# String Matching

- Substring search
  - Find a pattern of length M in a text of length N
    (typically N >> M)



## Brute-Force Substring Search

### Naive algorithm

- Check for pattern starting at each text position



- Order: O(MN)
- Any improvement?

```
int naiveStringMatch(char text[], char pattern[])
{
    int patLength, txtLength;

    patLength = strlen(pattern);
    txtLength = strlen(text);

    for(int i=0; i <= txtLength - patLength; i++)
    {
        for(int j=0; j < patLength; j++)
            if(text[i+j] != pattern[j])
                break;
        if(j == patLength)
            return i;
    }
    return -1;
}
```
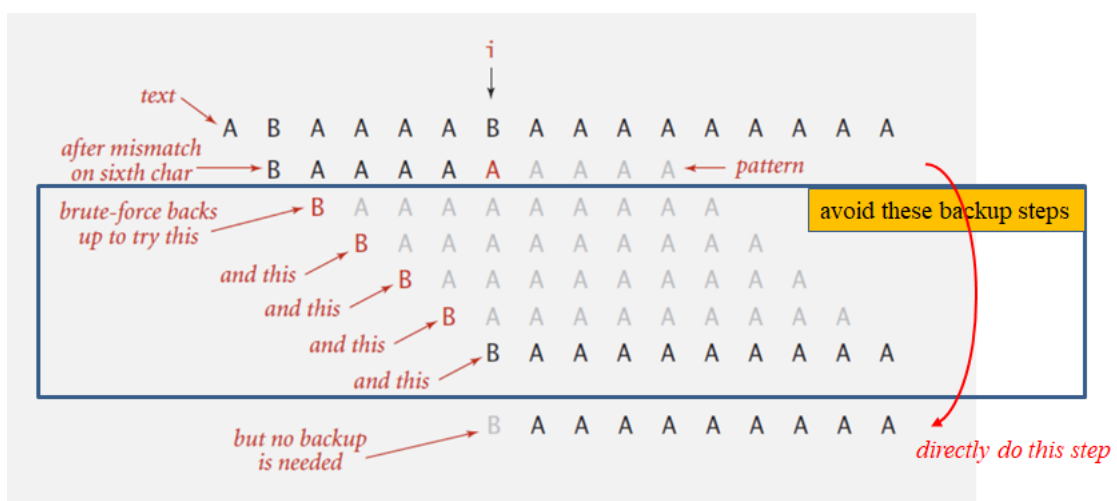
- - can be slow if text and pattern are repetitive
- Improvement
  - develop a linear time algorithm
  - avoid **backup**
    - naive algorithm needs backup for every mismatch
    - thus naive algorithm cannot be used when input text is a stream



# Knuth-Morris-Pratt(KMP) Algorithm

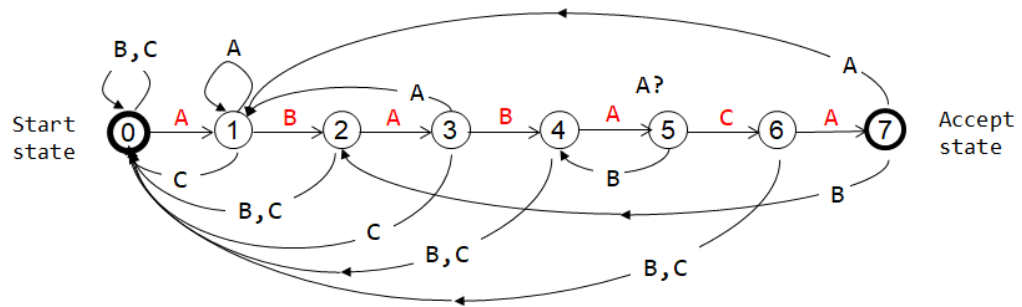- Clever method to always avoid **backup** problem



# Deterministic Finite Automaton

- DFA
  - Finite number of states (including **start** and **accept states**)
  - Exactly one transition for each char

- Accept if sequence of transitions leads to accept state
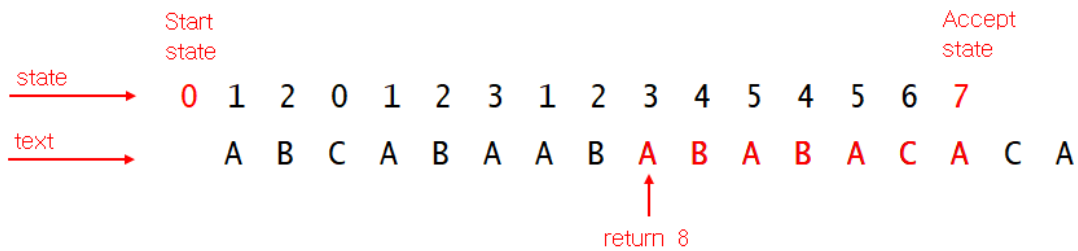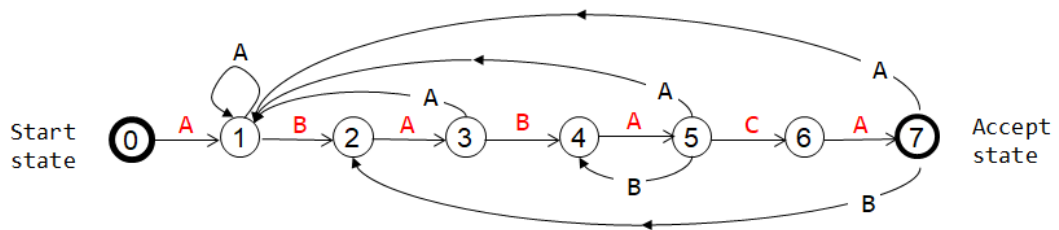
DFA for pattern ABABACA



DFA[][]

| char / state | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|---|
| A | 1 | 1 | 3 | 1 | 5 | 1 | 7 | 1 |
| B | 0 | 2 | 0 | 4 | 0 | 4 | 0 | 2 |
| C | 0 | 0 | 0 | 0 | 0 | 6 | 0 | 0 |
| others | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

```
if in state j reading char c:
    if j is 7, halt and accept
    else move to state DFA[c][j]

Example:
    text:   ABCABAABABABACACA
    state:01201231234545567
```

DFA for pattern ABABACA



```
state   0 1 2 0 1 2 3 1 2 3 4 5 4 5 6 7
text    A B C A B A A B A B A B A C A C A
```

return 8

Simplified Diagram: remove transitions to state 0



- Difference from naive algorithm
  - precomputation of DFA[][] from pattern
  - text pointer i never decrements (**no backup**)

```
// patLength = strlen(pattern);
int DFAmatching(char text[])
{
    int i, j, txtLength;

    txtLength = strlen(text);

    for(i=0, j=0; i <= txtLength && j < patLength; i++)
        j = DFA[text[i]][j];   // text[i] to be modified

    if(j == patLength)
        return i - patLength;
    else
        return -1;
}                        -  Order: O(N)
                   simulation of DFA on text with no backup
              -  How to build DFA efficiently?
```
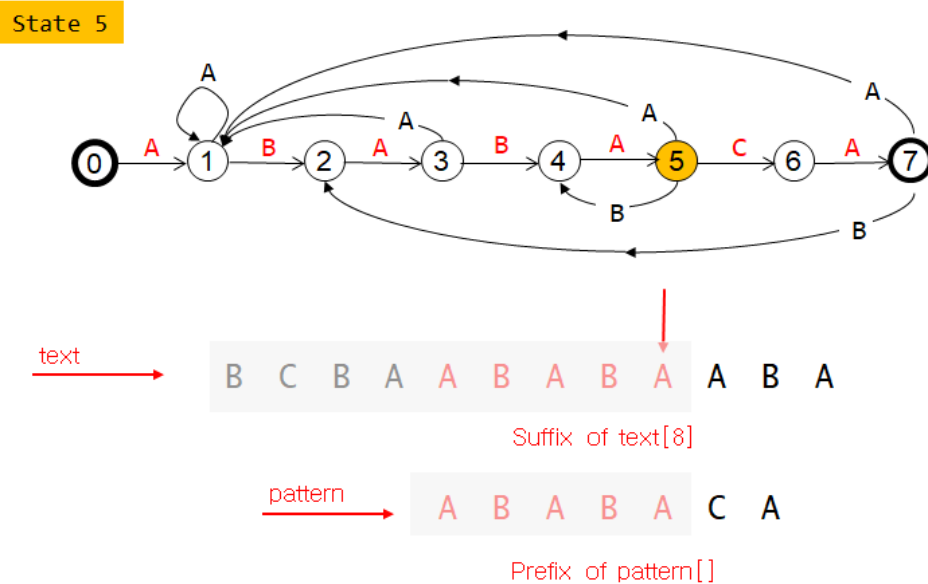
- The state of DFA represents
  - the number of characters in pattern that have been matched
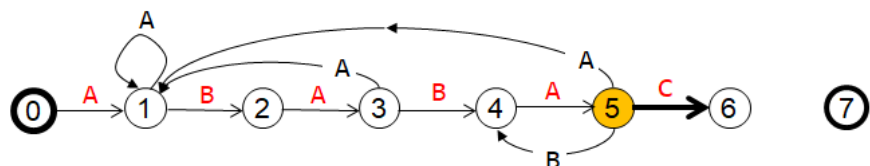


- Prefix / Sufix of a Text

| bananada | | |
| --- | --- | --- |

| | **Prefix** | **Suffix** |
| --- | --- | --- |
| NULL string → | bananada | bananada |
| | **b**ananada | bananad**a** |
| | **ba**nanada | banana**da** |
| | **ban**anada | banan**ada** |
| | **bana**nada | bana**nada** |
| | **banan**ada | ban**anada** |
| | **banana**da | ba**nanada** |
| | **bananad**a | b**ananada** |
| | **bananada** | **bananada** |

## DFA Construction

- Suppose that all transitions from state `0` to stat `j-1` are already computed
- Match transition
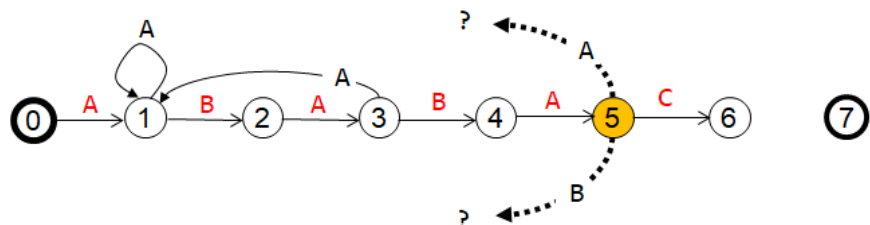  - If in state `j` and next char `char c = pattern[j]`, then transit to state `j+1`

Pattern: ABABA**C**A

State 5



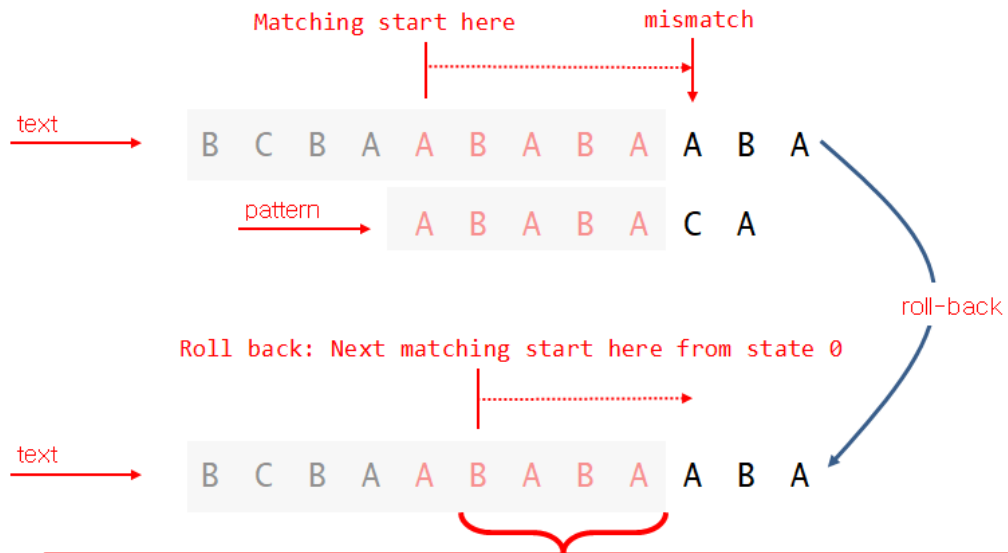- Mismatch transition
  - If in state `j` and next `char c != pattern[j]`, then which state to transit
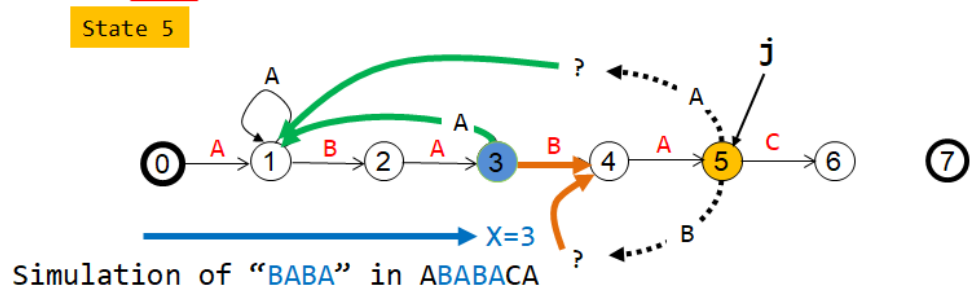
Pattern: ABABA**C**A

State 5

Matching start here      mismatch

text → B C B A A B A B A **A B A**

pattern → A B A B A C A

roll-back

Roll back: Next matching start here from state 0

text → B C B A A B A B A A B A

- *The same as* `pattern[1] ~ pattern[j-1]`
- Roll-back and transit to some state $X$ by matching
  `pattern[1] ~ pattern[j-1]` from state 0 on DFA.
- Transit to the next state `DFA['A'][X]` for the mismatched char `'A'`.

- then the last `j-1` characters of input text are
  `pattern[1] ~ pattern[j-1]`, followed by `c`

- to compute `DFA[c][j]`:

  - simulate `pattern[1] ~ pattern[j-1]` on DFA (still under construction) and let
    the current state `X`

  - Then `DFA[c][j] = DFA[c][X]`

$$DFA['A'][5] = DFA['A'][3] = 1$$
$$DFA['B'][5] = DFA['B'][3] = 4$$

Pattern: ABABA**CA**
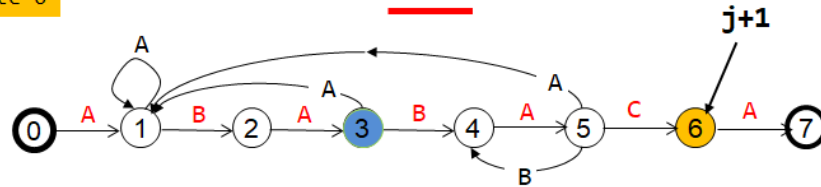
State 5



X=3

Simulation of "BABA" in ABABA**CA**

  - take a transition `c` from state `x`

  - Running time : require `j` steps

  - **But, if we maintain state X, it takes only constant time!**

- Maintaining state `x`:

  - finished computing transitions from state `j`

  - Now, move to next state `j+1`

  - then what the new state( `x'` ) of `x` be?

State 6

Pattern: ABABACA

j+1

X=3

Simulation of "BABA" in ABABACA

X' = ?

Simulation of "BABAC" in ABABACA

- Simulation requires j+1 steps
- But, X' = DFA['C'][X]

X' = DFA['C'][3] = 0

- A Linear Time Algorithm
  - for each state j
    - Match case : set DFA[pattern[j]][j] = j+1
    - Mismatch case : copy DFA[][X] to DFA[][j]
    - Update X

```
int DFA[MAX_SIZE][MAX_SIZE];  /* initially all elements are 0 */
// int R;              /* text character set size */

void constructDFA(char pattern[])
{

    int patLength = strlen(pattern);

    DFA[pattern[0]][0] = 1;

    for(int X=0, j=1; j<patLength; j++)
    {
        for(int c=0; c<R; c++)                 // copy mismatch cases
            DFA[c][j] = DFA[c][X];

        DFA[pattern[j]][j] = j+1;              // copy match case
        X = DFA[pattern[j]][X];                // update X
    }
}
```

- Example

DFA[][]

Pattern: ABABACA
0123456

| char | state | | | | | | | |
|------|-------|---|---|---|---|---|---|---|
|      | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| A    |   |   |   |   |   |   |   |   |
| B    |   |   |   |   |   |   |   |   |
| C    |   |   |   |   |   |   |   |   |
| others |  |   |   |   |   |   |   |   |

DFA[][]

Pattern: ABABACA (0123456)

j start from 1

| char \ state | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|---|
| A | 1 | 1 | 3 | 1 | 5 | 1 | 7 | 1 |
| B | 0 | 2 | 0 | 4 | 0 | 4 | 0 | 2 |
| C | 0 | 0 | 0 | 0 | 0 | 6 | 0 | 0 |
| others | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

X   0   0   1   2   3   0   1   1

DFA[B][0]   DFA[A][0]   DFA[B][1]   DFA[A][2]   DFA[C][3]   DFA[A][0]   ~~DFA[A][1]~~

```
int patLength = strlen(pattern);

DFA[pattern[0]][0] = 1;

for(int X=0, j=1; j<patLength; j++)
{
    for(int c=0; c<R; c++)
        DFA[c][j] = DFA[c][X];

    DFA[pattern[j]][j] = j+1;
    X = DFA[pattern[j]][X];
}
```

## Algorithm with DFA

- String matching algorithm with DFA accesses no more than M+N chars to search for a pattern of length M in a text of length N

- `DFA[][]` can be constructed in time and space of order `O(RM)`, where `R` is the number of characters used in a text

- Questions : Text에 나타나는 모든 pattern을 찾을 수 있는가?
    - Text : AAAAAAAAA
    - Pattern : AAAAA
    - Solution : 0, 1, 2, 3, 4, 5