

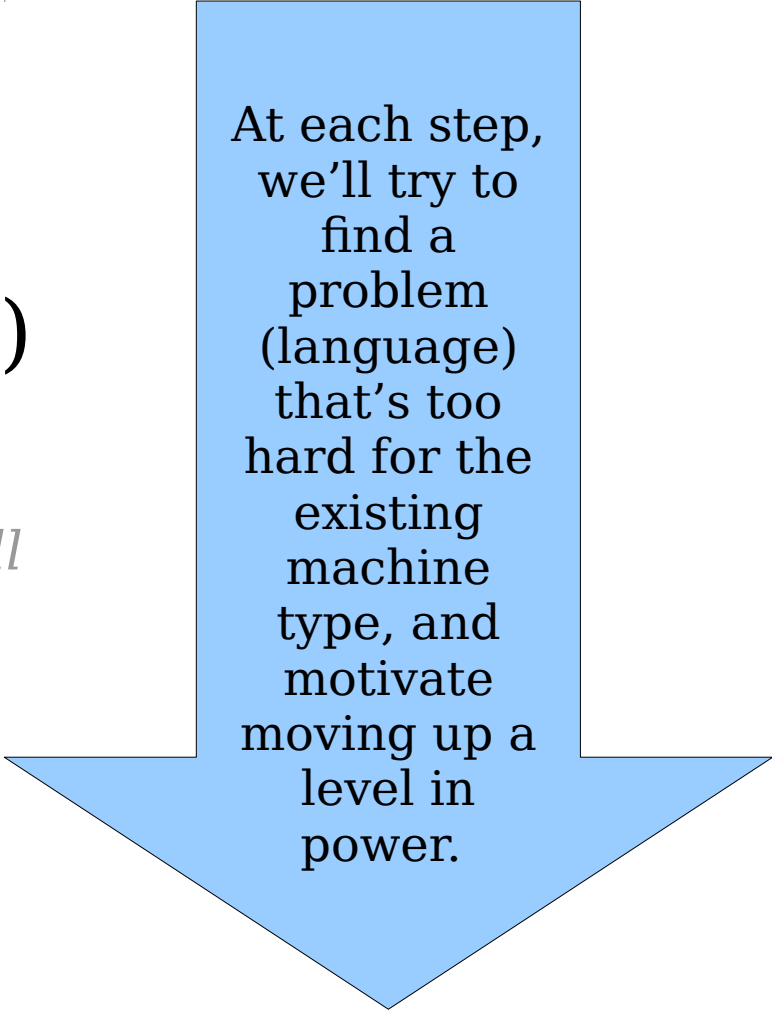
# Regular Expressions

Recap from Last Time

# Automata Progression

- **DFAs** (Deterministic Finite Automata)  
= **NFAs** (DFA + superpowers of guessing)
- **PDA**s (NFA + a stack to use as memory)  
*We'll actually explore an equivalent mechanism called CFGs, you'll see soon...*
- **Turing Machines** (DFA or NFA + an infinite array to use as memory)

**New**



At each step, we'll try to find a problem (language) that's too hard for the existing machine type, and motivate moving up a level in power.

# Language Concatenation

- If  $w \in \Sigma^*$  and  $x \in \Sigma^*$ , then  $wx$  is the **concatenation** of  $w$  and  $x$ .
- If  $L_1$  and  $L_2$  are languages over  $\Sigma$ , the **concatenation** of  $L_1$  and  $L_2$  is the language  $L_1L_2$  defined as

$$L_1L_2 = \{ wx \mid w \in L_1 \text{ and } x \in L_2 \}$$

- Example: if  $L_1 = \{ \mathbf{a}, \mathbf{ba}, \mathbf{bb} \}$  and  $L_2 = \{ \mathbf{aa}, \mathbf{bb} \}$ , then

$$L_1L_2 = \{ \mathbf{aaa}, \mathbf{abb}, \mathbf{baaa}, \mathbf{babb}, \mathbf{bbaa}, \mathbf{bbbb} \}$$

# Lots and Lots of Concatenation

- Consider the language  $L = \{ \mathbf{aa}, \mathbf{b} \}$
- $L^0 = \{\varepsilon\}$
- $LL=L^2$  is the set of strings formed by concatenating pairs of strings in  $L$ .

$\{ \mathbf{aaaa}, \mathbf{aab}, \mathbf{baa}, \mathbf{bb} \}$

- $LLL = L^3$  is the set of strings formed by concatenating triples of strings in  $L$ .

$\{ \mathbf{aaaaaa}, \mathbf{aaaab}, \mathbf{aabaa}, \mathbf{aabb}, \mathbf{baaaa}, \mathbf{baab}, \mathbf{bbaa}, \mathbf{bbb} \}$

- $LLLL = L^4$  is the set of strings formed by concatenating quadruples of strings in  $L$ .

$\{ \mathbf{aaaaaaaa}, \mathbf{aaaaaab}, \mathbf{aaaabaa}, \mathbf{aaaabb}, \mathbf{aabaaaa}, \mathbf{aabaab}, \mathbf{aabbaa}, \mathbf{aabbb}, \mathbf{baaaaaa}, \mathbf{baaaab}, \mathbf{baabaa}, \mathbf{baabb}, \mathbf{bbaaaa}, \mathbf{bbaab}, \mathbf{bbbaa}, \mathbf{bbbb} \}$

# The Kleene Closure

- An important operation on languages is the ***Kleene Closure***, which is defined as

$$L^* = \{ w \in \Sigma^* \mid \exists n \in \mathbb{N}. w \in L^n \}$$

# Closure Properties

- **Theorem:** If  $L_1$  and  $L_2$  are regular languages over an alphabet  $\Sigma$ , then so are the following languages:

- $\bar{L}_1$
- $L_1 \cup L_2$
- $L_1 \cap L_2$
- $L_1 L_2$
- $L_1^*$

- These properties are called ***closure properties of the regular languages.***

## Quick check 1:

Let  $\Sigma = \{1, 2, 3, a, b, c\}$ .  
Let  $L_1 = \{aa, b\}$ ,  $L_2 = \{33, 2\}$   
be languages over  $\Sigma$ .

Name two strings string **in**  $\bar{L}_1$ .

# Closure Properties

- **Theorem:** If  $L_1$  and  $L_2$  are regular languages over an alphabet  $\Sigma$ , then so are the following languages:

- $\bar{L}_1$
- $L_1 \cup L_2$
- $L_1 \cap L_2$
- $L_1 L_2$
- $L_1^*$

- These properties are called ***closure properties of the regular languages.***

## Quick check 2:

Let  $\Sigma = \{1, 2, 3, a, b, c\}$ .  
Let  $L_1 = \{aa, b\}$ ,  $L_2 = \{33, 2\}$   
be languages over  $\Sigma$ .

Name two strings **in**  $L_1 \cup L_2$ .

# Closure Properties

- **Theorem:** If  $L_1$  and  $L_2$  are regular languages over an alphabet  $\Sigma$ , then so are the following languages:

- $\bar{L}_1$
- $L_1 \cup L_2$
- $L_1 \cap L_2$
- $L_1L_2$
- $L_1^*$

- These properties are called ***closure properties of the regular languages.***

### Quick check 3:

Let  $\Sigma = \{1, 2, 3, a, b, c\}$ .  
Let  $L_1 = \{aa, b\}$ ,  $L_2 = \{33, 2\}$   
be languages over  $\Sigma$ .

Name two strings **in**  $L_1 \cap L_2$ .

# Closure Properties

- **Theorem:** If  $L_1$  and  $L_2$  are regular languages over an alphabet  $\Sigma$ , then so are the following languages:

- $\bar{L}_1$
- $L_1 \cup L_2$
- $L_1 \cap L_2$
- $L_1L_2$
- $L_1^*$

- These properties are called **closure properties of the regular languages**.

## Quick check 4:

Let  $\Sigma = \{1, 2, 3, a, b, c\}$ .  
Let  $L_1 = \{aa, b\}$ ,  $L_2 = \{33, 2\}$   
be languages over  $\Sigma$ .

Name two strings **in**  $L_1L_2$ .

# Closure Properties

- **Theorem:** If  $L_1$  and  $L_2$  are regular languages over an alphabet  $\Sigma$ , then so are the following languages:

- $\bar{L}_1$
- $L_1 \cup L_2$
- $L_1 \cap L_2$
- $L_1 L_2$
- $L_1^*$

- These properties are called ***closure properties of the regular languages.***

## Quick check 5:

Let  $\Sigma = \{1, 2, 3, a, b, c\}$ .  
Let  $L_1 = \{aa, b\}$ ,  $L_2 = \{33, 2\}$   
be languages over  $\Sigma$ .

Name two strings **in**  $L_1^*$ .

# Defining Regular Languages

- We currently have several tools for showing a language  $L$  is regular:
  - Construct a **DFA** for  $L$ .
    - *We used this for a variety of languages, as well as proving closure for complement.*
  - Construct an **NFA** for  $L$ .
    - *We used this for a variety of languages, as well as proving closure for union, concatenation, and Kleene  $*$ .*
  - Combine several simpler regular languages together via **closure properties** to form  $L$ .
    - *We used this for proving closure for intersection.*
    - *(Today we'll greatly expand on this last idea.)*

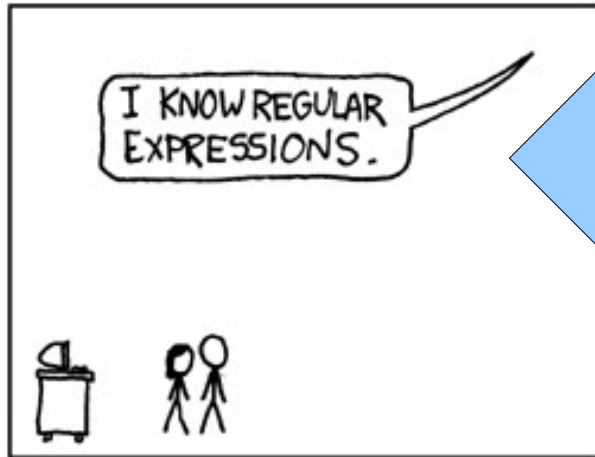
New Stuff!

Another View of Regular Languages:

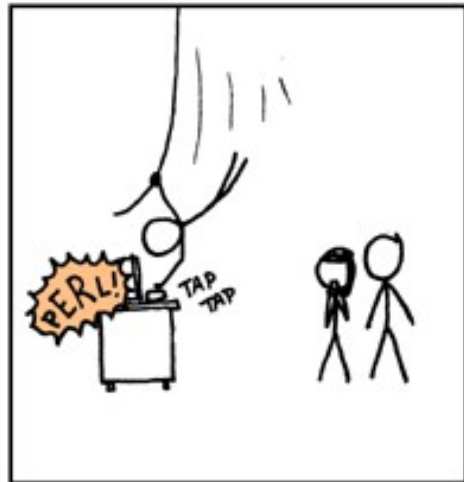
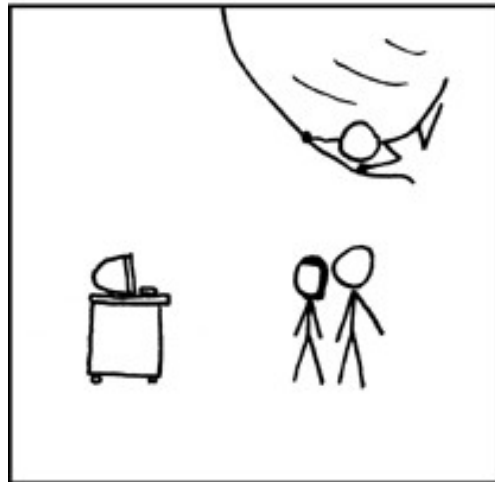
Regular Expressions (regex)

xkcd

WHENEVER I LEARN A NEW SKILL I CONCOCT ELABORATE FANTASY SCENARIOS WHERE IT LETS ME SAVE THE DAY.



80 minutes from now, you will be able to say this



# Regular Expressions

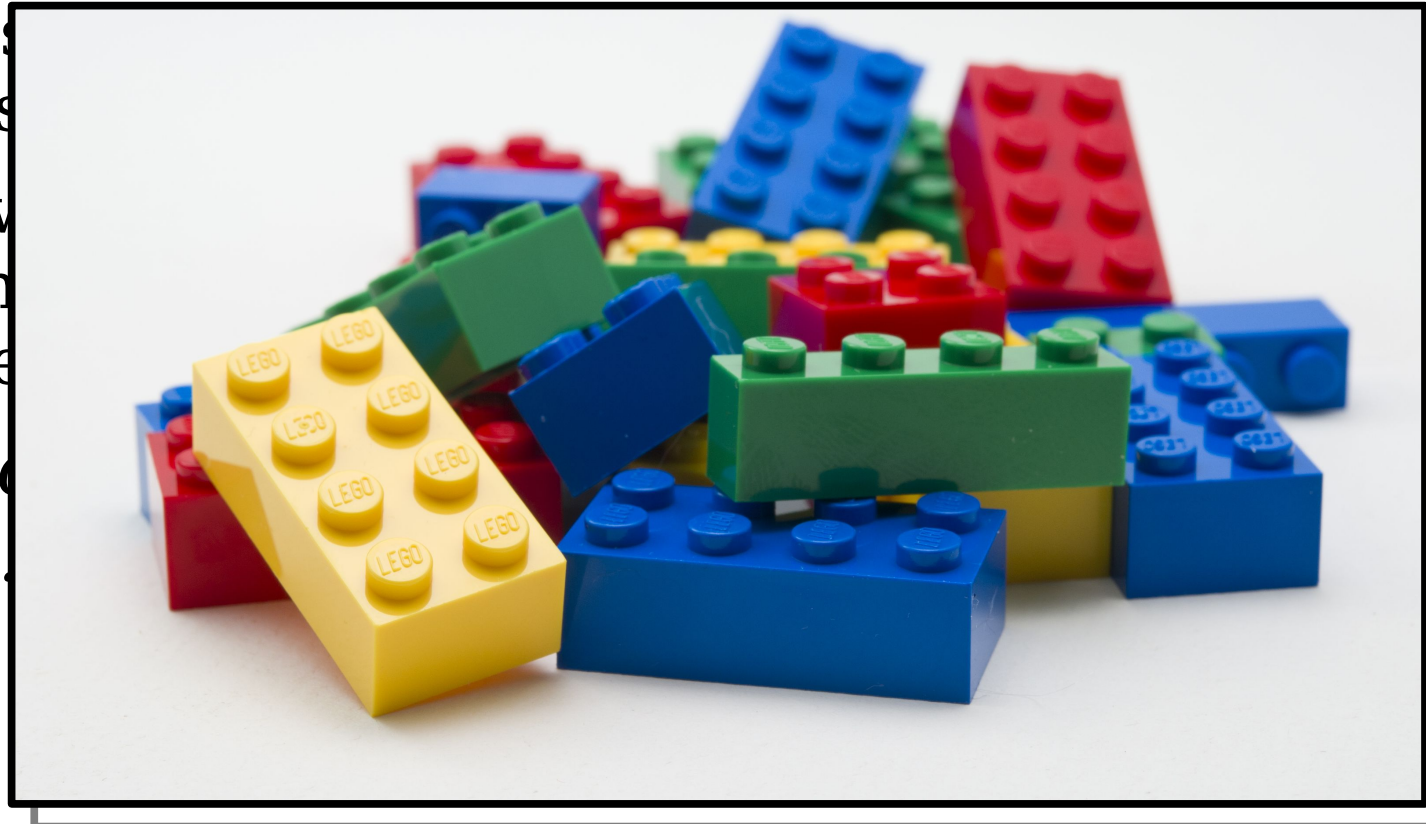
- ***Regular expressions*** are a way of describing a language via a string pattern-match representation.
  - They're used all over in practical day-to-day programming:
    - Data validation tools in just about every programming language (python, Javascript, C++, etc)
    - UNIX grep and flex tools to search files
    - Compilers

# Constructing Regular Languages with RegEx

- **Idea:** Build up all regular languages as follows:
  - Start with a small set of super-simple languages we already know to be regular.
  - Using closure properties, combine these simple languages together to form more languages.
  - Because we only use closure properties, the results of even the most elaborate combinations are guaranteed to be regular.
- *This is a bottom-up approach to the regular languages.*

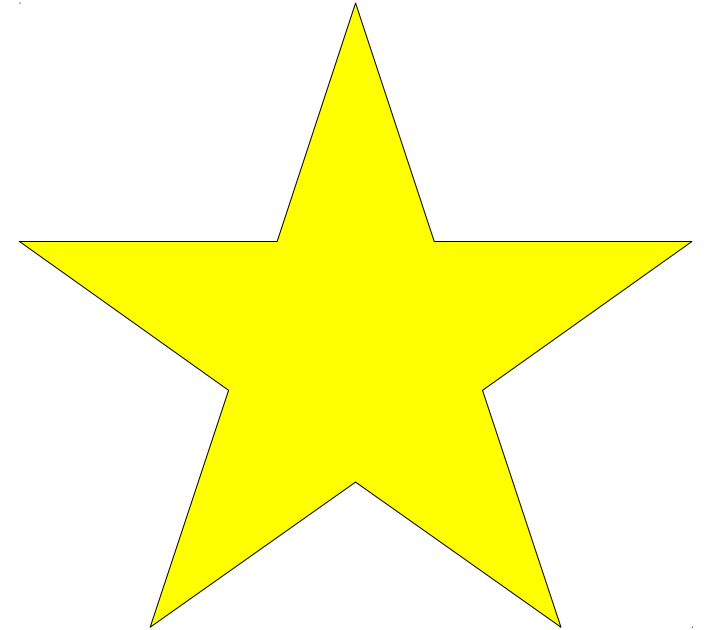
# Constructing Regular Languages with RegEx

- **Idea:** Build up all regular languages as follows:
  - Start with a small set of super-simple languages we already know to be regular.
  - Using closure operations on regular languages.
  - Because we can build up all regular languages from a small set of even the simplest regular languages, we are guaranteed that all regular languages are built up from a small set of languages.
- *This is a building block for regular languages.*



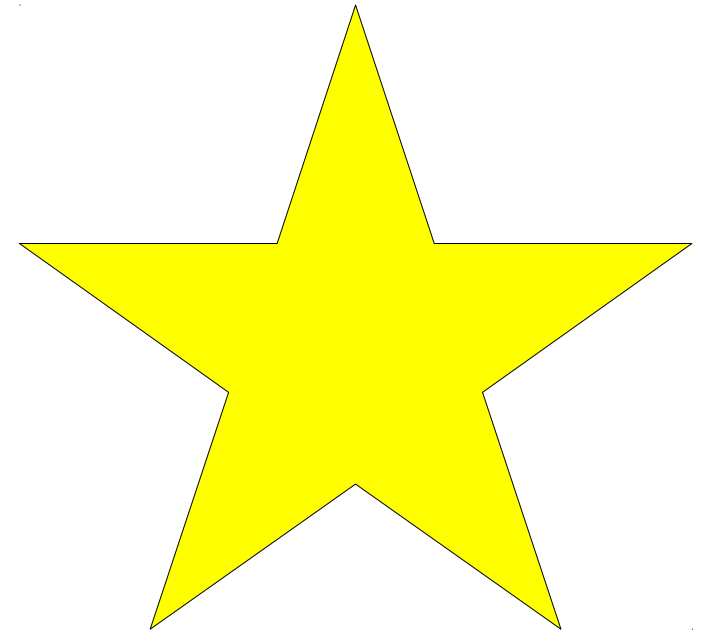
# RegEx “Starter Languages”

- The regular expressions begin with **three starter languages** (*a bit like base cases in induction!*)
  1. The symbol  $\emptyset$  is a regular expression that represents the empty language  $\emptyset$ .
  2. For any  $a \in \Sigma$ , the symbol  $a$  is a regular expression for the language  $\{a\}$ . (*So this “starter language” is actually many languages, one for each character in  $\Sigma$ .*)
  3. The symbol  $\epsilon$  is a regular expression that represents the language  $\{\epsilon\}$ .



# Combining RegEx

- Let  $R_1$  and  $R_2$  be regular expressions (either one of the starter languages, or a more complex regex). We can combine them in these ways:
- **Concatenation:**  $R_1R_2$
- **Union:**  $R_1 \cup R_2$
- **Kleene closure:**  $R_1^*$
- **Parentheses:**  $(R_1)$  may be used to control order of operations, and for reading clarity.



# Operator Precedence

- Here's the operator precedence for regular expressions:

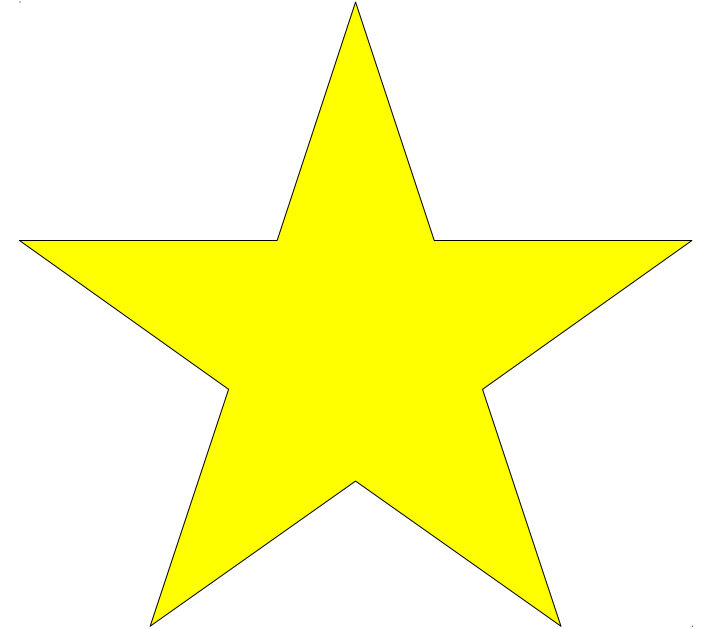
$(R_1)$

$R_1^*$

$R_1R_2$

$R_1 \cup R_2$

- So **ab\*cUd** is parsed as **((a(b\*))c)Ud**



# Regular Expression Examples

- The regular expression **trickUtreat** represents the language

{ **trick, treat** }.

- The regular expression **booo\*** represents the regular language

{ **boo, booo, boooo, ...** }.

- The regular expression **candy!(candy!)\*** represents the regular language

{ **candy!, candy!candy!, candy!candy!candy!,  
...**  }.

# Regular Expressions, Formally

- The **language of a regular expression** is the language described by that regular expression.
- Formally:
  - $\mathcal{L}(\epsilon) = \{\epsilon\}$
  - $\mathcal{L}(\emptyset) = \emptyset$
  - $\mathcal{L}(a) = \{a\}$
  - $\mathcal{L}(R_1R_2) = \mathcal{L}(R_1) \mathcal{L}(R_2)$
  - $\mathcal{L}(R_1 \cup R_2) = \mathcal{L}(R_1) \cup \mathcal{L}(R_2)$
  - $\mathcal{L}(R^*) = \mathcal{L}(R)^*$
  - $\mathcal{L}((R)) = \mathcal{L}(R)$

## Regex quick check:

Let  $\Sigma = \{a, b, c, d\}$ .

Let  $L_1 = \mathcal{L}(a(b \cup c)((d)))$  be a language over  $\Sigma$ .

Name two strings **in**  $L_1$ .

Name one string **not in**  $L_1$ .

# Designing Regular Expressions

- Let  $\Sigma = \{\mathbf{a}, \mathbf{b}\}$ .
- Let  $L = \{ w \in \Sigma^* \mid w \text{ contains } \mathbf{aa} \text{ as a substring} \}$ .

# Designing Regular Expressions

- Let  $\Sigma = \{\mathbf{a}, \mathbf{b}\}$ .
- Let  $L = \{ w \in \Sigma^* \mid w \text{ contains } \mathbf{aa} \text{ as a substring} \}$ .

**$(\mathbf{a} \cup \mathbf{b})^* \mathbf{aa} (\mathbf{a} \cup \mathbf{b})^*$**

# Designing Regular Expressions

- Let  $\Sigma = \{\mathbf{a}, \mathbf{b}\}$ .
- Let  $L = \{ w \in \Sigma^* \mid w \text{ contains } \mathbf{aa} \text{ as a substring} \}$ .

$(\mathbf{a} \cup \mathbf{b})^* \mathbf{aa} (\mathbf{a} \cup \mathbf{b})^*$

# Designing Regular Expressions

- Let  $\Sigma = \{a, b\}$ .
- Let  $L = \{ w \in \Sigma^* \mid w \text{ contains } aa \text{ as a substring} \}$ .

$(a \cup b)^*aa(a \cup b)^*$

**bbabbaabab**

**aaaa**

**bbbbabbbbaabbbb**

# Designing Regular Expressions

- Let  $\Sigma = \{\mathbf{a}, \mathbf{b}\}$ .
- Let  $L = \{ w \in \Sigma^* \mid w \text{ contains } \mathbf{aa} \text{ as a substring} \}$ .

$(\mathbf{a} \cup \mathbf{b})^* \mathbf{aa} (\mathbf{a} \cup \mathbf{b})^*$

**bbabbaabab**

**aaaa**

**bbbbabbbbaabbbb**

# Designing Regular Expressions

- Let  $\Sigma = \{\mathbf{a}, \mathbf{b}\}$ .
- Let  $L = \{ w \in \Sigma^* \mid w \text{ contains } \mathbf{aa} \text{ as a substring} \}$ .

$\Sigma^* \mathbf{aa} \Sigma^*$

$\mathbf{bbabbb} \mathbf{aa} \mathbf{bab}$

$\mathbf{aaaa}$

$\mathbf{bbbbabbbb} \mathbf{aa} \mathbf{bbbbbb}$

# Designing Regular Expressions

- Let  $\Sigma = \{\mathbf{a}, \mathbf{b}\}$ .
- Let  $L = \{ w \in \Sigma^* \mid |w| = 4 \}$ .

# Designing Regular Expressions

Let  $\Sigma = \{\mathbf{a}, \mathbf{b}\}$ .

Let  $L = \{ w \in \Sigma^* \mid |w| = 4 \}$ .

The length of a  
string  $w$  is  
denoted  $|w|$

# Designing Regular Expressions

- Let  $\Sigma = \{\mathbf{a}, \mathbf{b}\}$ .
- Let  $L = \{ w \in \Sigma^* \mid |w| = 4 \}$ .

# Designing Regular Expressions

- Let  $\Sigma = \{\mathbf{a}, \mathbf{b}\}$ .
- Let  $L = \{ w \in \Sigma^* \mid |w| = 4 \}$ .

**$\Sigma\Sigma\Sigma\Sigma$**

# Designing Regular Expressions

- Let  $\Sigma = \{\mathbf{a}, \mathbf{b}\}$ .
- Let  $L = \{ w \in \Sigma^* \mid |w| = 4 \}$ .

$\Sigma\Sigma\Sigma\Sigma$

# Designing Regular Expressions

- Let  $\Sigma = \{\mathbf{a}, \mathbf{b}\}$ .
- Let  $L = \{ w \in \Sigma^* \mid |w| = 4 \}$ .

$\Sigma\Sigma\Sigma\Sigma$

**aaaa**

**baba**

**bbbb**

**baaa**

# Designing Regular Expressions

- Let  $\Sigma = \{a, b\}$ .
- Let  $L = \{ w \in \Sigma^* \mid |w| = 4 \}$ .

$\Sigma\Sigma\Sigma\Sigma$

$aaaa$

$baba$

$bbbb$

$baaa$

# Designing Regular Expressions

- Let  $\Sigma = \{\mathbf{a}, \mathbf{b}\}$ .
- Let  $L = \{ w \in \Sigma^* \mid |w| = 4 \}$ .

$\Sigma^4$

**a****a****a****a**

**b****a****b****a**

**b****b****b****b**

**b****a****a****a**

# Designing Regular Expressions

- Let  $\Sigma = \{\mathbf{a}, \mathbf{b}\}$ .
- Let  $L = \{ w \in \Sigma^* \mid |w| = 4 \}$ .

$\Sigma^4$

**aaaa**

**baba**

**bbbb**

**baaa**

# Designing Regular Expressions

- Let  $\Sigma = \{\mathbf{a}, \mathbf{b}\}$ .
- Let  $L = \{ w \in \Sigma^* \mid w \text{ contains at most one } \mathbf{a} \}$ .

Here are some candidate regular expressions for the language  $L$ . **How many** of these are correct? (Discuss specifically which with your neighbors.)

$\Sigma^* \mathbf{a} \Sigma^*$

$\mathbf{b}^* \mathbf{a} \mathbf{b}^* \cup \mathbf{b}^*$

$\mathbf{b}^* (\mathbf{a} \cup \epsilon) \mathbf{b}^*$

$\mathbf{b}^* \mathbf{a}^* \mathbf{b}^* \cup \mathbf{b}^*$

$\mathbf{b}^* (\mathbf{a}^* \cup \epsilon) \mathbf{b}^*$

# Designing Regular Expressions

- Let  $\Sigma = \{\mathbf{a}, \mathbf{b}\}$ .
- Let  $L = \{ w \in \Sigma^* \mid w \text{ contains at most one } \mathbf{a} \}$ .

$$\mathbf{b^*(a \cup \epsilon)b^*}$$

# Designing Regular Expressions

- Let  $\Sigma = \{\mathbf{a}, \mathbf{b}\}$ .
- Let  $L = \{ w \in \Sigma^* \mid w \text{ contains at most one } \mathbf{a} \}$ .

$\mathbf{b}^* (\mathbf{a} \cup \epsilon) \mathbf{b}^*$

# Designing Regular Expressions

- Let  $\Sigma = \{\mathbf{a}, \mathbf{b}\}$ .
- Let  $L = \{ w \in \Sigma^* \mid w \text{ contains at most one } \mathbf{a} \}$ .

$\mathbf{b^* (a \cup \epsilon) b^*}$

**bbbbabbb**

**bbbbbb**

**abbb**

**a**

# Designing Regular Expressions

- Let  $\Sigma = \{a, b\}$ .
- Let  $L = \{ w \in \Sigma^* \mid w \text{ contains at most one } a \}$ .

$b^*(a \cup \epsilon)b^*$

bbbbabb

bbbbbb

abbb

a

# Designing Regular Expressions

- Let  $\Sigma = \{\mathbf{a}, \mathbf{b}\}$ .
- Let  $L = \{ w \in \Sigma^* \mid w \text{ contains at most one } \mathbf{a} \}$ .

**b\*a?b\***

**bbbbabbb**

**bbbbbb**

**abbb**

**a**

# A More Elaborate Design

- Let  $\Sigma = \{ \mathbf{a}, \mathbf{.}, \mathbf{@} \}$ , where  $\mathbf{a}$  represents “some letter.”
- Let's make a regex for email addresses.

# A More Elaborate Design

- Let  $\Sigma = \{ \mathbf{a}, \mathbf{.}, \mathbf{@} \}$ , where  $\mathbf{a}$  represents “some letter.”
- Let's make a regex for email addresses.

**cs103@cs.stanford.edu**

**first.middle.last@mail.site.org**

**dot.at@dot.com**

# A More Elaborate Design

- Let  $\Sigma = \{ \mathbf{a}, \cdot, @ \}$ , where **a** represents “some letter.”
- Let's make a regex for email addresses.

**cs103**@cs.stanford.edu

**first**.middle.last@mail.site.org

**dot**.at@dot.com

# A More Elaborate Design

- Let  $\Sigma = \{ \mathbf{a}, \cdot, @ \}$ , where **a** represents “some letter.”
- Let's make a regex for email addresses.

**aa\***

**cs103**@cs.stanford.edu

**first**.middle.last@mail.site.org

**dot**.at@dot.com

# A More Elaborate Design

- Let  $\Sigma = \{ \mathbf{a}, \mathbf{.}, \mathbf{@} \}$ , where **a** represents “some letter.”
- Let's make a regex for email addresses.

**aa\***

**cs103**@cs.stanford.edu

**first**.**middle**.**last**@mail.site.org

**dot**.**at**@dot.com

# A More Elaborate Design

- Let  $\Sigma = \{ \mathbf{a}, \mathbf{.}, \mathbf{@} \}$ , where **a** represents “some letter.”
- Let's make a regex for email addresses.

**aa\*(.aa\*)\***

**cs103**@cs.stanford.edu

**first**.**middle**.**last**@mail.site.org

**dot**.**at**@dot.com

# A More Elaborate Design

- Let  $\Sigma = \{ \mathbf{a}, \mathbf{.}, \mathbf{@} \}$ , where **a** represents “some letter.”
- Let's make a regex for email addresses.

**aa\*(.aa\*)\***

**cs103@cs.stanford.edu**

**first.middle.last@mail.site.org**

**dot.at@dot.com**

# A More Elaborate Design

- Let  $\Sigma = \{ \mathbf{a}, \mathbf{.}, \mathbf{@} \}$ , where **a** represents “some letter.”
- Let's make a regex for email addresses.

**aa\*(.aa\*)\*@**

**cs103@cs.stanford.edu**

**first.middle.last@mail.site.org**

**dot.at@dot.com**

# A More Elaborate Design

- Let  $\Sigma = \{ \mathbf{a}, \mathbf{.}, \mathbf{@} \}$ , where **a** represents “some letter.”
- Let's make a regex for email addresses.

**aa\*(.aa\*)\*@**

**cs103@cs.stanford.edu**

**first.middle.last@mail.site.org**

**dot.at@dot.com**

# A More Elaborate Design

- Let  $\Sigma = \{ \mathbf{a}, \mathbf{.}, \mathbf{@} \}$ , where **a** represents “some letter.”
- Let's make a regex for email addresses.

**aa\*(.aa\*)\*@aa\*.aa\***

**cs103@cs.stanford.edu**

**first.middle.last@mail.site.org**

**dot.at@dot.com**

# A More Elaborate Design

- Let  $\Sigma = \{ \mathbf{a}, \mathbf{.}, \mathbf{@} \}$ , where **a** represents “some letter.”
- Let's make a regex for email addresses.

**aa\*(.aa\*)\*@aa\*.aa\***

**cs103@cs.stanford.edu**

**first.middle.last@mail.site.org**

**dot.at@dot.com**

# A More Elaborate Design

- Let  $\Sigma = \{ \mathbf{a}, \mathbf{.}, \mathbf{@} \}$ , where  $\mathbf{a}$  represents “some letter.”
- Let's make a regex for email addresses.

**$\mathbf{aa^*(.aa^*)^*@aa*.aa*(.aa^*)^*$**

**$\mathbf{cs103@cs.stanford.edu}$**

**$\mathbf{first.middle.last@mail.site.org}$**

**$\mathbf{dot.at@dot.com}$**

# A More Elaborate Design

- Let  $\Sigma = \{ \mathbf{a}, \mathbf{.}, \mathbf{@} \}$ , where **a** represents “some letter.”
- Let's make a regex for email addresses.

```
aa*(.aa*)*@aa*.aa*(.aa*)*
```

**cs103@cs.stanford.edu**

**first.middle.last@mail.site.org**

**dot.at@dot.com**

# A More Elaborate Design

- Let  $\Sigma = \{ \mathbf{a}, \mathbf{.}, \mathbf{@} \}$ , where  $\mathbf{a}$  represents “some letter.”
- Let's make a regex for email addresses.

```
a+ (.aa*)*@aa*.aa*(.aa*)*
```

**cs103@cs.stanford.edu**

**first.middle.last@mail.site.org**

**dot.at@dot.com**

# A More Elaborate Design

- Let  $\Sigma = \{ \mathbf{a}, \mathbf{.}, \mathbf{@} \}$ , where **a** represents “some letter.”
- Let's make a regex for email addresses.

**a<sup>+</sup> (.aa\*)\* @aa\*.aa\*(.aa\*)\***

**cs103@cs.stanford.edu**

**first.middle.last@mail.site.org**

**dot.at@dot.com**

# A More Elaborate Design

- Let  $\Sigma = \{ \mathbf{a}, \mathbf{.}, \mathbf{@} \}$ , where  $\mathbf{a}$  represents “some letter.”
- Let's make a regex for email addresses.

$\mathbf{a^+ (.a^+)^* @ a^+ .a^+ (.a^+)^*}$

**cs103@cs.stanford.edu**

**first.middle.last@mail.site.org**

**dot.at@dot.com**

# A More Elaborate Design

- Let  $\Sigma = \{ \mathbf{a}, \mathbf{.}, \mathbf{@} \}$ , where  $\mathbf{a}$  represents “some letter.”
- Let's make a regex for email addresses.

$\mathbf{a}^+ (\mathbf{.a}^+)^* \mathbf{@} \mathbf{a}^+ \boxed{\mathbf{.a}^+ (\mathbf{.a}^+)^*}$

**cs103@cs.stanford.edu**

**first.middle.last@mail.site.org**

**dot.at@dot.com**

# A More Elaborate Design

- Let  $\Sigma = \{ \mathbf{a}, \mathbf{.}, \mathbf{@} \}$ , where  $\mathbf{a}$  represents “some letter.”
- Let's make a regex for email addresses.

$\mathbf{a}^+ (\mathbf{.a}^+)^* \mathbf{@} \mathbf{a}^+ \boxed{\mathbf{.a}^+ (\mathbf{.a}^+)^*}$

**cs103@cs.stanford.edu**

**first.middle.last@mail.site.org**

**dot.at@dot.com**

# A More Elaborate Design

- Let  $\Sigma = \{ \mathbf{a}, \mathbf{.}, \mathbf{@} \}$ , where  $\mathbf{a}$  represents “some letter.”
- Let's make a regex for email addresses.

$\mathbf{a}^+ (\mathbf{.a}^+)^* \mathbf{@} \mathbf{a}^+ \boxed{\mathbf{.a}^+ (\mathbf{.a}^+)^*}$

**cs103@cs.stanford.edu**

**first.middle.last@mail.site.org**

**dot.at@dot.com**

# A More Elaborate Design

- Let  $\Sigma = \{ \mathbf{a}, \mathbf{.}, \mathbf{@} \}$ , where  $\mathbf{a}$  represents “some letter.”
- Let's make a regex for email addresses.

$\mathbf{a}^+ (\mathbf{.a}^+)^* \mathbf{@} \mathbf{a}^+ (\mathbf{.a}^+)^+$

**cs103@cs.stanford.edu**

**first.middle.last@mail.site.org**

**dot.at@dot.com**

# A More Elaborate Design

- Let  $\Sigma = \{ \mathbf{a}, \mathbf{.}, \mathbf{@} \}$ , where  $\mathbf{a}$  represents “some letter.”
- Let's make a regex for email addresses.

$\mathbf{a}^+ (\mathbf{.a}^+)^* \mathbf{@} \mathbf{a}^+ (\mathbf{.a}^+)^+$

**cs103@cs.stanford.edu**

**first.middle.last@mail.site.org**

**dot.at@dot.com**

# A More Elaborate Design

- Let  $\Sigma = \{ \mathbf{a}, \mathbf{.}, \mathbf{@} \}$ , where  $\mathbf{a}$  represents “some letter.”
- Let's make a regex for email addresses.

$\mathbf{a^+}(\mathbf{.a^+})^*\mathbf{@a^+}(\mathbf{.a^+})^+$

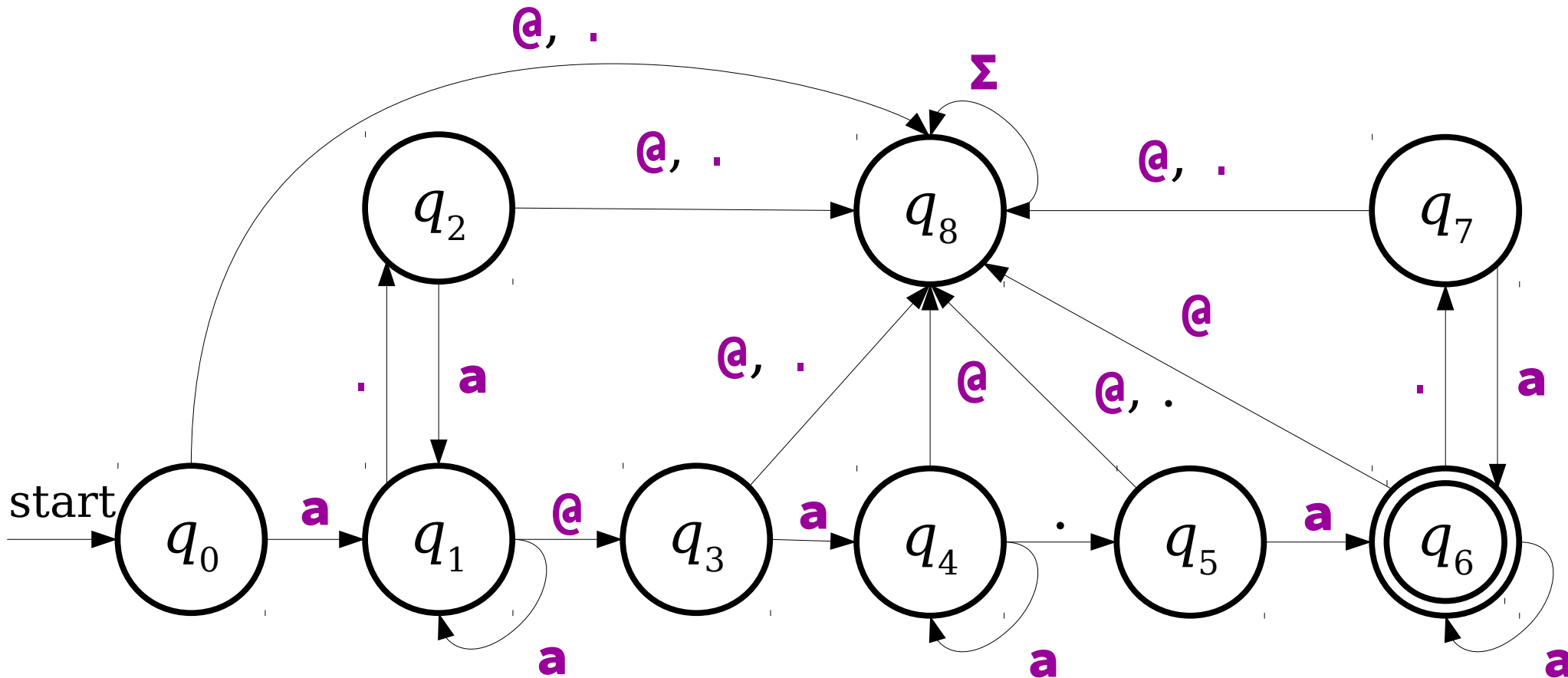
**cs103@cs.stanford.edu**

**first.middle.last@mail.site.org**

**dot.at@dot.com**

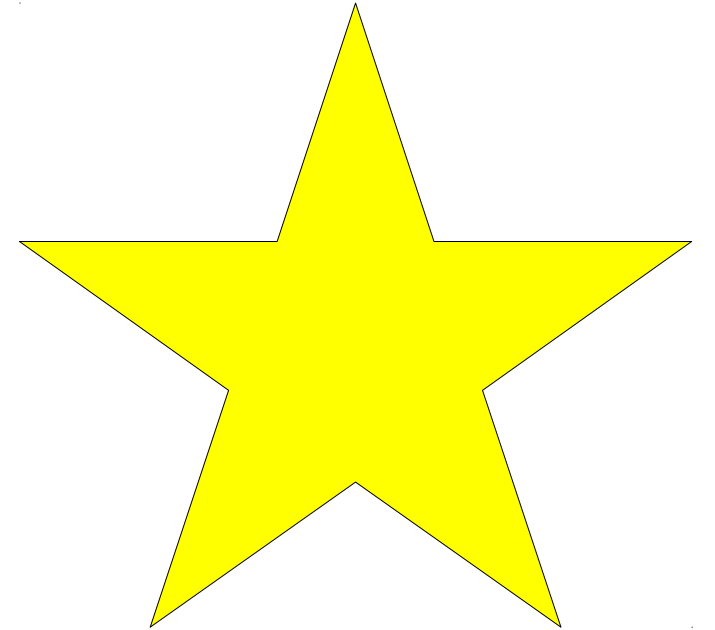
# For Comparison

$a^+ (\cdot a^+) * @ a^+ (\cdot a^+)^+$



# RegEx Shorthand Summary

- $R^n$  is shorthand for  $RR \dots R$  ( $n$  times).
  - Edge case: define  $R^0 = \epsilon$ .
- $\Sigma$  is shorthand for “any character in  $\Sigma$ .”
- $R?$  is shorthand for  $(R \cup \epsilon)$ , meaning “zero or one copies of  $R$ .”
- $R^+$  is shorthand for  $RR^*$ , meaning “one or more copies of  $R$ .”



# The Lay of the Land

Languages you can  
build a DFA for.

Languages you can  
build an NFA for.

***Regular  
Languages***

**Where do Regex  
languages fit?**

# The Power of Regular Expressions

***Theorem:*** If  $R$  is a regular expression, then  $\mathcal{L}(R)$  is regular.

***Proof idea:*** Use induction!

- **Base cases:** The atomic regular expressions all represent regular languages. (*Remember I said they were like induction base cases? Well they literally are!*)
- The combination steps represent closure properties.
- So anything you can make from them must be regular!

# Thompson's Algorithm

- In practice, many regex matchers use an algorithm called ***Thompson's algorithm*** to convert regular expressions into NFAs (and, from there, to DFAs).
  - Read Sipser if you're curious!
- ***Fun fact:*** the “Thompson” here is Ken Thompson, one of the co-inventors of Unix!

Languages you can  
build a DFA for.

Languages you can  
build an NFA for.

***Regular  
Languages***

Languages You Can  
Write a Regex For



Languages you can  
build a DFA for.

Languages you can  
build an NFA for.

***Regular  
Languages***

Languages You Can  
Write a Regex For



Languages you can  
build a DFA for.

Languages you can  
build an NFA for.

***Regular  
Languages***

Languages You Can  
Write a Regex For



Languages you can  
build a DFA for.

Languages you can  
build an NFA for.

***Regular  
Languages***

Languages You Can  
Write a Regex For

**This  
one?  
or**

Languages you can  
build a DFA for.

Languages you can  
build an NFA for.

***Regular  
Languages***

Languages You Can  
Write a Regex For

this  
one?

# The Power of Regular Expressions

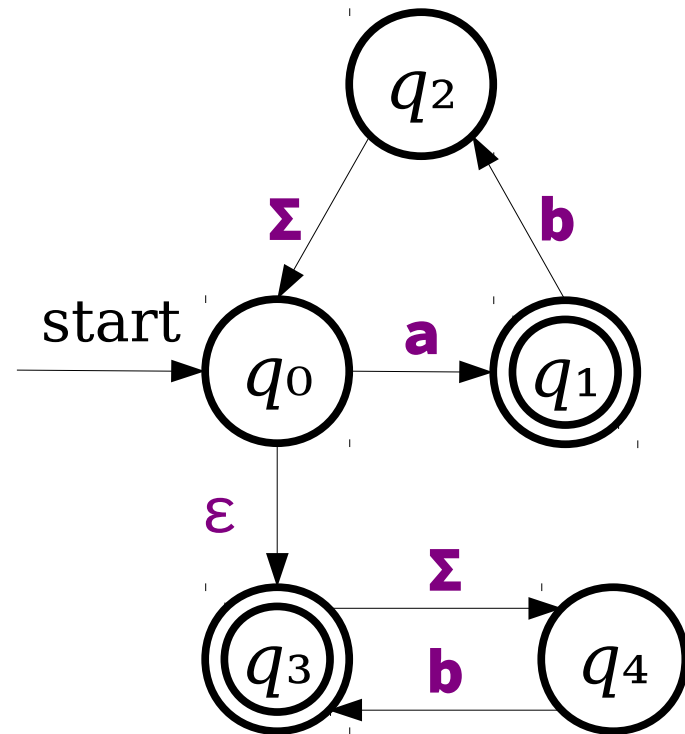
***Theorem:*** For all languages  $L$ , if  $L$  is a regular language, then there is a regular expression for  $L$ .

***This is not obvious!***

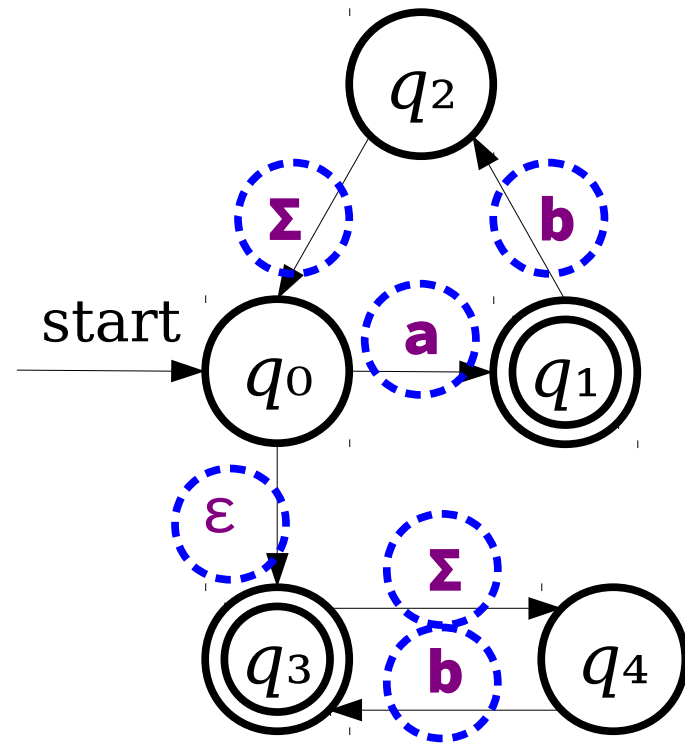
***Proof plan:***

- ***Assume:*** Pick an arbitrary regular language  $L$ .
- ***Want to show:*** We want to show there exists a regex for  $L$ .
- Since  $L$  is regular, there exists an NFA  $N$  for  $L$ .
- Describe how to transform  $N$  into a Regex, in a generalizable way. *For the next few slides, we'll explore how to do that.*

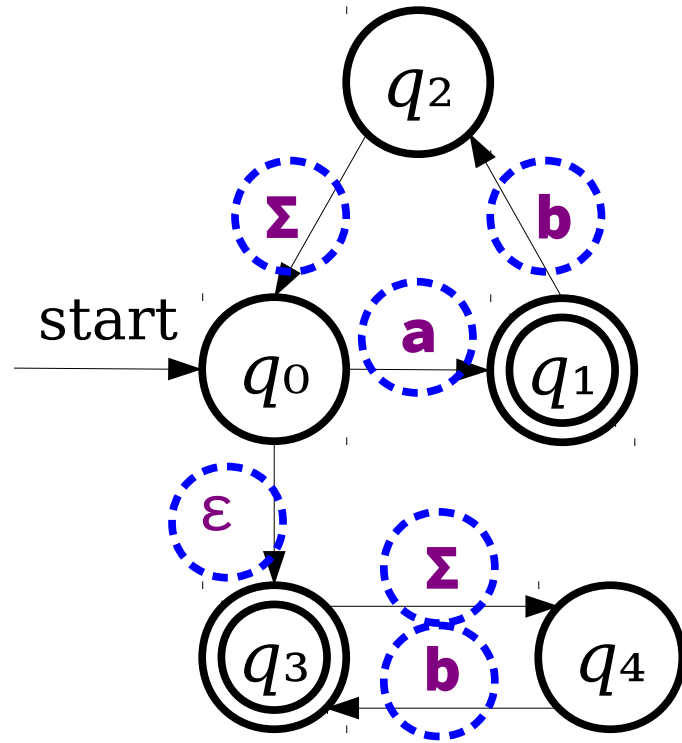
# Generalizing NFAs



# Generalizing NFAs

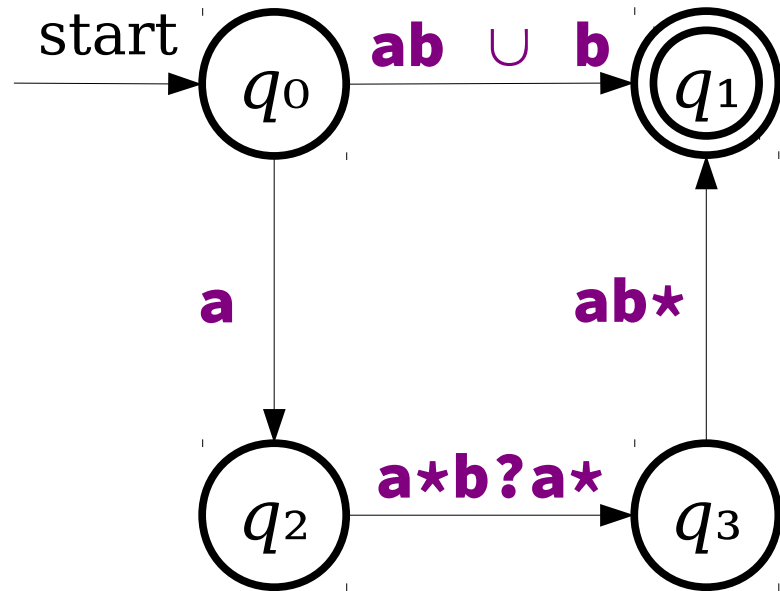


# Generalizing NFAs

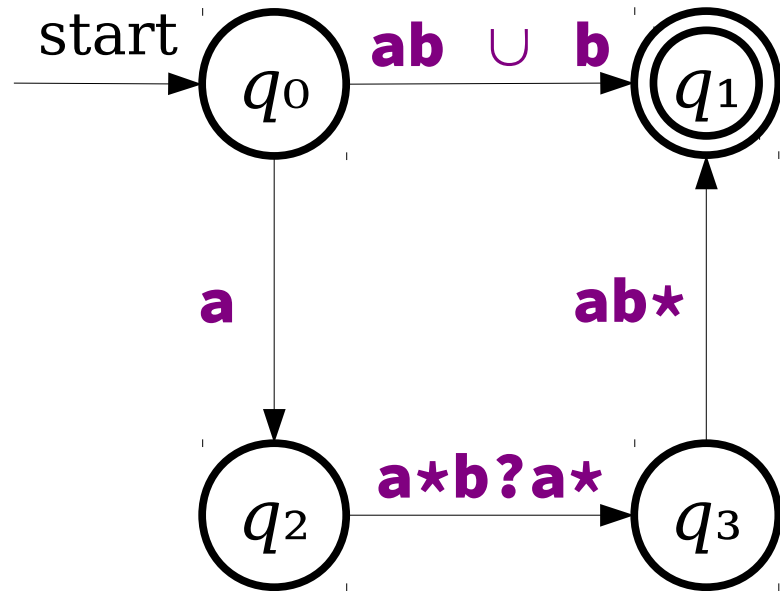


These are all regular expressions!

# Generalizing NFAs

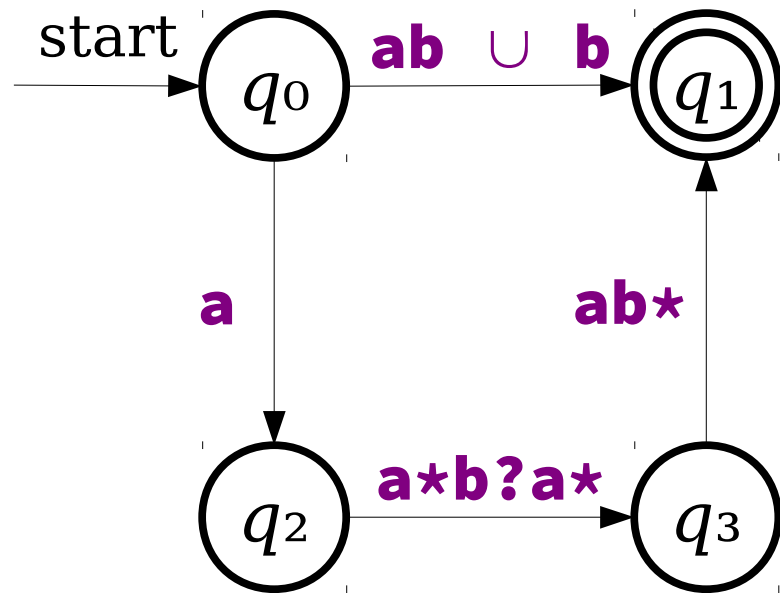


# Generalizing NFAs



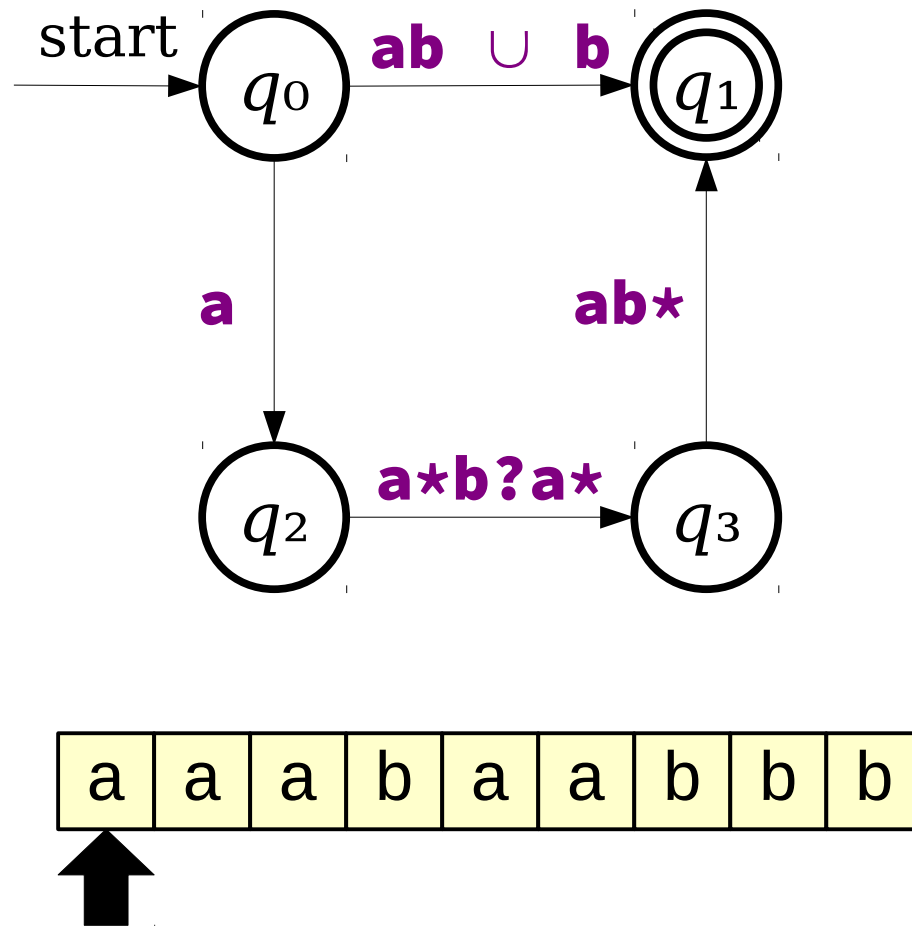
Note: Actual NFAs aren't allowed to have full regex transitions like these. This is just a thought experiment.

# Generalizing NFAs

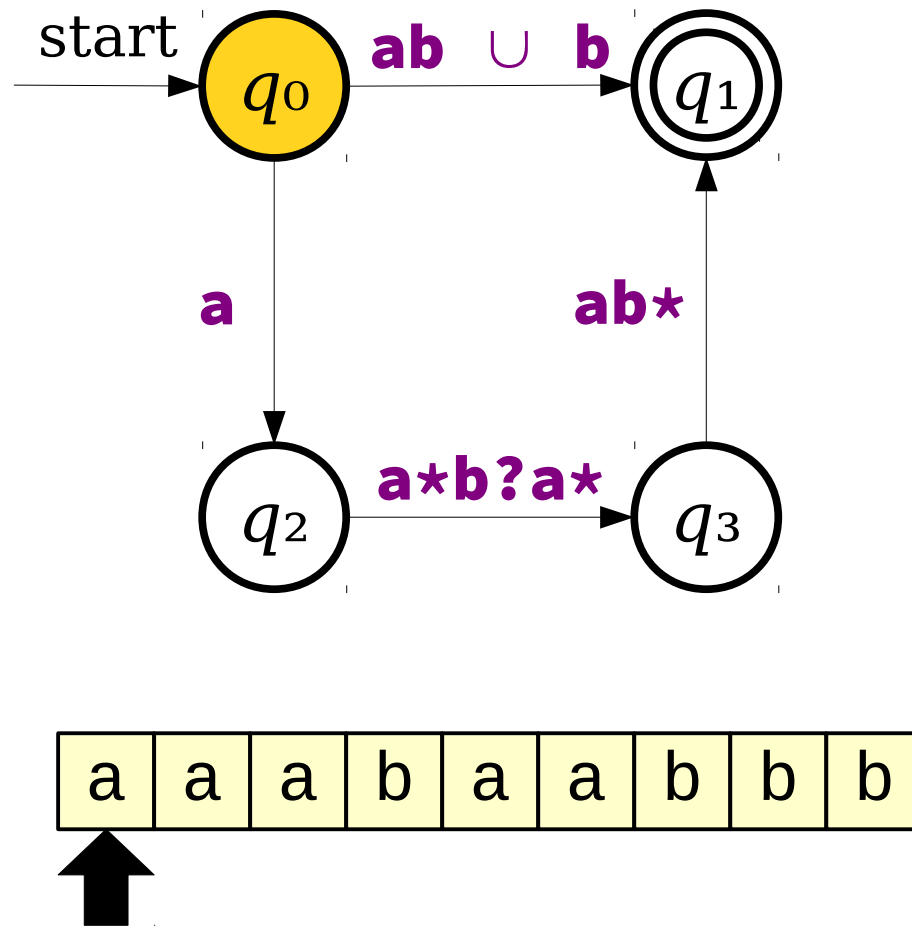


a	a	a	b	a	a	b	b	b
---	---	---	---	---	---	---	---	---

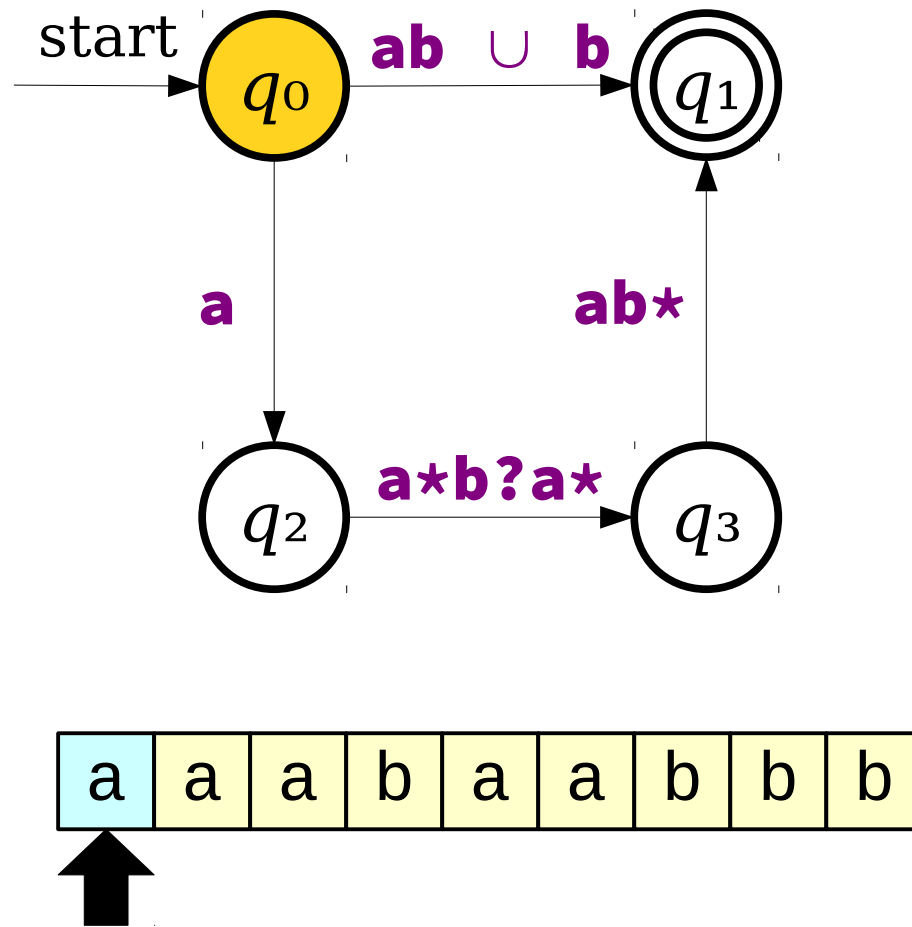
# Generalizing NFAs



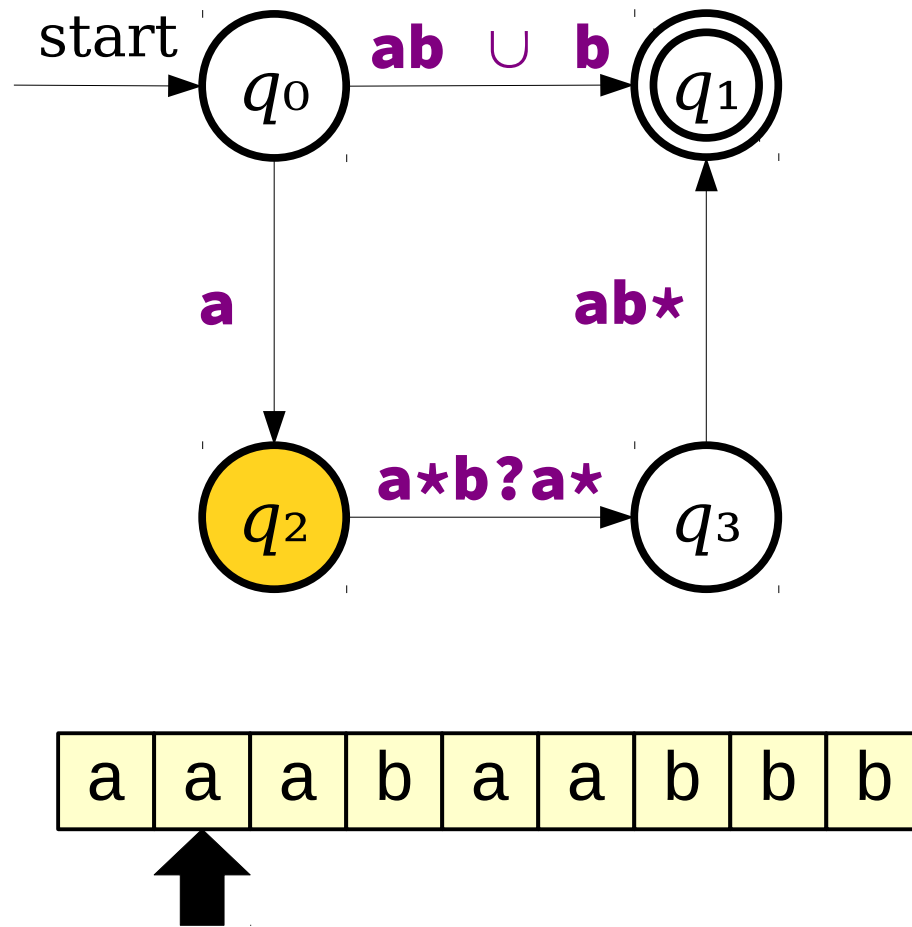
# Generalizing NFAs



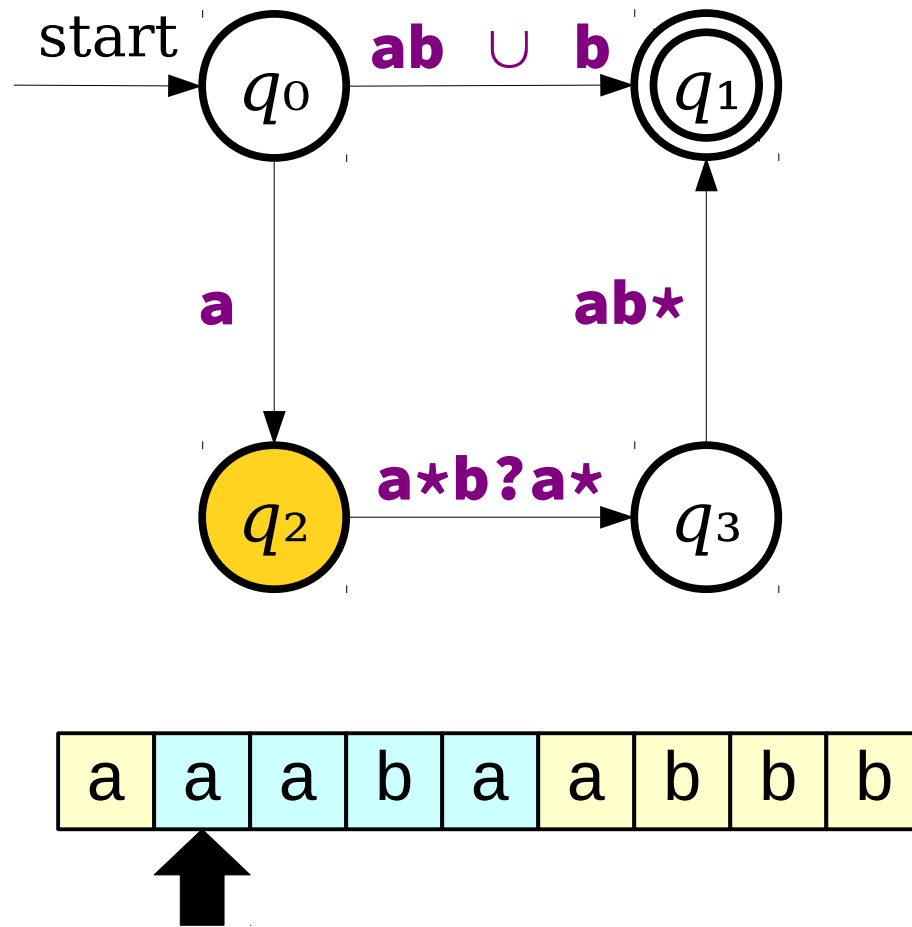
# Generalizing NFAs



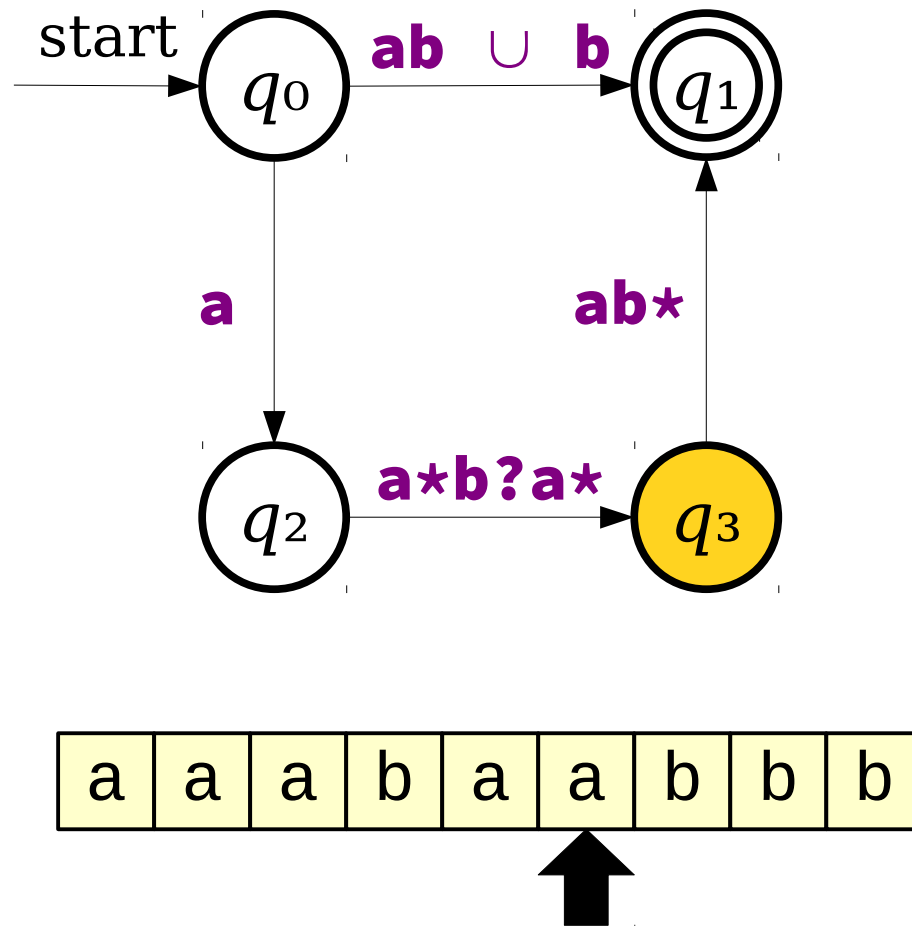
# Generalizing NFAs



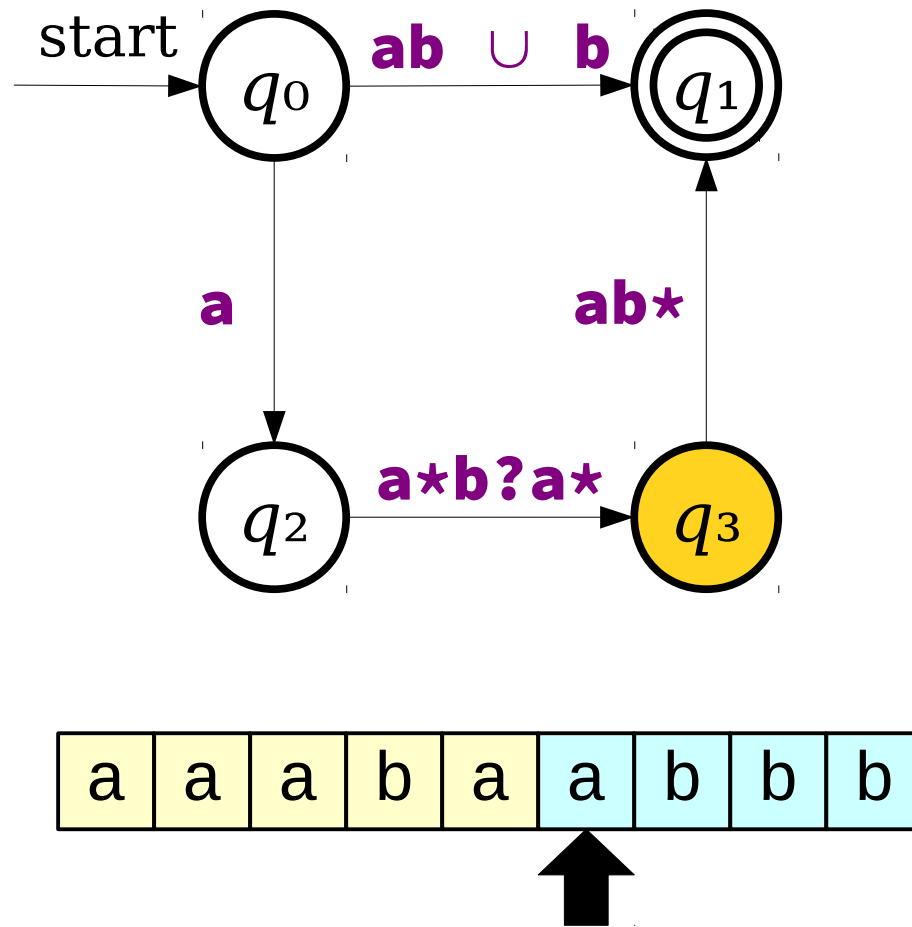
# Generalizing NFAs



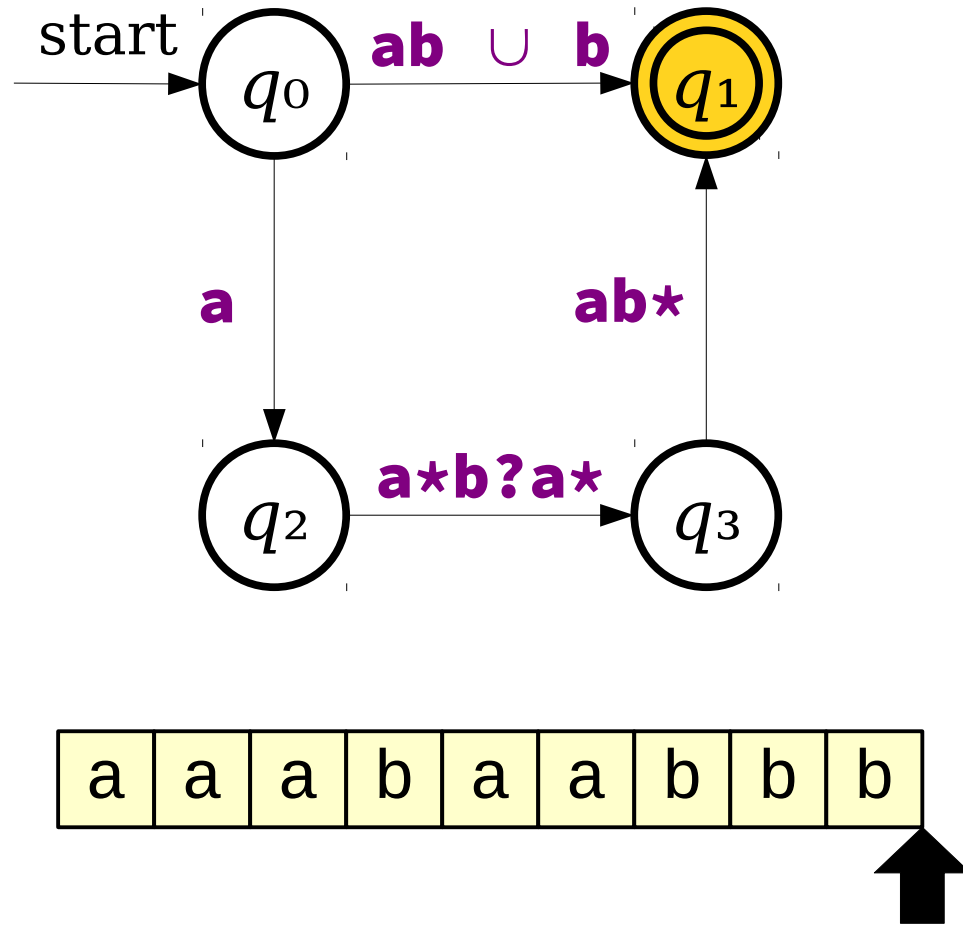
# Generalizing NFAs



# Generalizing NFAs

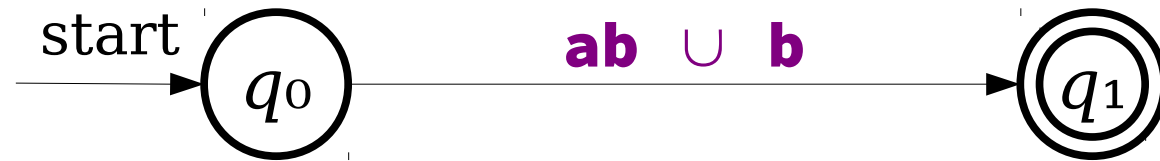


# Generalizing NFAs

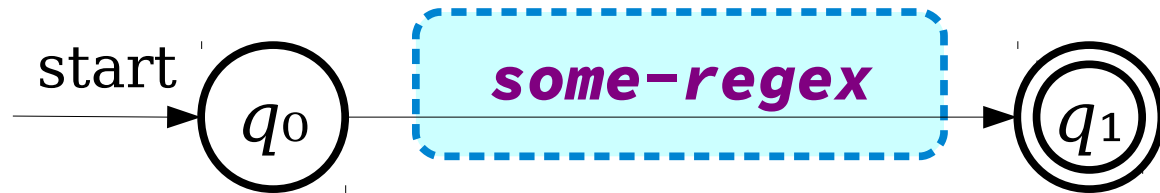


***Key Idea 1:*** Imagine that we can label transitions in an NFA with arbitrary regular expressions.

# Generalizing NFAs



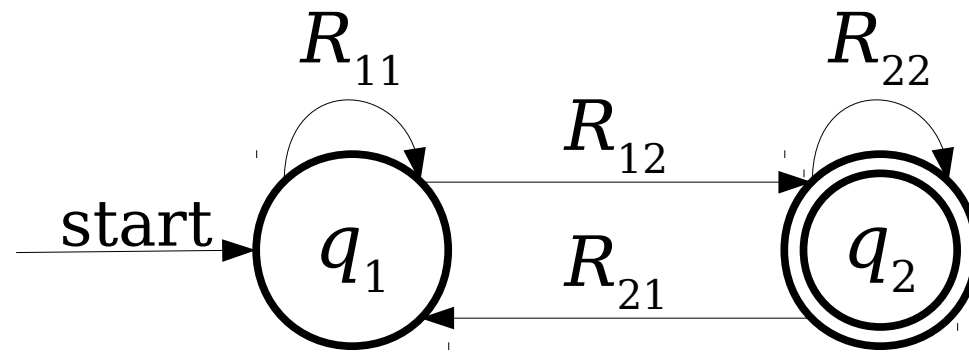
**Key Idea 2:** If we can convert an NFA into a generalized NFA that looks like this...



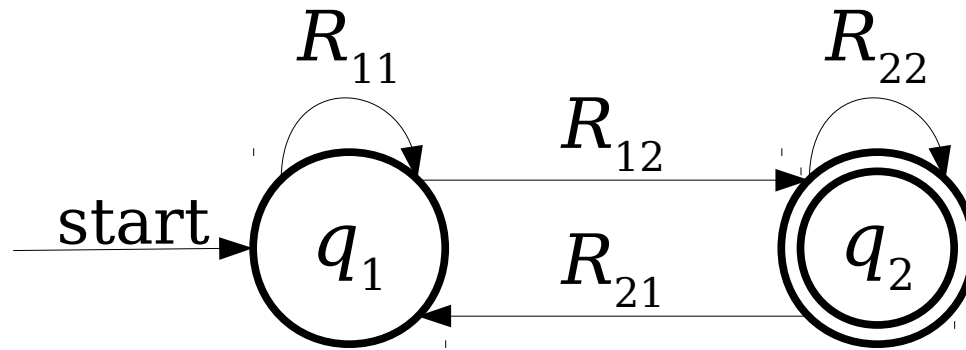
...then we can easily read off a regular expression for the original NFA.

*...but not all NFAs look like this :-(  
but maybe we can fix that?*

# From NFAs to Regular Expressions

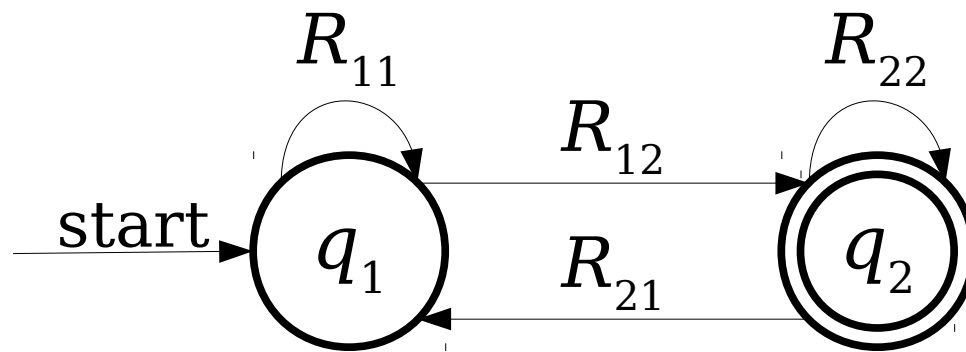


# From NFAs to Regular Expressions



Here,  $R_{11}$ ,  $R_{12}$ ,  $R_{21}$ , and  $R_{22}$  are arbitrary regular expressions.

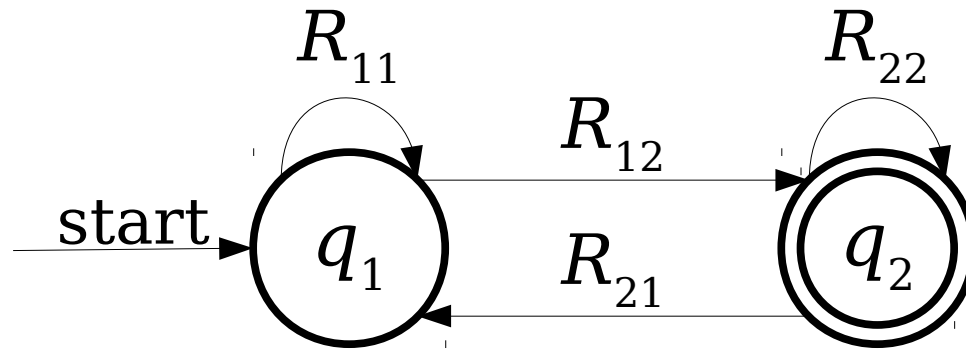
# From NFAs to Regular Expressions



**Key Idea 2 Goal:** Somehow transform this NFA so that it looks like this:

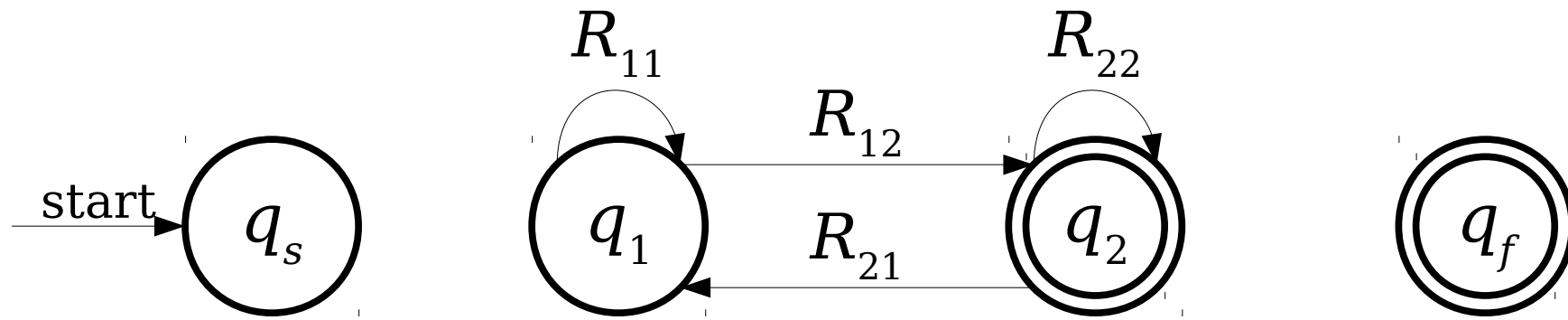


# From NFAs to Regular Expressions

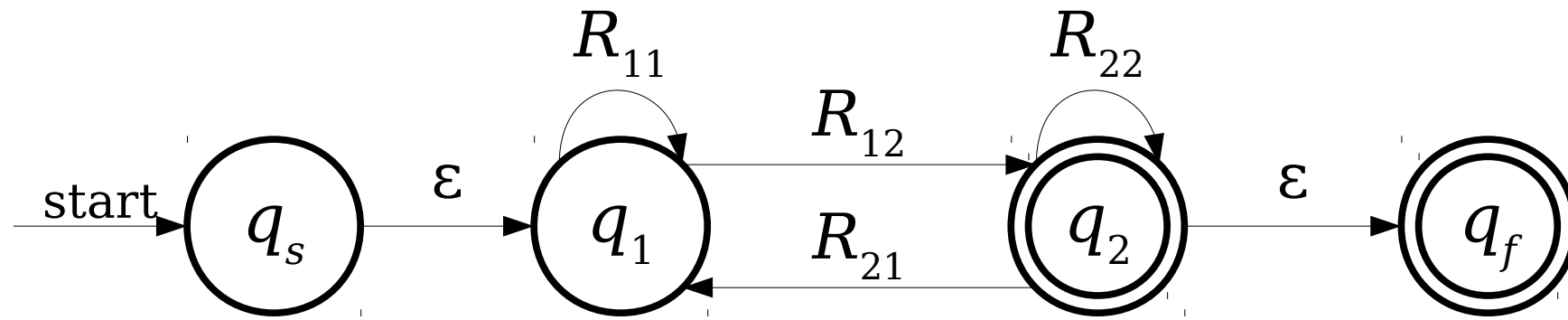


The first step is going to be a simple one...

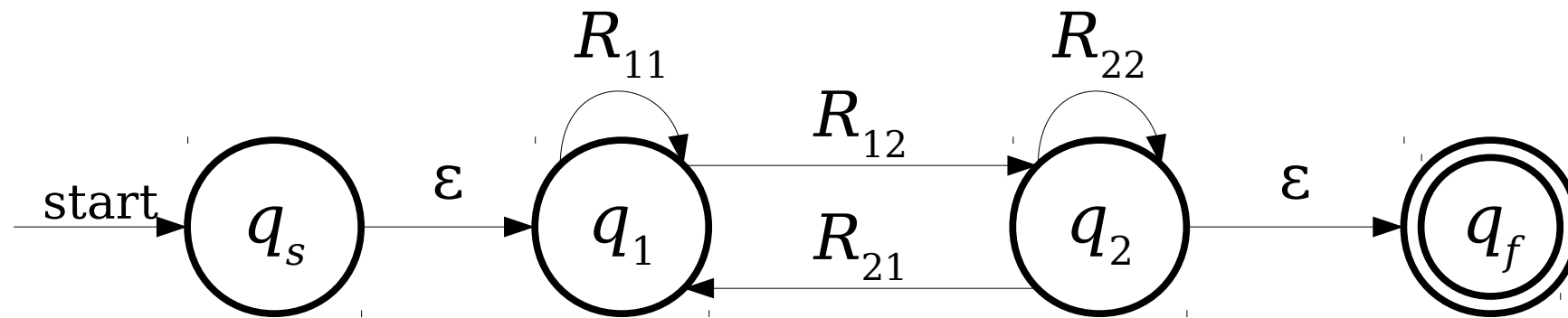
# From NFAs to Regular Expressions



# From NFAs to Regular Expressions



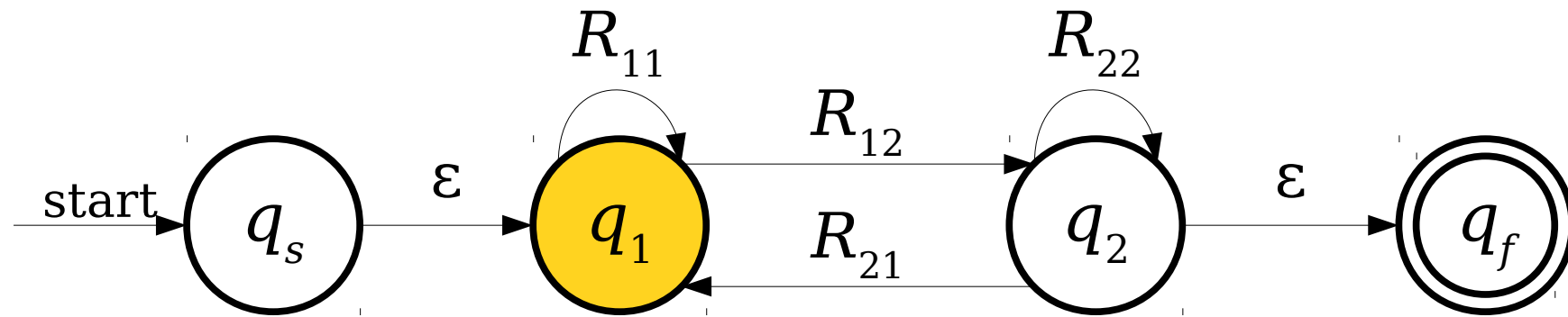
# From NFAs to Regular Expressions



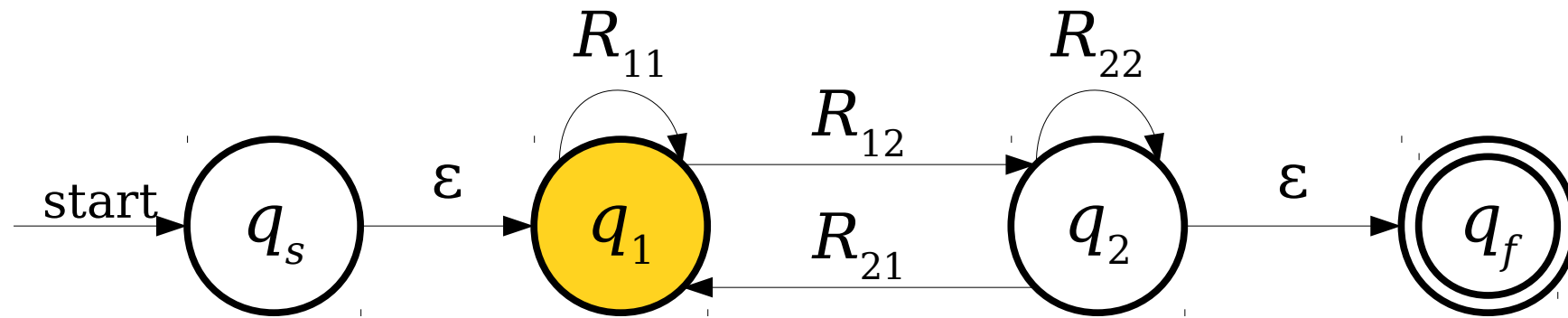
**Key Idea 2 Goal Progress:** Now it has clean start and end, but middle is messy. Need to get rid of those states.



# From NFAs to Regular Expressions

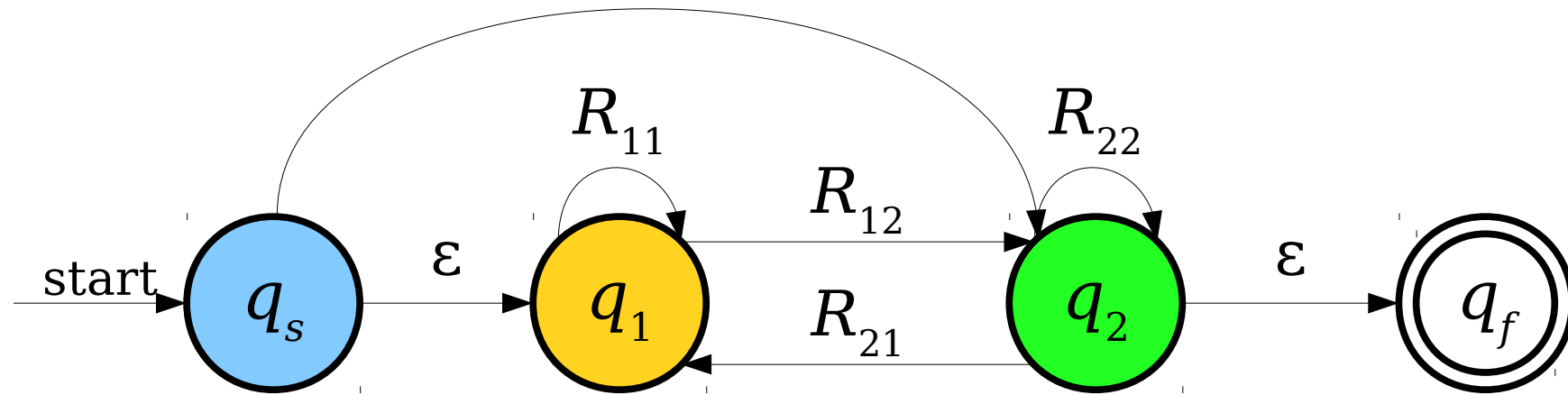


# From NFAs to Regular Expressions

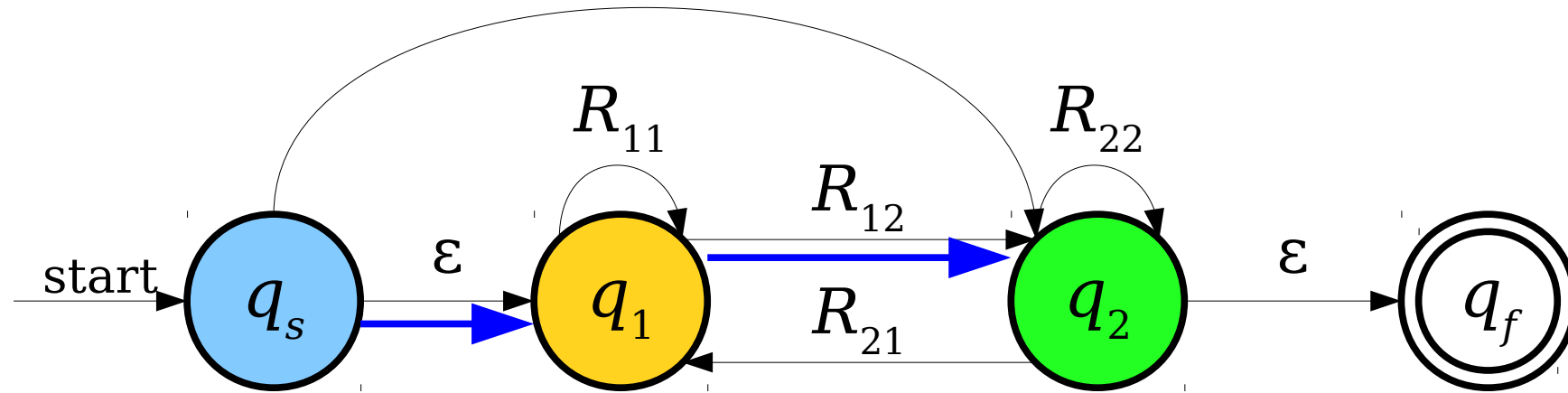


Could we eliminate  
this state from the  
NFA?

# From NFAs to Regular Expressions

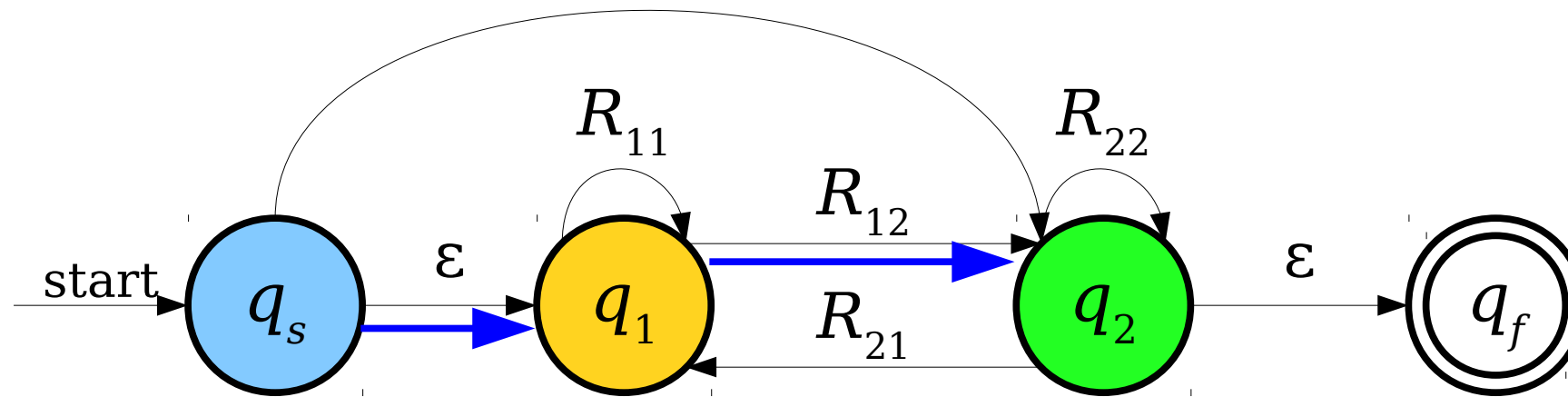


# From NFAs to Regular Expressions



Here is a pattern that we might process  
when going from  $q_s$  to  $q_2$ :  $\epsilon R_{12}$

# From NFAs to Regular Expressions



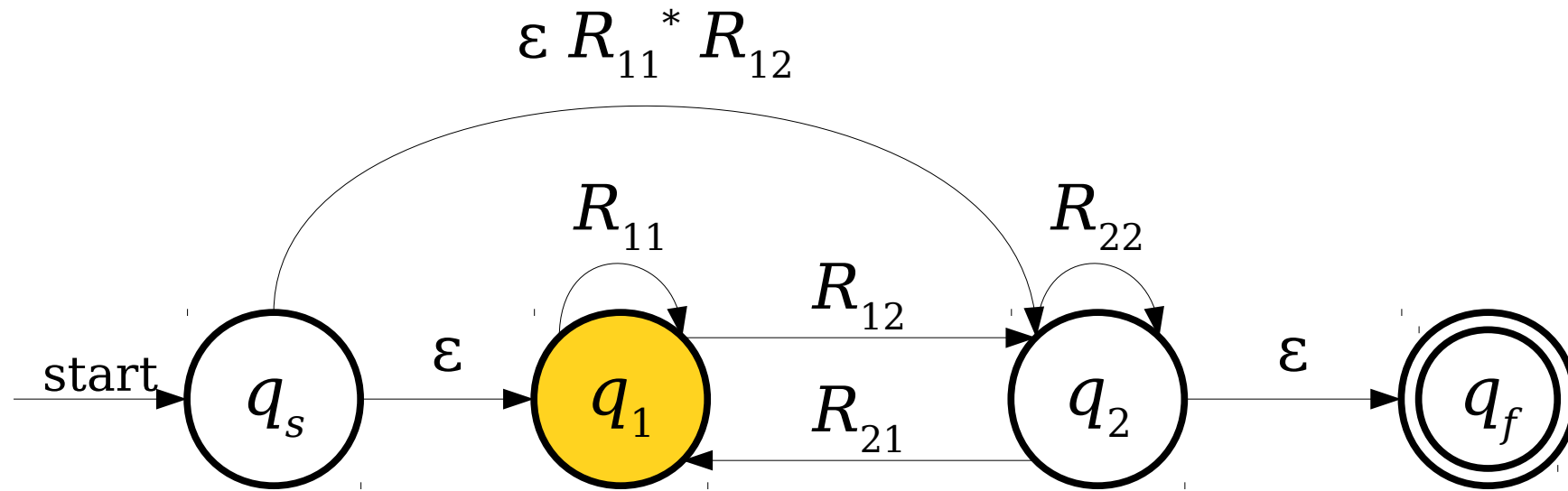
Here is a pattern that we might process when going from  $q_s$  to  $q_2$ :  $\epsilon R_{12}$

## State elimination question:

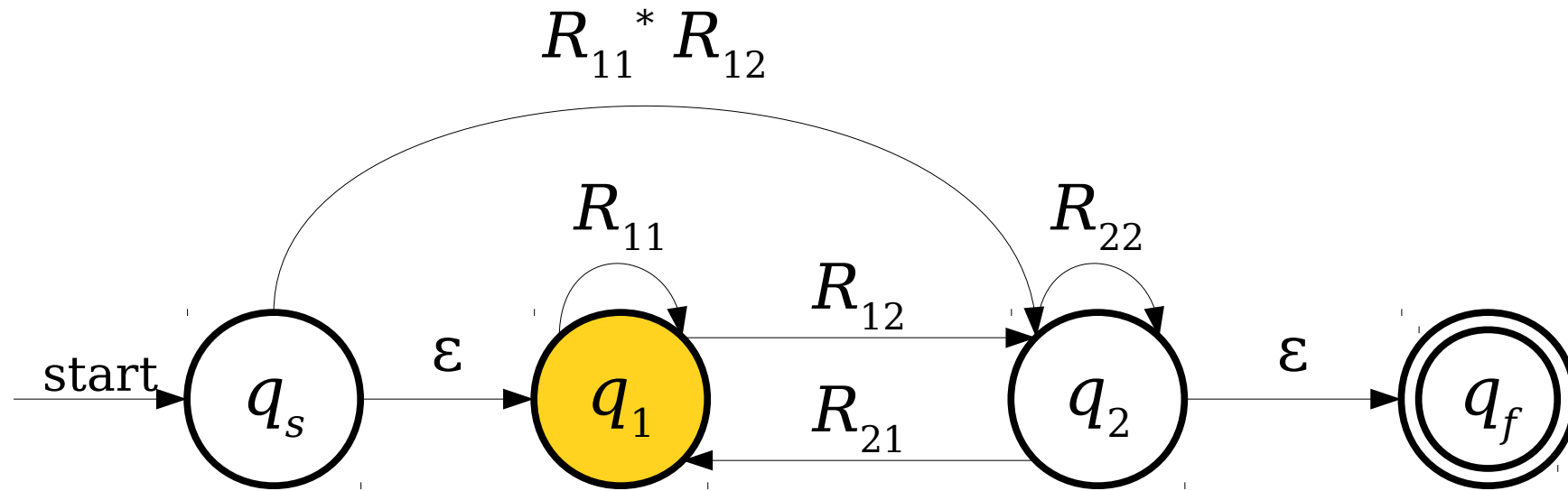
How many of the following are also patterns might we process when going from  $q_s$  to  $q_2$ ?

- $\epsilon R_{11} R_{12}$
- $\epsilon R_{11} R_{11} R_{12}$
- $\epsilon R_{11} R_{12} R_{11}$
- $\epsilon R_{11} R_{11} R_{12}$
- $\epsilon R_{11} R_{11} R_{11} R_{11} R_{11} R_{12}$

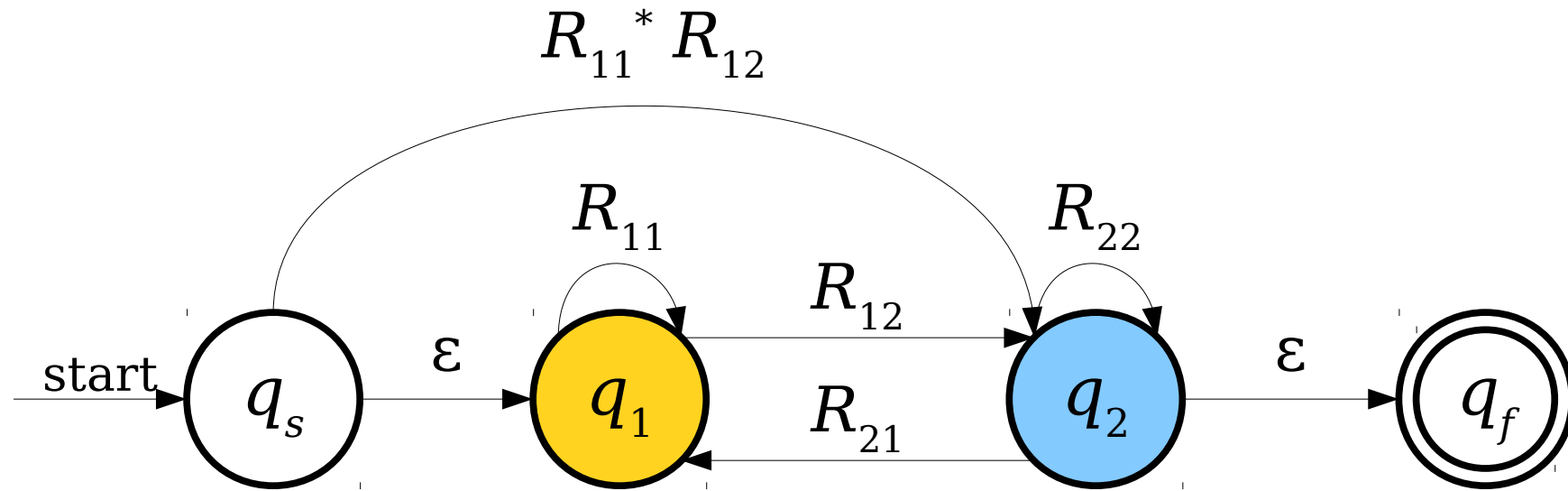
# From NFAs to Regular Expressions



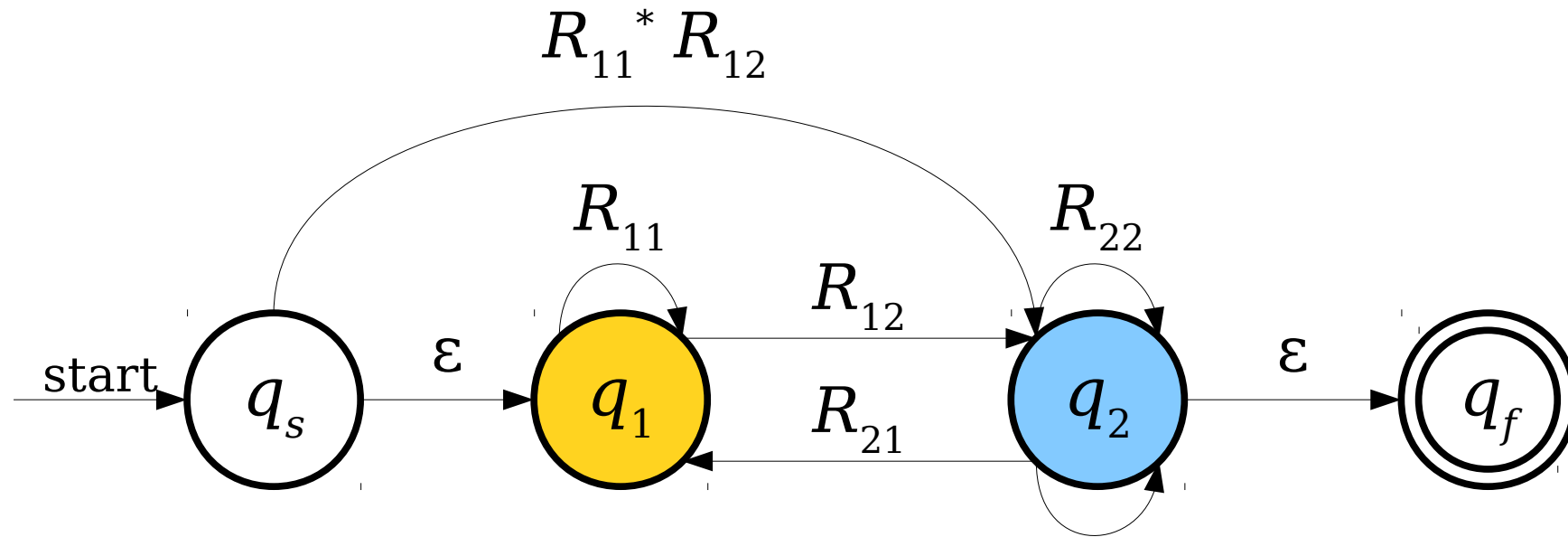
# From NFAs to Regular Expressions



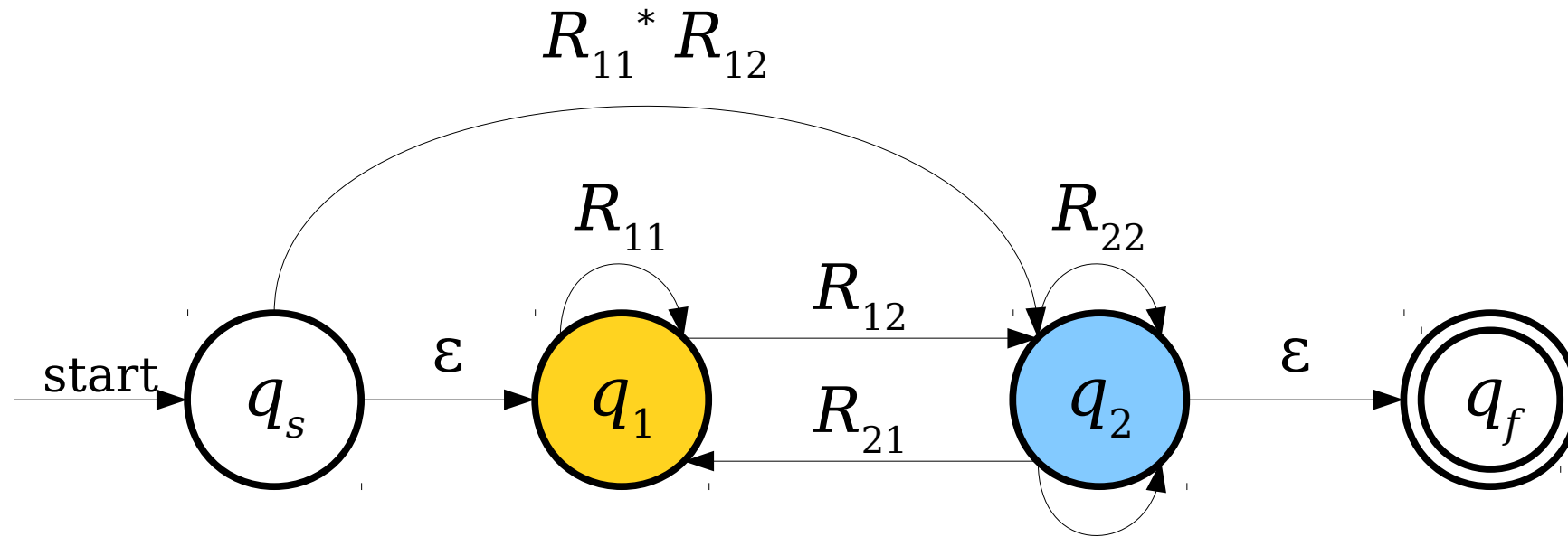
# From NFAs to Regular Expressions



# From NFAs to Regular Expressions

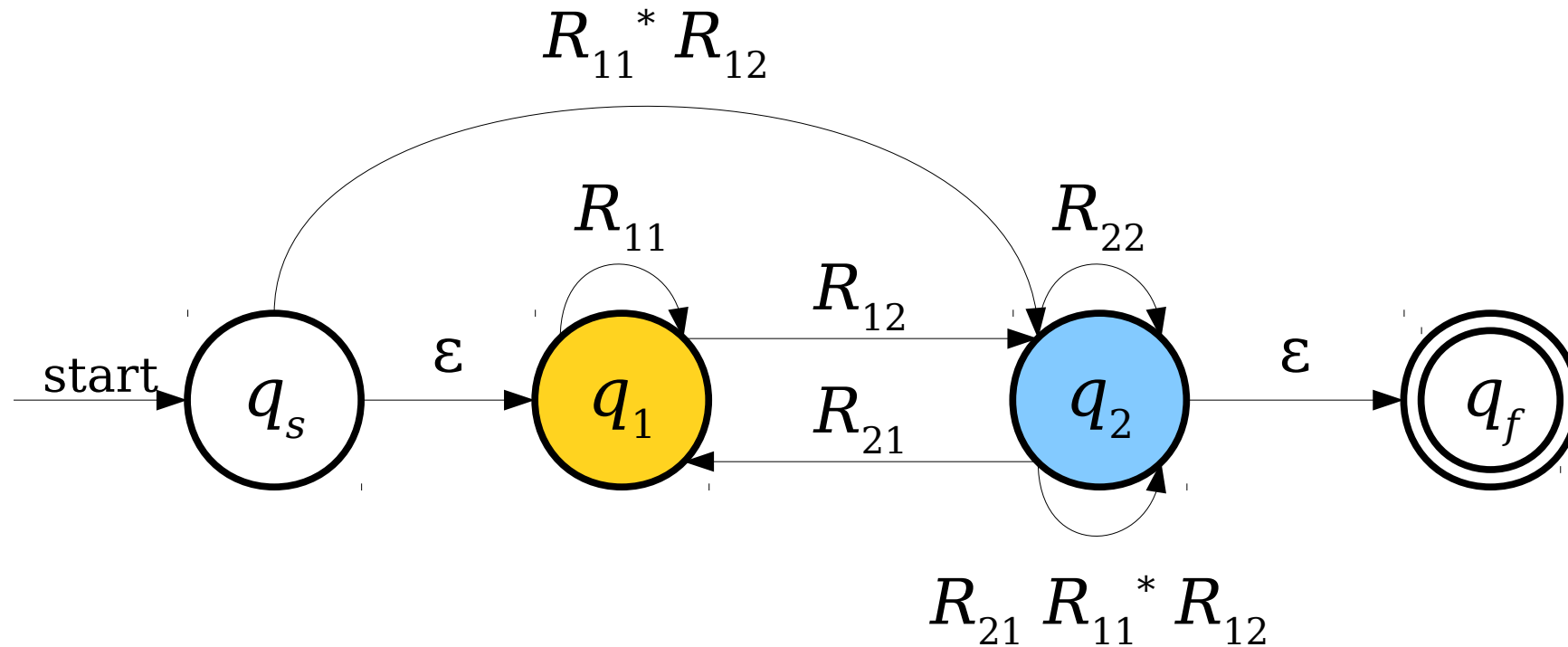


# From NFAs to Regular Expressions



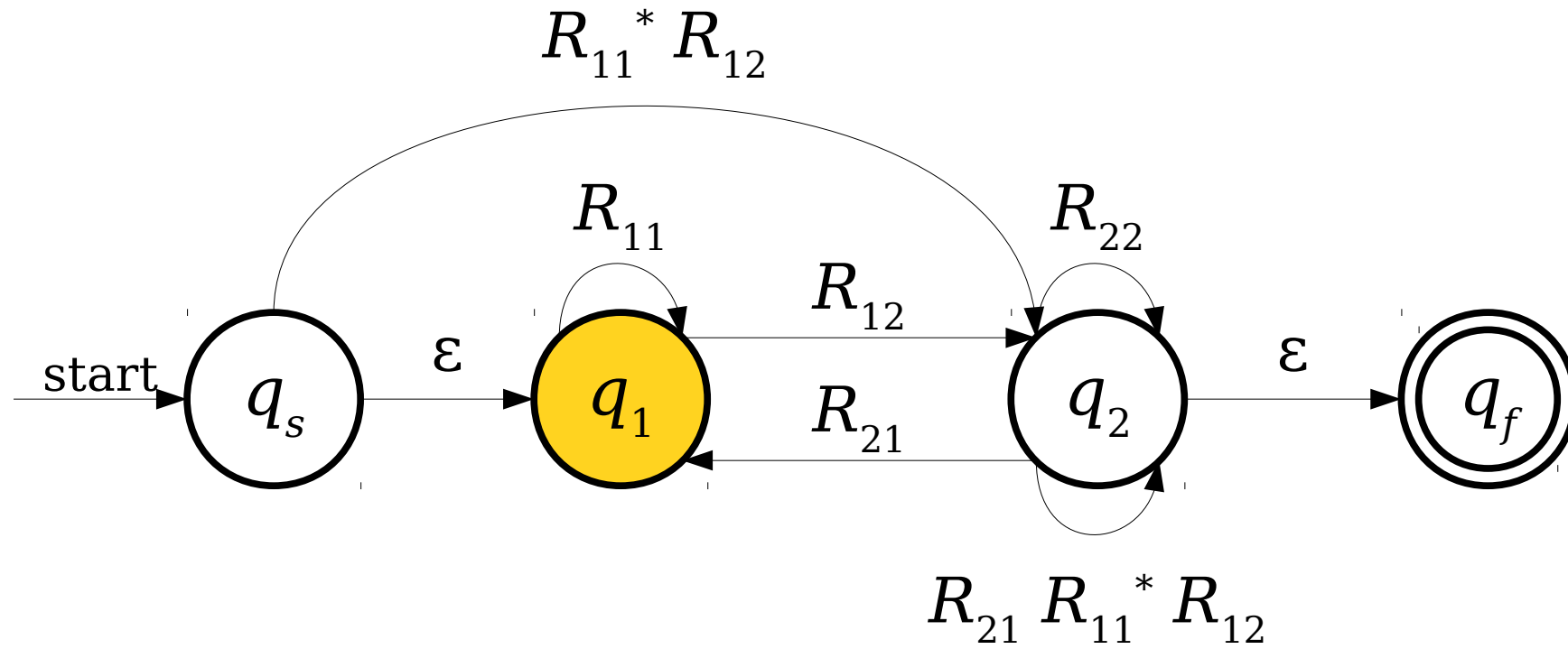
Here is a pattern that we might process  
when going from  $q_2$  to  $q_2$ :  $R_{21} R_{11} R_{12}$

# From NFAs to Regular Expressions

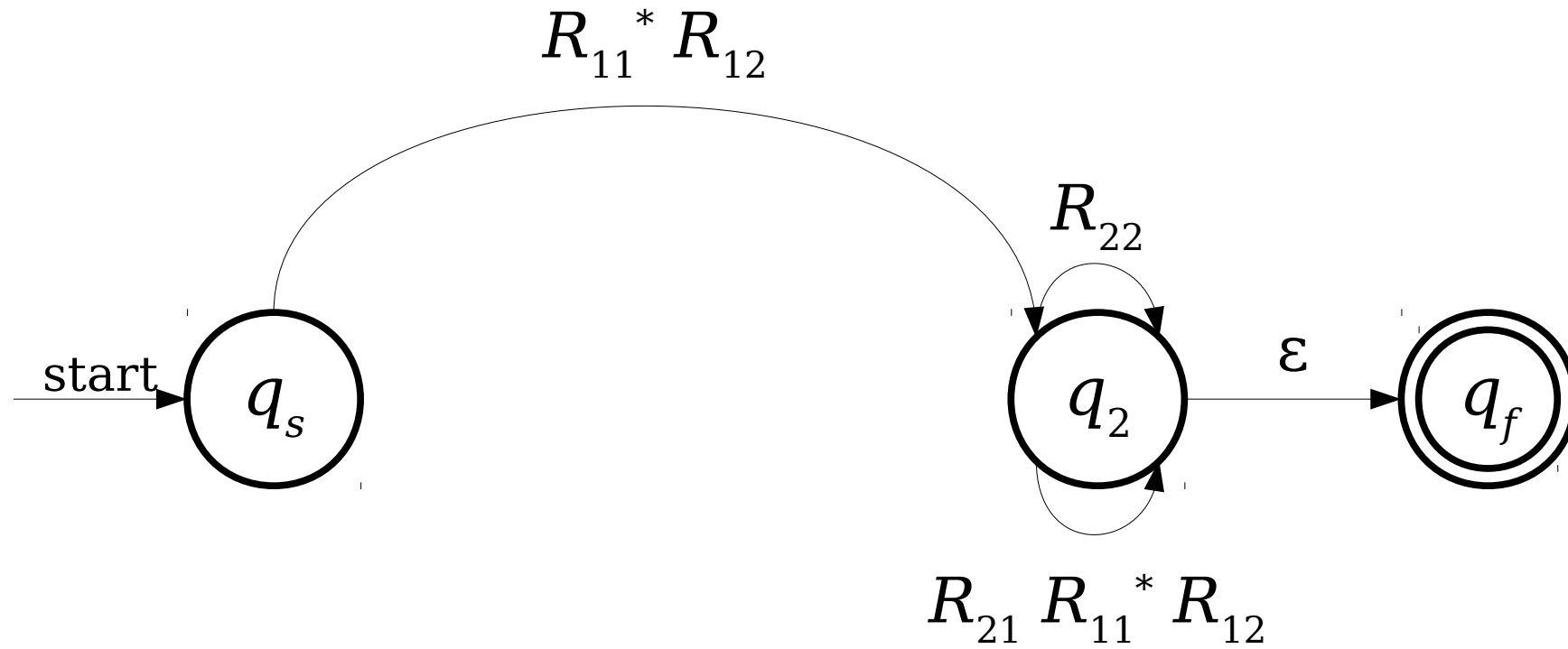


Here is a pattern that we might process  
when going from  $q_2$  to  $q_2$ :  $R_{21} R_{11}^* R_{12}$

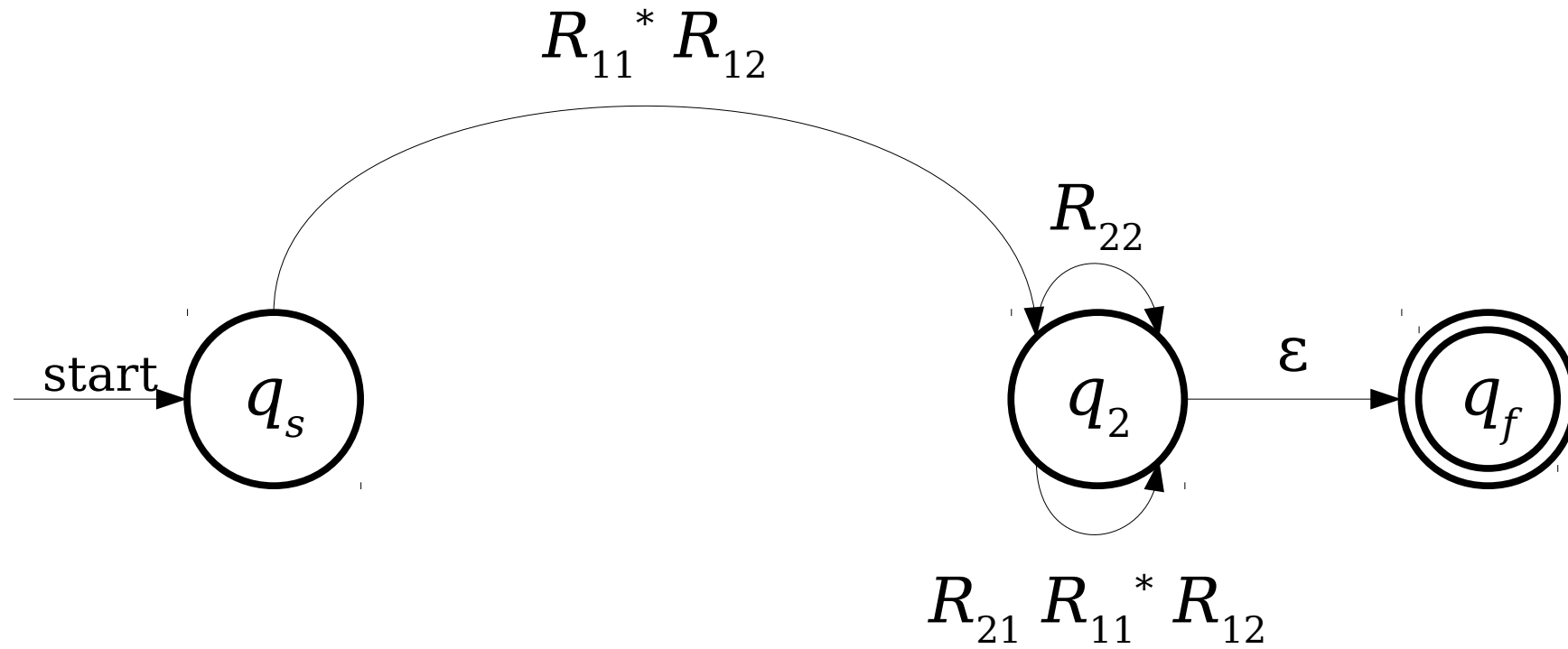
# From NFAs to Regular Expressions



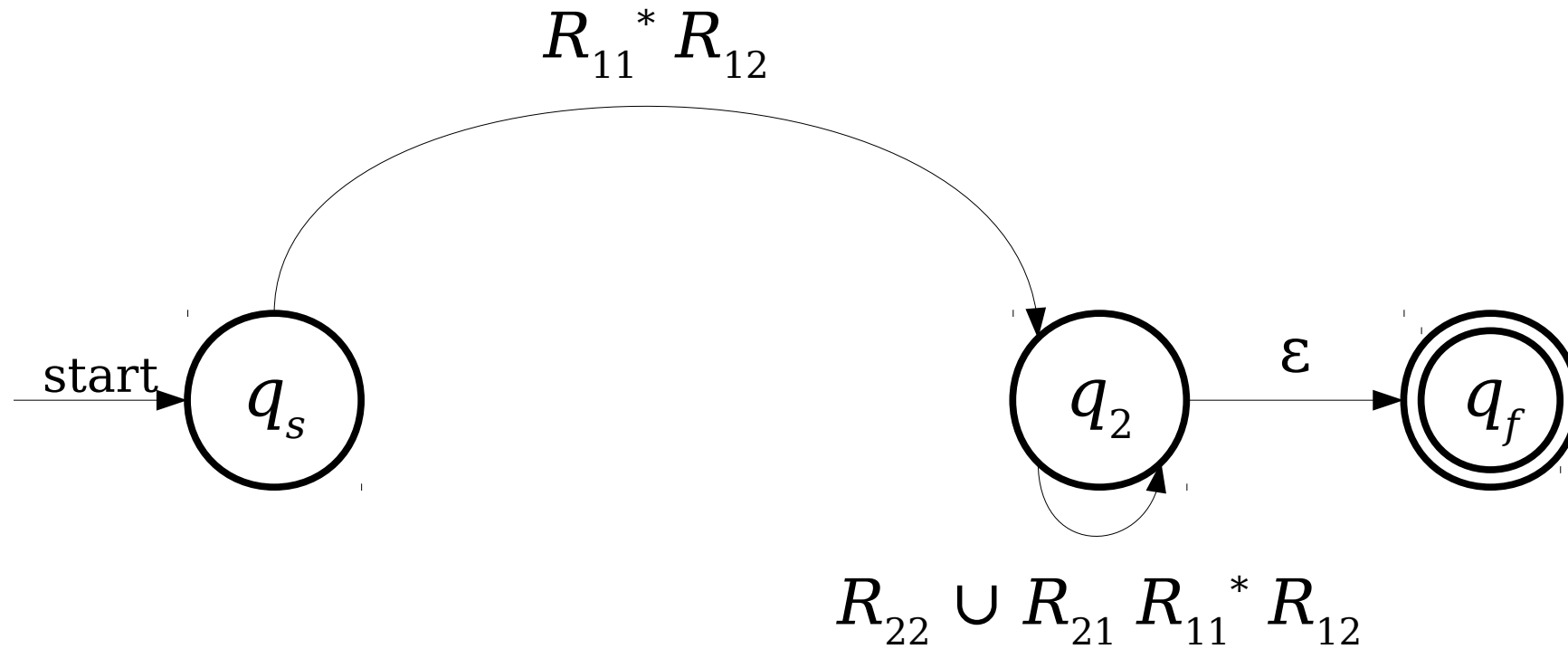
# From NFAs to Regular Expressions



# From NFAs to Regular Expressions

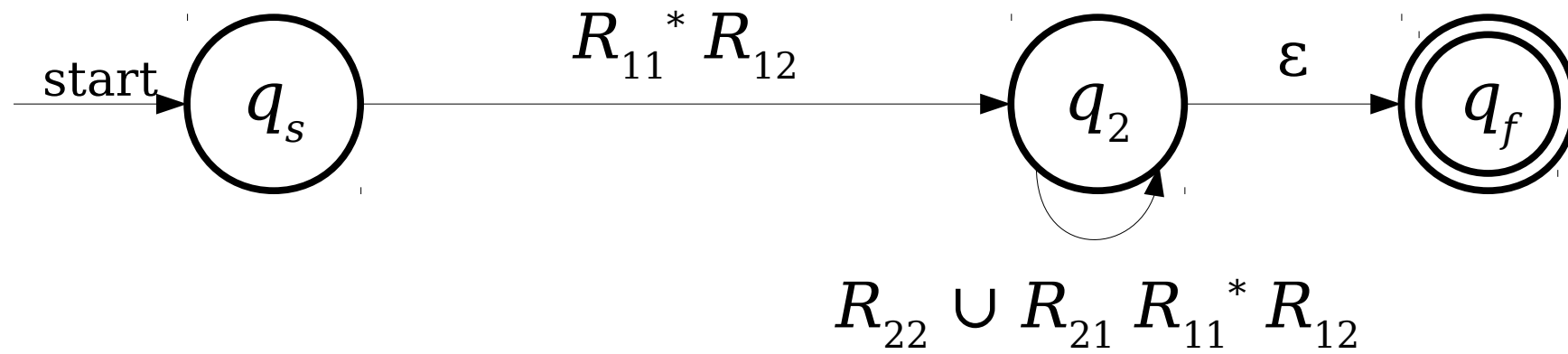


# From NFAs to Regular Expressions

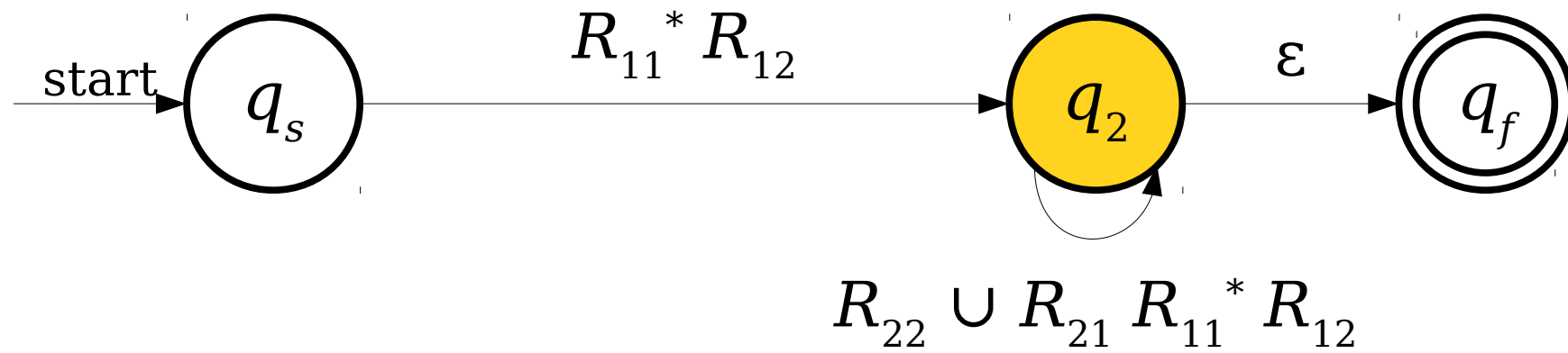


Note: We're using **union** to combine these transitions together.

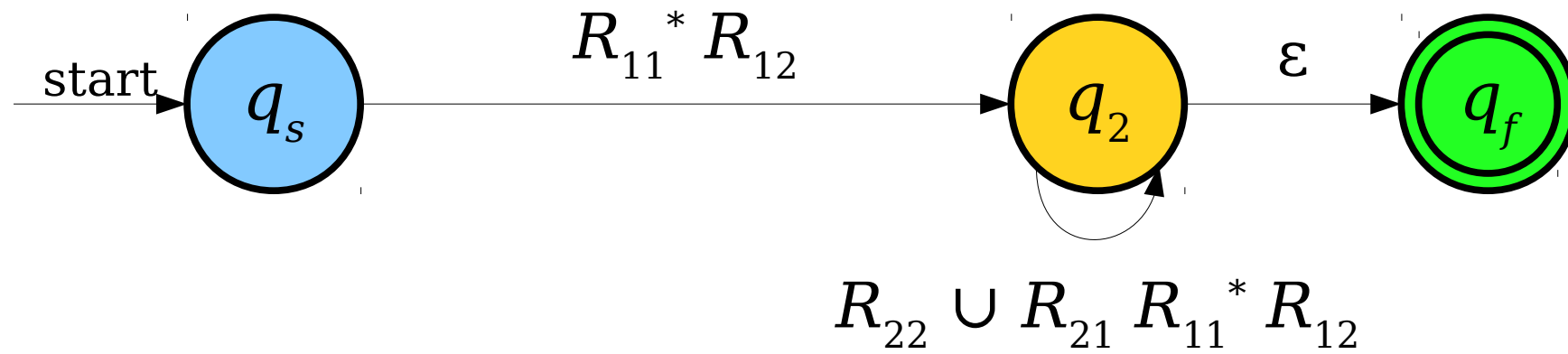
# From NFAs to Regular Expressions



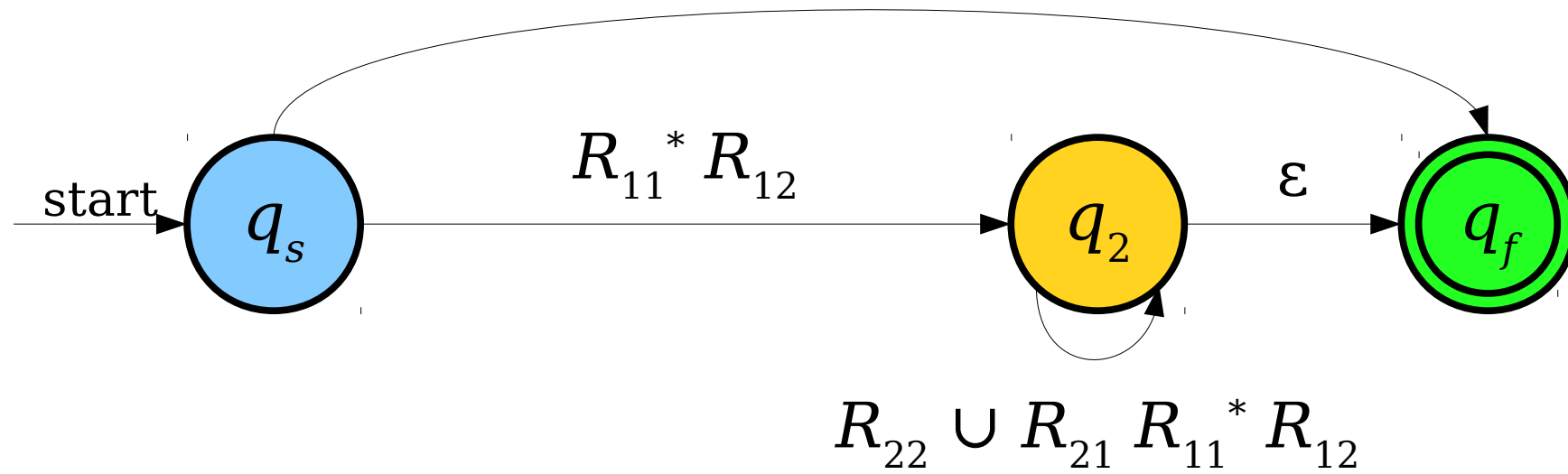
# From NFAs to Regular Expressions



# From NFAs to Regular Expressions

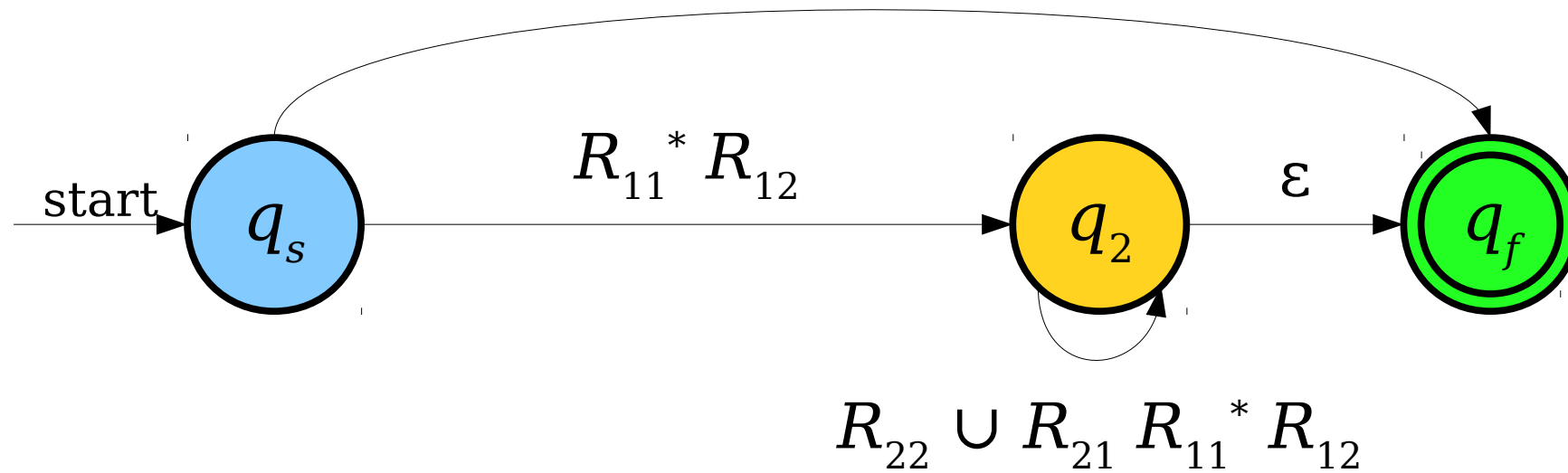


# From NFAs to Regular Expressions

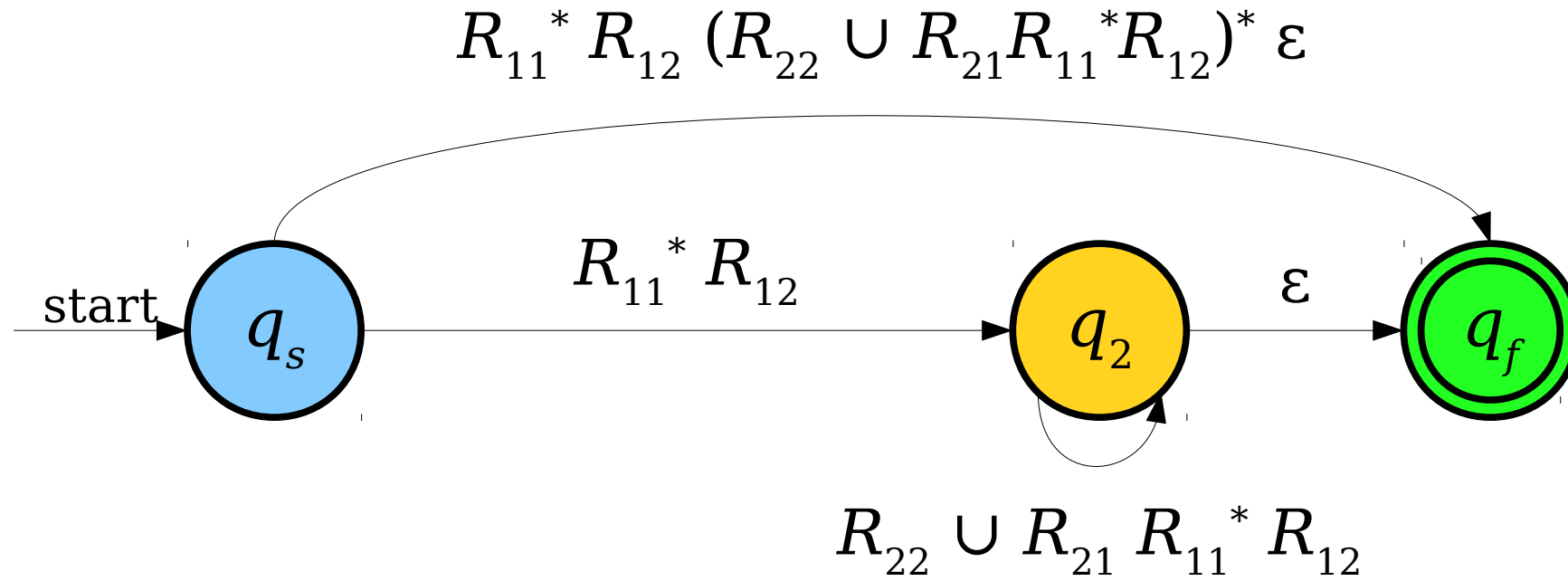


# From NFAs to Regular Expressions

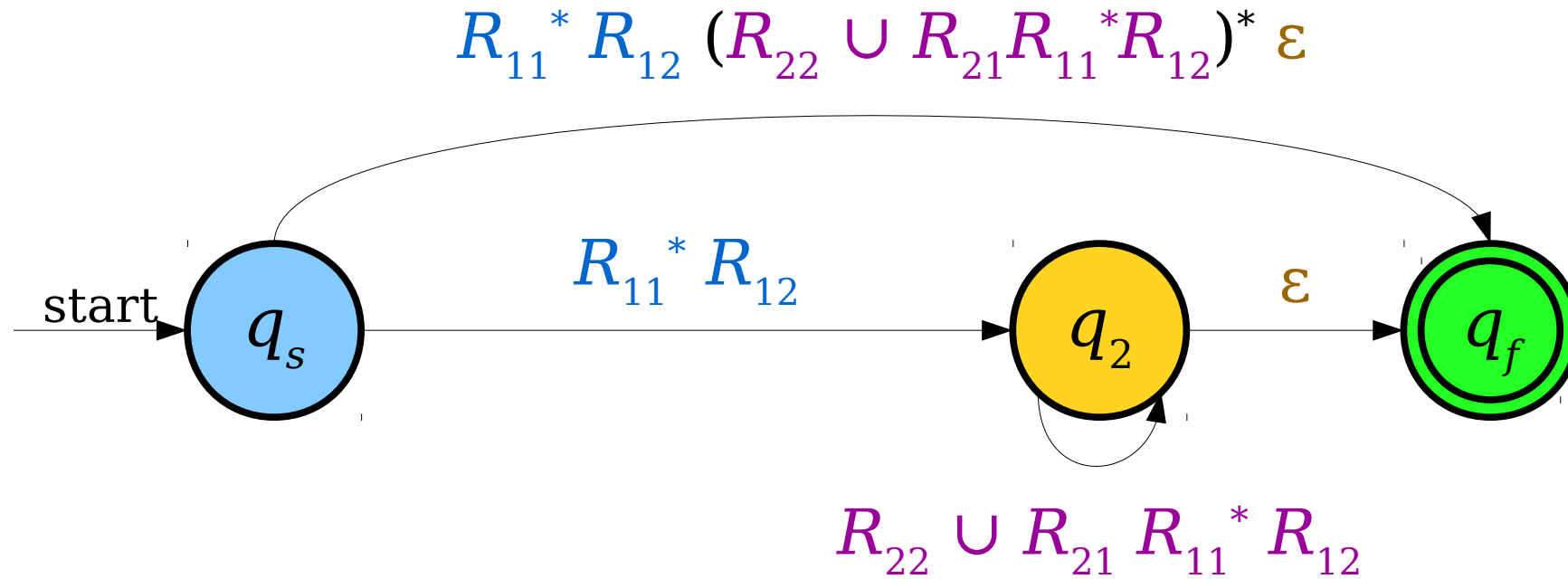
**Quick check:** what goes on this transition?



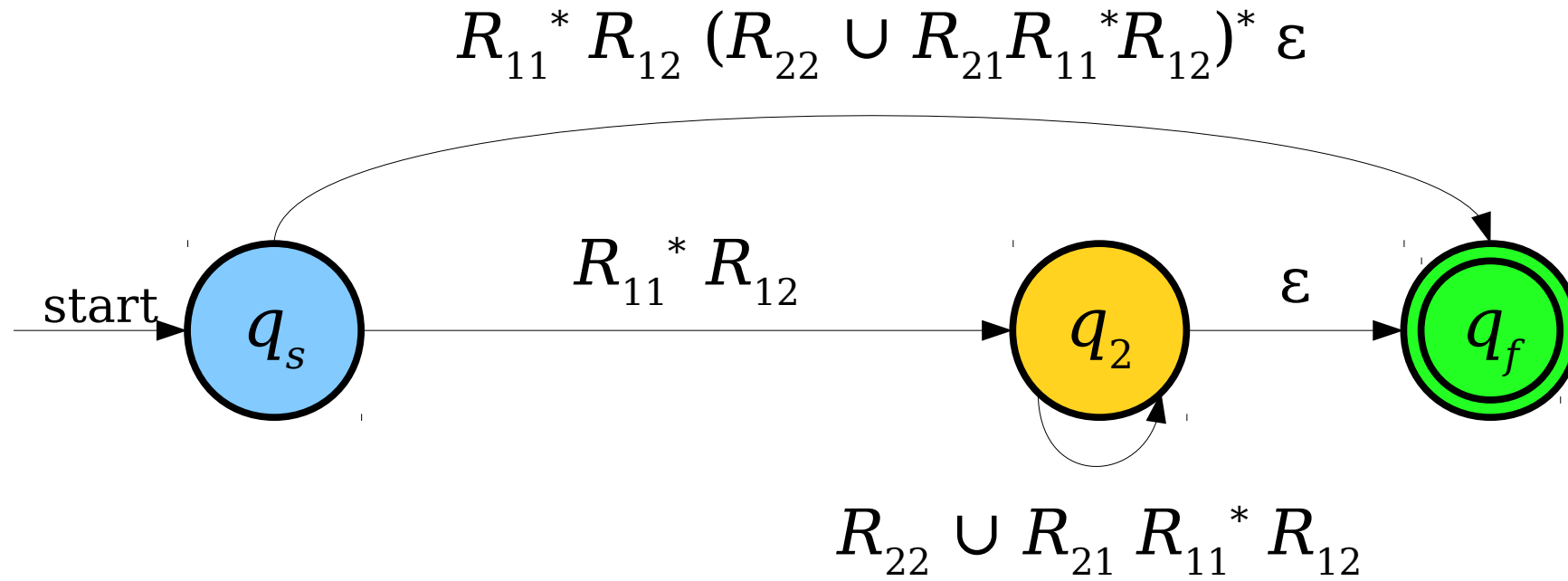
# From NFAs to Regular Expressions



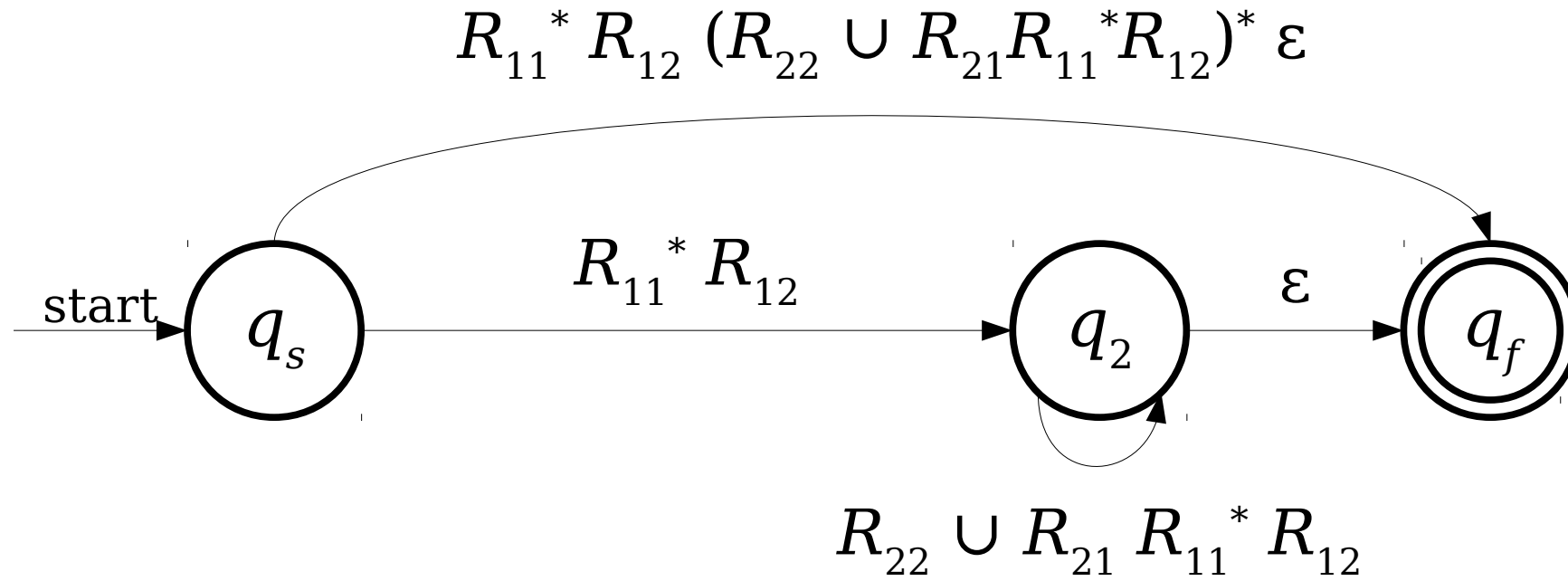
# From NFAs to Regular Expressions



# From NFAs to Regular Expressions

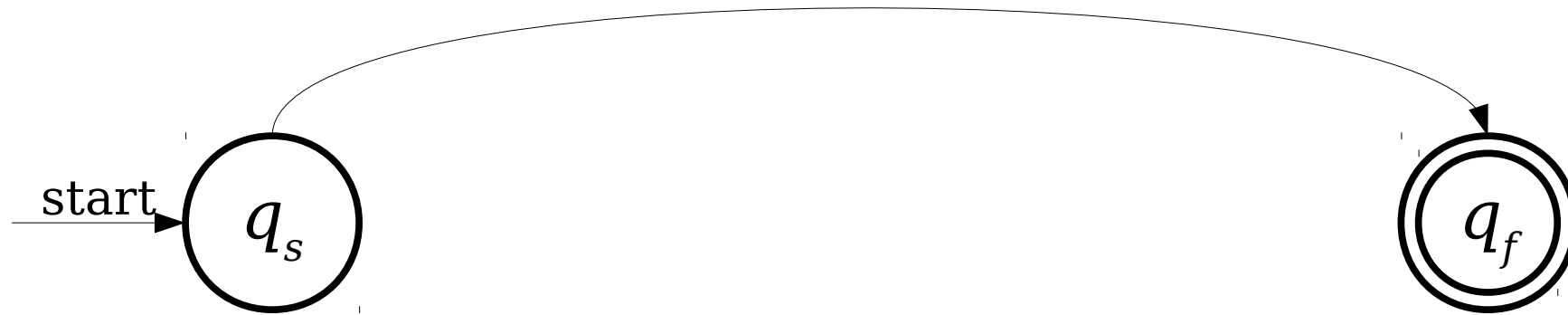


# From NFAs to Regular Expressions

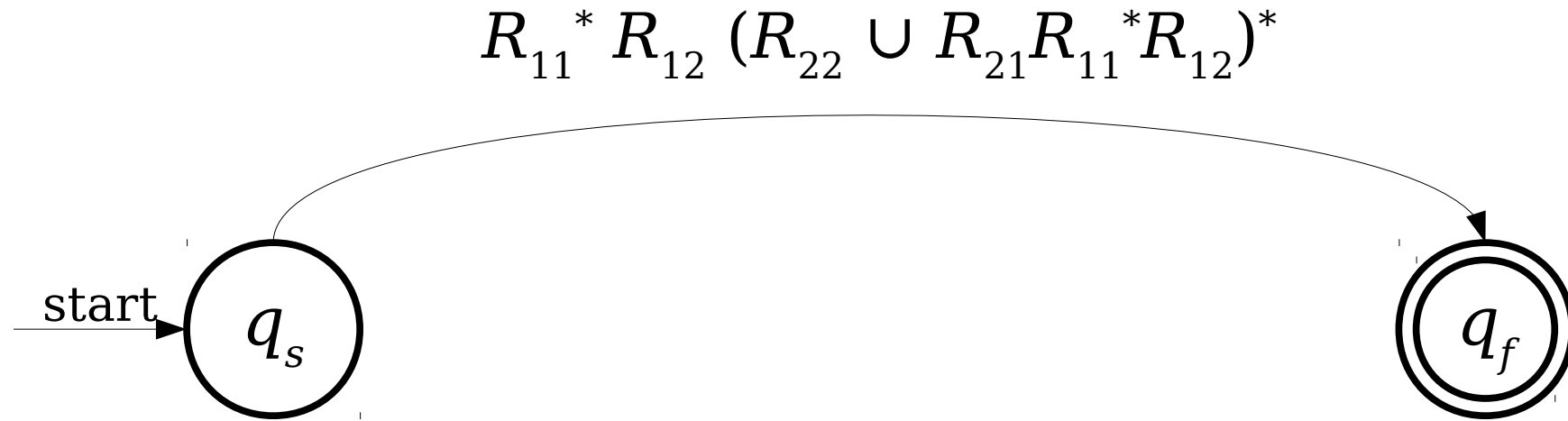


# From NFAs to Regular Expressions

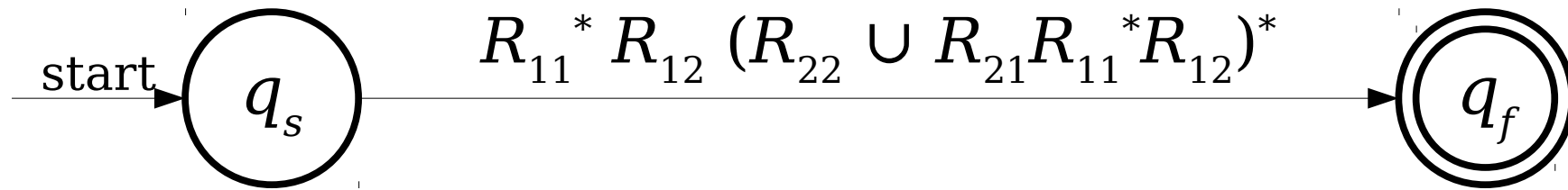
$$R_{11}^* R_{12} (R_{22} \cup R_{21} R_{11}^* R_{12})^* \varepsilon$$



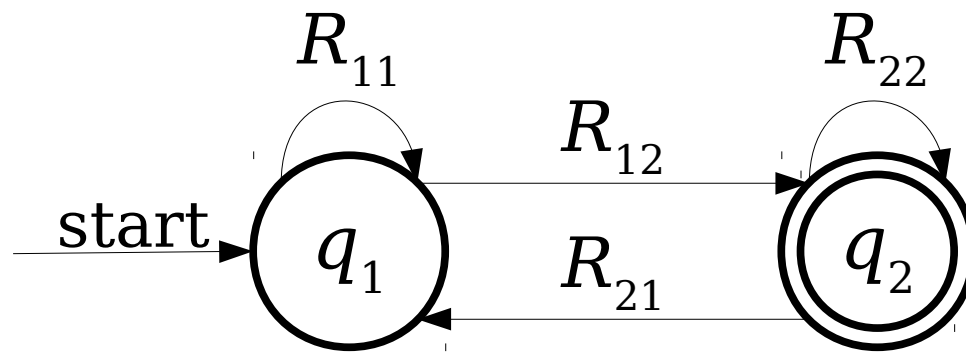
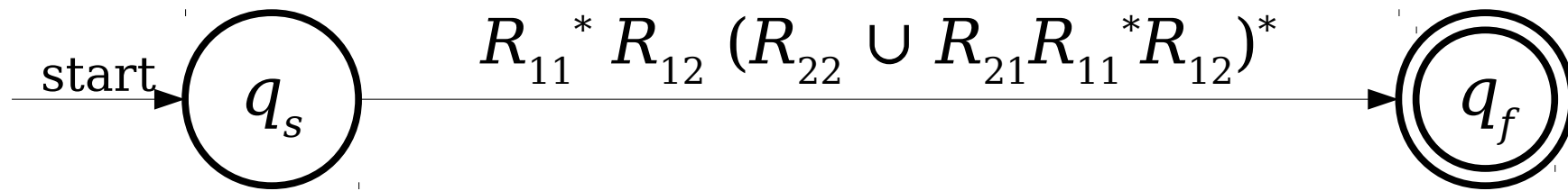
# From NFAs to Regular Expressions



# From NFAs to Regular Expressions

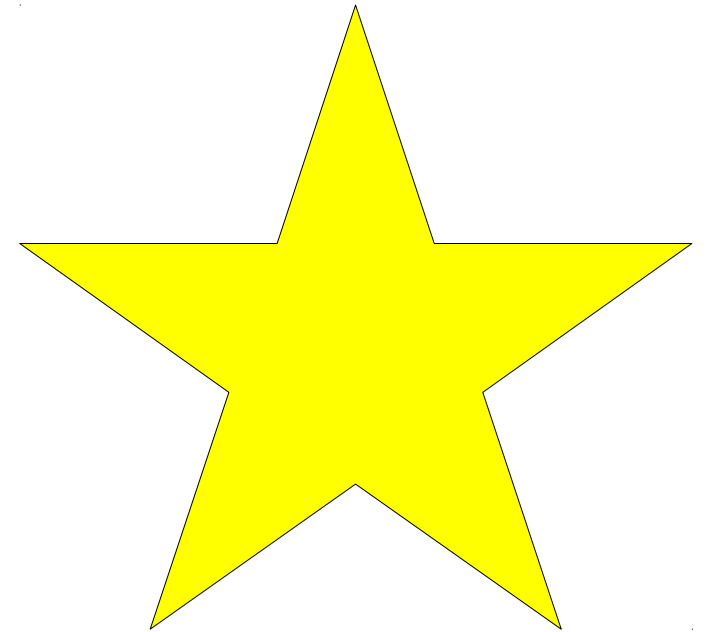


# From NFAs to Regular Expressions



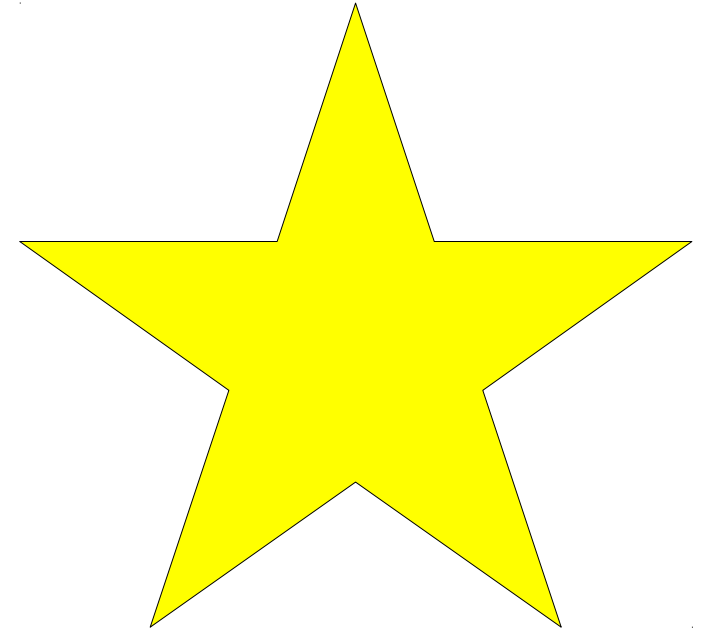
# The State-Elimination Algorithm

- Start with an NFA  $N$  for the language  $L$ .
- Add a new start state  $q_s$  and accept state  $q_f$  to the NFA.
  - Add an  $\varepsilon$ -transition from  $q_s$  to the old start state of  $N$ .
  - Add  $\varepsilon$ -transitions from each accepting state of  $N$  to  $q_f$ , then mark them as not accepting.
- Repeatedly remove states other than  $q_s$  and  $q_f$  from the NFA by “shortcutting” them until only two states remain:  $q_s$  and  $q_f$ .
- The transition from  $q_s$  to  $q_f$  is then a regular expression for the NFA.



# The State-Elimination Algorithm

- To eliminate a state  $q$  from the automaton, do the following for each pair of states  $q_0$  and  $q_1$ , where there's a transition from  $q_0$  into  $q$  and a transition from  $q$  into  $q_1$ :
  - Let  $R_{in}$  be the regex on the transition from  $q_0$  to  $q$ .
  - Let  $R_{out}$  be the regex on the transition from  $q$  to  $q_1$ .
  - If there is a regular expression  $R_{stay}$  on a transition from  $q$  to itself, add a new transition from  $q_0$  to  $q_1$  labeled  $((R_{in})(R_{stay})^*(R_{out}))$ .
  - If there isn't, add a new transition from  $q_0$  to  $q_1$  labeled  $((R_{in})(R_{out}))$
- If a pair of states has multiple transitions between them labeled  $R_1, R_2, \dots, R_k$ , replace them with a single transition labeled  $R_1 \cup R_2 \cup \dots \cup R_k$ .



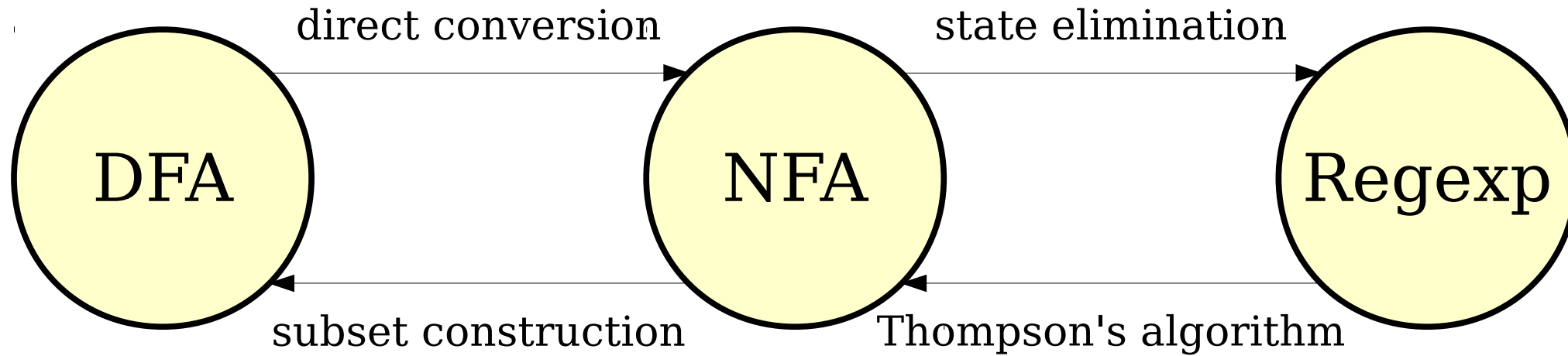
Languages you can  
build a DFA for.

Languages you can  
build an NFA for.

***Regular  
Languages***

Languages you can  
write a Regex for.

# Our Transformations



**Theorem:** The following are all equivalent:

- $L$  is a regular language.
- There is a DFA  $D$  such that  $\mathcal{L}(D) = L$ .
- There is an NFA  $N$  such that  $\mathcal{L}(N) = L$ .
- There is a regular expression  $R$  such that  $\mathcal{L}(R) = L$ .

# Next Time

- ***On Beyond Regular!***
  - What kinds of languages *aren't* regular?
- ***The Myhill-Nerode Theorem***
  - Proof technique to show non-regular