

CS420

Introduction to the Theory of Computation

Lecture 12: Regular expressions & NFAs

Tiago Cogumbreiro

Today we will learn...

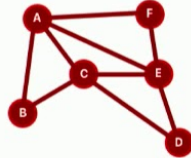
- Define a data type that represents Regular Expressions
- Define inductively acceptance for Regular Expressions
- Define equivalence of Regular Expressions



Sept 13-14, 2019
thestrangeloop.com

Typical solutions

- Collect all milestones and wildcard matching
- Analytical window functions
- Graph data model



V: {A,B,C,D,E,F}
E: {AB,AC,AF,BC,CD,CE,DE,EF}

9/14/2019

=

	AB	AC	AF	BC	CD	CE	DE	EF
A	1	1	1	0	0	0	0	0
B	1	0	0	1	0	0	0	0
C	0	1	0	1	1	1	0	0
D	0	0	0	0	1	0	1	0
E	0	0	0	0	0	1	1	1
F	0	0	1	0	0	0	0	1

N

[www.youtube.com/watch?](http://www.youtube.com/watch?v=WykSdgtLDD0)

[v=WykSdgtLDD0](http://www.youtube.com/watch?v=WykSdgtLDD0)

Pattern Matching @ Scale Using Finite State Machine, by Ajit Koti and Rashmi Shamprasad

Learn how Netflix engineers use Regular Expressions to explore their data.



Implementing Regular Expressions

- We identified a set of language-operators
- We want to explore their expressiveness:
What kind of questions can we pose using that set of operators?
- How do we implement such theory in Coq?

Regular expressions

Inductive definition

```
Inductive regex :=  
  | r_void: regex  
  | r_nil: regex  
  | r_char: Ascii.ascii → regex  
  | r_app: regex → regex → regex  
  | r_union: regex → regex → regex  
  | r_star: regex → regex.
```

Informal description

- **r_void**: represents the **Void** language
- **r_nil**: the empty string **Nil** language
- **r_char**: the **Char** language
- **r_union**: represents the union of two languages
- **r_app**: represents concatenation of languages
- **r_star**: represents zero-or-more copies of an expression

Regular expression Coq notation

Informal description

- `r_void`: the `Void` language
- `r_nil`: the `Nil` language
- `r_char`: the `Char` language
- `r_union`: the `Union` operator
(notation `r1 || r2`)
- `r_app`: the `Append` operator
(notation `r1 ;; r2`)
- `r_star`: the `Star` operator

Exercises

1. Strings with a's and b's that end with "aa"
"aa", "aaa", "baaa", "bbbbbaa"
2. Strings that have at an even number of a's
"aa", "", "aaaa", "aaaaaa"
3. Nonempty strings that only contain any number of a's and b's
"a", "b", "ab", "aaaaa", "bbbbbb", "abaaa"
4. Strings that interleave one "a" with one "b"
"a", "b", "ab", "ba", "aba", "bab", "abab", "baba"

Exercise

Strings with a's and b's that end with "aa"

Examples: "aa", "aaa", "baaa", "bbbbbaa"

Exercise

Strings with a's and b's that end with "aa"

Examples: "aa", "aaa", "baaa", "bbbbbaa"

Solution

$(a||b)^* \cdot aa$

Exercise

Strings that have at an even number of a's

Examples: "aa", "", "aaaa", "aaaaaa"

Exercise

Strings that have at an even number of a's

Examples: "aa", "", "aaaa", "aaaaaa"

Solution

$(aa)^*$

Exercise

Nonempty strings that only contain any number of a's and b's

Examples: "a", "b", "ab", "aaaaa", "bbbbbb", "abaaa"

Exercise

Nonempty strings that only contain any number of a's and b's

Examples: "a", "b", "ab", "aaaaa", "bbbbbb", "abaaa"

Solution

$$(a||b)^* \cdot (a||b)$$

Exercise

Strings that interleave one "a" with one "b"

Examples: "a", "b", "ab", "ba", "aba", "bab", "abab", "baba"

Exercise

Strings that interleave one "a" with one "b"

Examples: "a", "b", "ab", "ba", "aba", "bab", "abab", "baba"

Solution

$$(ab)^* \mid (ab)^* a \mid (ba)^* \mid (ba)^* b$$

Inductive propositions: acceptance

Rules `accept_nil` and `accept_char`

$$\epsilon \in \mathbf{r_nil} \qquad [c] \in \mathbf{c}$$

Rule `accept_app`

$$\frac{w_1 \in R_1 \quad w_2 \in R_2}{w_1 \cdot w_2 \in R_1 ; ; R_2}$$

Rules `accept_union_l` and `accept_union_r`

$$\frac{w \in R_1}{w \in R_1 \parallel R_2} \qquad \frac{w \in R_2}{w \in R_1 \parallel R_2}$$

Rules `accept_star_nil` and `accept_star_cons_neq`

$$\epsilon \in R^* \qquad \frac{w_1 \neq \epsilon \quad w_1 \in R \quad w_2 \in R^*}{w_1 \cdot w_2 \in R^*}$$



Regex.v

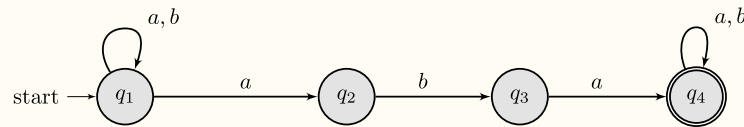
Nondeterministic Finite Automata (NFA)

NFA by example

Strings with a's and b's that end with "aa"

Examples: "aa", "aaa", "baaa", "bbbbbbbaa"

State diagram



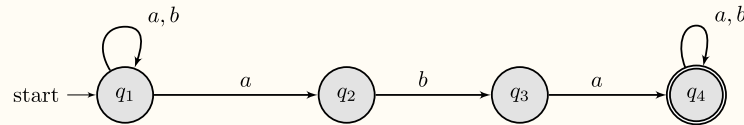
About

- The diagram is a **graph**
- Nodes are called **states**
- Edges are called **transition**
- Accepting a word: a path in the graph
- Initial state, identified **start** →
- Accepting state, double edge
- Consume one character per transition
- Comma in transitions means OR

NFA by example

Strings with a's and b's that end with "aa"

Examples: "aa", "aaa", "baaa", "bbbbbbbaa"



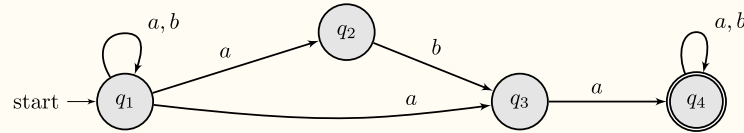
1. In q_1 read as many a's and b's as needed
2. Eventually, read one a and move to q_2
3. Finally, if we are able to read two a's, then we can accept the string (q_3)

As long as we can find **one** path, we can accept the input. There may exist multiple paths in the same state diagram (nondeterminism).

Exercise

Strings that have at an even number of a's

Examples: "aa", "", "aaaa", "aaaaaa"



1. State q_1 accepts the empty string
2. If we consume an **a**, then we have read an odd-number of a's.
Thus, q_2 is non-accepting
3. If we read another **a**, we have read an even-number of a's
Thus, we go back to q_1 , which is an accepting state.

Exercise

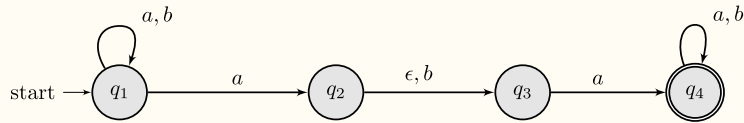
Nonempty strings that only contain any number of a's and b's

Examples: "a", "b", "ab", "aaaaa", "bbbbbb", "abaaa"

Exercise

Nonempty strings that only contain any number of a's and b's

Examples: "a", "b", "ab", "aaaaa", "bbbbbb", "abaaa"



- In state q_1 we can read as many a's as we want
- Eventually, we read at least one a or b and proceed to q_2

Exercise

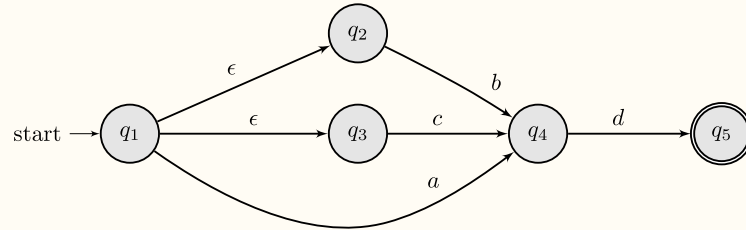
Strings that interleave one "a" with one "b"

Examples: "a", "b", "ab", "ba", "aba", "bab", "abab", "baba"

Exercise

Strings that interleave one "a" with one "b"

Examples: "a", "b", "ab", "ba", "aba", "bab", "abab", "baba"



- We start in an accepting state
- Reading an **a** moves us to q_2 which expects a **b**
- Reading a **b** moves us to q_3 which expects an **a**
- All states are accepting. **However, not all strings are accepted.**

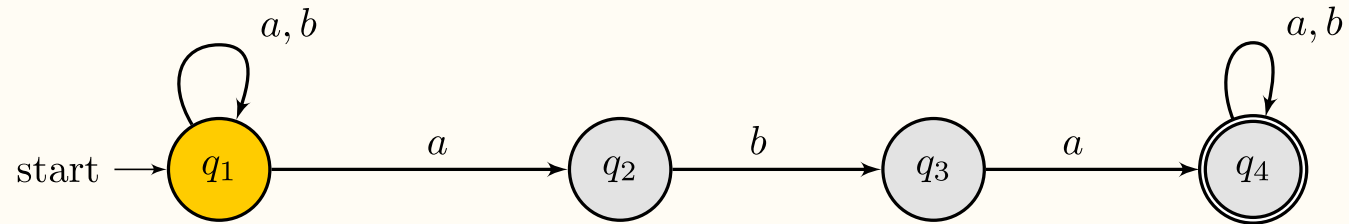
Acceptance in an NFA

Acceptance is path finding

The given string must be a path from the starting node into the accepting node.

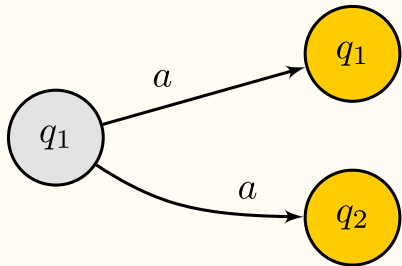
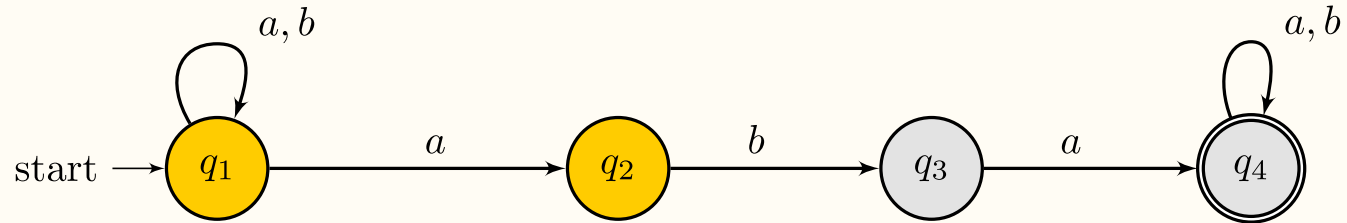
Acceptance in an NFA

Acceptance of **a**bbaba



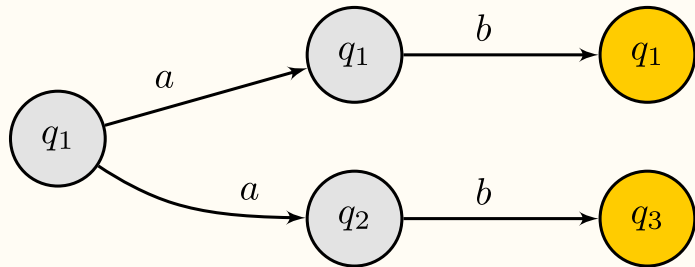
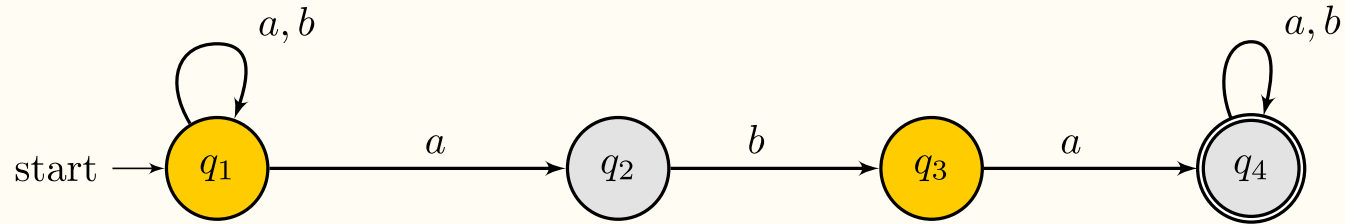
Acceptance in an NFA

Acceptance of **ab**baba



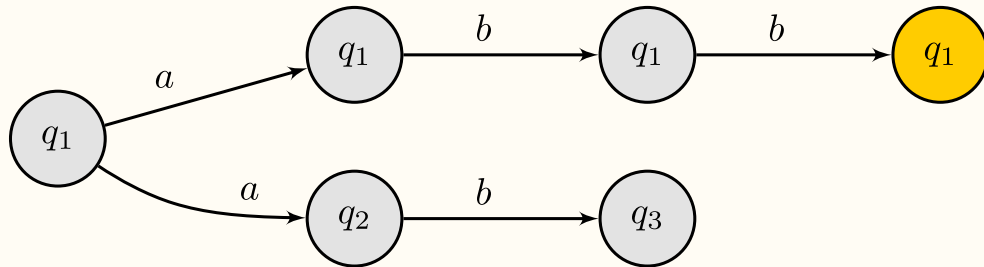
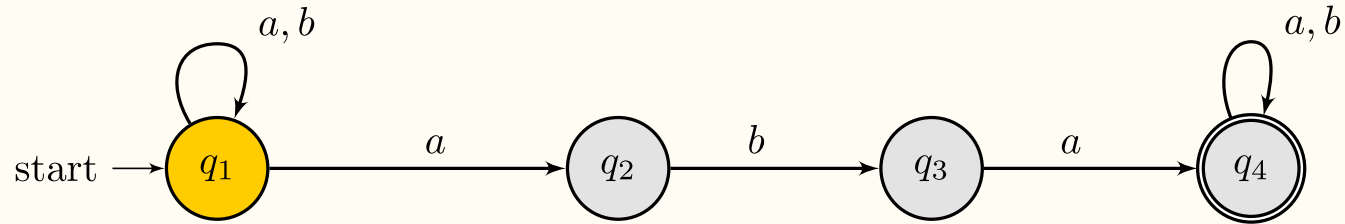
Acceptance in an NFA

Acceptance of **ab**aba



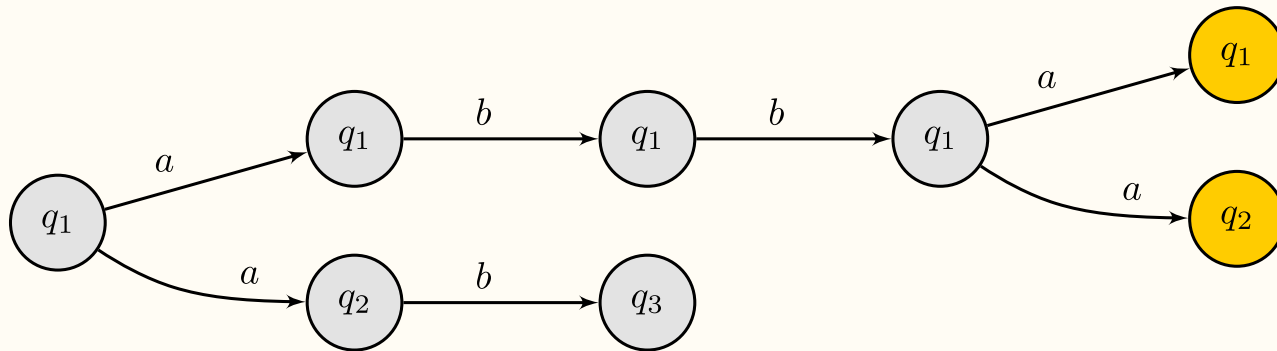
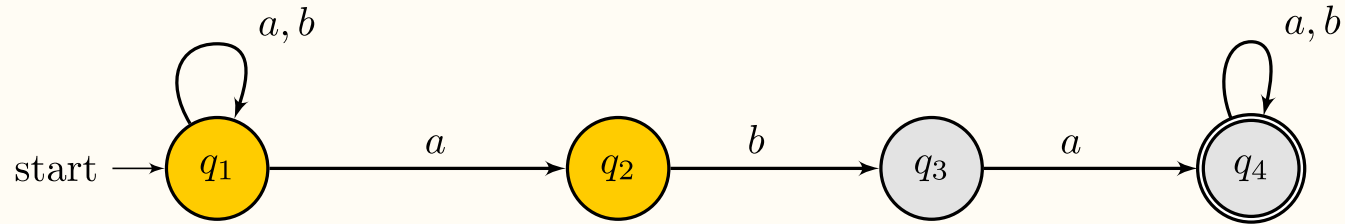
Acceptance in an NFA

Acceptance of **abbaba**



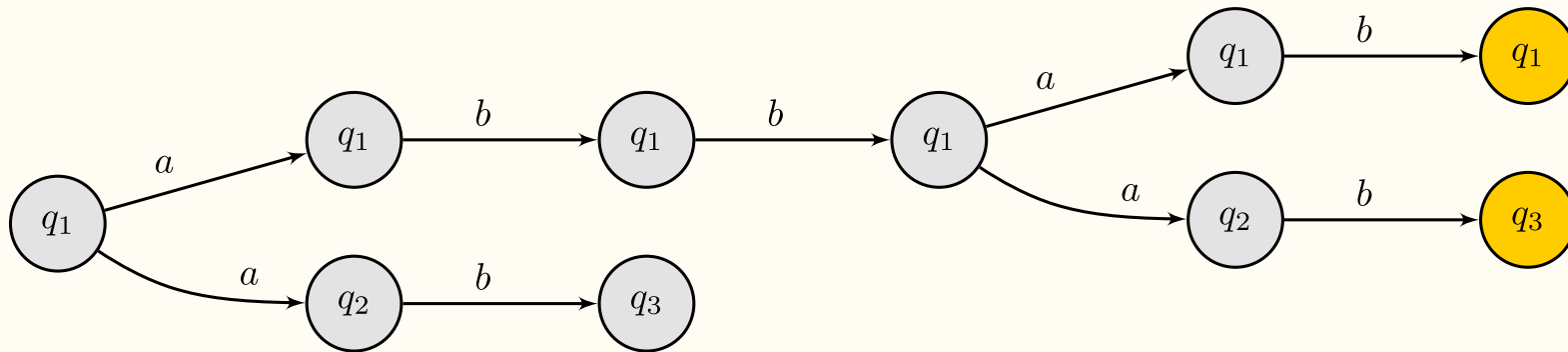
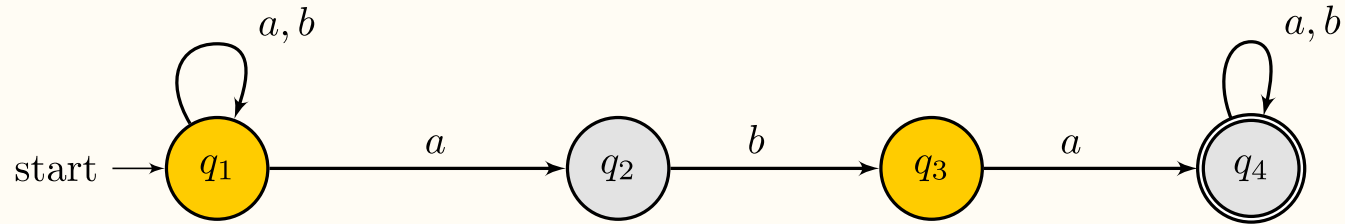
Acceptance in an NFA

Acceptance of **abbaba**



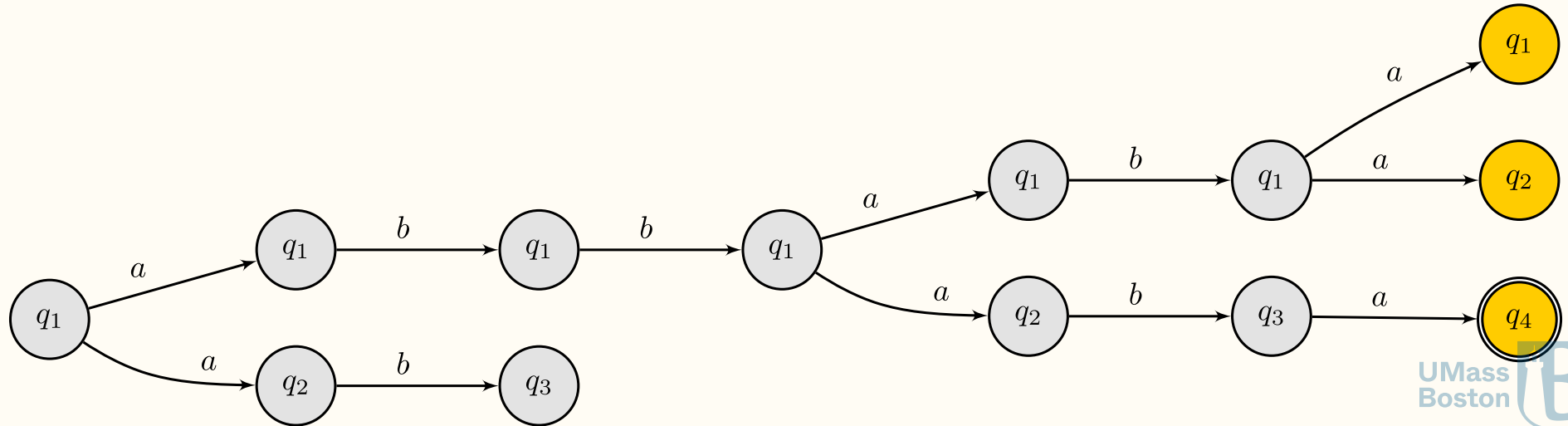
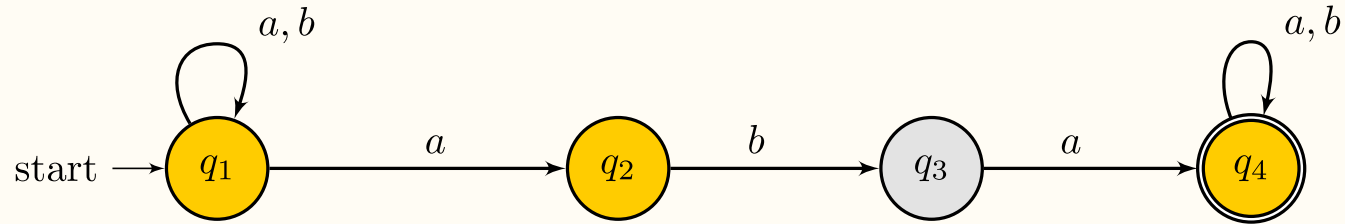
Acceptance in an NFA

Acceptance of **abbaba**



Acceptance in an NFA

Acceptance of `abbaba`



Acceptance in an NFA

- There are multiple concurrent possible paths and a current state
- Given a current state, if there are no transitions for a given input, the path ends
- Once we reach the final path, we check if there are accepting states

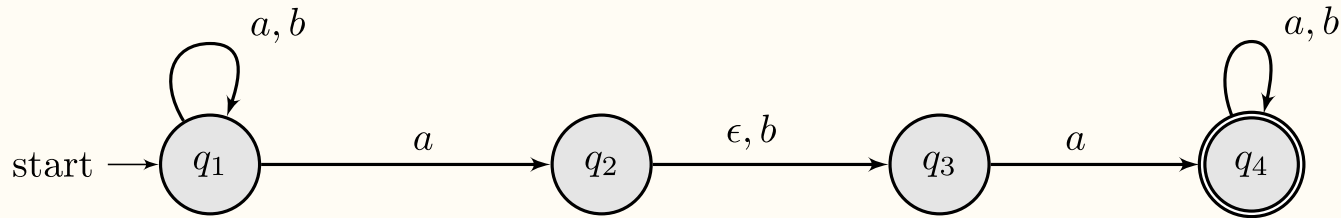
Epsilon transitions

Epsilon transitions

Exercise 2

Let $\Sigma = \{a, b\}$. Give an NFA with four states that recognizes the following language

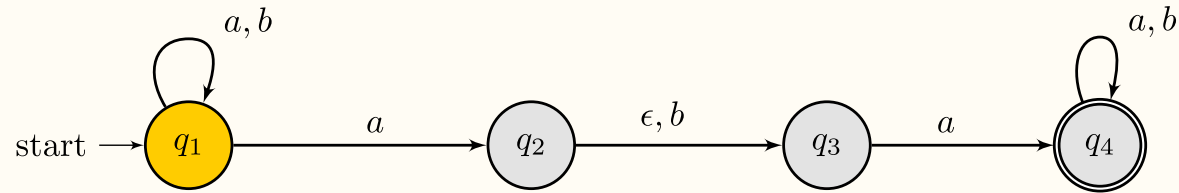
$\{w \mid w \text{ contains the strings } aba \text{ or } aa\}$



Note

- NFAs can also include ϵ transitions, which may be taken without consuming an input

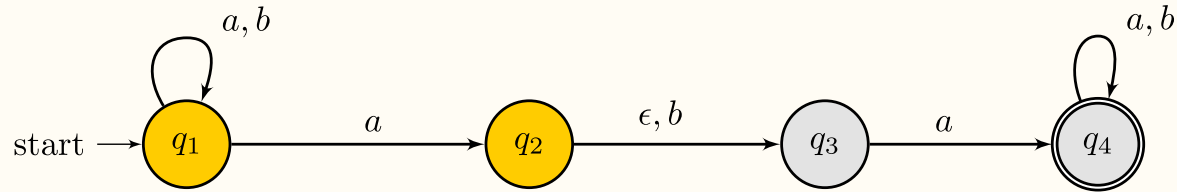
Exercise 2: acceptance of **a**aba



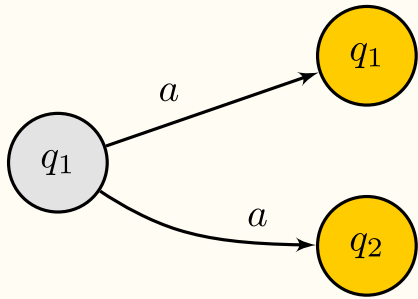
Interleave
input with ϵ .
Read a



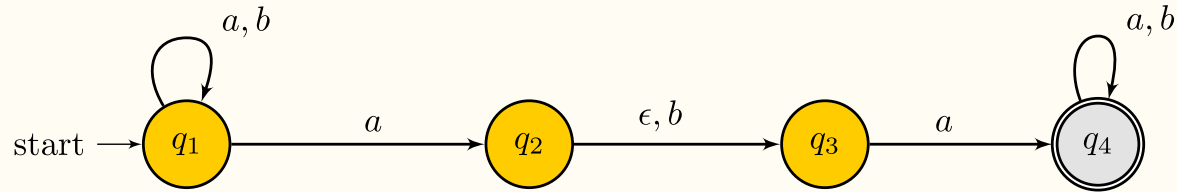
Exercise 2: acceptance of $a\epsilon aba$



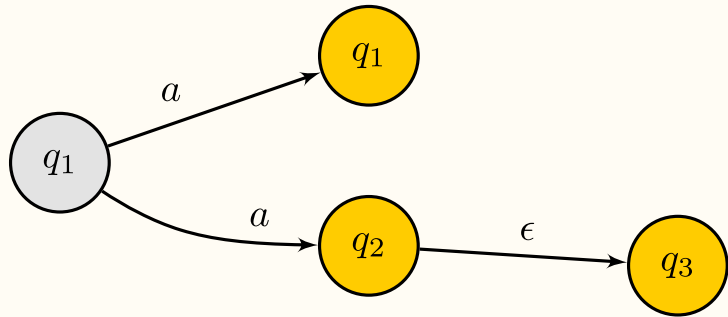
Interleave
input with ϵ .
Read ϵ



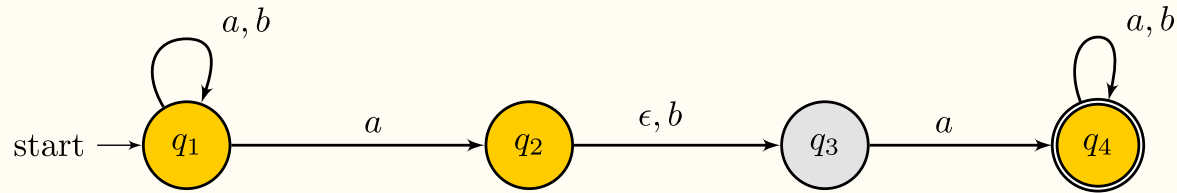
Exercise 2: acceptance of **aa**ba



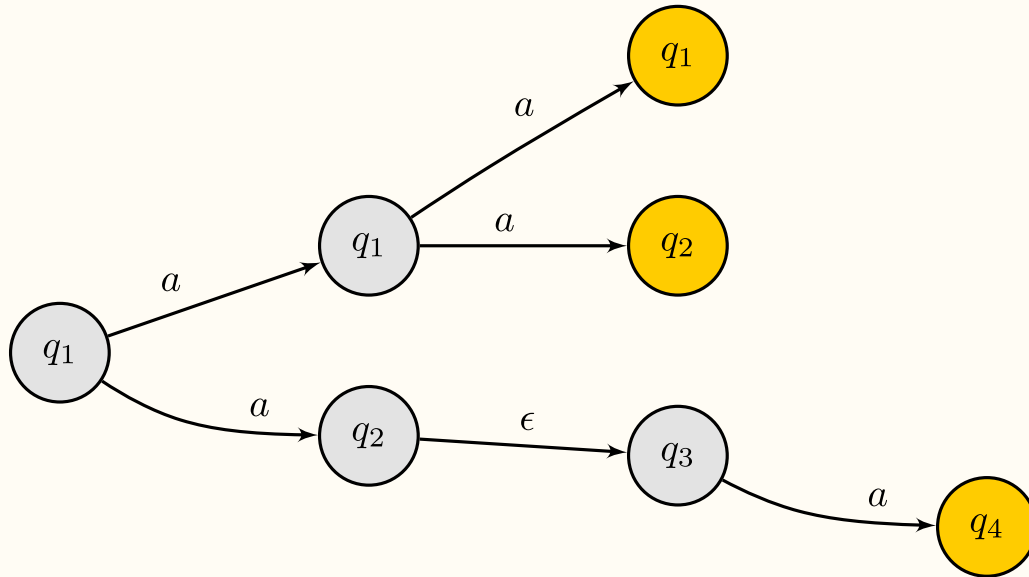
Interleave
input with ϵ .
Read a



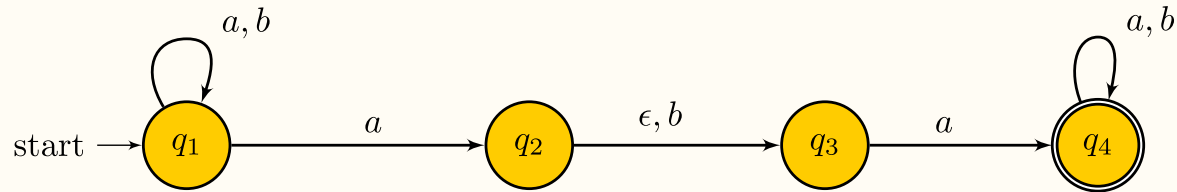
Exercise 2: acceptance of $aab\epsilon a$



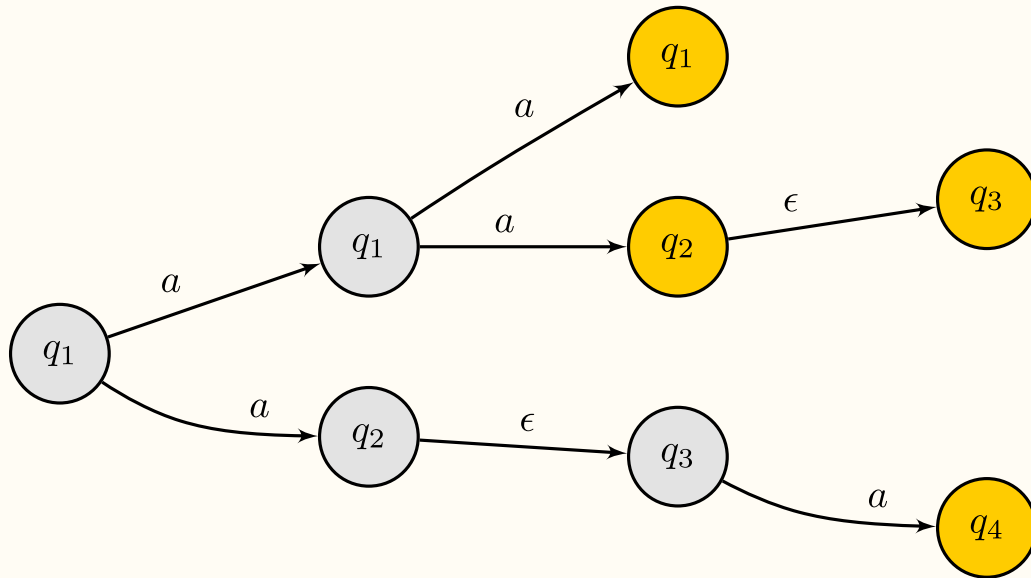
Interleave
input with ϵ .
Read ϵ



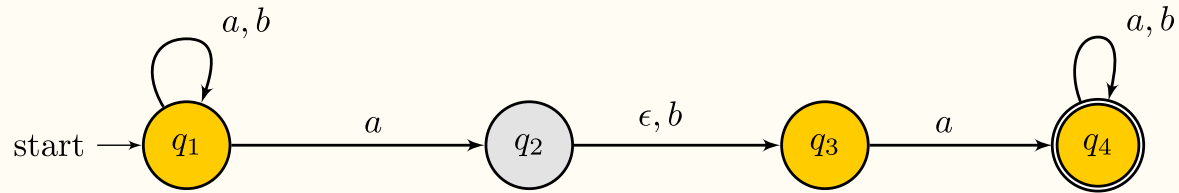
Exercise 2: acceptance of **aa**ba



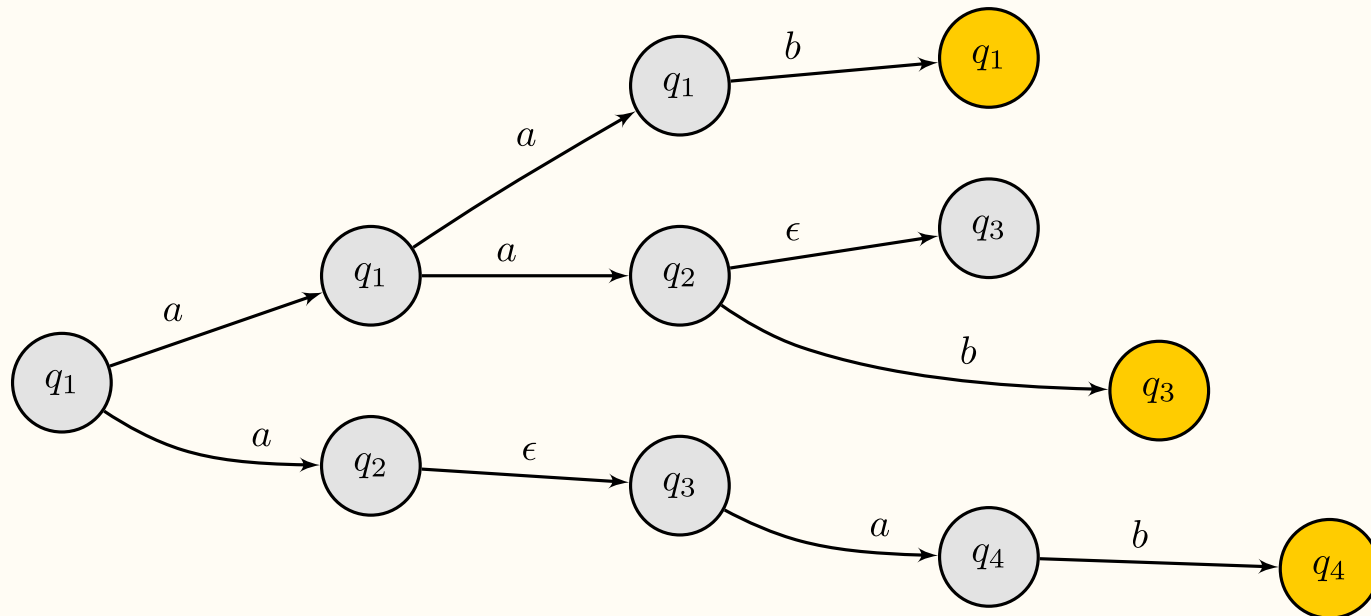
Interleave
input with ϵ .
Read b



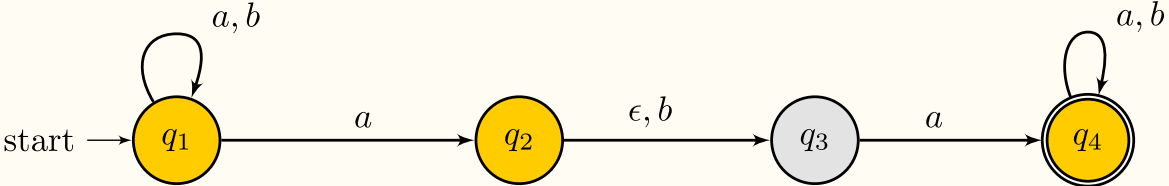
Exercise 2: acceptance of **aa**ba****



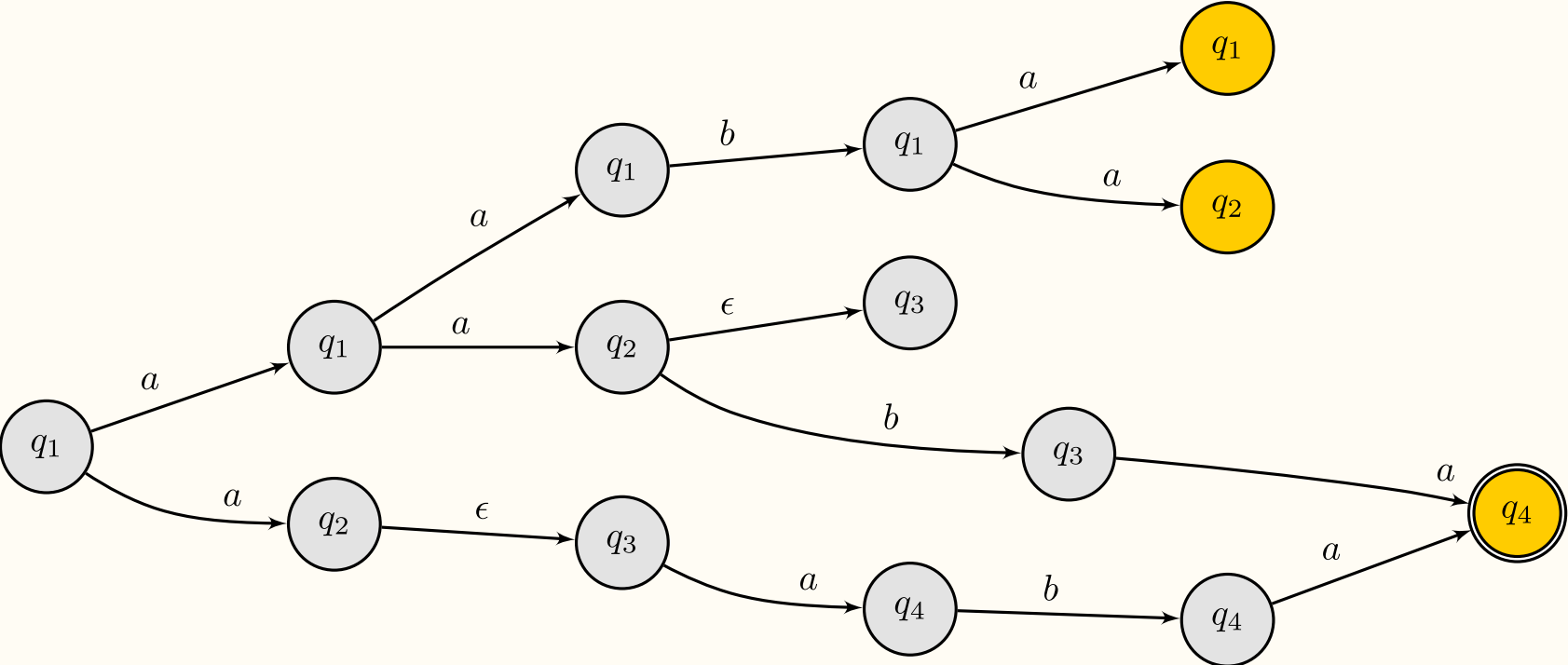
Interleave
input with ϵ .
Read a



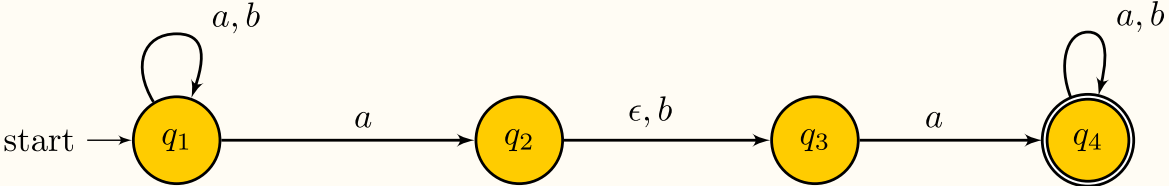
Exercise 2: acceptance of **aaba** ϵ



Interleave
input with ϵ .
Read ϵ

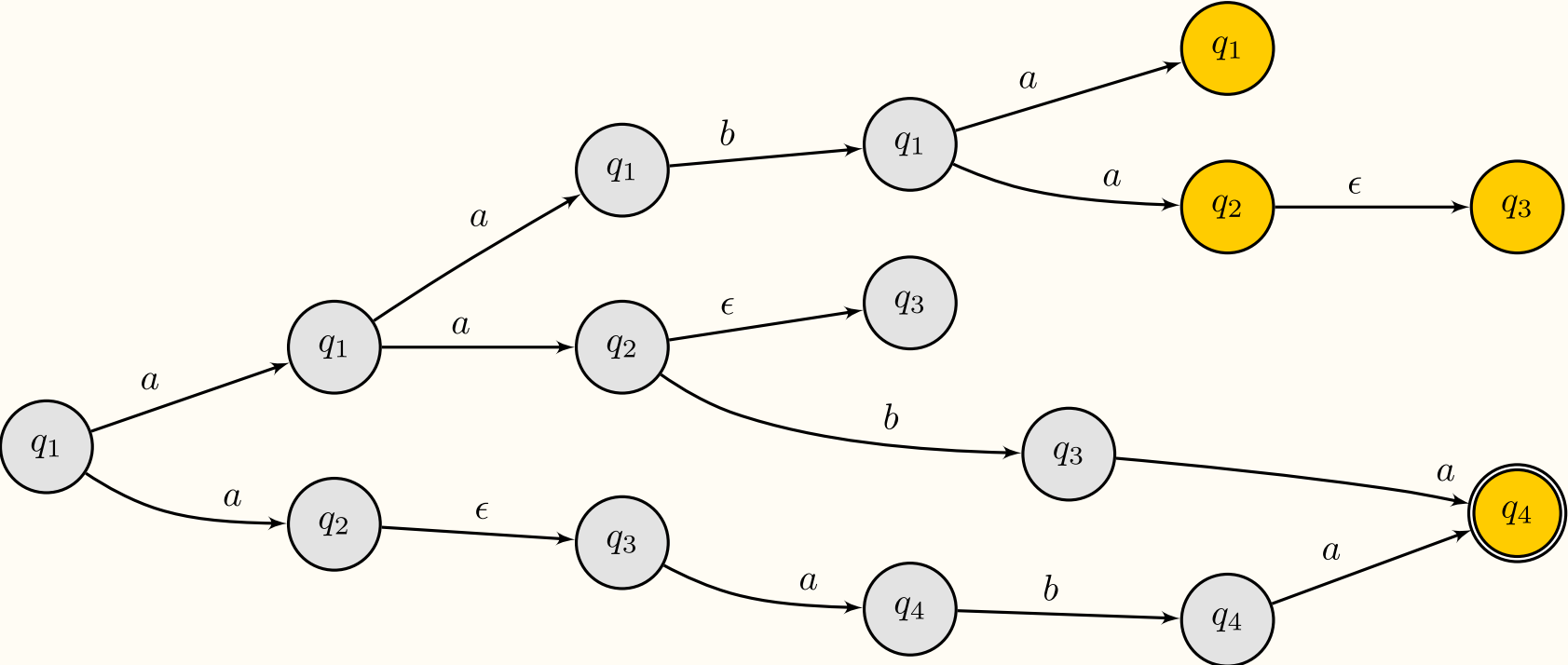


Exercise 2: acceptance of aaba



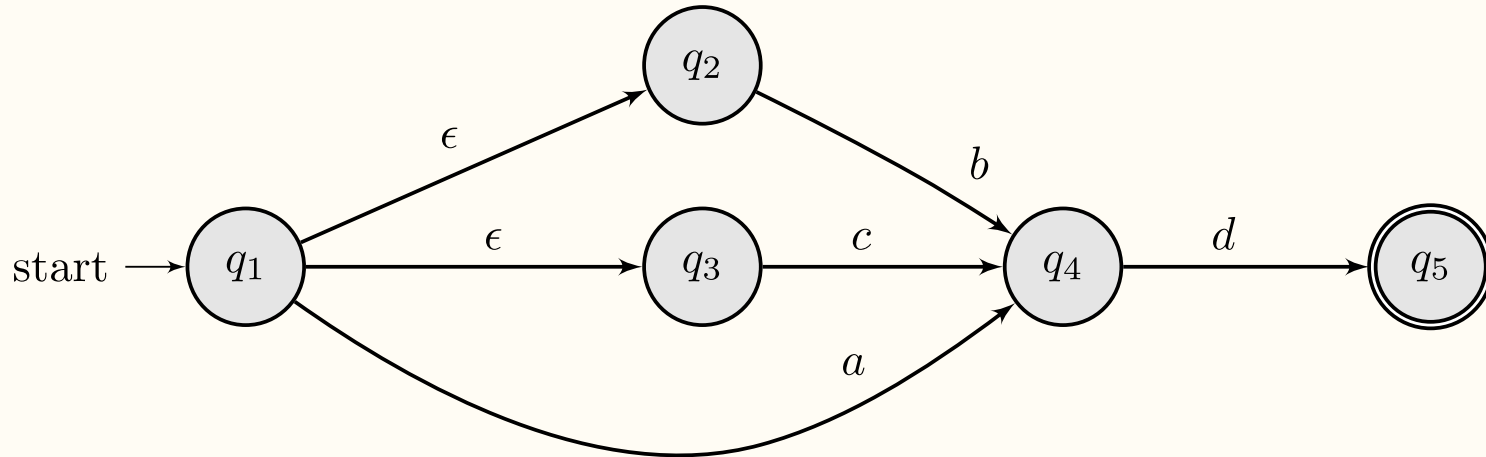
Interleave
input with ϵ .

Read ϵ

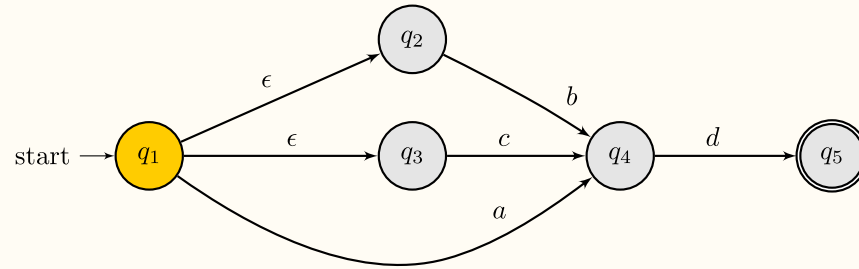


Note ϵ transitions in the initial state

We looked at ϵ in the middle of the state diagram. Let us observe their effect in the initial state.



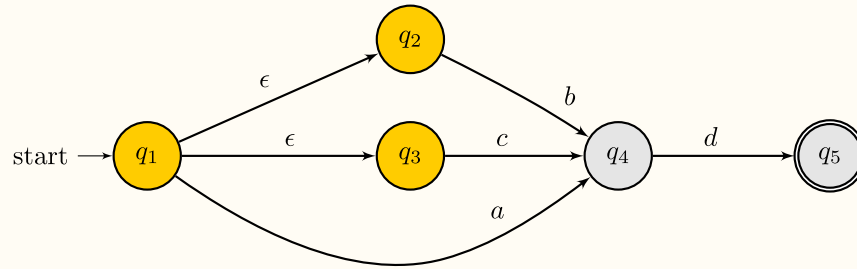
Exercise 2: acceptance of `bd`



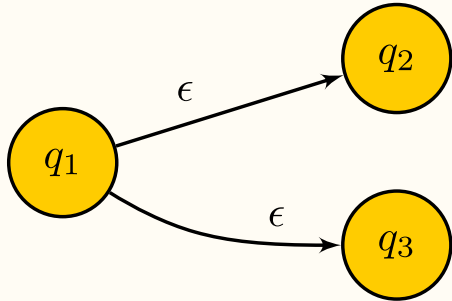
Read ϵ



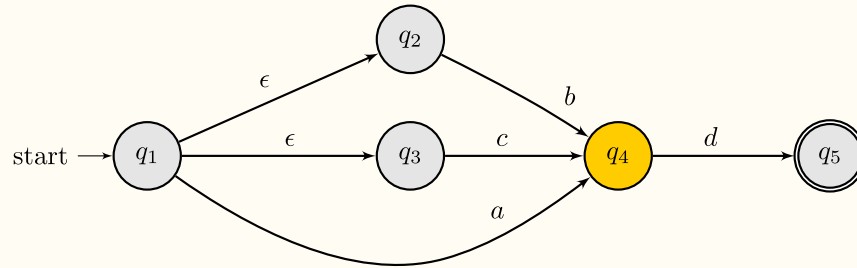
Exercise 2: acceptance of `bd`



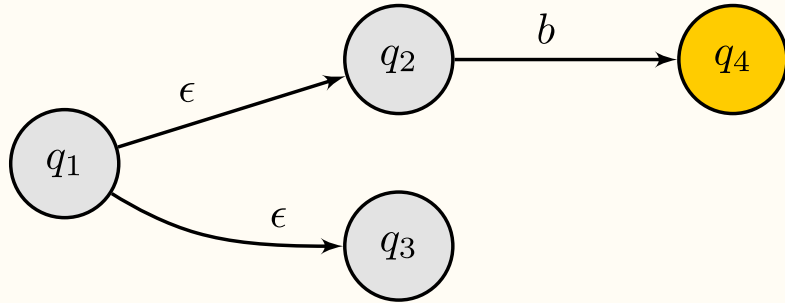
Read b



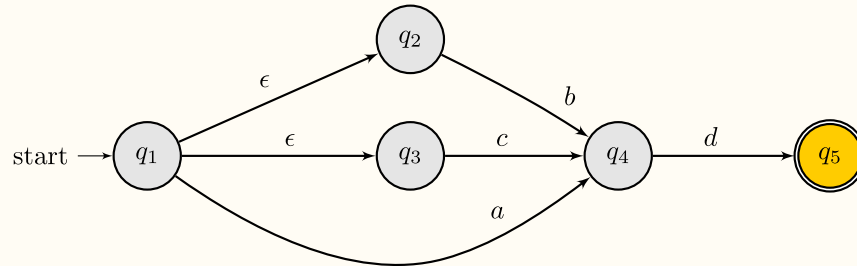
Exercise 2: acceptance of `bd`



**Read ϵ and
then read `d`**



Exercise 2: acceptance of `bd`



Accepted!

