



<http://algs4.cs.princeton.edu>

## 5.3 SUBSTRING SEARCH

---

- ▶ *introduction*
- ▶ *brute force*
- ▶ *Knuth–Morris–Pratt*



<http://algs4.cs.princeton.edu>

## 5.3 SUBSTRING SEARCH

---

- ▶ *introduction*
- ▶ *brute force*
- ▶ *Knuth–Morris–Pratt*

# Substring search

---

**Goal.** Find pattern of length  $M$  in a text of length  $N$ .

typically  $N \gg M$

*pattern* → N E E D L E

*text* → I N A H A Y S T A C K N E E D L E I N A

↑  
*match*

# Substring search applications

---

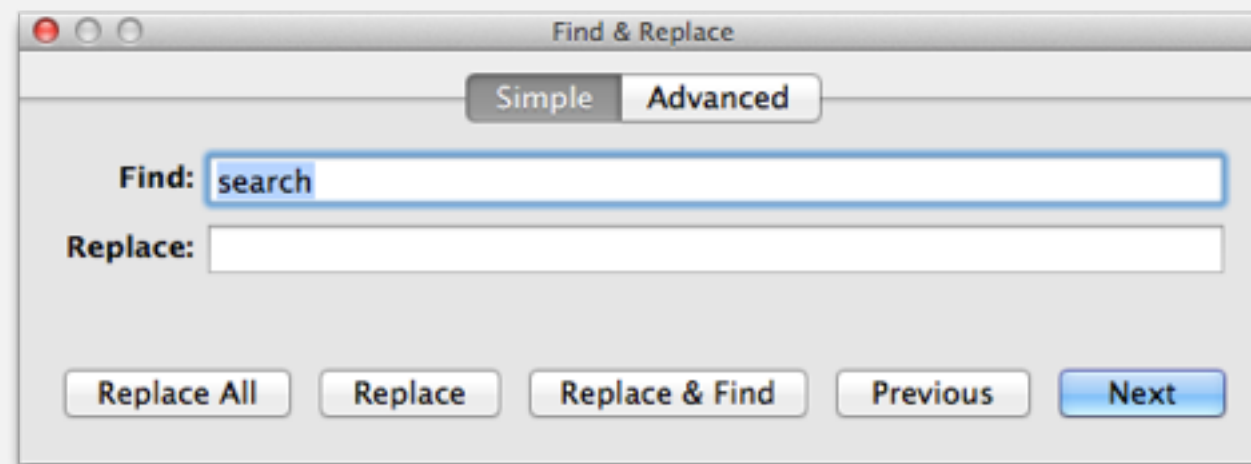
**Goal.** Find pattern of length  $M$  in a text of length  $N$ .

typically  $N \gg M$

*pattern* → N E E D L E

*text* → I N A H A Y S T A C K N E E D L E I N A

match



# Substring search applications

**Goal.** Find pattern of length  $M$  in a text of length  $N$ .

typically  $N \gg M$

*pattern* → N E E D L E

*text* → I N A H A Y S T A C K N E E D L E I N A

match

**Computer forensics.** Search memory or disk for signatures, e.g., all URLs or RSA keys that the user has entered.



<http://citp.princeton.edu/memory>

# Substring search applications

**Goal.** Find pattern of length  $M$  in a text of length  $N$ .

typically  $N \gg M$

*pattern* → N E E D L E

*text* → I N A H A Y S T A C K N E E D L E I N A

match

Identify patterns indicative of spam.

- PROFITS
- LOSE WE1GHT
- herbal Viagra
- There is no catch.
- This is a one-time mailing.
- This message is sent in compliance with spam regulations.



# Substring search applications

## Electronic surveillance.



Need to monitor all  
internet traffic.  
(security)

No way!  
(privacy)



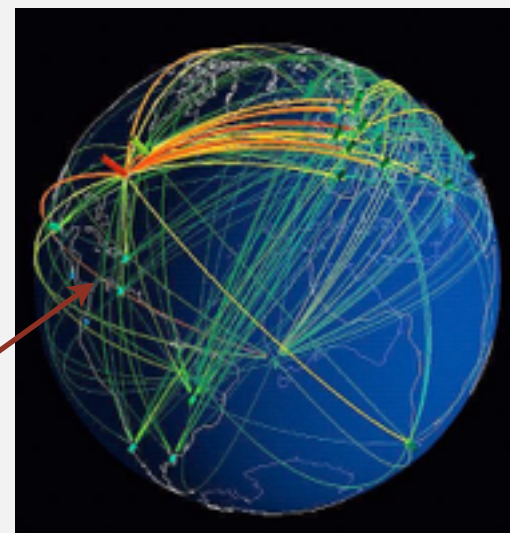
Well, we're mainly  
interested in  
"ATTACK AT DAWN"

OK. Build a  
machine that just  
looks for that.



"ATTACK AT DAWN"  
substring search  
machine

found

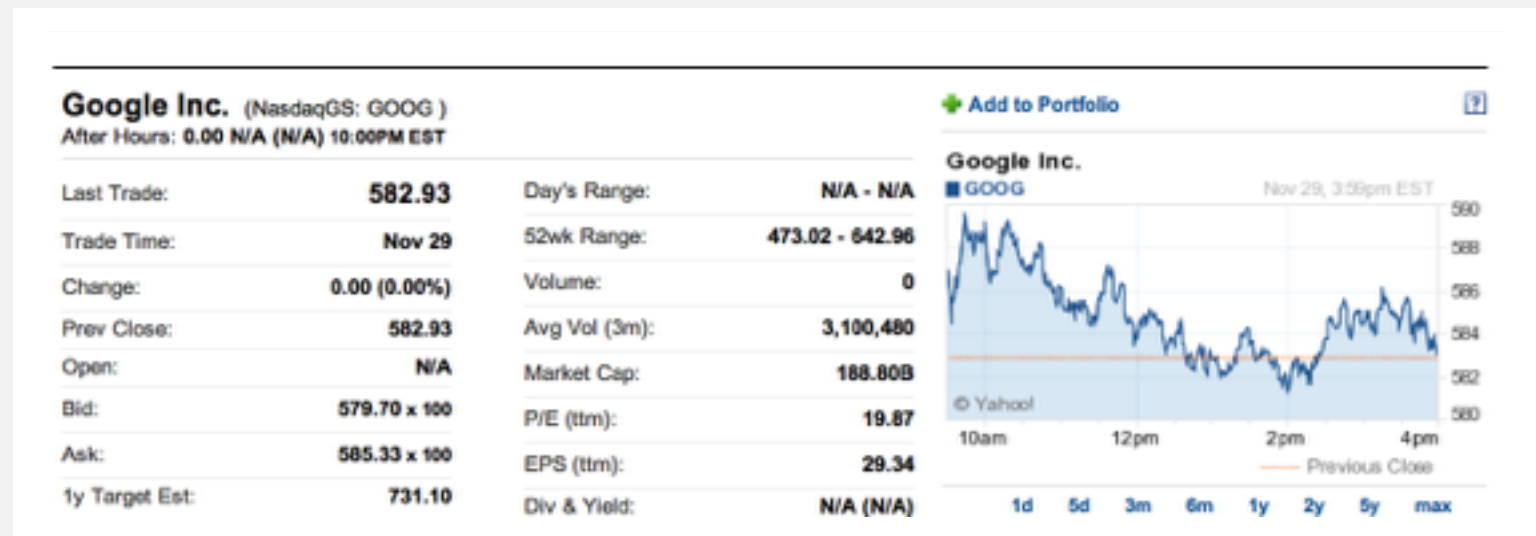




# Substring search applications

**Screen scraping.** Extract relevant data from web page.

**Ex.** Find string delimited by `<b>` and `</b>` after first occurrence of pattern Last Trade:.



<http://finance.yahoo.com/q?s=goog>

```
...
<tr>
<td class= "yfnc_tablehead1"
width= "48%">
Last Trade:
</td>
<td class= "yfnc_tabledata1">
<big><b>582.93</b></big>
</td></tr>
<td class= "yfnc_tablehead1"
width= "48%">
Trade Time:
</td>
<td class= "yfnc_tabledata1">
...
```



# Screen scraping: Java implementation

---

**Java library.** The `indexOf()` method in Java's `String` data type returns the index of the first occurrence of a given string, starting at a given offset.

```
public class StockQuote
{
    public static void main(String[] args)
    {
        String name = "http://finance.yahoo.com/q?s=";
        In in = new In(name + args[0]);
        String text = in.readAll();
        int start    = text.indexOf("Last Trade:", 0);
        int from     = text.indexOf("<b>", start);
        int to       = text.indexOf("</b>", from);
        String price = text.substring(from + 3, to);
        StdOut.println(price);
    }
}
```

```
% java StockQuote goog
582.93
```

**Caveat.** Must update program if Yahoo format changes.



<http://algs4.cs.princeton.edu>

## 5.3 SUBSTRING SEARCH

---

- ▶ *introduction*
- ▶ *brute force*
- ▶ *Knuth–Morris–Pratt*
- ▶ *Boyer–Moore*
- ▶ *Rabin–Karp*

# Brute-force substring search

Check for pattern starting at each text position.

i	j	i+j	0	1	2	3	4	5	6	7	8	9	10
txt →			A	B	A	C	A	D	A	B	R	A	C
0	2	2	A	B	R	A	← pat						
1	0	1		A	B	R	A	entries in red are mismatches					
2	1	3			A	B	R	A	entries in gray are for reference only				
3	0	3				A	B	R	A	entries in black match the text			
4	1	5					A	B	R	A	entries in gray are for reference only		
5	0	5						A	B	R	A	entries in black match the text	
6	4	10							A	B	R	A	
return i when j is M			match										

# Brute-force substring search: Java implementation

---

Check for pattern starting at each text position.

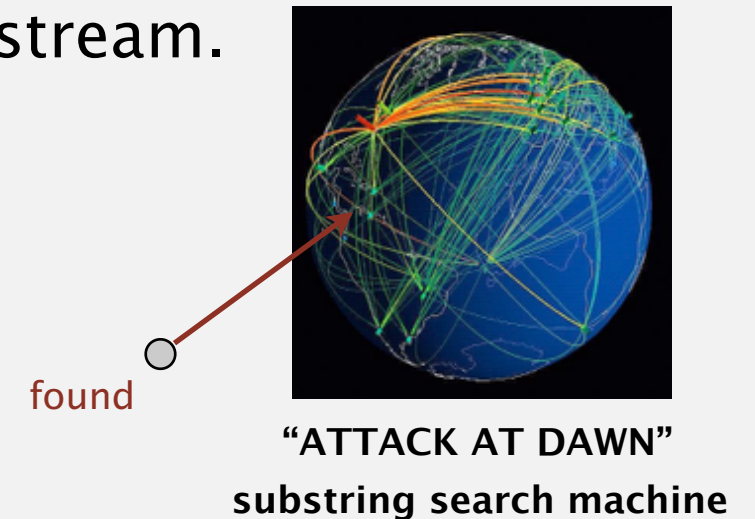
i	j	i + j	0	1	2	3	4	5	6	7	8	9	10
			A	B	A	C	A	D	A	B	R	A	C
4	3	7					A	D	A	C	R		
5	0	5						A	D	A	C	R	

```
public static int search(String pat, String txt)
{
    int M = pat.length();
    int N = txt.length();
    for (int i = 0; i <= N - M; i++)
    {
        int j;
        for (j = 0; j < M; j++)
            if (txt.charAt(i+j) != pat.charAt(j))
                break;
        if (j == M) return i; ← index in text where
                                pattern starts
    }
    return N; ← not found
}
```

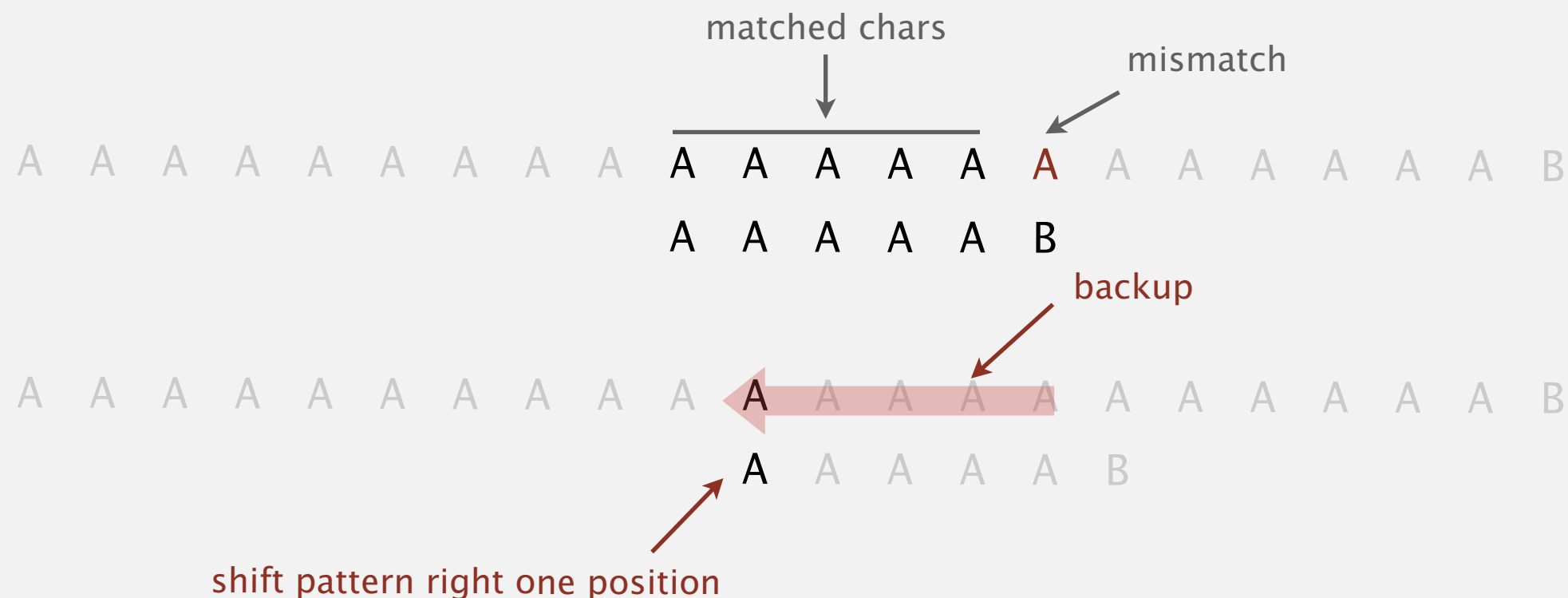
# Backup

In many applications, we want to avoid **backup** in text stream.

- Treat input as stream of data.
- Abstract model: standard input.



Brute-force algorithm needs backup for every mismatch.



**Approach 1.** Maintain buffer of last  $M$  characters.

**Approach 2.** Stay tuned.

# Brute-force substring search: alternate implementation

Same sequence of character compares as previous implementation.

- $i$  points to end of sequence of already-matched characters in text.
- $j$  stores # of already-matched characters (end of sequence in pattern).

$i$	$j$	0	1	2	3	4	5	6	7	8	9	10
		A	B	A	C	A	D	A	B	R	A	C
7	3					A	D	A	C	R		
5	0					A	D	A	C	R		

```
public static int search(String pat, String txt)
{
    int i, N = txt.length();
    int j, M = pat.length();
    for (i = 0, j = 0; i < N && j < M; i++)
    {
        if (txt.charAt(i) == pat.charAt(j)) j++;
        else { i -= j; j = 0; }
    }
    if (j == M) return i - M;
    else return N;
}
```

← explicit backup

# Algorithmic challenges in substring search

---

Brute-force is not always good enough.

**Theoretical challenge.** Linear-time guarantee. ← fundamental algorithmic problem

**Practical challenge.** Avoid backup in text stream. ← often no space (or time) to save text

Now is the time for all people to come to the aid of their party. Now is the time for all good people to come to the aid of their party. Now is the time for many good people to come to the aid of their party. Now is the time for all good people to come to the aid of their party. Now is the time for a lot of good people to come to the aid of their party. Now is the time for all of the good people to come to the aid of their party. Now is the time for all good people to come to the aid of their party. Now is the time for each good person to come to the aid of their party. Now is the time for all good people to come to the aid of their party. Now is the time for all good Republicans to come to the aid of their party. Now is the time for all good people to come to the aid of their party. Now is the time for many or all good people to come to the aid of their party. Now is the time for all good people to come to the aid of their party. Now is the time for all good Democrats to come to the aid of their party. Now is the time for all people to come to the aid of their party. Now is the time for all good people to come to the aid of their party. Now is the time for many good people to come to the aid of their party. Now is the time for all good people to come to the aid of their party. Now is the time for a lot of good people to come to the aid of their party. Now is the time for all of the good people to come to the aid of their party. Now is the time for all good people to come to the aid of their **attack at dawn** party. Now is the time for each person to come to the aid of their party. Now is the time for all good people to come to the aid of their party. Now is the time for all good Republicans to come to the aid of their party. Now is the time for all good people to come to the aid of their party. Now is the time for many or all good people to come to the aid of their party. Now is the time for all good people to come to the aid of their party. Now is the time for all good Democrats to come to the aid of their party.





<http://algs4.cs.princeton.edu>

## 5.3 SUBSTRING SEARCH

---

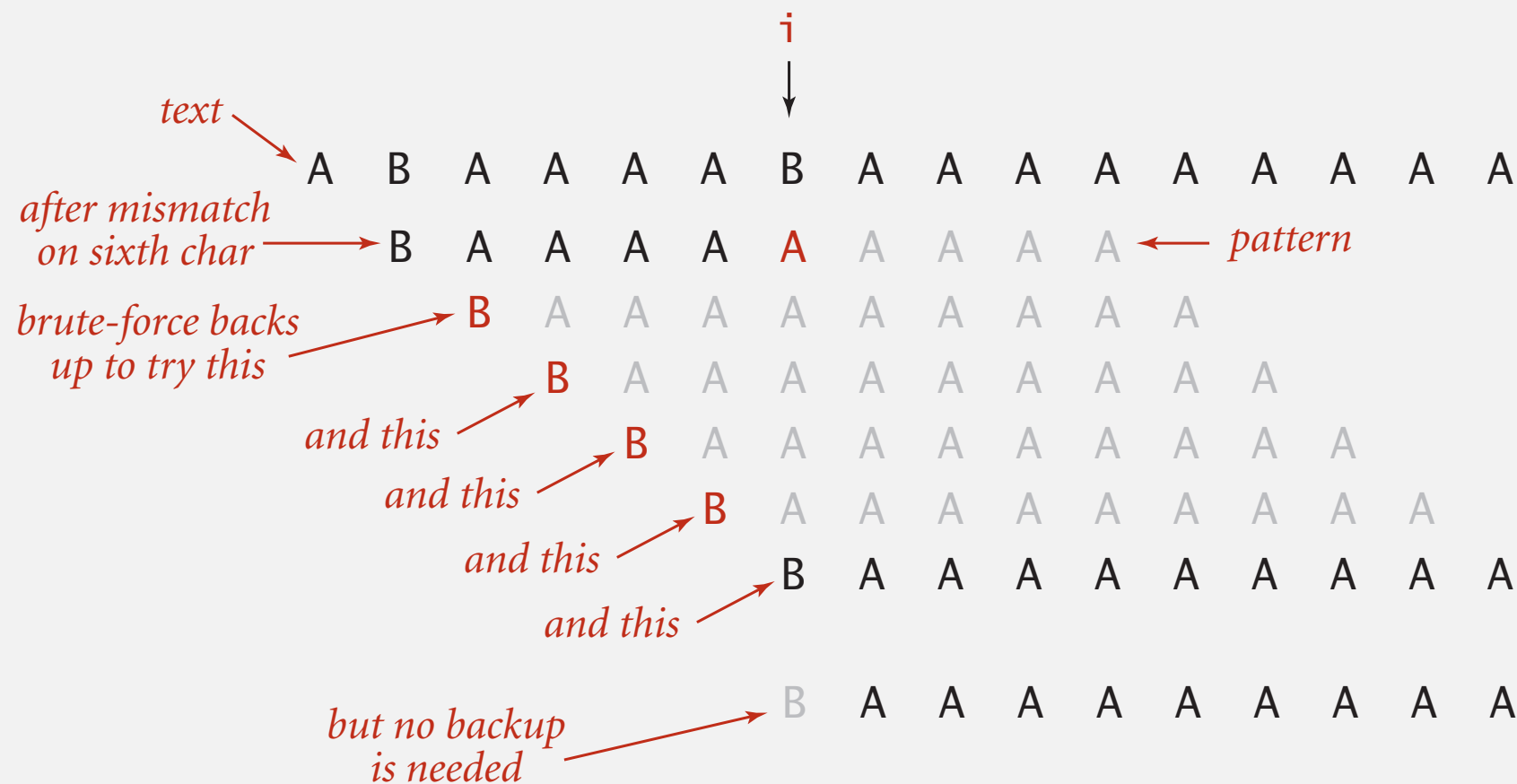
- ▶ *introduction*
- ▶ *brute force*
- ▶ *Knuth–Morris–Pratt*

# Knuth–Morris–Pratt substring search

**Intuition.** Suppose we are searching in text for pattern BAAAAAAAAA.

- Suppose we match 5 chars in pattern, with mismatch on 6<sup>th</sup> char.
- We know previous 6 chars in text are BAAAAB.
- Don't need to back up text pointer!

assuming { A, B } alphabet



**Knuth–Morris–Pratt algorithm.** Clever method to always avoid backup!

# Deterministic finite state automaton (DFA)

DFA is abstract string-searching machine.

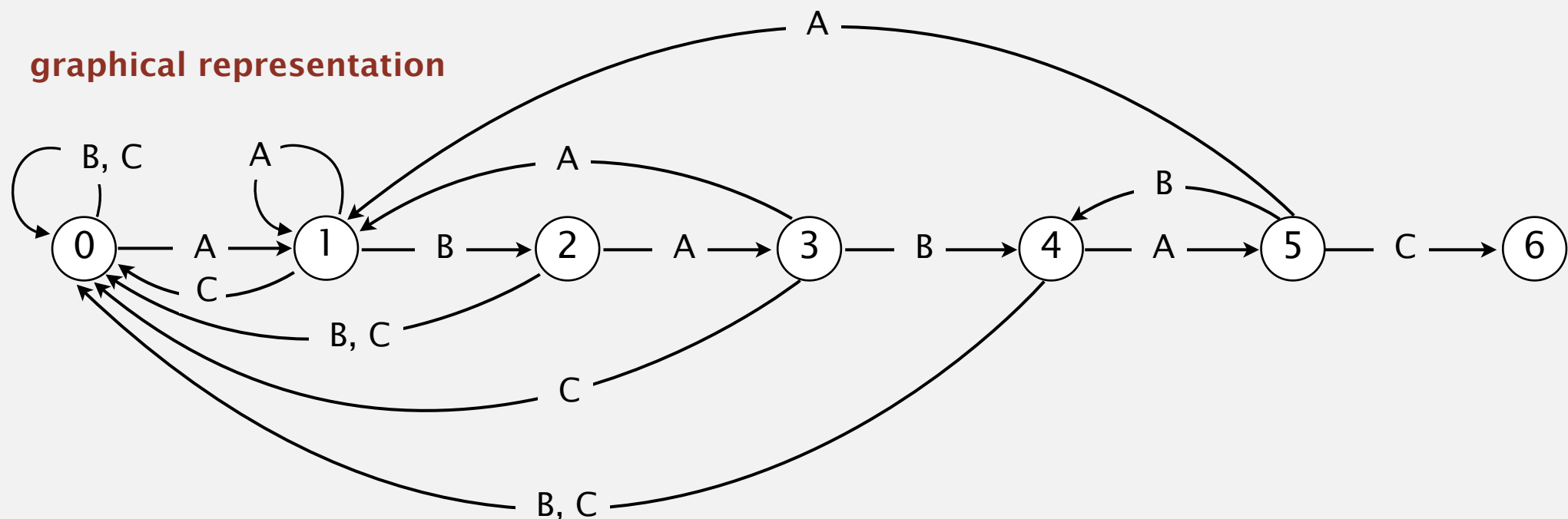
- Finite number of states (including start and halt).
- Exactly one state transition for each char in alphabet.
- Accept if sequence of state transitions leads to halt state.

internal representation

j	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]	A	1	1	3	1	5
	B	0	2	0	4	0
	C	0	0	0	0	6

If in state  $j$  reading char  $c$ :  
if  $j$  is 6 halt and accept  
else move to state  $dfa[c][j]$

graphical representation

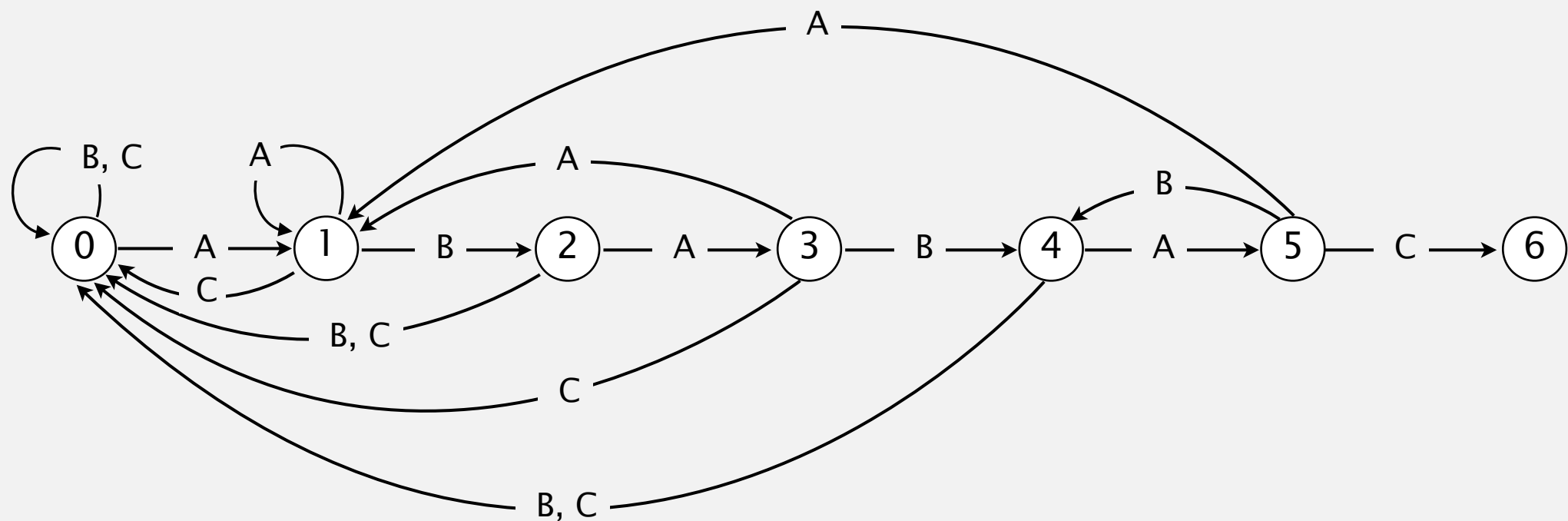


# Knuth–Morris–Pratt demo: DFA simulation

A A B A C A A B A B A C A A



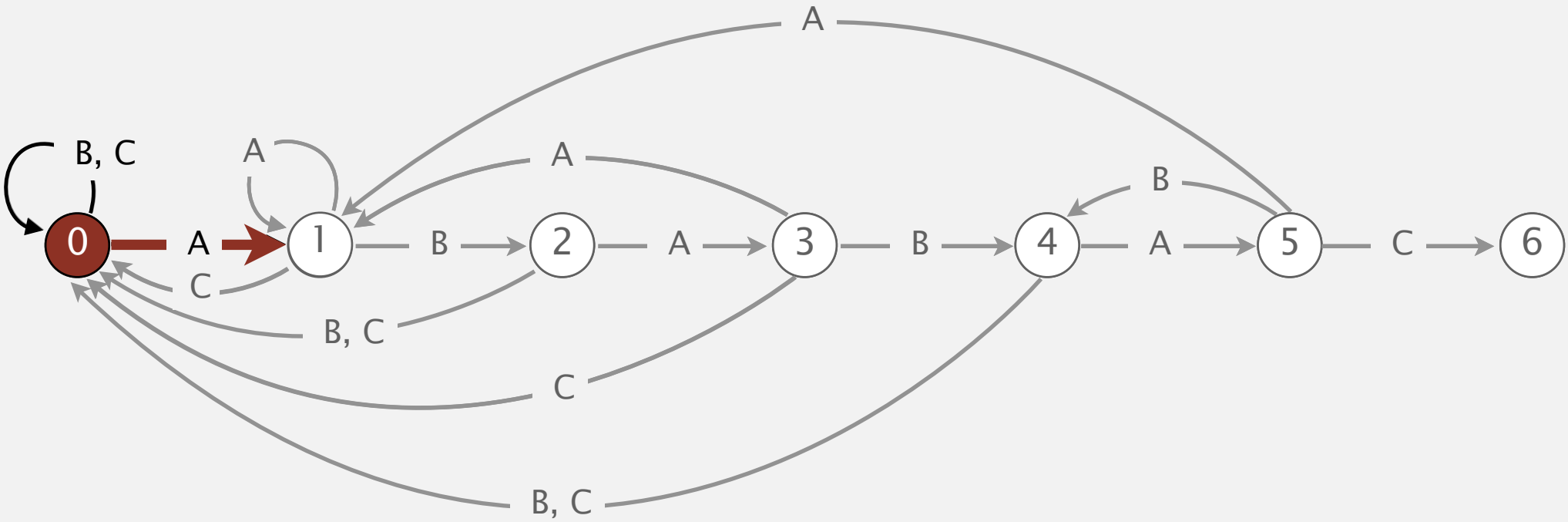
		0	1	2	3	4	5
pat.charAt(j)		A	B	A	B	A	C
dfa[][j]	A	1	1	3	1	5	1
	B	0	2	0	4	0	4
	C	0	0	0	0	0	6



# Knuth–Morris–Pratt demo: DFA simulation

A A B A C A A B A B A C A A

		0	1	2	3	4	5
pat.charAt(j)		A	B	A	B	A	C
dfa[][j]	A	1	1	3	1	5	1
	B	0	2	0	4	0	4
	C	0	0	0	0	0	6

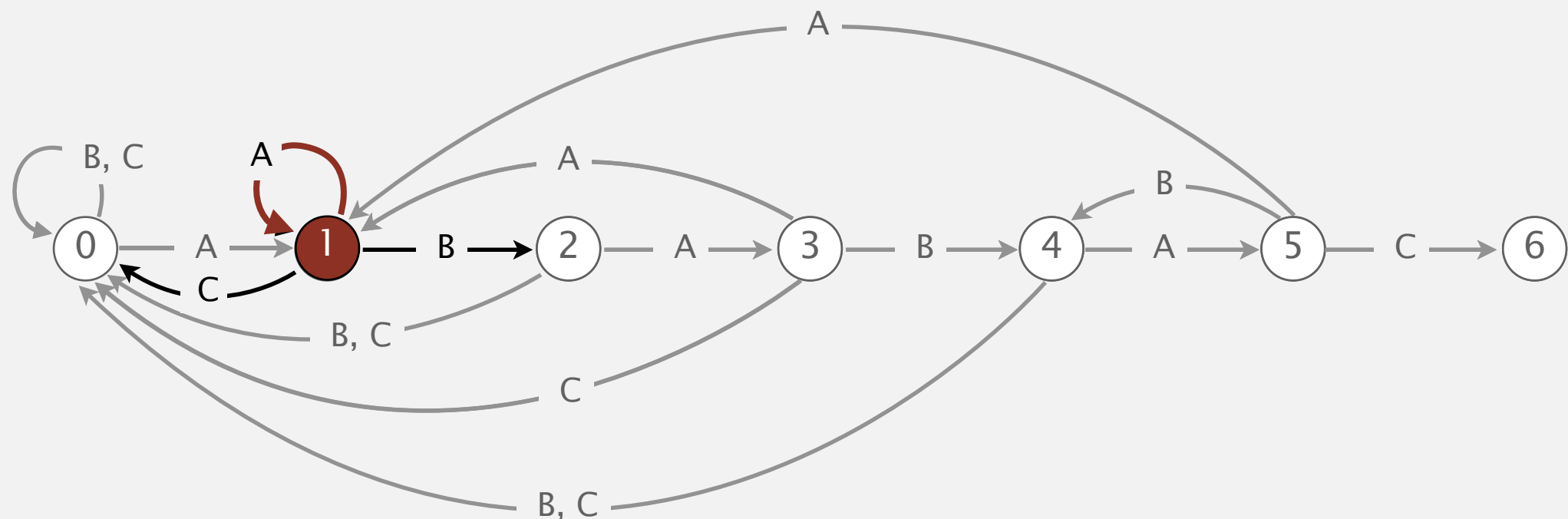


# Knuth–Morris–Pratt demo: DFA simulation

A A B A C A A B A B A C A A

↑

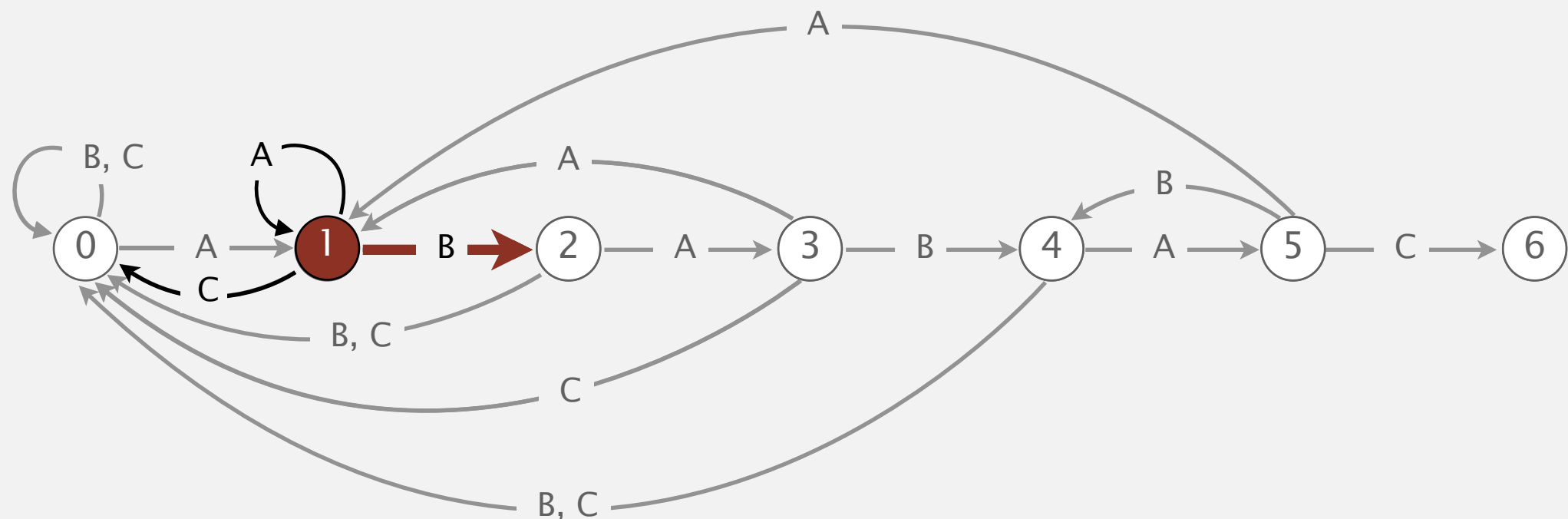
	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]						
A	1	1	3	1	5	1
B	0	2	0	4	0	4
C	0	0	0	0	0	6



# Knuth–Morris–Pratt demo: DFA simulation

A A B A C A A B A B A C A A

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]						
A	1	1	3	1	5	1
B	0	2	0	4	0	4
C	0	0	0	0	0	6



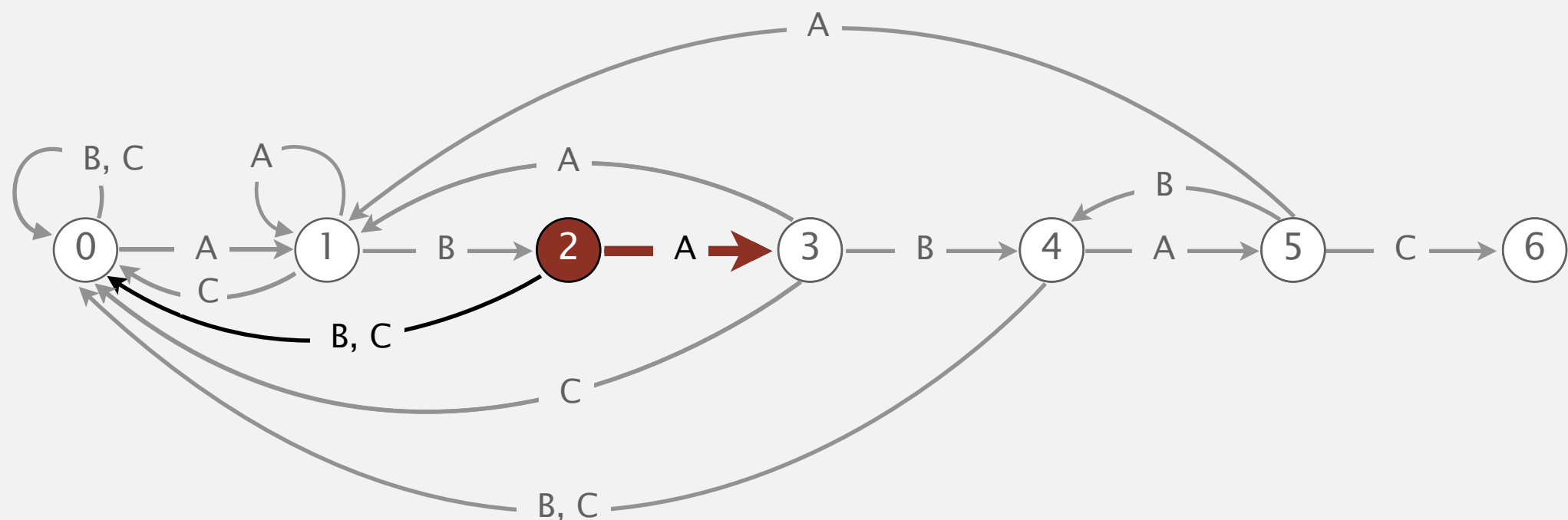


# Knuth–Morris–Pratt demo: DFA simulation

A A B A C A A B A B A C A A

↑

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]						
A	1	1	3	1	5	1
B	0	2	0	4	0	4
C	0	0	0	0	0	6

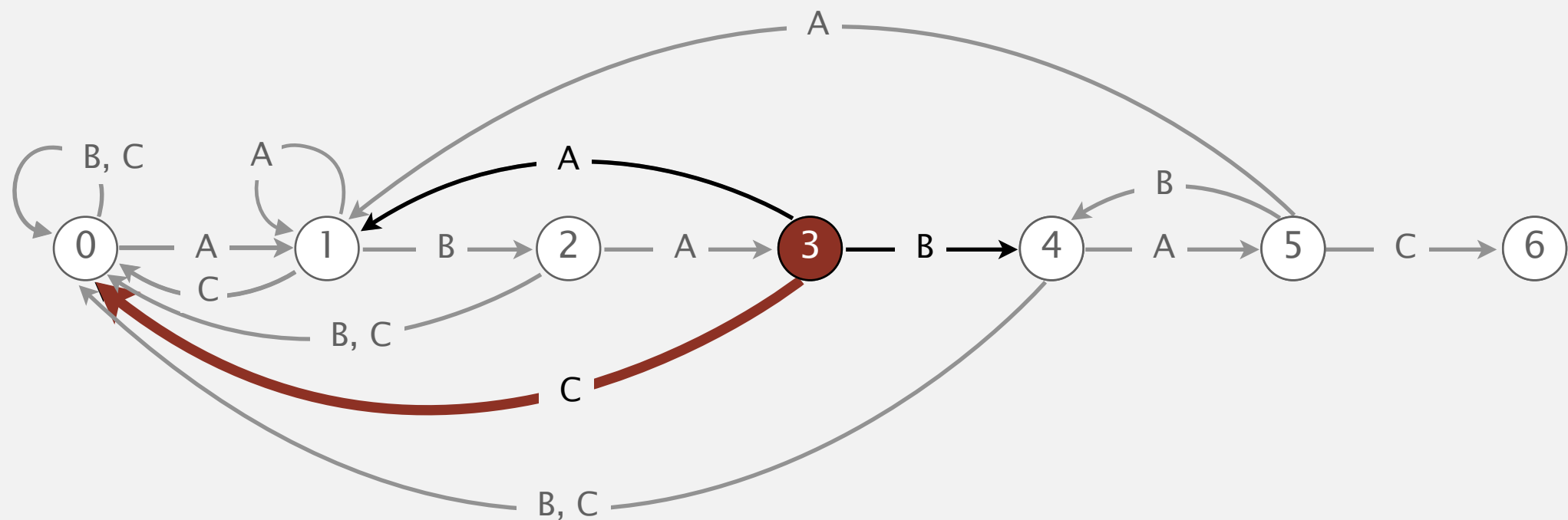


# Knuth–Morris–Pratt demo: DFA simulation

A A B A C A A B A B A C A A

↑

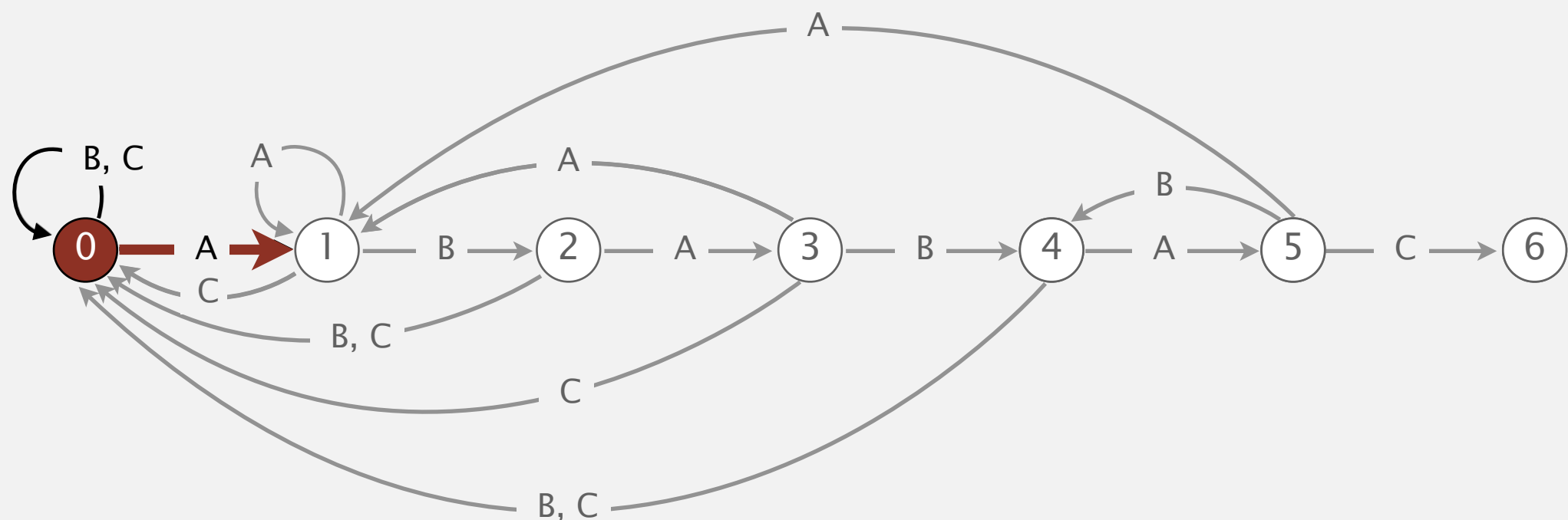
	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]						
A	1	1	3	1	5	1
B	0	2	0	4	0	4
C	0	0	0	0	0	6



# Knuth–Morris–Pratt demo: DFA simulation

A A B A C A A B A B A C A A

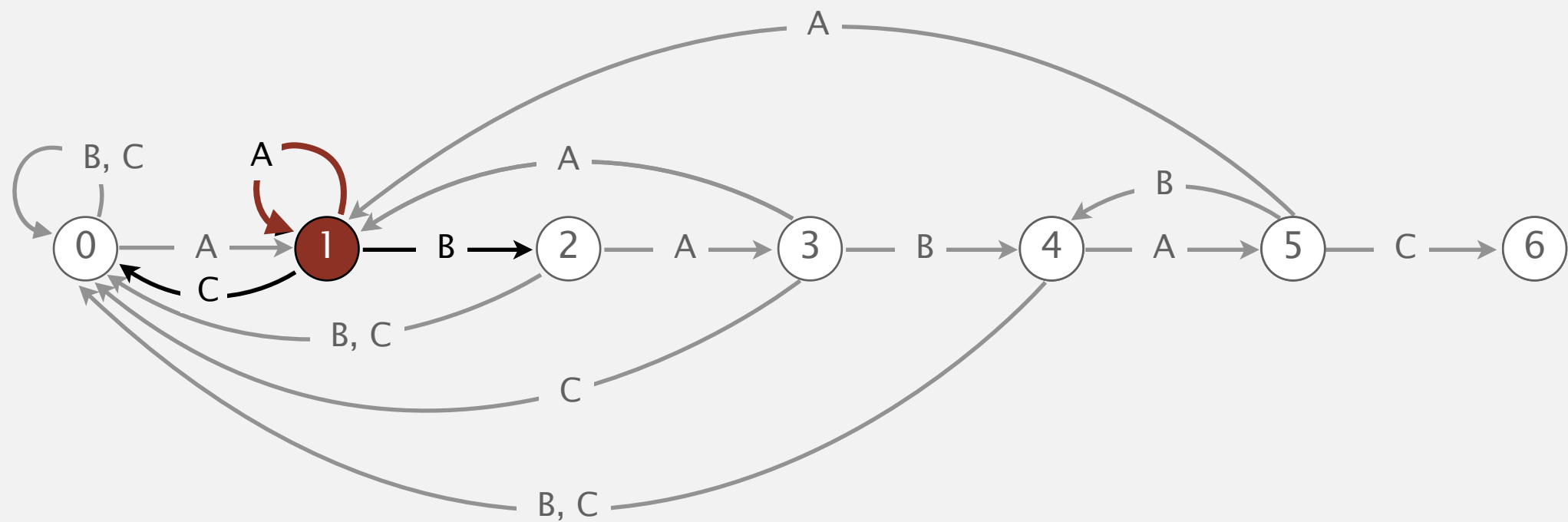
		0	1	2	3	4	5
pat.charAt(j)		A	B	A	B	A	C
dfa[][j]	A	1	1	3	1	5	1
	B	0	2	0	4	0	4
	C	0	0	0	0	0	6



# Knuth–Morris–Pratt demo: DFA simulation

A A B A C **A** A B A B A C A A

	0	<b>1</b>	2	3	4	5
pat.charAt(j)	A	<b>B</b>	A	B	A	C
dfa[][j]						
A	1	<b>1</b>	3	1	5	1
B	0	<b>2</b>	0	4	0	4
C	0	<b>0</b>	0	0	0	6

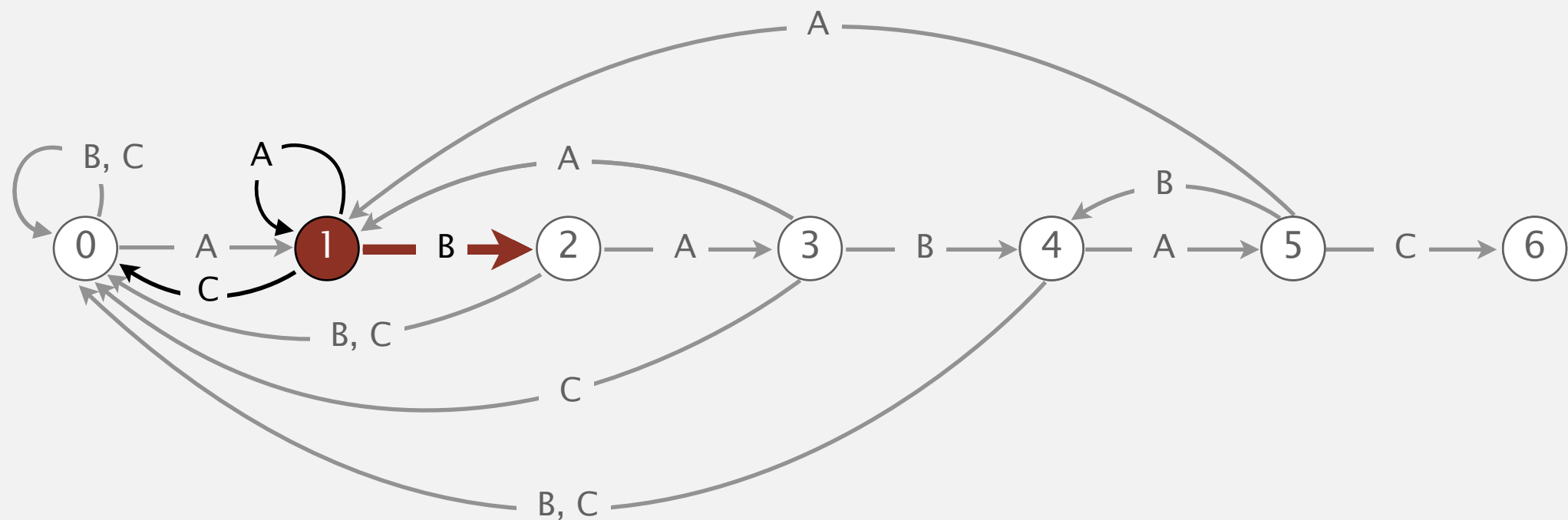


# Knuth–Morris–Pratt demo: DFA simulation

A A B A C A **A** B A B A C A A

↑

	0	<b>1</b>	2	3	4	5
pat.charAt(j)	A	<b>B</b>	A	B	A	C
dfa[][j]						
A	1	<b>1</b>	3	1	5	1
B	0	<b>2</b>	0	4	0	4
C	0	<b>0</b>	0	0	0	6

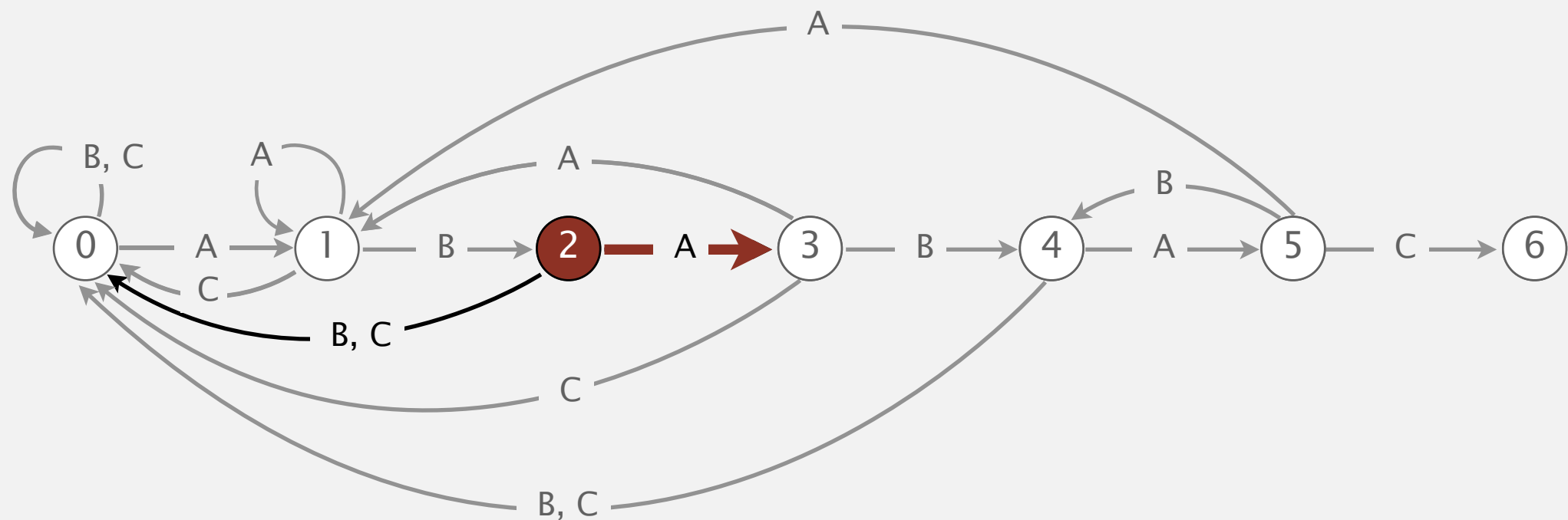


# Knuth–Morris–Pratt demo: DFA simulation

A A B A C A **A** **B** A B A C A A

↑

	0	1	<b>2</b>	3	4	5
pat.charAt(j)	A	B	<b>A</b>	B	A	C
dfa[][j]						
A	1	1	<b>3</b>	1	5	1
B	0	2	<b>0</b>	4	0	4
C	0	0	<b>0</b>	0	0	6

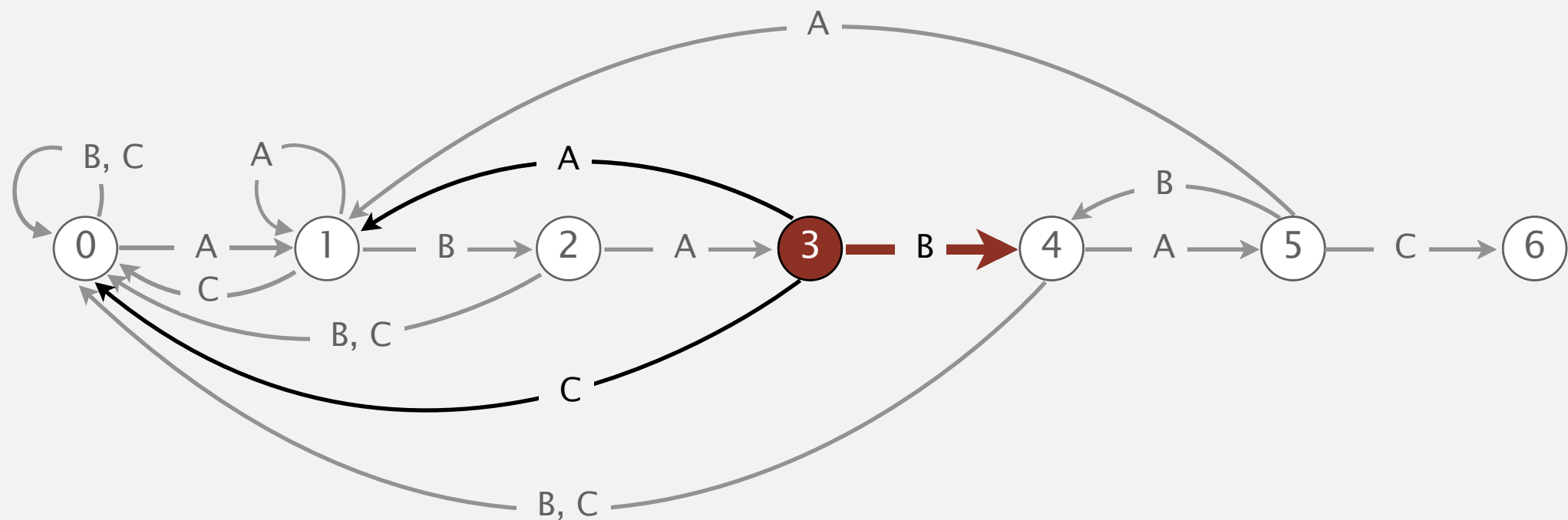


# Knuth–Morris–Pratt demo: DFA simulation

A A B A C A **A** **B** **A** B A C A A

↑

	0	1	2	<b>3</b>	4	5
pat.charAt(j)	A	B	A	<b>B</b>	A	C
dfa[][j] A	1	1	3	<b>1</b>	5	1
dfa[][j] B	0	2	0	<b>4</b>	0	4
dfa[][j] C	0	0	0	<b>0</b>	0	6

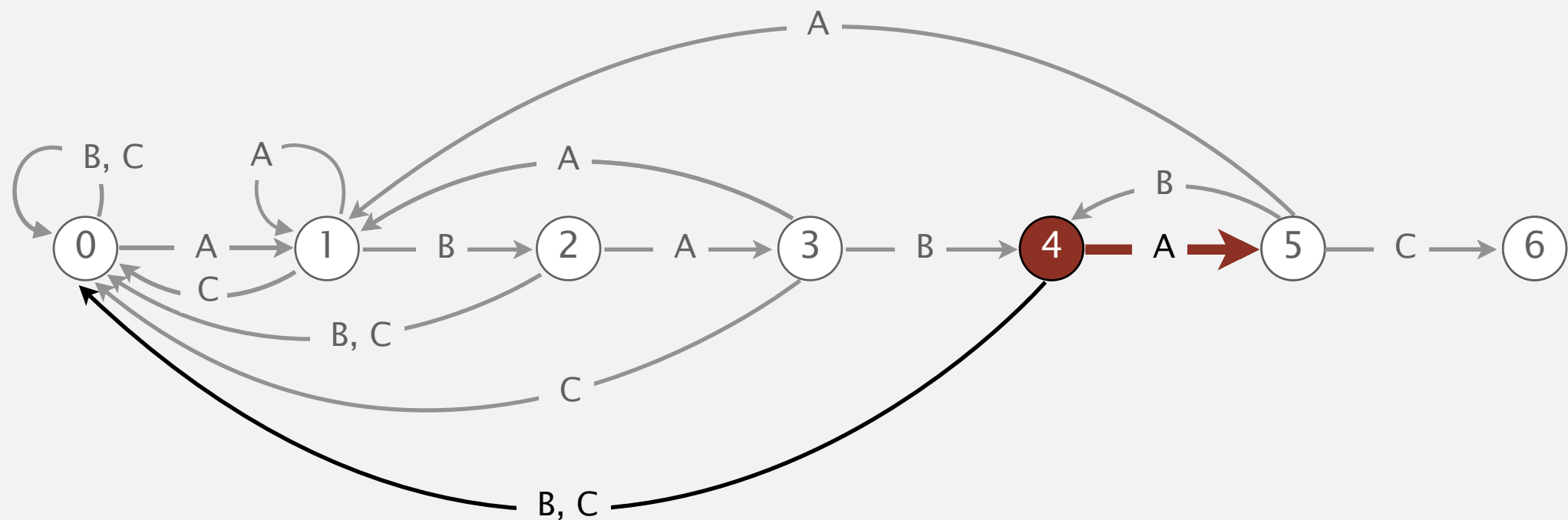




# Knuth–Morris–Pratt demo: DFA simulation

A A B A C A A B A B A C A A

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]						
A	1	1	3	1	5	1
B	0	2	0	4	0	4
C	0	0	0	0	0	6

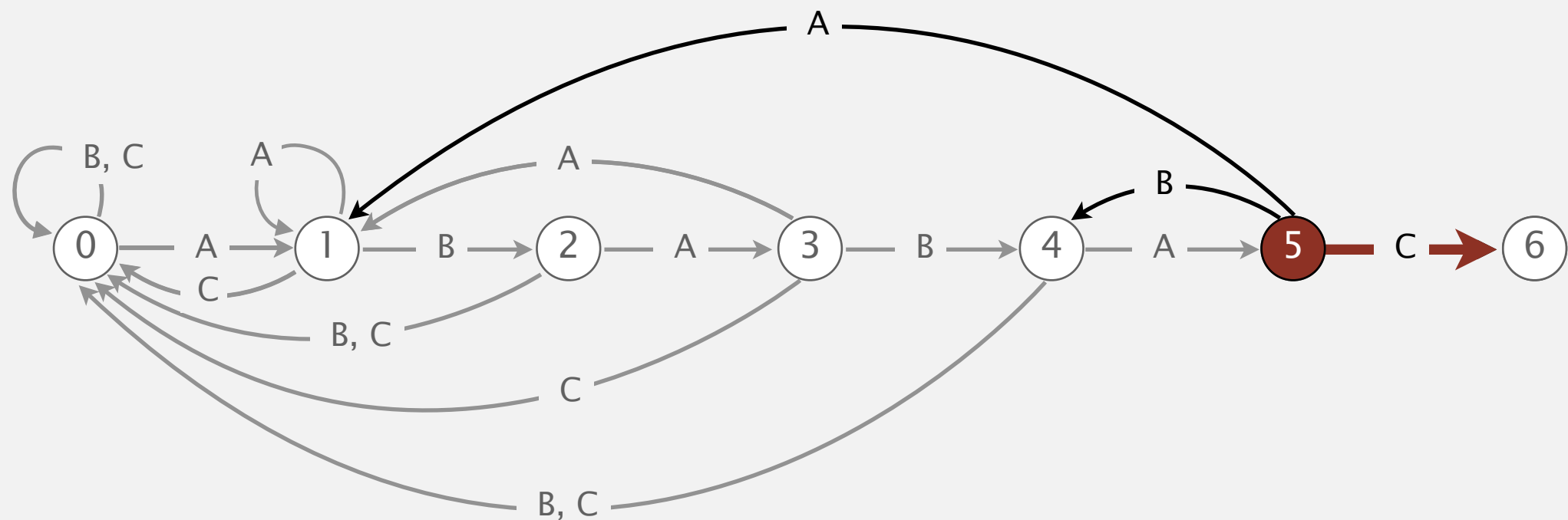


# Knuth–Morris–Pratt demo: DFA simulation

A A B A C A A B A B A C A A

↑

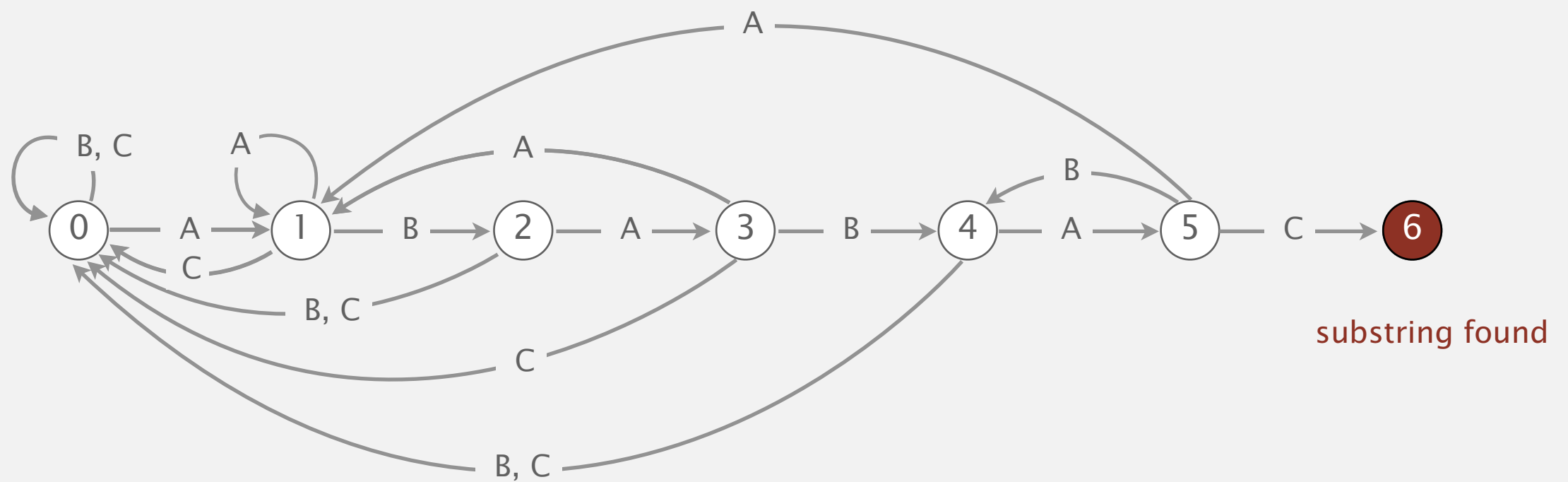
	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j] A	1	1	3	1	5	1
dfa[][j] B	0	2	0	4	0	4
dfa[][j] C	0	0	0	0	0	6



# Knuth-Morris-Pratt demo: DFA simulation

A A B A C A A B A B A C A A

		0	1	2	3	4	5
pat.charAt(j)		<hr/>					
		A	B	A	B	A	C
dfa[][j]	A	1	1	3	1	5	1
	B	0	2	0	4	0	4
	C	0	0	0	0	0	6



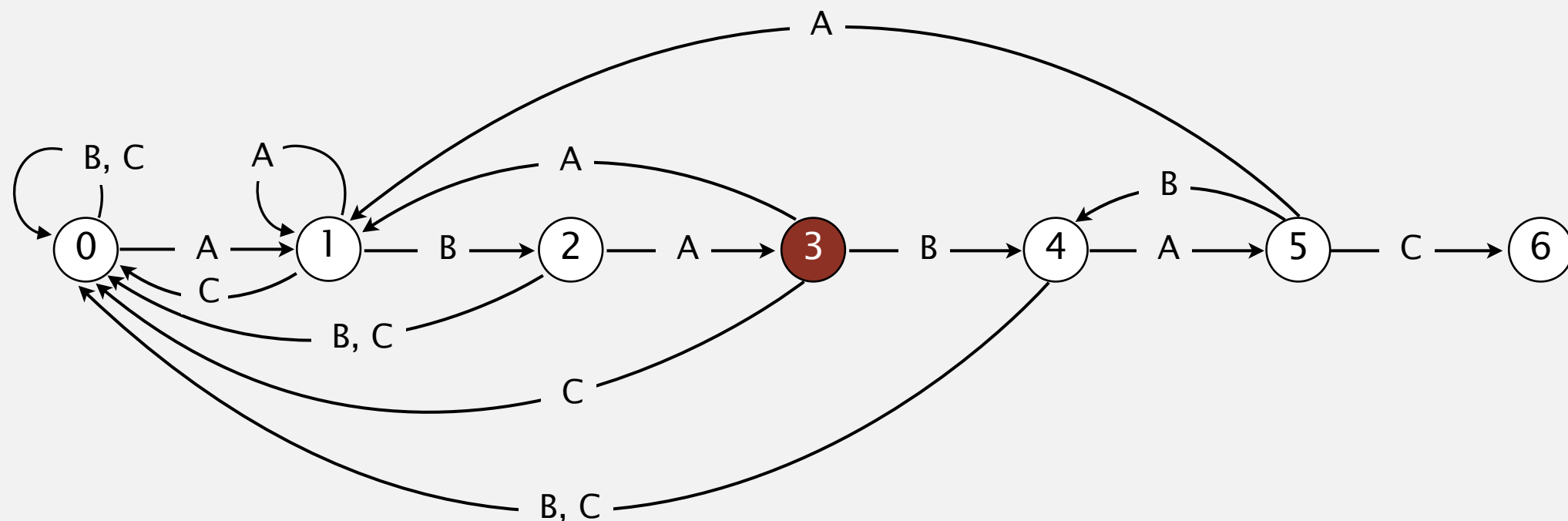
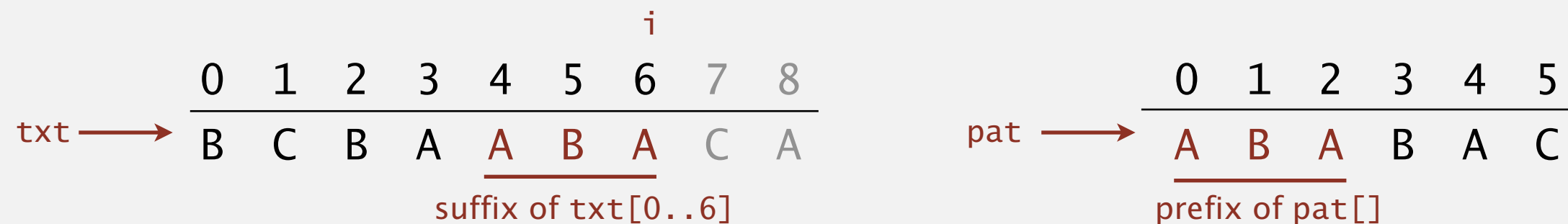
# Interpretation of Knuth–Morris–Pratt DFA

Q. What is interpretation of DFA state after reading in `txt[i]`?

A. State = number of characters in pattern that have been matched.

length of longest prefix of `pat[]`  
that is a suffix of `txt[0..i]`

Ex. DFA is in state 3 after reading in `txt[0..6]`.



# Knuth–Morris–Pratt substring search: Java implementation

---

## Key differences from brute-force implementation.

- Need to precompute `dfa[][]` from pattern.
- Text pointer `i` never decrements.

```
public int search(String txt)
{
    int i, j, N = txt.length();
    for (i = 0, j = 0; i < N && j < M; i++)
        j = dfa[txt.charAt(i)][j];
    if (j == M) return i - M;
    else      return N;
}
```

← no backup

## Running time.

- Simulate DFA on text: at most  $N$  character accesses.
- Build DFA: how to do efficiently? [warning: tricky algorithm ahead]



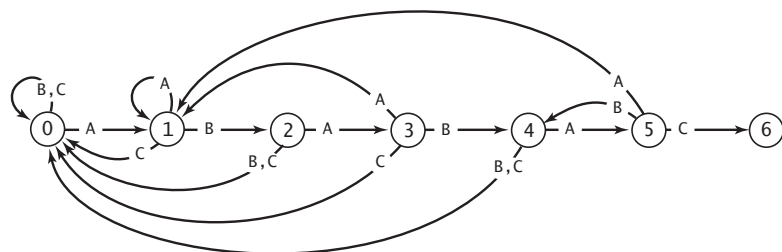
# Knuth–Morris–Pratt substring search: Java implementation

## Key differences from brute-force implementation.

- Need to precompute `dfa[][]` from pattern.
- Text pointer `i` never decrements.
- Could use **input stream**.

```
public int search(In in)
{
    int i, j;
    for (i = 0, j = 0; !in.isEmpty() && j < M; i++)
        j = dfa[in.readChar()][j];
    if (j == M) return i - M;
    else      return NOT_FOUND;
}
```

no backup

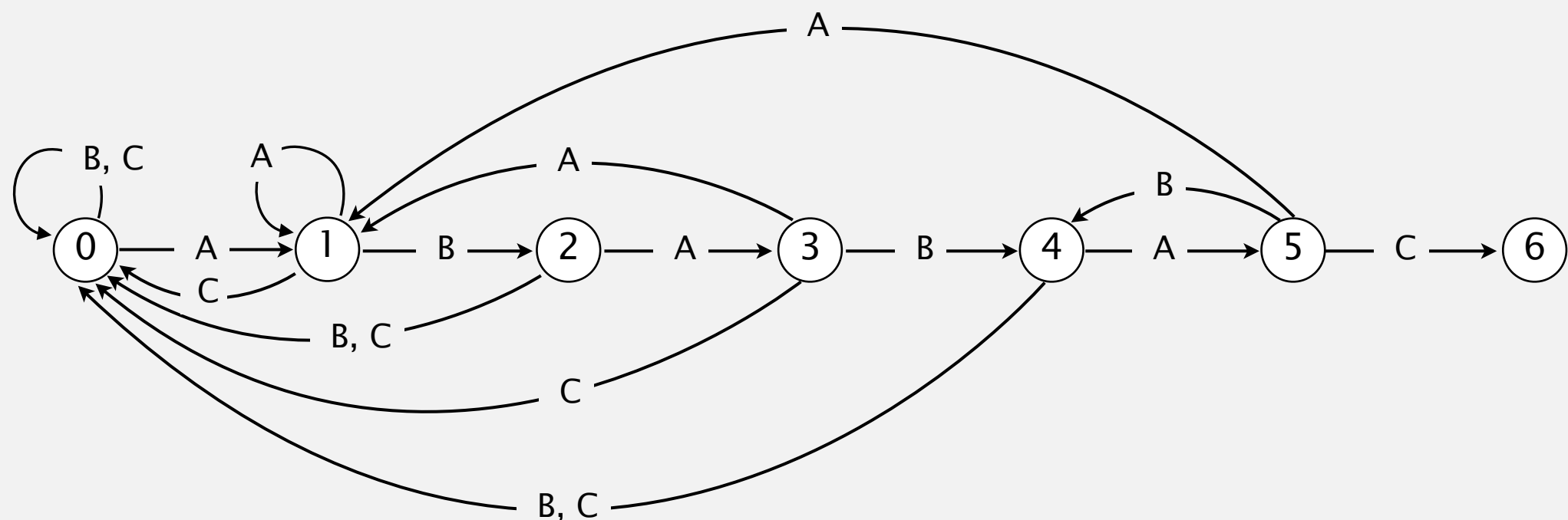


# Knuth–Morris–Pratt demo: DFA construction



		0	1	2	3	4	5
pat.charAt(j)		A	B	A	B	A	C
dfa[][j]	A	1	1	3	1	5	1
	B	0	2	0	4	0	4
	C	0	0	0	0	0	6

Constructing the DFA for KMP substring search for A B A B A C





# Knuth–Morris–Pratt demo: DFA construction

---

Include one state for each character in pattern (plus accept state).

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]	A					
	B					
	C					

Constructing the DFA for KMP substring search for A B A B A C

0

1

2

3

4

5

6

# Knuth–Morris–Pratt demo: DFA construction

**Match transition.** If in state  $j$  and next char  $c == \text{pat.charAt}(j)$ , go to  $j+1$ .

↑ first  $j$  characters of pattern have already been matched      ↑ next char matches      ↑ now first  $j+1$  characters of pattern have been matched

		0	1	2	3	4	5
pat.charAt(j)		A	B	A	B	A	C
dfa[][j]	A	1		3		5	
	B		2		4		
	C						6

Constructing the DFA for KMP substring search for A B A B A C



# Knuth–Morris–Pratt demo: DFA construction

Mismatch transition: back up if  $c \neq \text{pat.charAt}(j)$ .

		0	1	2	3	4	5
pat.charAt(j)		A	B	A	B	A	C
dfa[][j]	A	1		3		5	
	B	0	2		4		
	C	0					6

Constructing the DFA for KMP substring search for A B A B A C

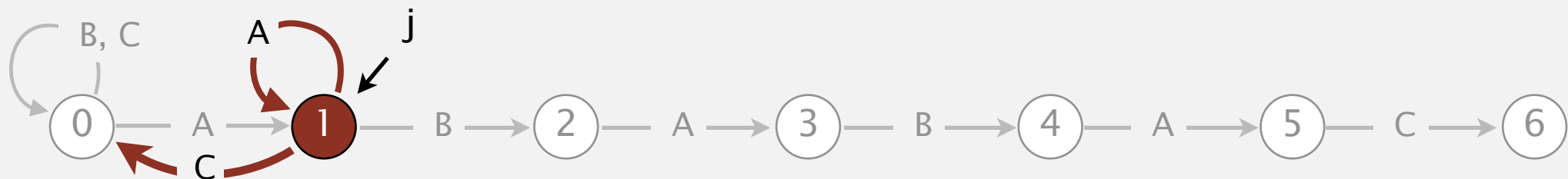


# Knuth–Morris–Pratt demo: DFA construction

Mismatch transition: back up if  $c \neq \text{pat.charAt}(j)$ .

		0	1	2	3	4	5
pat.charAt(j)		A	B	A	B	A	C
dfa[][j]	A	1	1	3		5	
	B	0	2		4		
	C	0	0				6

Constructing the DFA for KMP substring search for A B A B A C

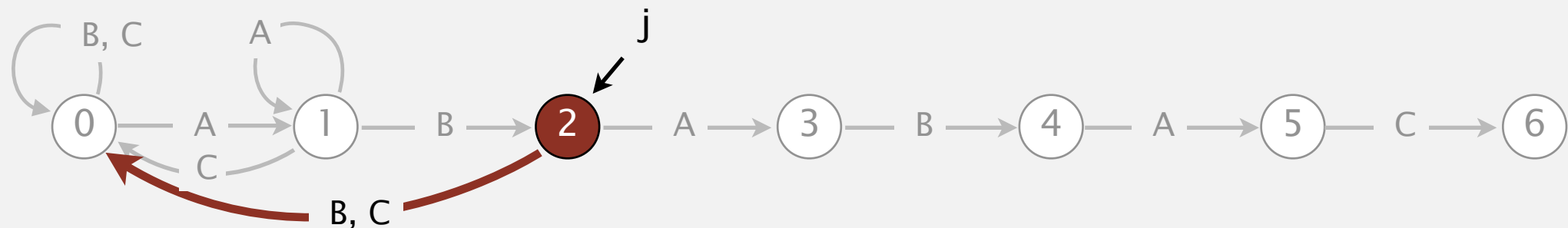


# Knuth–Morris–Pratt demo: DFA construction

**Mismatch transition:** back up if  $c \neq \text{pat.charAt}(j)$ .

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]						
A	1	1	3		5	
B	0	2	0	4		
C	0	0	0			6

Constructing the DFA for KMP substring search for A B A B A C

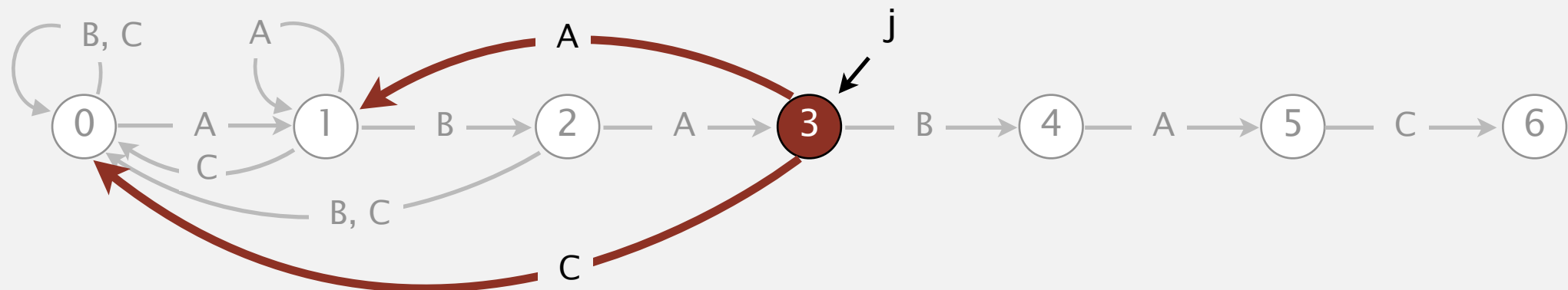


# Knuth–Morris–Pratt demo: DFA construction

Mismatch transition: back up if  $c \neq \text{pat.charAt}(j)$ .

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]						
A	1	1	3	1	5	
B	0	2	0	4		
C	0	0	0	0		6

Constructing the DFA for KMP substring search for A B A B A C

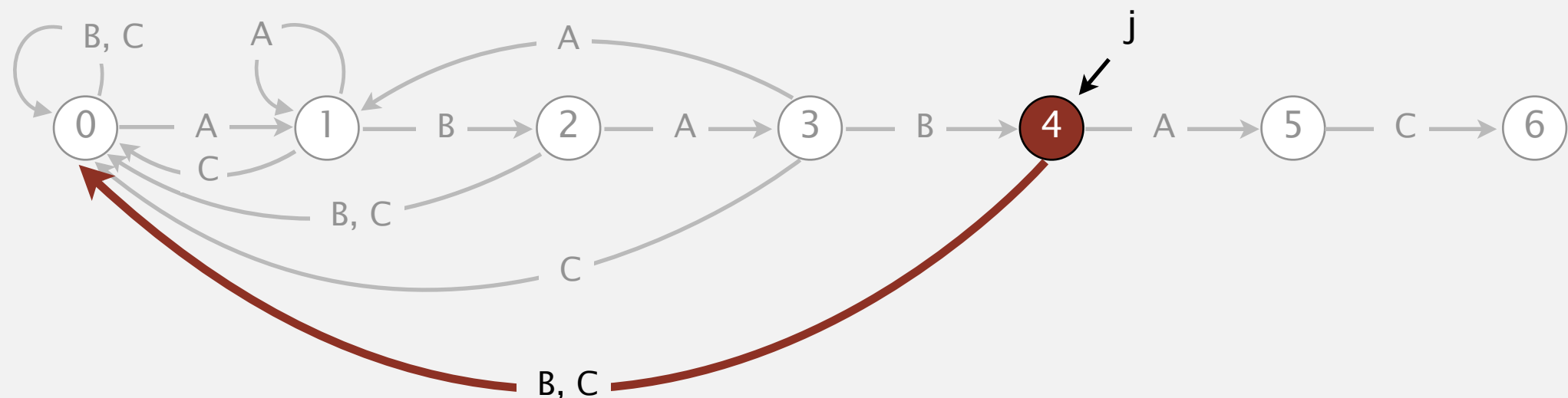


# Knuth–Morris–Pratt demo: DFA construction

**Mismatch transition:** back up if  $c \neq \text{pat.charAt}(j)$ .

		0	1	2	3	4	5
pat.charAt(j)		A	B	A	B	A	C
dfa[][j]	A	1	1	3	1	5	
	B	0	2	0	4	0	
	C	0	0	0	0	0	6

Constructing the DFA for KMP substring search for A B A B A C

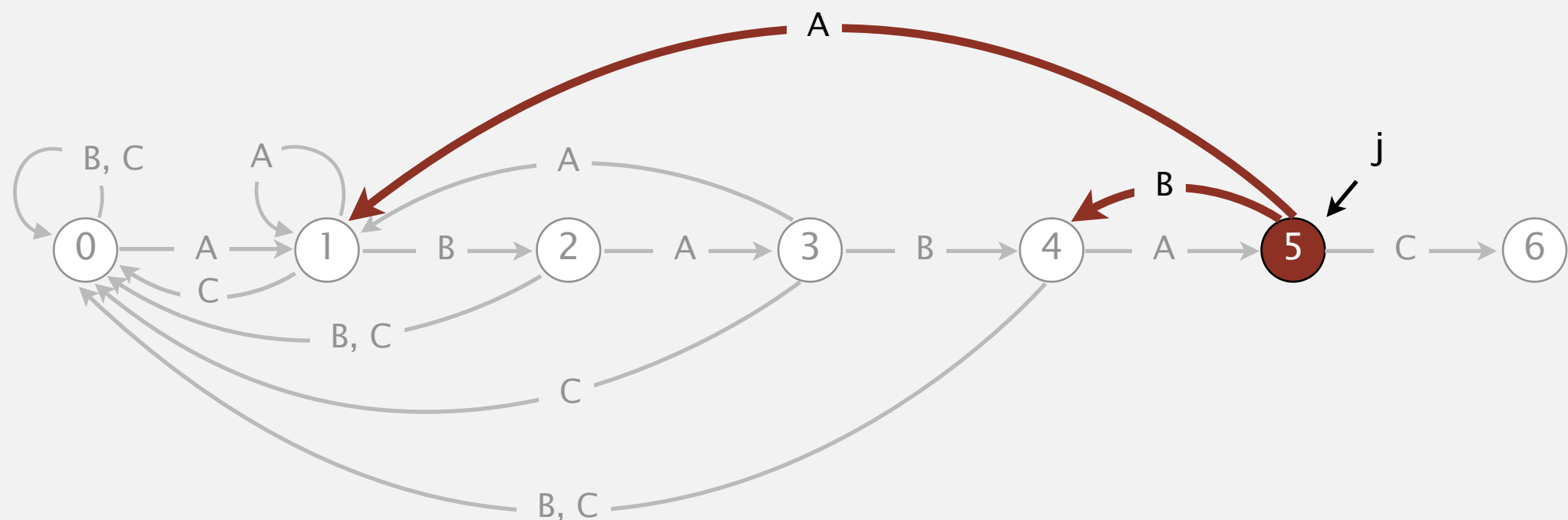


# Knuth–Morris–Pratt demo: DFA construction

Mismatch transition: back up if  $c \neq \text{pat.charAt}(j)$ .

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]						
A	1	1	3	1	5	1
B	0	2	0	4	0	4
C	0	0	0	0	0	6

Constructing the DFA for KMP substring search for A B A B A C

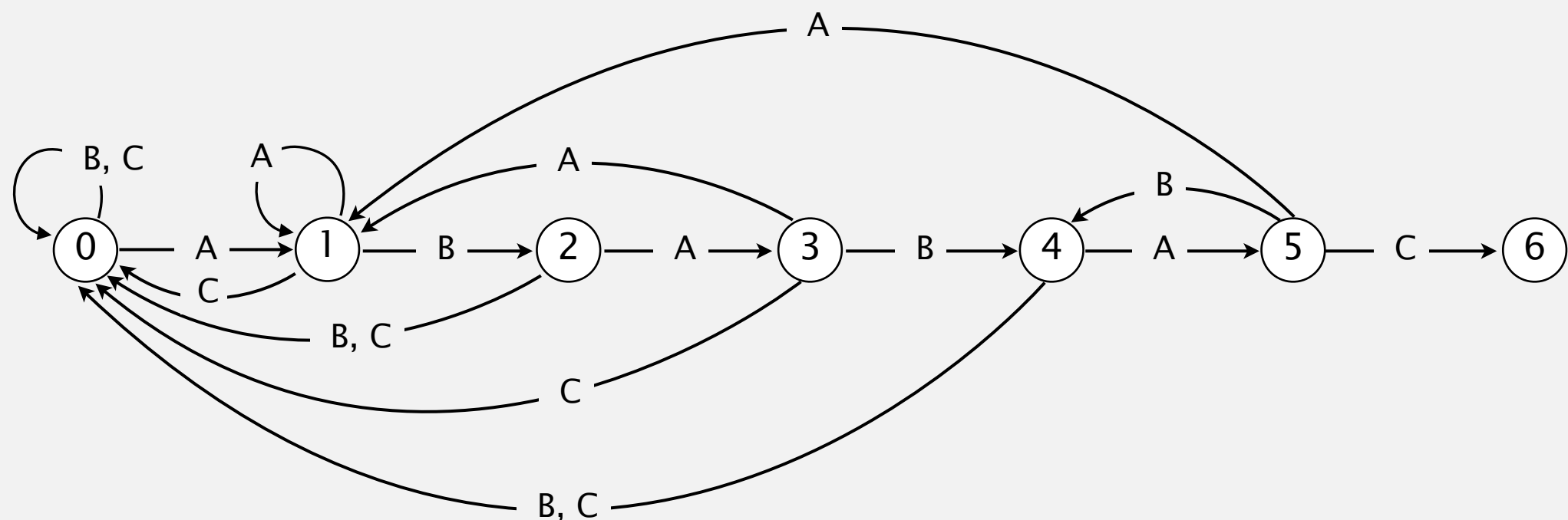




# Knuth–Morris–Pratt demo: DFA construction

		0	1	2	3	4	5
pat.charAt(j)		A	B	A	B	A	C
dfa[][j]	A	1	1	3	1	5	1
	B	0	2	0	4	0	4
	C	0	0	0	0	0	6

Constructing the DFA for KMP substring search for A B A B A C



# How to build DFA from pattern?

---

Include one state for each character in pattern (plus accept state).

		0	1	2	3	4	5
pat.charAt(j)	A	A	B	A	B	A	C
	B						
	C						

0

1

2

3

4

5

6

# How to build DFA from pattern?

**Match transition.** If in state  $j$  and next char  $c == \text{pat.charAt}(j)$ , go to  $j+1$ .

↑  
first  $j$  characters of pattern  
have already been matched

↑  
next char matches

↑  
now first  $j + 1$  characters of  
pattern have been matched

		0	1	2	3	4	5
pat.charAt(j)		A	B	A	B	A	C
dfa[][j]	A	1		3		5	
	B		2		4		
	C						6



# How to build DFA from pattern?

**Mismatch transition.** If in state  $j$  and next char  $c \neq \text{pat.charAt}(j)$ , then the last  $j-1$  characters of input are  $\text{pat}[1..j-1]$ , followed by  $c$ .

**To compute  $\text{dfa}[c][j]$ :** Simulate  $\text{pat}[1..j-1]$  on DFA and take transition  $c$ .

**Running time.** Seems to require  $j$  steps.

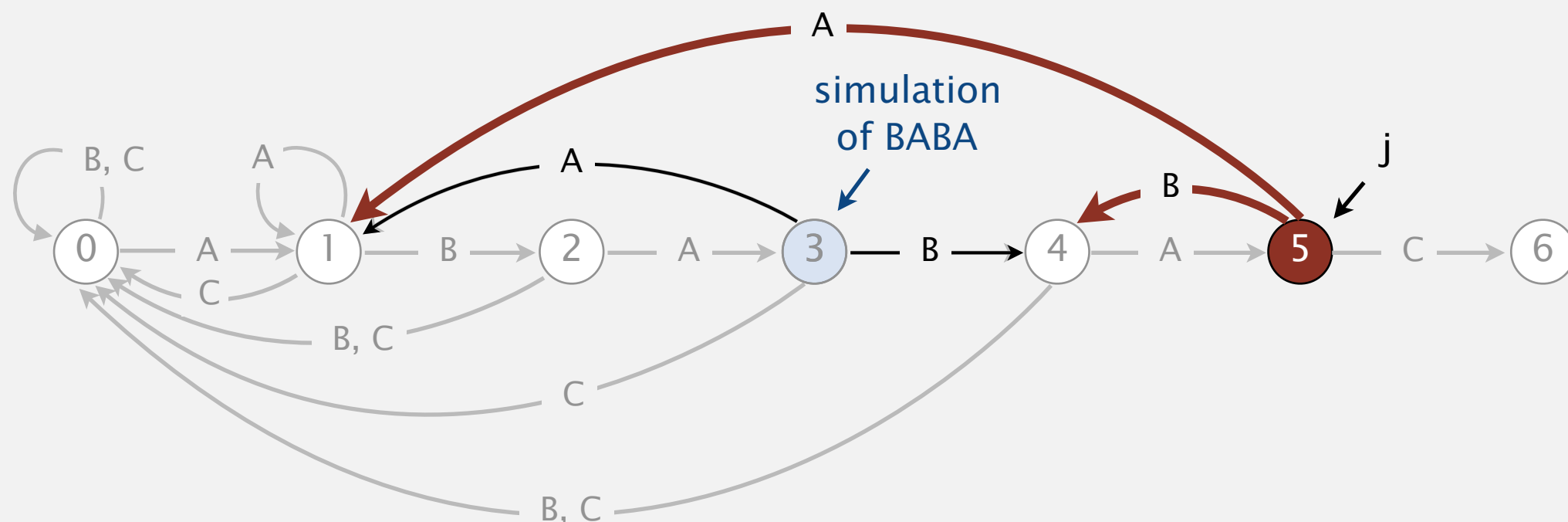
still under construction (!)

**Ex.**  $\text{dfa}['A'][5] = 1$        $\text{dfa}['B'][5] = 4$

simulate BABAA

simulate BABAB

$j$	0	1	2	3	4	5
$\text{pat.charAt}(j)$	A	B	A	B	A	C



# How to build DFA from pattern?

**Mismatch transition.** If in state  $j$  and next char  $c \neq \text{pat.charAt}(j)$ , then the last  $j-1$  characters of input are  $\text{pat}[1..j-1]$ , followed by  $c$ .

state  $x$

To compute  $\text{dfa}[c][j]$ : Simulate  $\text{pat}[1..j-1]$  on DFA and take transition  $c$ .

Running time. Takes only **constant time** if we maintain state  $x$ .

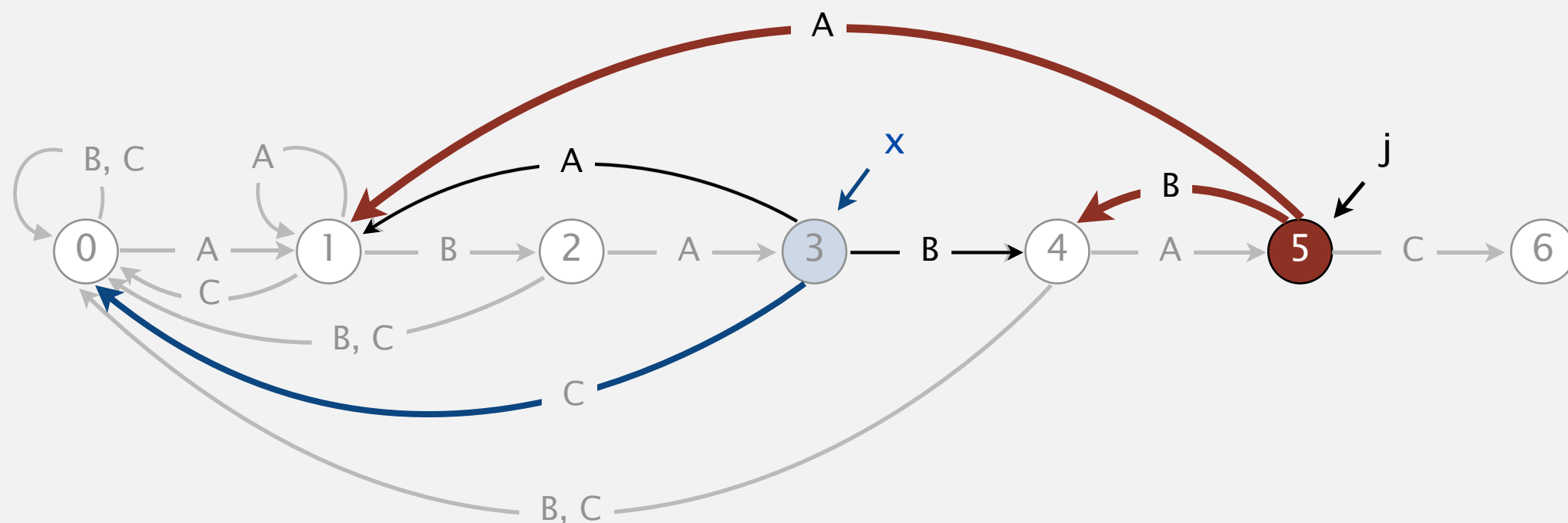
Ex.  $\text{dfa}['A'][5] = 1$      $\text{dfa}['B'][5] = 4$      $x' = 0$

from state  $x$ ,  
take transition 'A'  
 $= \text{dfa}['A'][x]$

from state  $x$ ,  
take transition 'B'  
 $= \text{dfa}['B'][x]$

from state  $x$ ,  
take transition 'C'  
 $= \text{dfa}['C'][x]$

0	1	2	3	4	5
A	B	A	B	A	C

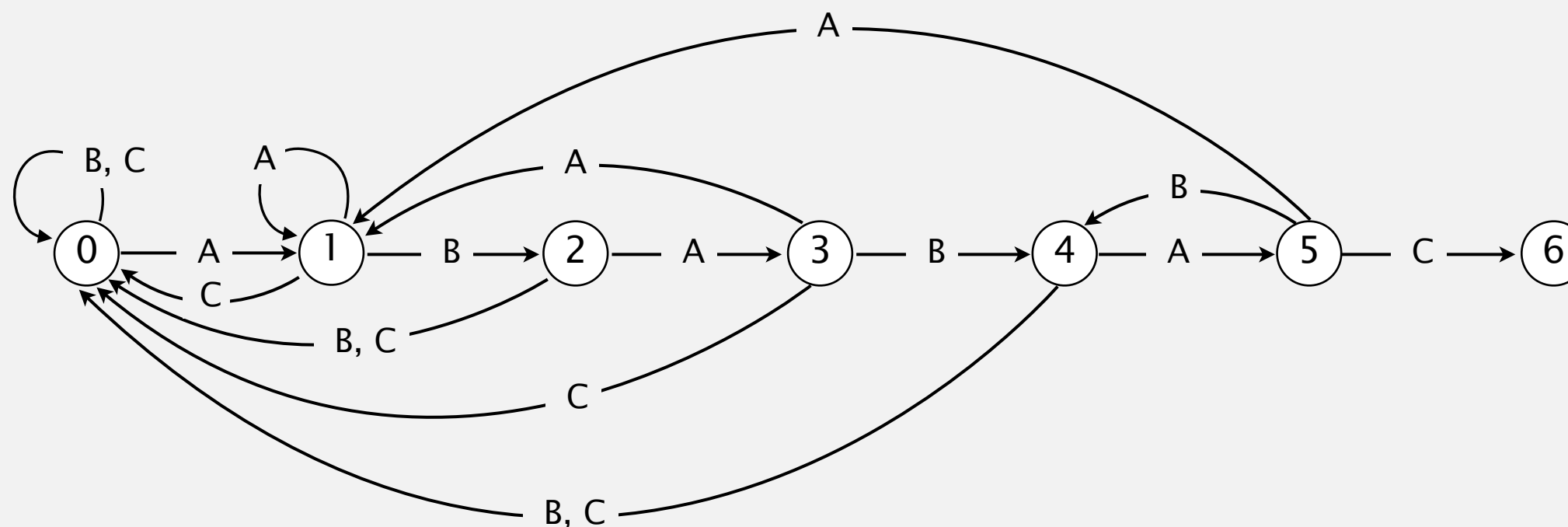


# Knuth–Morris–Pratt demo: DFA construction in linear time



		0	1	2	3	4	5
pat.charAt(j)		A	B	A	B	A	C
dfa[][j]	A	1	1	3	1	5	1
	B	0	2	0	4	0	4
	C	0	0	0	0	0	6

Constructing the DFA for KMP substring search for A B A B A C



# Knuth–Morris–Pratt demo: DFA construction in linear time

---

Include one state for each character in pattern (plus accept state).

		0	1	2	3	4	5
pat.charAt(j)		A	B	A	B	A	C
dfa[][j]	A						
	B						
	C						

Constructing the DFA for KMP substring search for A B A B A C



# Knuth–Morris–Pratt demo: DFA construction in linear time

**Match transition.** For each state  $j$ ,  $\text{dfa}[\text{pat.charAt}(j)][j] = j+1$ .

↑  
first  $j$  characters of pattern  
have already been matched

↑  
now first  $j+1$  characters of  
pattern have been matched

		0	1	2	3	4	5
<code>pat.charAt(j)</code>		A	B	A	B	A	C
<code>dfa[][j]</code>	A	1		3		5	
	B		2		4		
	C						6

Constructing the DFA for KMP substring search for A B A B A C





# Knuth–Morris–Pratt demo: DFA construction in linear time

**Mismatch transition.** For state 0 and char  $c \neq \text{pat.charAt}(j)$ ,  
set  $\text{dfa}[c][0] = 0$ .

		0	1	2	3	4	5
pat.charAt(j)		A	B	A	B	A	C
dfa[][j]	A	1		3		5	
	B	0	2		4		
	C	0					6

Constructing the DFA for KMP substring search for A B A B A C



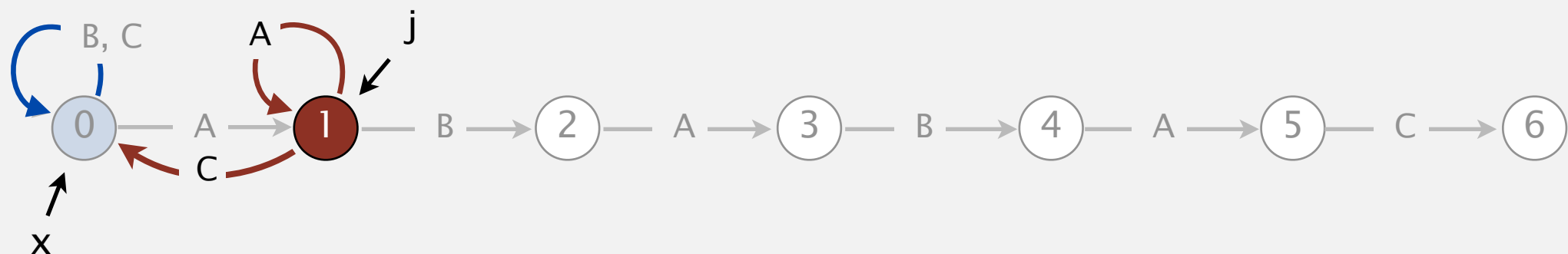
# Knuth–Morris–Pratt demo: DFA construction in linear time

**Mismatch transition.** For each state  $j$  and char  $c \neq \text{pat.charAt}(j)$ , set  $\text{dfa}[c][j] = \text{dfa}[c][x]$ ; then update  $x = \text{dfa}[\text{pat.charAt}(j)][x]$ .

$x = \text{simulation of empty string}$   
↓

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]	A	1	3		5	
B	0	2		4		
C	0					6

Constructing the DFA for KMP substring search for A B A B A C



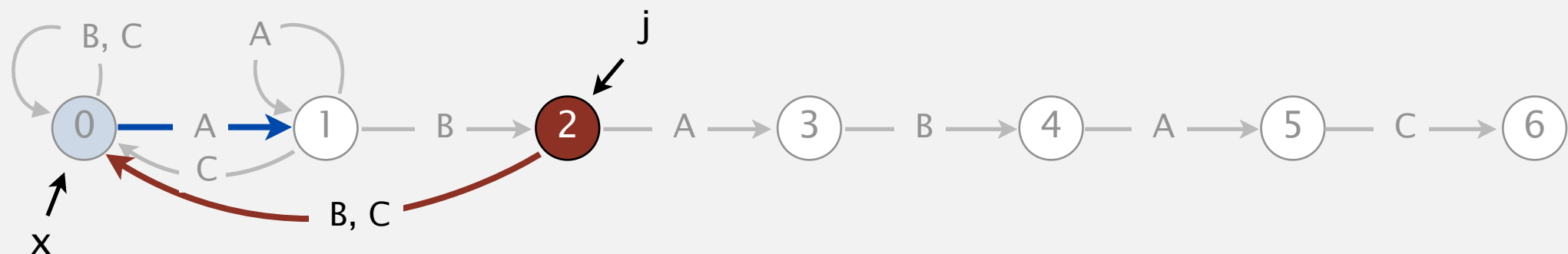
# Knuth–Morris–Pratt demo: DFA construction in linear time

**Mismatch transition.** For each state  $j$  and char  $c \neq \text{pat.charAt}(j)$ , set  $\text{dfa}[c][j] = \text{dfa}[c][x]$ ; then update  $x = \text{dfa}[\text{pat.charAt}(j)][x]$ .

$x = \text{simulation of B}$   
↓

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]	A	1	3		5	
B	0	2		4		
C	0	0				6

Constructing the DFA for KMP substring search for A B A B A C



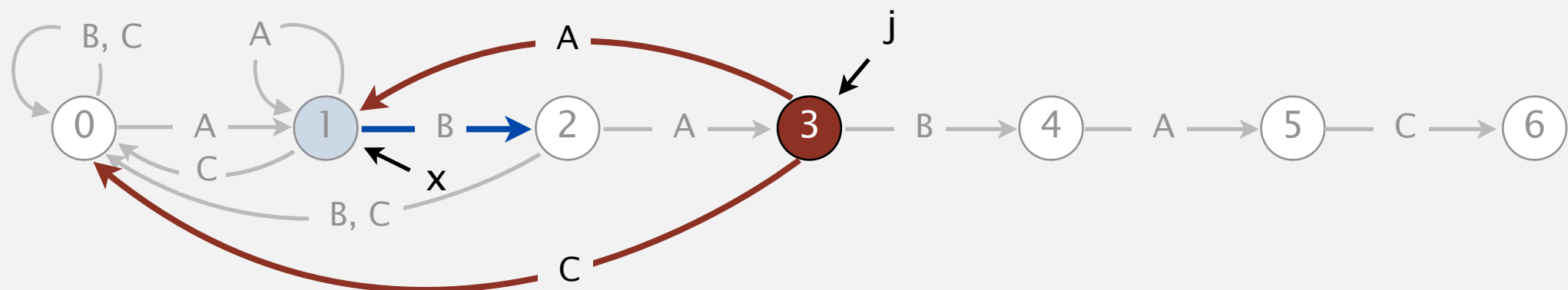
# Knuth–Morris–Pratt demo: DFA construction in linear time

**Mismatch transition.** For each state  $j$  and char  $c \neq \text{pat.charAt}(j)$ , set  $\text{dfa}[c][j] = \text{dfa}[c][x]$ ; then update  $x = \text{dfa}[\text{pat.charAt}(j)][x]$ .

$x = \text{simulation of B A}$   
↓

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]						
A	1	1	3		5	
B	0	2	0	4		
C	0	0	0			6

Constructing the DFA for KMP substring search for A B A B A C



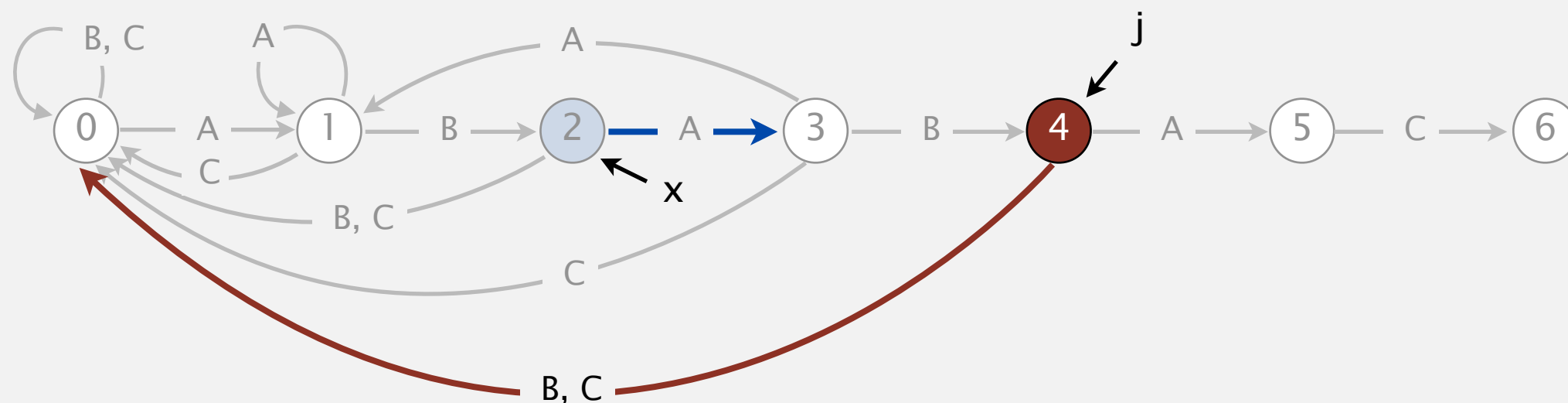
# Knuth–Morris–Pratt demo: DFA construction in linear time

**Mismatch transition.** For each state  $j$  and char  $c \neq \text{pat.charAt}(j)$ , set  $\text{dfa}[c][j] = \text{dfa}[c][x]$ ; then update  $x = \text{dfa}[\text{pat.charAt}(j)][x]$ .

$x = \text{simulation of B A B}$   
↓

	0	1	2	3	4	5
$\text{pat.charAt}(j)$	A	B	A	B	A	C
$\text{dfa[]} [j]$						
A	1	1	3	1	5	
B	0	2	0	4		
C	0	0	0	0		6

Constructing the DFA for KMP substring search for A B A B A C



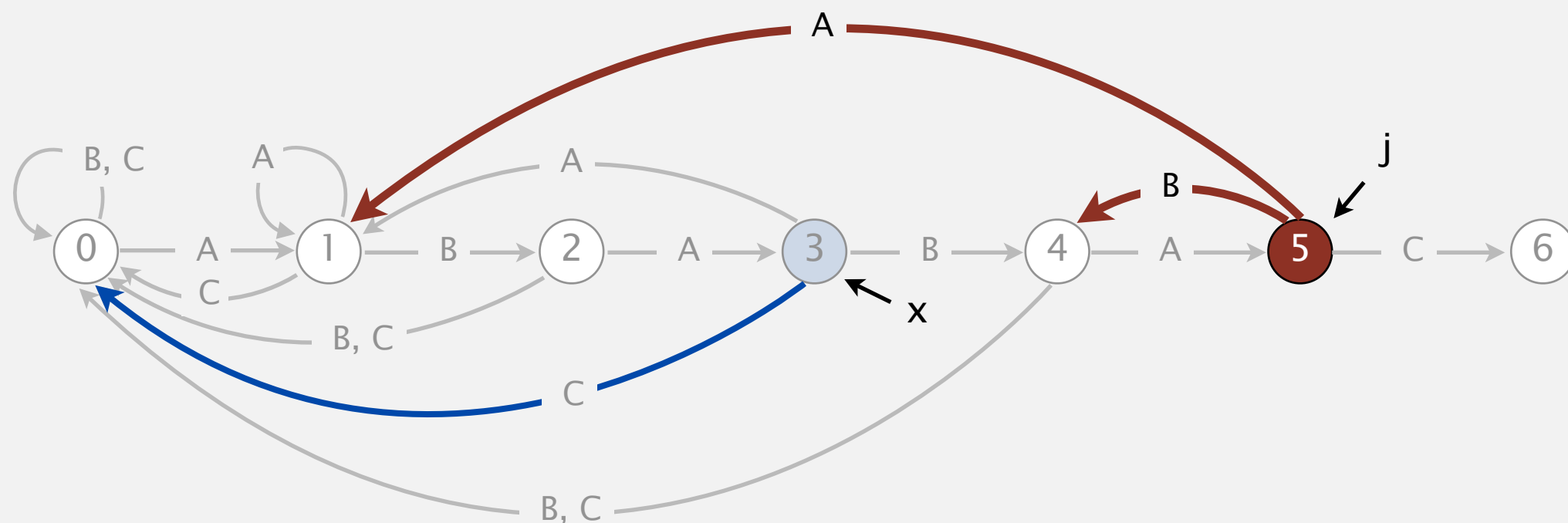
# Knuth–Morris–Pratt demo: DFA construction in linear time

**Mismatch transition.** For each state  $j$  and char  $c \neq \text{pat.charAt}(j)$ , set  $\text{dfa}[c][j] = \text{dfa}[c][x]$ ; then update  $x = \text{dfa}[\text{pat.charAt}(j)][x]$ .

x = simulation of B A B A  
↓

	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C
dfa[][j]						
A	1	1	3	1	5	
B	0	2	0	4	0	
C	0	0	0	0	0	6

Constructing the DFA for KMP substring search for A B A B A C



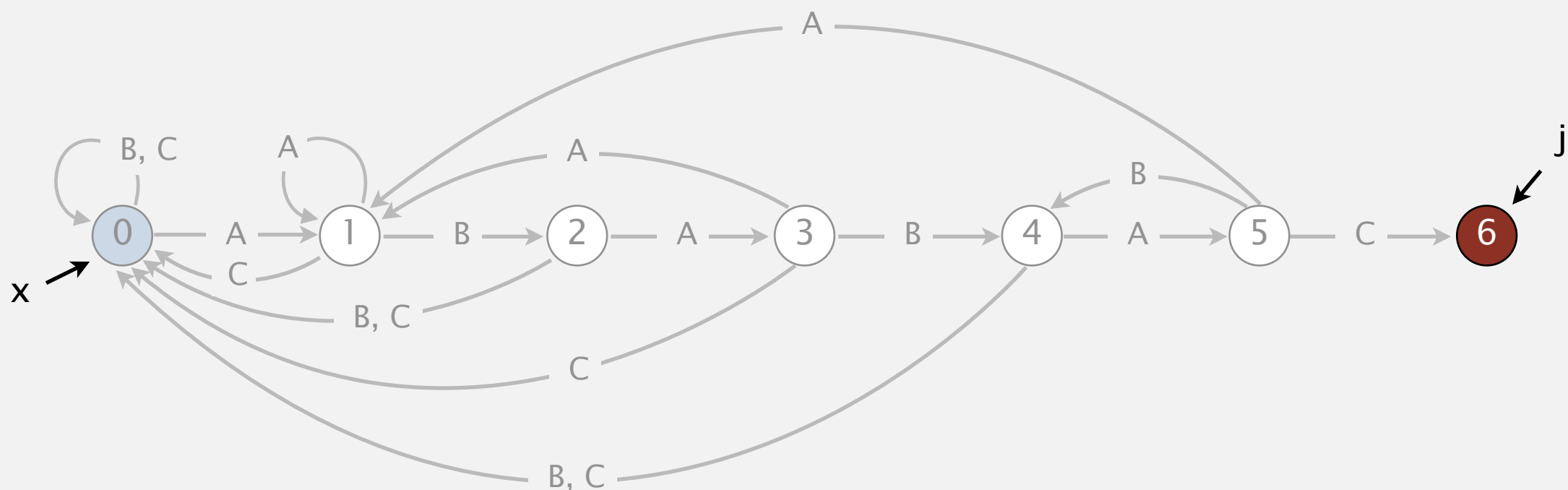
# Knuth–Morris–Pratt demo: DFA construction in linear time

**Mismatch transition.** For each state  $j$  and char  $c \neq \text{pat.charAt}(j)$ , set  $\text{dfa}[c][j] = \text{dfa}[c][x]$ ; then update  $x = \text{dfa}[\text{pat.charAt}(j)][x]$ .

$x = \text{simulation of B A B A C}$   
↓

	0	1	2	3	4	5
$\text{pat.charAt}(j)$	A	B	A	B	A	C
$\text{dfa}[][j]$	A	1	1	3	1	5
B	0	2	0	4	0	4
C	0	0	0	0	0	6

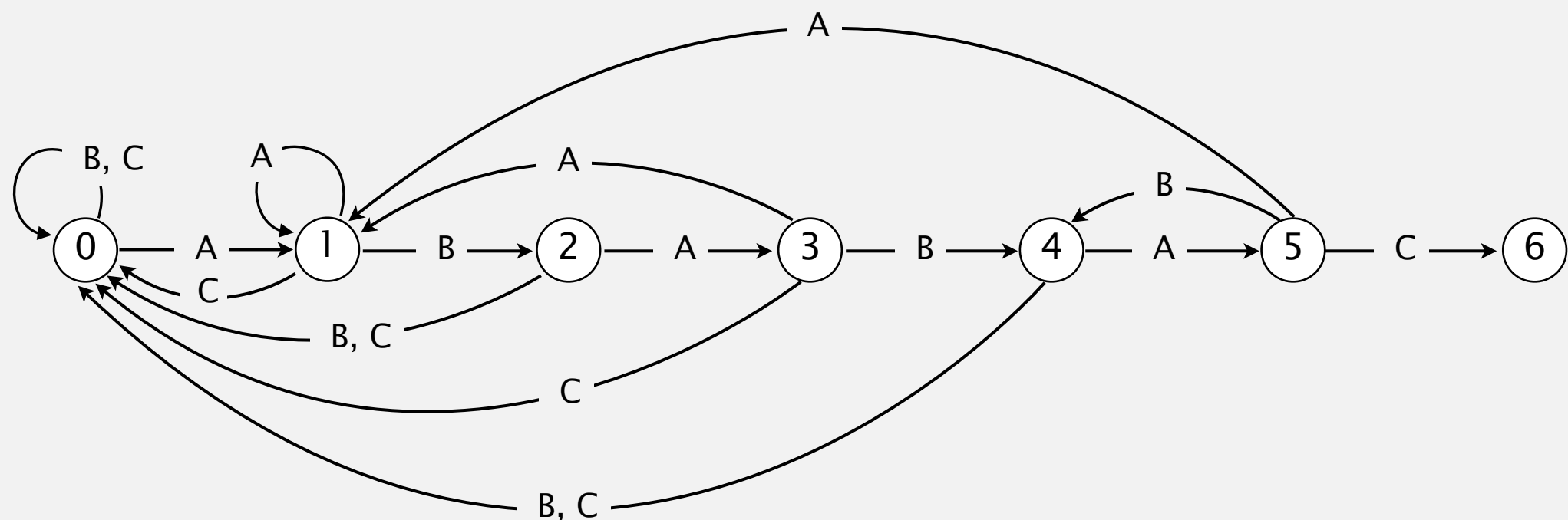
Constructing the DFA for KMP substring search for A B A B A C



# Knuth–Morris–Pratt demo: DFA construction in linear time

		0	1	2	3	4	5
pat.charAt(j)		A	B	A	B	A	C
dfa[][j]	A	1	1	3	1	5	1
	B	0	2	0	4	0	4
	C	0	0	0	0	0	6

Constructing the DFA for KMP substring search for A B A B A C





# Constructing the DFA for KMP substring search: Java implementation

---

For each state  $j$ :

- Copy `dfa[][x]` to `dfa[][j]` for mismatch case.
- Set `dfa[pat.charAt(j)][j]` to  $j+1$  for match case.
- Update  $x$ .

```
public KMP(String pat)
{
    this.pat = pat;
    M = pat.length();
    dfa = new int[R][M];
    dfa[pat.charAt(0)][0] = 1;
    for (int x = 0, j = 1; j < M; j++)
    {
        for (int c = 0; c < R; c++)
            dfa[c][j] = dfa[c][x];
        dfa[pat.charAt(j)][j] = j+1;
        x = dfa[pat.charAt(j)][x];
    }
}
```

← copy mismatch cases

← set match case

← update restart state

**Running time.**  $M$  character accesses (but space/time proportional to  $R M$ ).

# Knuth–Morris–Pratt: brief history

---

- Independently discovered by two theoreticians and a hacker.
  - Knuth: inspired by esoteric theorem, discovered linear algorithm
  - Pratt: made running time independent of alphabet size
  - Morris: built a text editor for the CDC 6400 computer
- Theory meets practice.

SIAM J. COMPUT.  
Vol. 6, No. 2, June 1977

## FAST PATTERN MATCHING IN STRINGS\*

DONALD E. KNUTH<sup>†</sup>, JAMES H. MORRIS, JR.<sup>‡</sup> AND VAUGHAN R. PRATT<sup>¶</sup>

**Abstract.** An algorithm is presented which finds all occurrences of one given string within another, in running time proportional to the sum of the lengths of the strings. The constant of proportionality is low enough to make this algorithm of practical use, and the procedure can also be extended to deal with some more general pattern-matching problems. A theoretical application of the algorithm shows that the set of concatenations of even palindromes, i.e., the language  $\{a a^R\}^*$ , can be recognized in linear time. Other algorithms which run even faster on the average are also considered.



Don Knuth



Jim Morris



Vaughan Pratt