed	Remaining	S
$\epsilon$	$abbba$	$\{q_0\}$
$a$	$bbba$	$\{q_0\}$
$ab$	$bba$	$\{q_0, q_1\}$
$abb$	$ba$	$\{q_0, q_1, q_2\}$
$abbba$	$a$	$\{q_0, q_1, q_2\}$
$abbba$	$\epsilon$	$\{q_0, q_3\}$

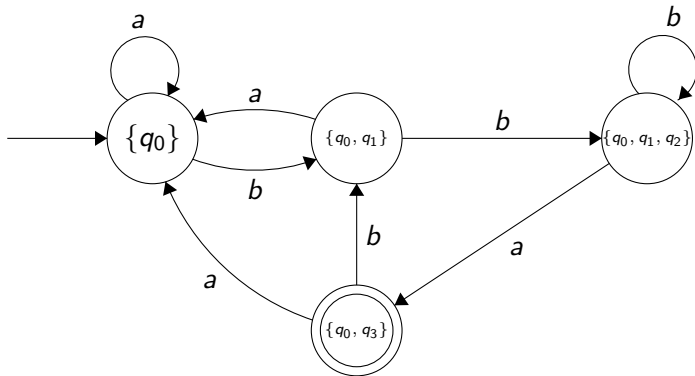
Since  $\{q_0, q_3\} \cap \{q_3\} \neq \emptyset$ , accept.

## NFA to DFA

To convert an NFA to a DFA, one could write down all of the  $2^Q$  possible states and then connect them one by one based on  $\delta$  and each letter in  $\Sigma$ . This however leads to a lot of extra states and a lot of unnecessary work. Instead,

- Start with the state  $S = \{q_0\}$
- From this state, go to the NFA and determine what happens on each  $a \in \Sigma$  for each  $q \in S$ . The set of resulting states should become its own state in your DFA.
- Repeat the previous step for each new state created until you have exhausted every possibility.
- Accepting states are any states that included an accepting state of the original NFA.

## Previous NFA as a DFA





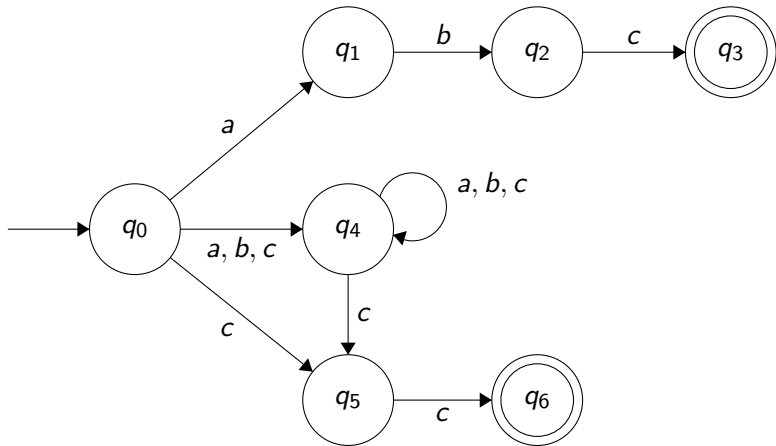
## More Examples In Class

Let  $\Sigma = \{a, b, c\}$ . Write an NFA and the associated DFA for the following examples:

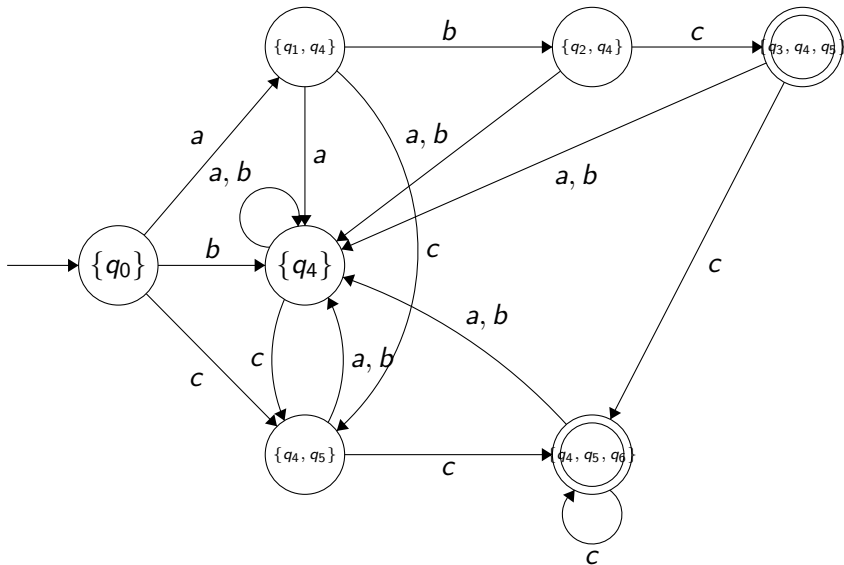
- $L = \{abc\} \cup \{w : w \text{ ends with } cc\}$
- $L = \{abc\} \cup \{w : w \text{ contains } cc\}$
- $L = \{w : w \text{ contains } cab\} \cup \{w : w \text{ contains an even number of } bs\}$
- $L = \{w : w \text{ contains exactly one } abb\} \cup \{w : w \text{ does not contain } ac\}$

## First Example (Non-Optimal)

$$L = \{abc\} \cup \{w : w \text{ ends with } cc\}$$

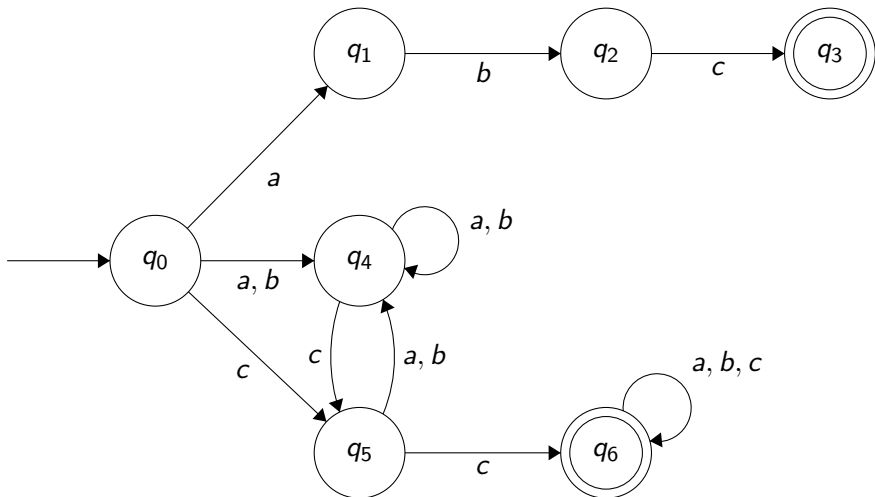


# DFA Associated with Previous Example

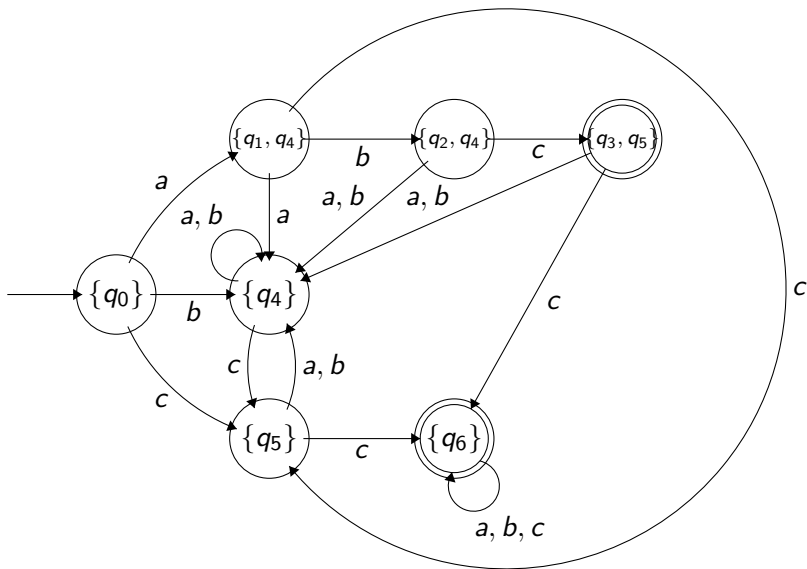


## Second Example (Non-Optimal)

$L = \{abc\} \cup \{w : w \text{ contains a copy of } cc\}$



# DFA Associated with Previous Example

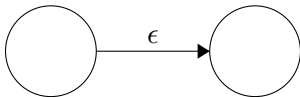


# Summary

- From Kleene's theorem, the set of languages accepted by a DFA are the regular languages
- The set of languages accepted by DFAs are the same as those accepted by NFAs
- Therefore, the set of languages accepted by an NFA are precisely the regular languages!

## One More Generalization

- We gained no new computing powers from an NFA.
- What about if we permitted state changes without reading a character?
- These are known as  $\epsilon$  transitions:



# $\epsilon$ -Non-Deterministic Finite Automata

## Definition

An  $\epsilon$ -NFA is a 5-tuple  $(\Sigma, Q, q_0, A, \delta)$ :

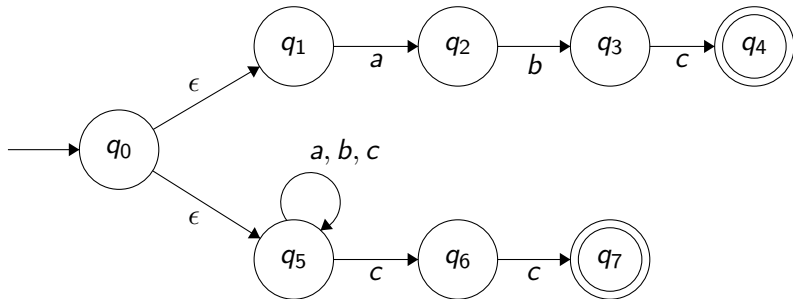
- $\Sigma$  is a finite non-empty set (alphabet) **that does not contain the symbol  $\epsilon$** .
- $Q$  is a finite non-empty set of states.
- $q_0 \in Q$  is a start state
- $A \subseteq Q$  is a set of accepting states
- $\delta : (Q \times \Sigma \cup \{\epsilon\}) \rightarrow 2^Q$  is our [total] transition function. Note that  $2^Q$  denotes the *power set* of  $Q$ , that is, the set of all subsets of  $Q$ . This allows us to go to multiple states at once!



## Why Do This?

These  $\epsilon$ -transitions make it trivial to take the union of two NFAs:

Example:  $L = \{abc\} \cup \{w : w \text{ ends with } cc\}$



## Simulating an $\epsilon$ -NFA

Define  $E(S)$  to be the epsilon closure of a set of states  $S$ , that is, the set of all states reachable from  $S$  in 0 or more  $\epsilon$ -transitions.

Note this implies that  $S \subset E(S)$ .

---

### Algorithm 2 Algorithm to Simulate an $\epsilon$ -NFA

---

```
1:  $S = E(\{q_0\})$ 
2: while not EOF do
3:    $c = \text{read\_char}()$ 
4:    $S = E(\cup_{q \in S} \delta(q, c))$ 
5: end while
6: if  $S \cap A \neq \emptyset$  then
7:   Accept
8: else
9:   Reject
10: end if
```

---

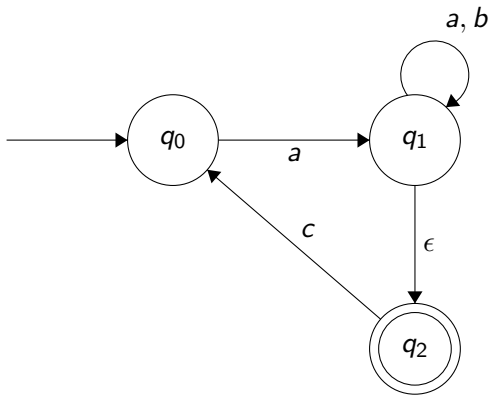
## Tracing a Word with Previous Example

Processed	Remaining	S
$\epsilon$	<i>abcaccc</i>	$\{q_0, q_1, q_5\}$
<i>a</i>	<i>bcaccc</i>	$\{q_2, q_5\}$
<i>ab</i>	<i>caccc</i>	$\{q_3, q_5\}$
<i>abc</i>	<i>accc</i>	$\{q_4, q_5, q_6\}$
<i>abca</i>	<i>ccc</i>	$\{q_5\}$
<i>abcac</i>	<i>cc</i>	$\{q_5, q_6\}$
<i>abcacc</i>	<i>c</i>	$\{q_5, q_6, q_7\}$
<i>abcaccc</i>	$\epsilon$	$\{q_5, q_6, q_7\}$

Since  $\{q_5, q_6, q_7\} \cap \{q_4, q_7\} \neq \emptyset$ , accept.

## Second Example

Exercise: What regular language does this machine represent?



## Second Tracing Example

Processed	Remaining	$S$
$\epsilon$	$abca$	$\{q_0\}$
$a$	$bca$	$\{q_1, q_2\}$
$ab$	$ca$	$\{q_1, q_2\}$
$abc$	$a$	$\{q_0\}$
$abca$	$\epsilon$	$\{q_1, q_2\}$

Since  $\{q_1, q_2\} \cap \{q_2\} \neq \emptyset$ , accept.

## Third Tracing Example

Processed	Remaining	S
$\epsilon$	cc	$\{q_0\}$
c	c	$\{\}$
cc	$\epsilon$	$\{\}$

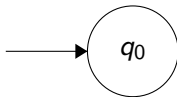
Since  $\{\} \cap \{q_2\} = \emptyset$ , reject.

# Equivalences

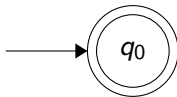
- Using the same technique as for an NFA, every  $\epsilon$ -NFA has a corresponding DFA.
- In both cases, this technique can be automated! It does not necessarily give the smallest DFA but it will give a DFA that is valid.
- This combined with Kleene's Theorem implies that every language recognized by an  $\epsilon$ -NFA is regular.
- In fact, we can show that an  $\epsilon$ -NFA exists for every regular expression, then we have proved one direction of Kleene's Theorem (See next slides). We do this by **structural induction**.

# $\epsilon$ -NFAs that Recognize Regular Languages

1.  $\emptyset$



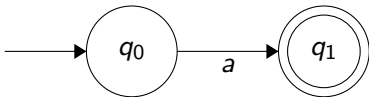
2.  $\{\epsilon\}$





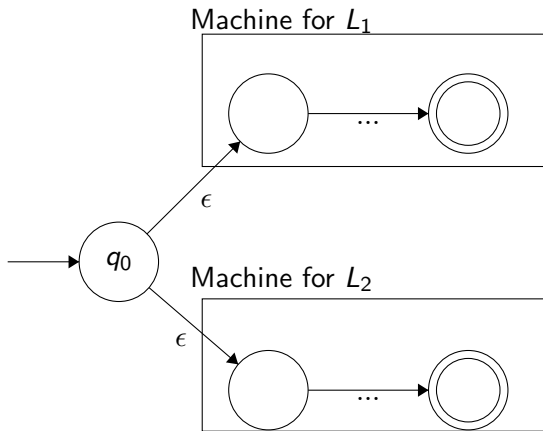
## $\epsilon$ -NFAs that Recognize Regular Languages

3.  $\{a\}$



## $\epsilon$ -NFAs that Recognize Regular Languages

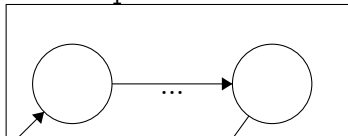
4.  $L_1 \cup L_2$  (that is, given  $\epsilon$ -NFAs that recognize  $L_1$  and  $L_2$  already, construct one for this language)



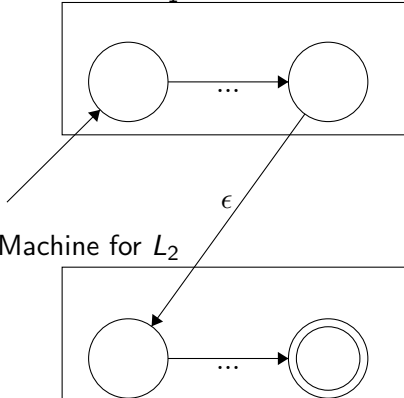
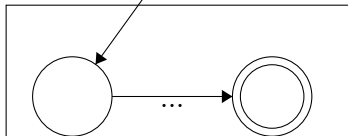
## $\epsilon$ -NFAs that Recognize Regular Languages

5.  $L_1L_2$  (that is, given  $\epsilon$ -NFAs that recognize  $L_1$  and  $L_2$  already, construct one for this language). In the diagram below, change all accepting states in the  $\epsilon$ -NFA for  $L_1$  to non-accepting states and for each such state, add an  $\epsilon$  transition to the start state of  $L_2$ .

Machine for  $L_1$

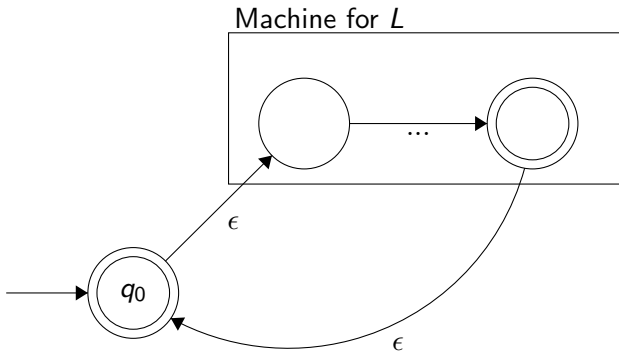


Machine for  $L_2$



## $\epsilon$ -NFAs that Recognize Regular Languages

6.  $L^*$ . Assume we have a  $\epsilon$ -NFA for  $L$  already. Below, from each accepting state, add an  $\epsilon$  transition back to the newly created start state.



## Summary of Previous Slides

- For each regular language, we can construct an  $\epsilon$ -NFA that recognizes the language.
- We can also convert each  $\epsilon$ -NFA into a DFA.
- Both of the above processes can be automated.
- There are several ways to go from a DFA to a regular language which we will not discuss here (see the **State Removal Method** for example).

# Scanning

Is C a regular language? The following are regular:

- C keywords
- C identifiers
- C literals
- C operators
- C comments

Sequences of these are hence also regular. Finite automata can do our tokenization (ie. our scanning).

What about punctuation? Even simpler, set  $\Sigma = \{ ( , ) \}$  and  $L = \{ \text{set of balanced parentheses} \}$ . Is  $L$  regular?

# Scanning

Is C a regular language? The following are regular:

- C keywords
- C identifiers
- C literals
- C operators
- C comments

Sequences of these are hence also regular. Finite automata can do our tokenization (ie. our scanning).

What about punctuation? Even simpler, set  $\Sigma = \{ ( , ) \}$  and  $L = \{ \text{set of balanced parentheses} \}$ . Is  $L$  regular? More on this later.

## Scanning Continued

How does our scanner work?

- Recall our goal is given some text, break the text up into tokens (eg. (ID, div), (REG, \$1), (COMMA, ','), (REG, \$2))
- Problem: Some tokens can be recognized in multiple different ways!
- Eg. 0x12cc. Could be a single HEXINT or could be an INT (0) followed by an ID (x) followed by another INT (1) followed by another INT (2) and followed by an ID (cc). (There are lots of other interpretations).

Which interpretation is correct?



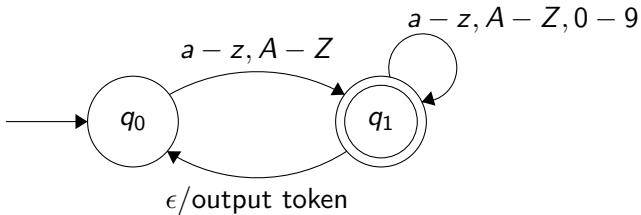
## Formalization of our Problem

Given a regular language  $L$  (say  $L$  is all valid MIPS or C tokens), determine if a given word  $w$  is in  $LL^*$  (or in other words, is  $w \in L^* \setminus \{\epsilon\}$ ?)

(We don't want to consider the empty program as being valid hence why we drop  $\epsilon$ ).

## Concrete Example for C

Consider the language  $L$  of just ID tokens in C:



With the input  $abcde$ , this could be considered as anywhere from 1 to 5 different tokens! What do we do?

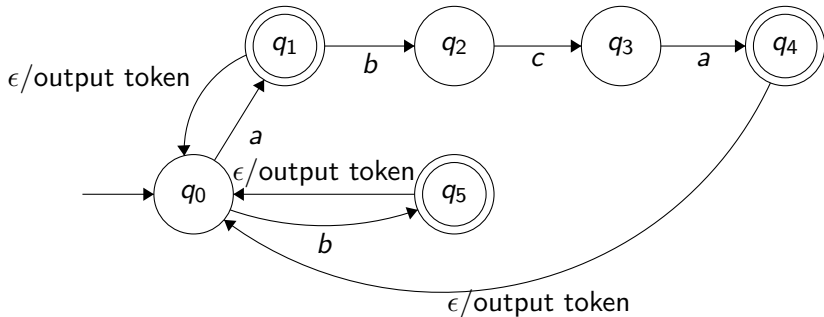
# What do to?

- We will discuss two algorithms called *maximal munch* and *simplified maximal munch*.
- General idea: Consume the largest possible token that makes sense. Produce the token and then proceed.
- Difference:
  - Maximal Munch: Consume characters until you no longer have a valid transition. If you have characters left to consume, backtrack to the last valid accepting state and resume.
  - Simplified Maximal Munch: Consume characters until you no longer have a valid transition. If you are currently in an accepting state, produce the token and proceed. Otherwise go to an error state.

## DFA for next two slides

$\Sigma = \{a, b, c\}$ ,  $L = \{a, b, abca\}$ , consider  $w = ababca$ .

Note that  $w \in LL^*$ . What follows is an  $\epsilon$ -NFA for  $LL^*$  based on our algorithm:



## Examples

**Maximal Munch:**  $\Sigma = \{a, b, c\}$ ,  $L = \{a, b, abca\}$ ,  $w = ababca$ .

Note that  $w \in LL^*$ .

- Algorithm consumes  $a$  **and flags this state as it is accepting**, then  $b$  then tries to consume  $a$  but ends up in an error state.
- Algorithm then backtracks to first  $a$  since that was the last accepting state. Token  $a$  is output.
- Algorithm then resumes consuming  $b$  **and flags this state as it is accepting**, then tries to consume  $a$  but ends up in an error state.
- Algorithm then backtracks to first  $b$  since that was the last accepting state. Token  $b$  is output.
- Algorithm then consumes the second  $a$ , the second  $b$ , the first  $c$ , the third  $a$  and runs out of input. This last state is accepting so output our last token  $abca$  and accept.

# Examples

**Simplified Maximal Munch:**  $\Sigma = \{a, b, c\}$ ,  $L = \{a, b, abba\}$ ,  
 $w = ababca$ . Note that  $w \in LL^*$ .

- Algorithm consumes  $a$ , then  $b$  then tries to consume  $a$  but ends up in an error state. **Note there is no keeping track of the first accepting state!**
- Algorithm then checks to see if  $ab$  is accepting. It is not (as  $ab \notin L$ ).
- Algorithm rejects  $ababca$ .
- Note: This gave the wrong answer! However this algorithm is usually good enough and is used in practice.

## Practical Implications

Consider the following C++ line:

```
vector<pair<string,int>> v;
```

Notice that at the end, there is the token >>! This on its own is a valid token! With Simplified Maximal Munch, we would actually reject this declaration. This would have been the case in versions of C++ compilers before C++11. How did they avoid this problem?

## Practical Implications

Consider the following C++ line:

```
vector<pair<string,int>> v;
```

Notice that at the end, there is the token >>! This on its own is a valid token! With Simplified Maximal Munch, we would actually reject this declaration. This would have been the case in versions of C++ compilers before C++11. How did they avoid this problem? Requiring an extra space after the first > character!

```
vector<pair<string,int> > v;
```

The above code would be valid in C++03 or C++11 (for example).