

# Regular Expressions

Recap from Last Time

# Regular Languages

- A language  $L$  is a ***regular language*** iff there is a DFA  $D$  such that  $\mathcal{L}(D) = L$ .
- ***Theorem:*** The following are equivalent:
  - $L$  is a regular language.
  - There is a **DFA** for  $L$ .
  - There is an **NFA** for  $L$ .

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

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

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

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

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

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

# 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 one string **in**  $\bar{L}_1$ .

Name one string **not 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:

- $\overline{L_1}$
- $L_1 \cup L_2$
- $L_1 \cap L_2$
- $L_1 L_2$
- $L_1^*$

- 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 one string **in**  $L_1 \cup L_2$ .

Name one string **not 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:

- $\overline{L_1}$
- $L_1 \cup L_2$
- $L_1 \cap L_2$
- $L_1 L_2$
- $L_1^*$

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

## Quick check 3:

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

Name one string **in**  $L_1 \cap L_2$ .

Name one string **not 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:

- $\overline{L_1}$
- $L_1 \cup L_2$
- $L_1 \cap L_2$
- $L_1 L_2$
- $L_1^*$

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

## Quick check 4:

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

Name one string **in**  $L_1 L_2$ .

Name one string **not in**  $L_1 L_2$ .

# Closure Properties

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

- $\overline{L_1}$
- $L_1 \cup L_2$
- $L_1 \cap L_2$
- $L_1 L_2$
- $L_1^*$

- 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 one string **in**  $L_1^*$ .

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

New Stuff!

# Another View of Regular Languages

# Rethinking Regular Languages

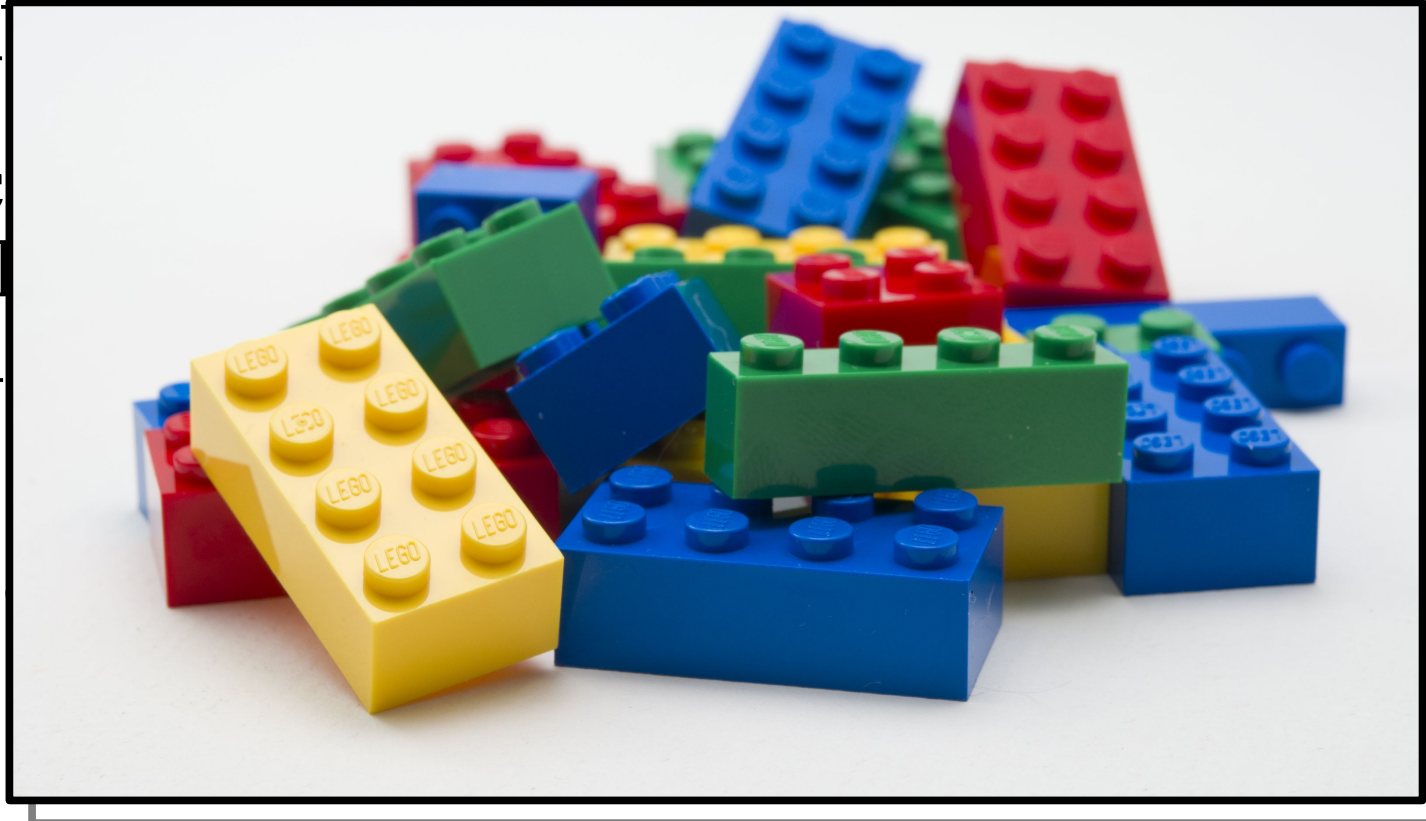
- We currently have several tools for showing a language  $L$  is regular:
  - Construct a **DFA** for  $L$ .
  - Construct an **NFA** for  $L$ .
  - Combine several simpler regular languages together via **closure properties** to form  $L$ .
- Today we expand on this last idea.

# Constructing Regular Languages

- **Idea:** Build up all regular languages as follows:
  - Start with a small set of simple languages we already know to be regular.
  - Using closure properties, combine these simple languages together to form more elaborate languages.
- *This is a bottom-up approach to the regular languages.*

# Constructing Regular Languages

- **Idea:** Build up all regular languages as follows:
  - Start with a small set of simple languages we already
  - Using c  
simple l  
elabora
- *This is a regular l*





# Regular Expressions

- ***Regular expressions*** are a way of describing a language via a string representation.
- They're used just about everywhere:
  - They're built into the JavaScript language and used for data validation.
  - They're used in the UNIX grep and flex tools to search files and build compilers.
  - They're employed to clean and scrape data for large-scale analysis projects.
- Conceptually, regular expressions are strings describing how to assemble a larger language out of smaller pieces.

# Atomic Regular Expressions

- The regular expressions begin with three simple building blocks.
- The symbol  $\emptyset$  is a regular expression that represents the empty language  $\emptyset$ .
- For any  $a \in \Sigma$ , the symbol  $a$  is a regular expression for the language  $\{a\}$ .
- The symbol  $\epsilon$  is a regular expression that represents the language  $\{\epsilon\}$ .
  - **Remember:**  $\{\epsilon\} \neq \emptyset!$
  - **Remember:**  $\{\epsilon\} \neq \epsilon!$

# Compound Regular Expressions

- If  $R_1$  and  $R_2$  are regular expressions,  $R_1R_2$  is a regular expression for the **concatenation** of the languages of  $R_1$  and  $R_2$ .
- If  $R_1$  and  $R_2$  are regular expressions,  $R_1 \cup R_2$  is a regular expression for the **union** of the languages of  $R_1$  and  $R_2$ .
- If  $R$  is a regular expression,  $R^*$  is a regular expression for the **Kleene closure** of the language of  $R$ .
- If  $R$  is a regular expression,  $(R)$  is a regular expression with the same meaning as  $R$ .

# Operator Precedence

- Here's the operator precedence for regular expressions:

$(R)$

$R^*$

$R_1R_2$

$R_1 \cup R_2$

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

# Regular Expression Examples

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

{ **trick**, **trear** }.

- 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:
  - $\mathcal{L}(\epsilon) = \{\epsilon\}$
  - $\mathcal{L}(\emptyset) = \emptyset$
  - $\mathcal{L}(a) = \{a\}$
  - $\mathcal{L}(R_1R_2) = \mathcal{L}(R_1) \mathcal{L}(R_2)$
  - $\mathcal{L}(R_1 \cup R_2) = \mathcal{L}(R_1) \cup \mathcal{L}(R_2)$
  - $\mathcal{L}(R^*) = \mathcal{L}(R)^*$
  - $\mathcal{L}((R)) = \mathcal{L}(R)$

Worthwhile activity: Apply this recursive definition to

**$a(b \cup c)((d))$**

and see what you get.

# Regular Expressions, Formally

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

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

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

**bbabbbaabab**

**aaaa**

**bbbbbabbbbaabbbb**

# Designing Regular Expressions

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

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

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

$\mathbf{aaaa}$

$\mathbf{bbbbbabbbb} \mathbf{aa} \mathbf{bbbbbb}$

# Designing Regular Expressions

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

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

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

$\mathbf{aaaa}$

$\mathbf{bbbbbabbbb} \mathbf{aa} \mathbf{bbbbbb}$

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

$\Sigma\Sigma\Sigma\Sigma$

$\mathbf{a}\mathbf{a}\mathbf{a}\mathbf{a}$

$\mathbf{b}\mathbf{a}\mathbf{b}\mathbf{a}$

$\mathbf{b}\mathbf{b}\mathbf{b}\mathbf{b}$

$\mathbf{b}\mathbf{a}\mathbf{a}\mathbf{a}$

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

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

$\mathbf{b} \mathbf{b} \mathbf{b} \mathbf{b} \mathbf{a} \mathbf{b} \mathbf{b} \mathbf{b}$

$\mathbf{b} \mathbf{b} \mathbf{b} \mathbf{b} \mathbf{b} \mathbf{b}$

$\mathbf{a} \mathbf{b} \mathbf{b} \mathbf{b}$

$\mathbf{a}$

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

**$\textcolor{teal}{aa}^*$**

**$\textcolor{teal}{cs103}$**  $\textcolor{gray}{@cs.stanford.edu}$

**$\textcolor{teal}{first}.\textcolor{violet}{middle}.\textcolor{violet}{last}$**  $\textcolor{gray}{@mail.site.org}$

**$\textcolor{teal}{dot}.\textcolor{violet}{at}$**  $\textcolor{gray}{@dot.com}$

# 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  $\Sigma = \{ \text{a}, \text{.}, \text{@} \}$ , 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.

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

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

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

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

# A More Elaborate Design

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

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

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

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

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

# A More Elaborate Design

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

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

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

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

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



# 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 **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 **a** represents “some letter.”
- Let's make a regex for email addresses.

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

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

**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.

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

**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 **a** represents “some letter.”
- Let's make a regex for email addresses.

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

**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.

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

**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 **a** represents “some letter.”
- Let's make a regex for email addresses.

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

**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.

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

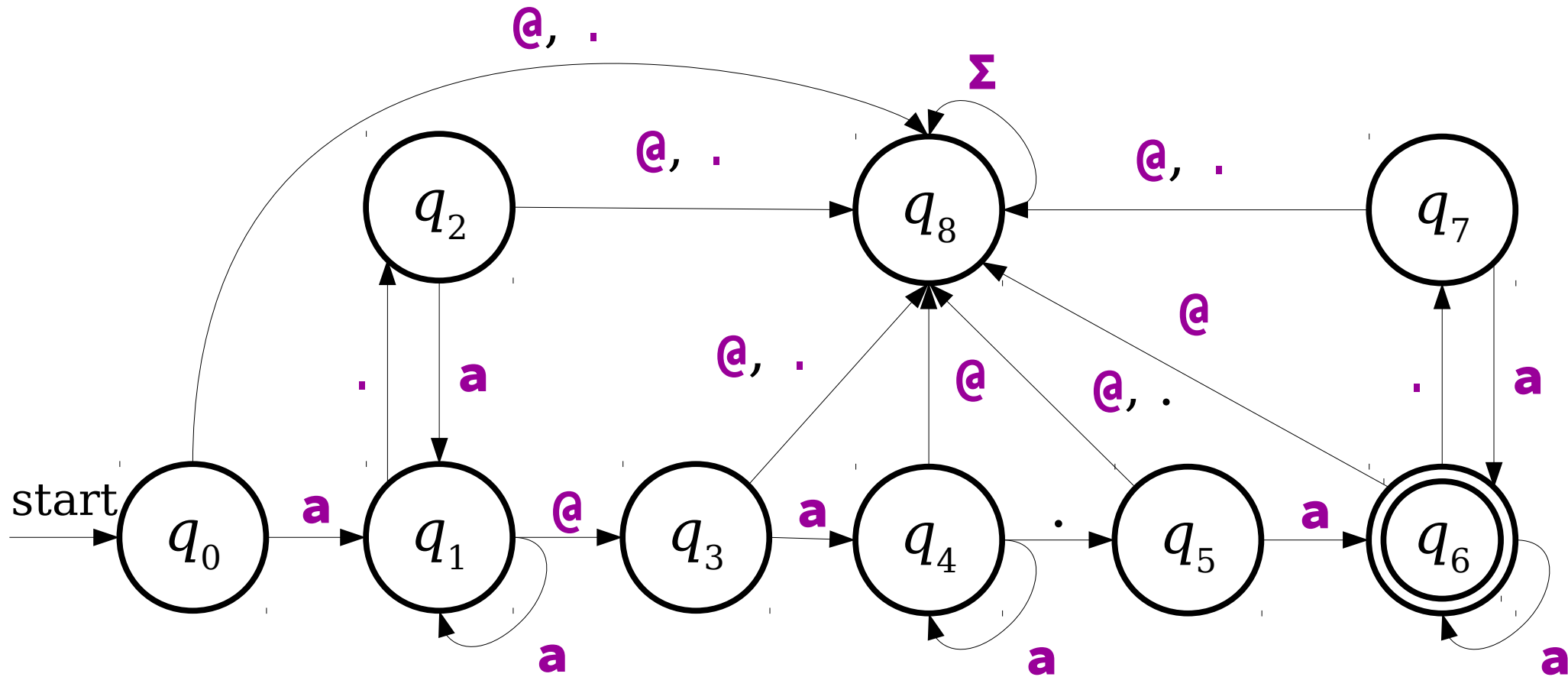
**cs103@cs.stanford.edu**

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

**dot.at@dot.com**

# For Comparison

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



# Shorthand Summary

- $R^n$  is shorthand for  $RR \dots R$  ( $n$  times).
  - Edge case: define  $R^0 = \varepsilon$ .
- $\Sigma$  is shorthand for “any character in  $\Sigma$ .”
- $R?$  is shorthand for  $(R \cup \varepsilon)$ , 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***

```
graph TD; A[Languages you can build a DFA for.] --> C((Regular Languages)); B[Languages you can build an NFA for.] --> C;
```

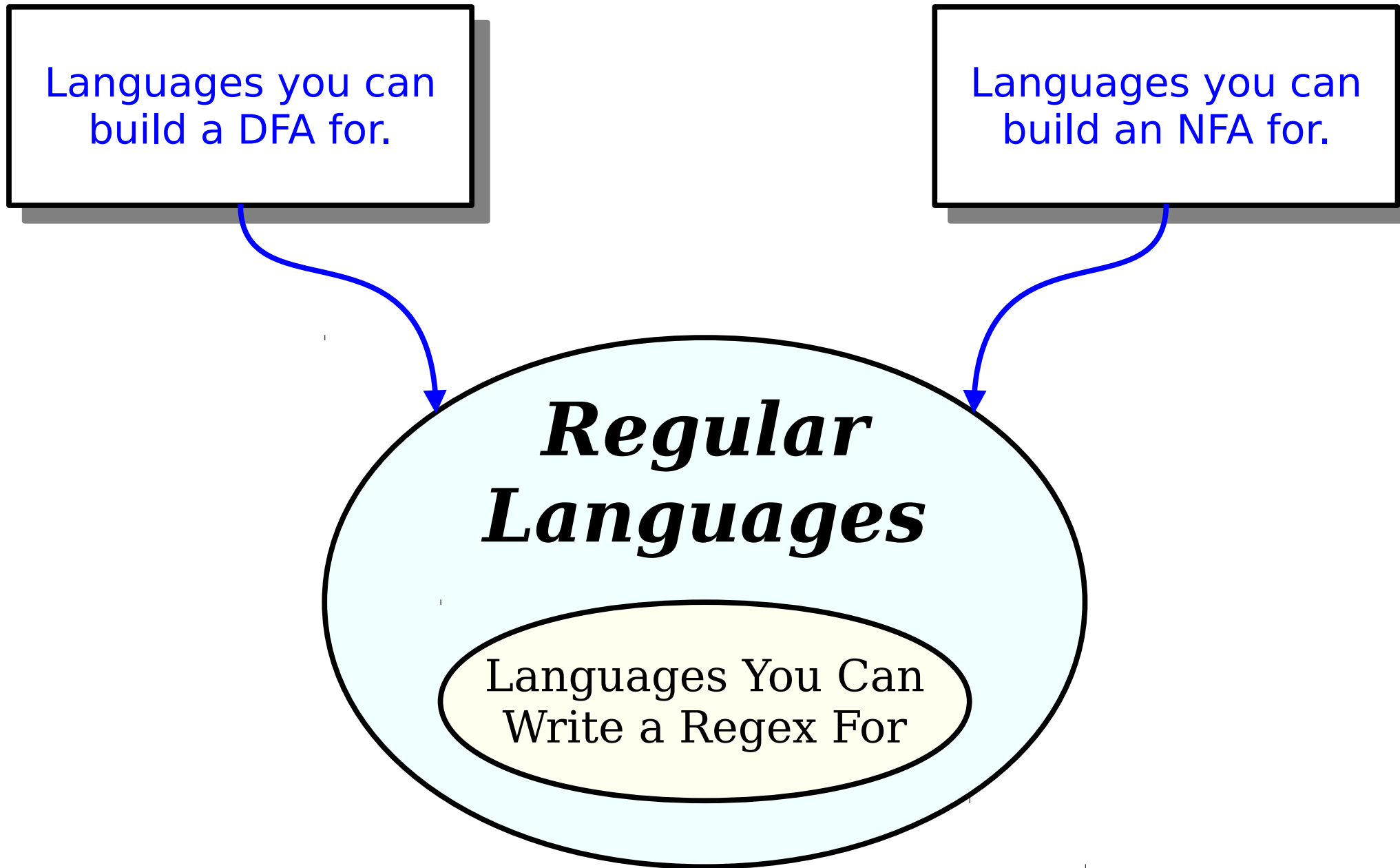
The diagram illustrates the relationship between different types of languages and regular languages. At the top, two rectangular boxes are positioned. The left box contains the text 'Languages you can build a DFA for.' and the right box contains 'Languages you can build an NFA for.'. Both boxes have a black border and a light gray drop shadow. Two blue curved arrows originate from the bottom of these boxes and point towards a central light blue oval. The oval has a black border and contains the text '***Regular Languages***' in a bold, italicized black font. This visualizes that both DFA-recognizable and NFA-recognizable languages are subsets of the class of regular languages.

Languages you can  
build a DFA for.

Languages you can  
build an NFA for.

***Regular  
Languages***

Languages You Can  
Write a Regex For

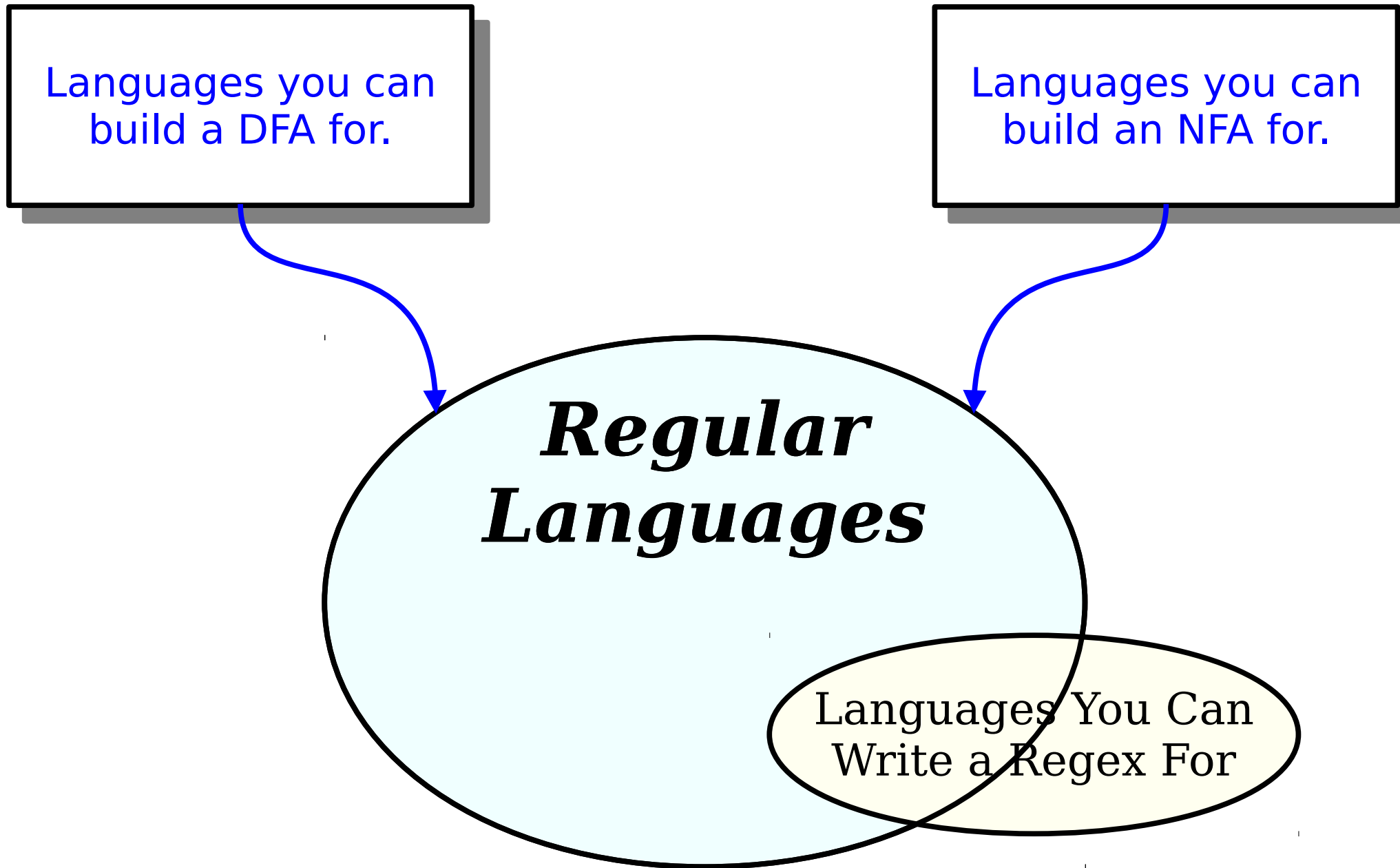


Languages you can  
build a DFA for.

Languages you can  
build an NFA for.

***Regular  
Languages***

Languages You Can  
Write a Regex For



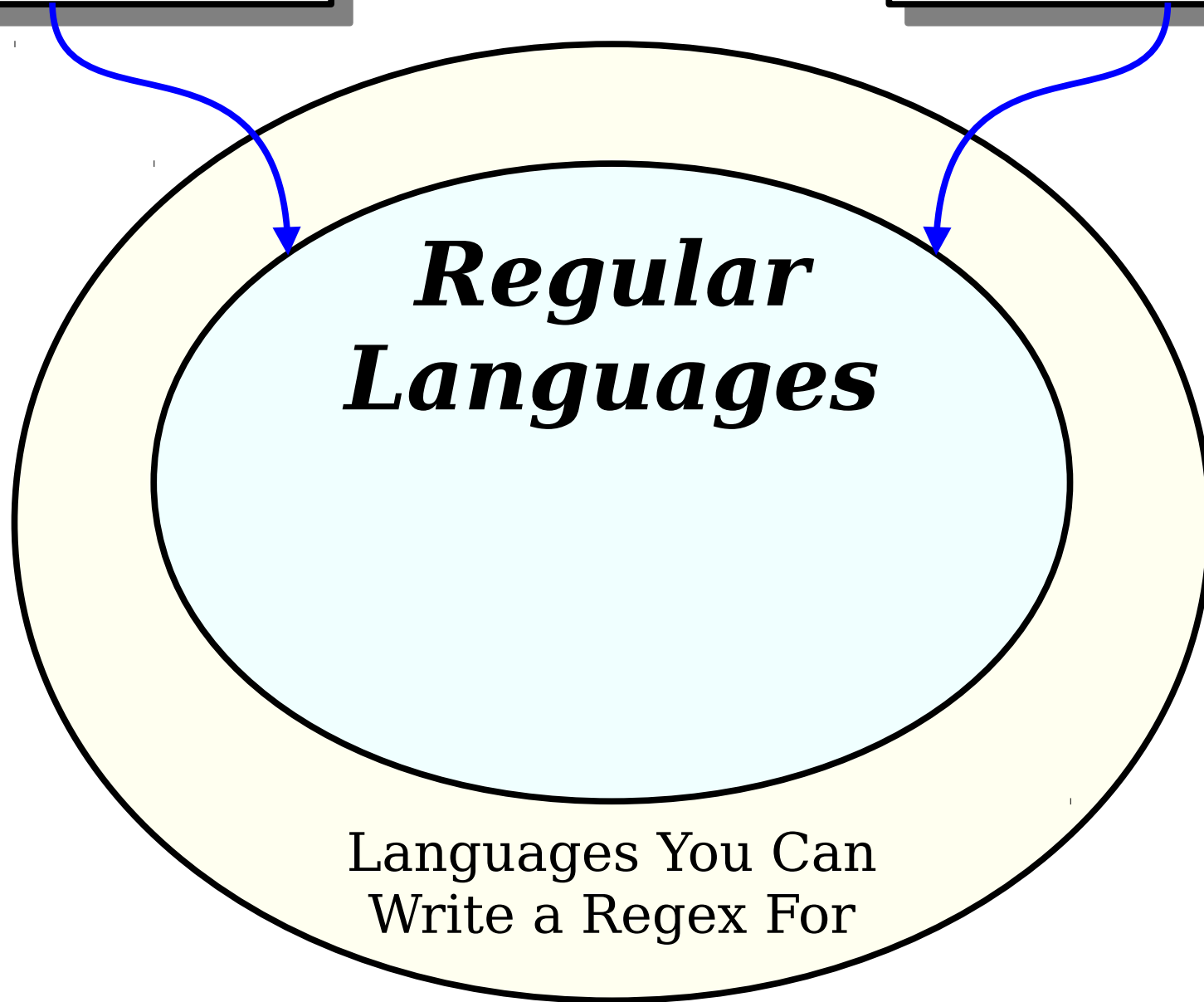


Languages you can  
build a DFA for.

Languages you can  
build an NFA for.

***Regular  
Languages***

Languages You Can  
Write a Regex For

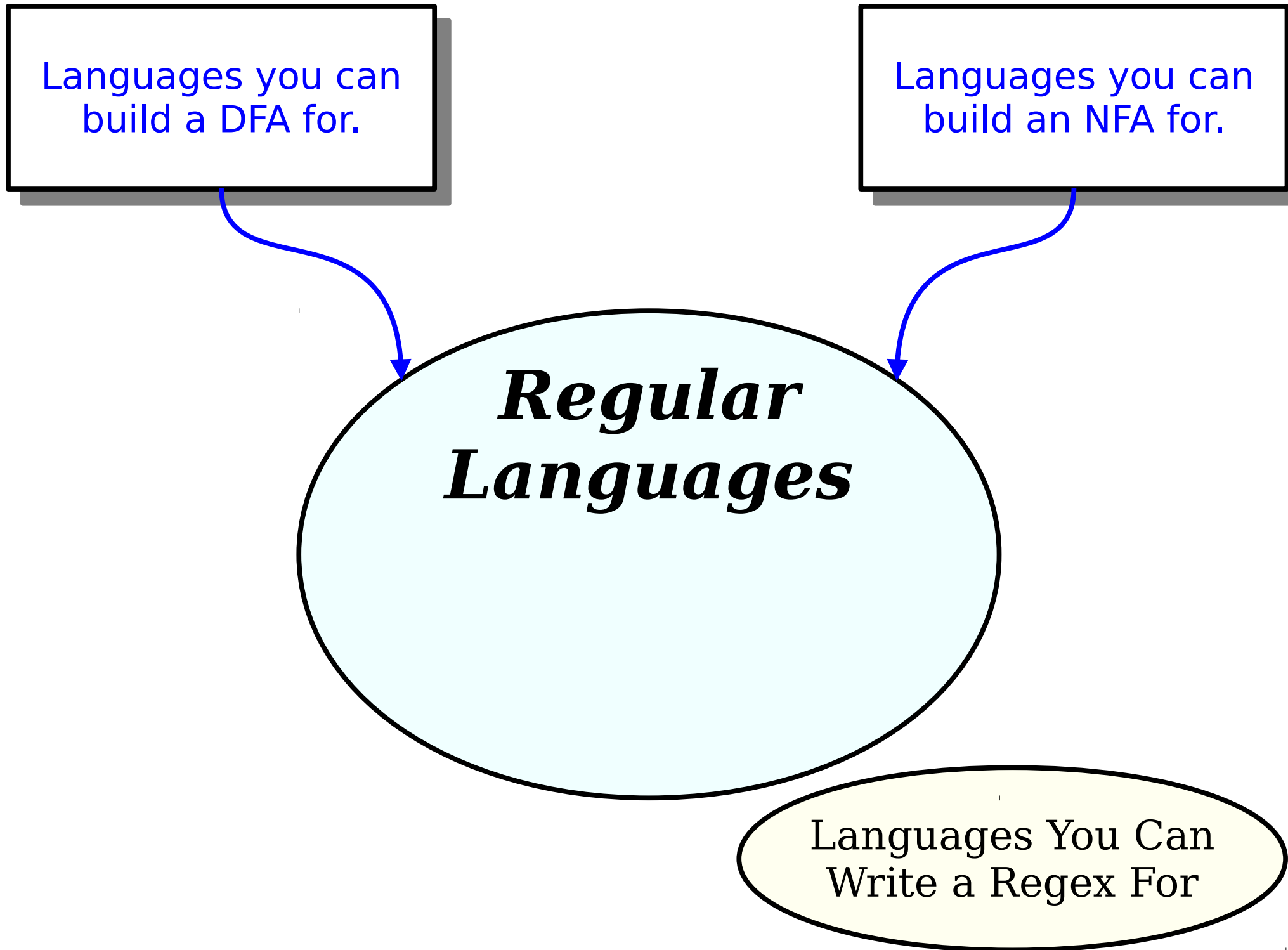


Languages you can  
build a DFA for.

Languages you can  
build an NFA for.

***Regular  
Languages***

Languages You Can  
Write a Regex For

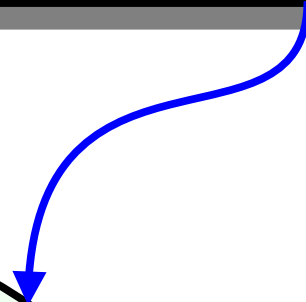
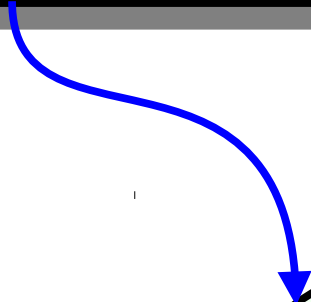


Languages you can  
build a DFA for.

Languages you can  
build an NFA for.

# ***Regular Languages***

Languages You Can  
Write a Regex For



# The Power of Regular Expressions

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

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

- The atomic regular expressions all represent regular languages.
- 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***

```
graph TD; A[Languages you can build a DFA for.] --> C((Regular Languages)); B[Languages you can build an NFA for.] --> C;
```

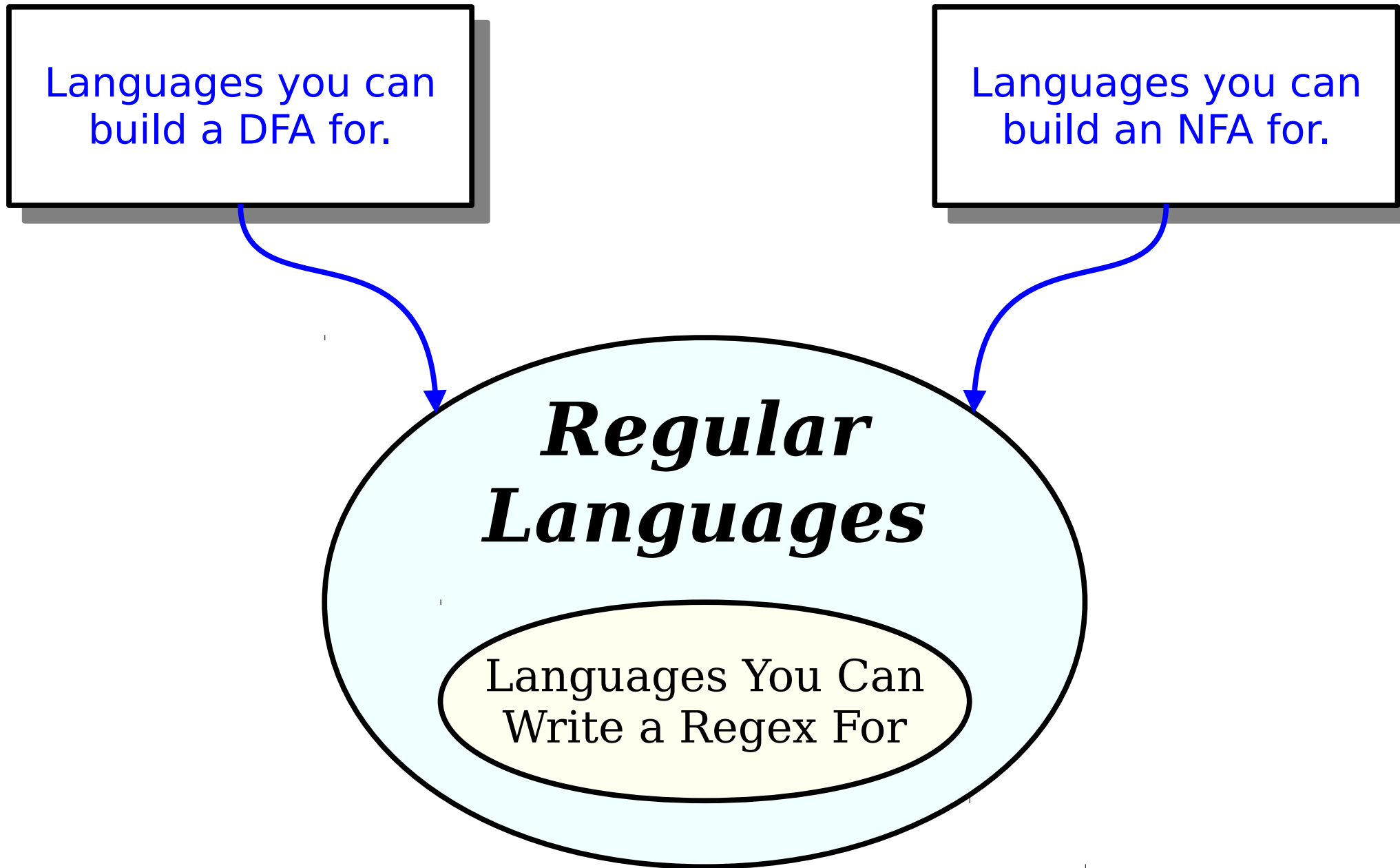
The diagram illustrates the relationship between different types of languages and the concept of Regular Languages. At the top, there are two rectangular boxes. The left box contains the text 'Languages you can build a DFA for.' and the right box contains 'Languages you can build an NFA for.'. Both boxes have a black border and a light gray drop shadow. Two blue curved arrows originate from the bottom of these boxes and point towards a central light blue oval. The oval has a black border and contains the text '***Regular Languages***' in a bold, italicized black font. This visualizes that both DFA-recognizable and NFA-recognizable languages are subsets of Regular Languages.

Languages you can  
build a DFA for.

Languages you can  
build an NFA for.

***Regular  
Languages***

Languages You Can  
Write a Regex For

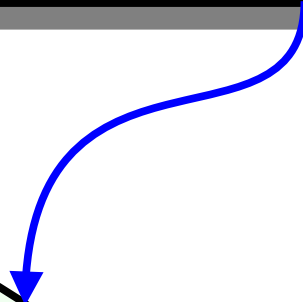
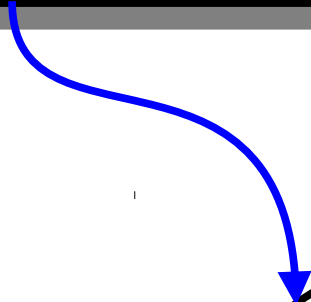


Languages you can  
build a DFA for.

Languages you can  
build an NFA for.

# ***Regular Languages***

Languages You Can  
Write a Regex For





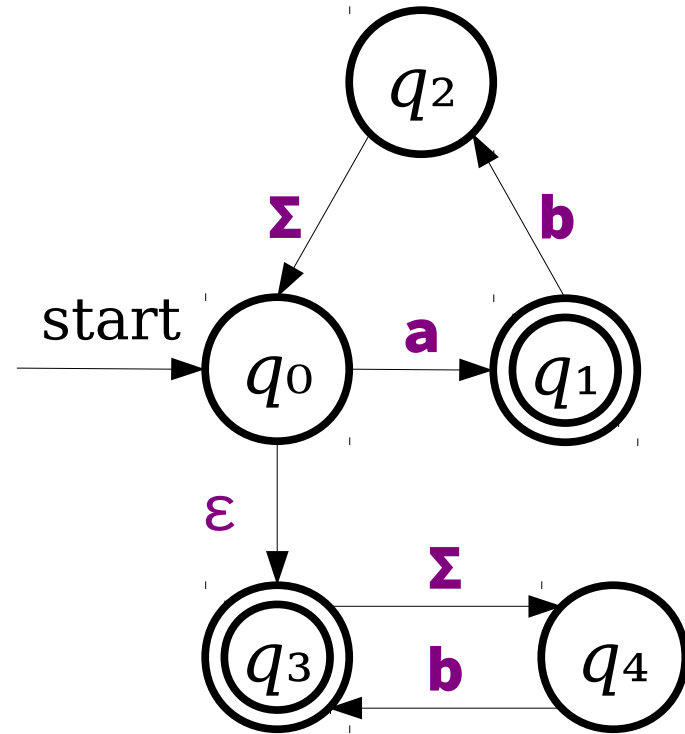
# The Power of Regular Expressions

***Theorem:*** If  $L$  is a regular language, then there is a regular expression for  $L$ .

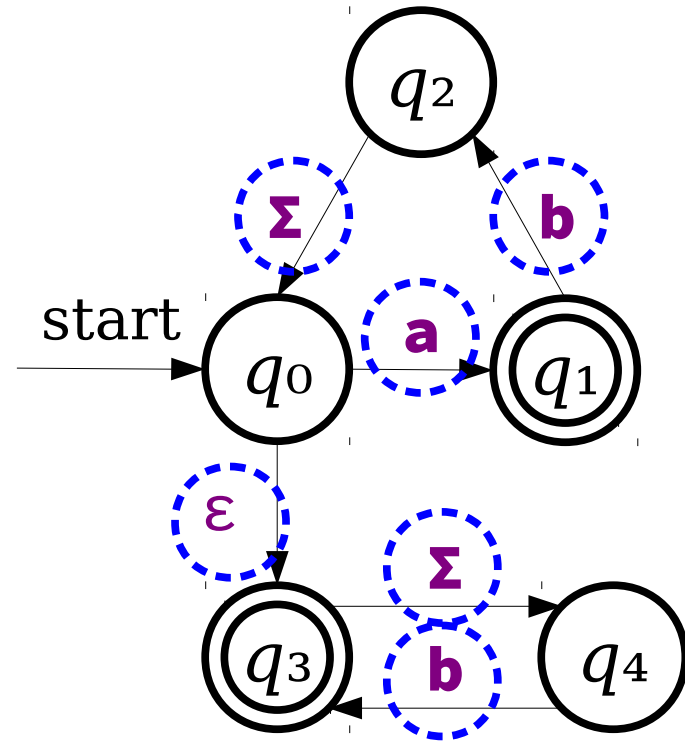
***This is not obvious!***

***Proof idea:*** Show how to convert an arbitrary NFA into a regular expression.

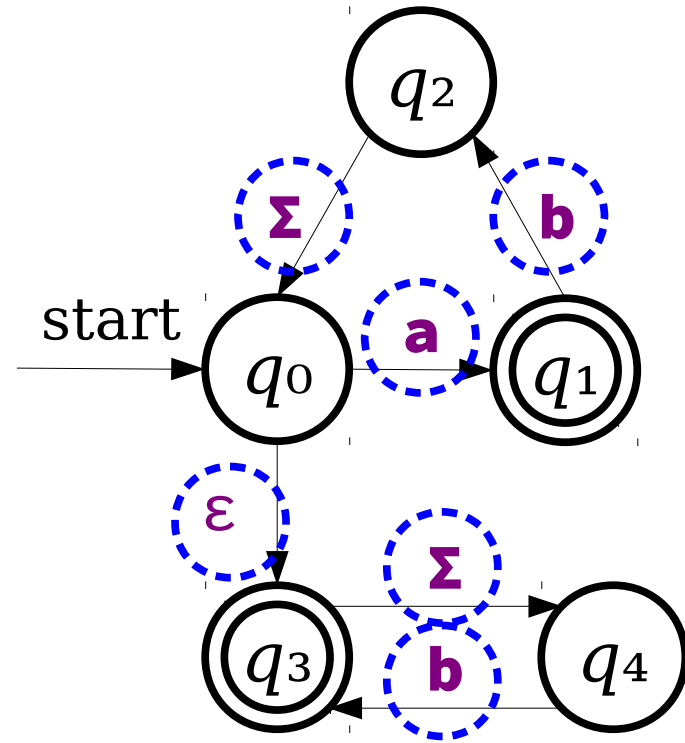
# Generalizing NFAs



# Generalizing NFAs

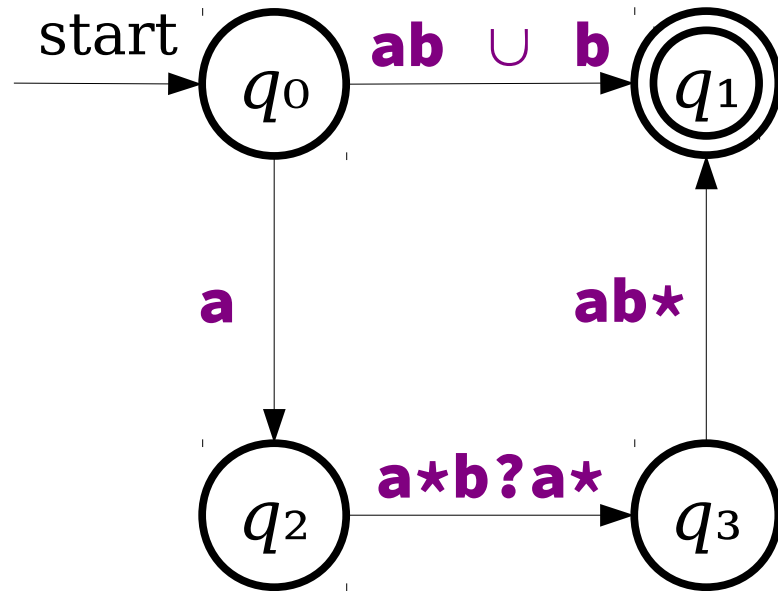


# Generalizing NFAs

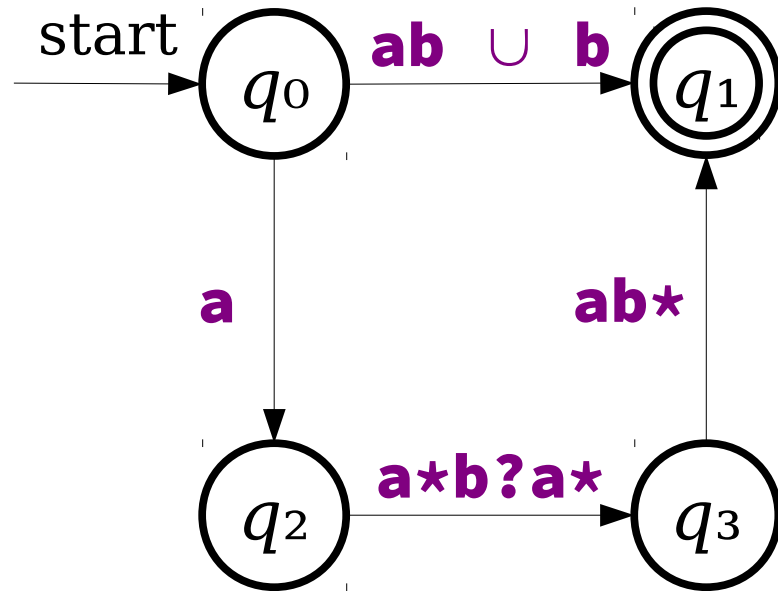


These are all regular expressions!

# Generalizing NFAs

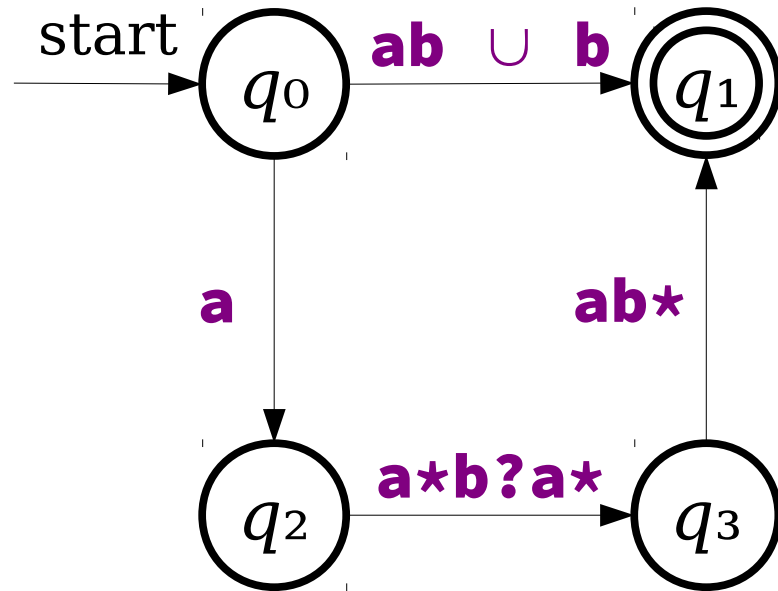


# Generalizing NFAs



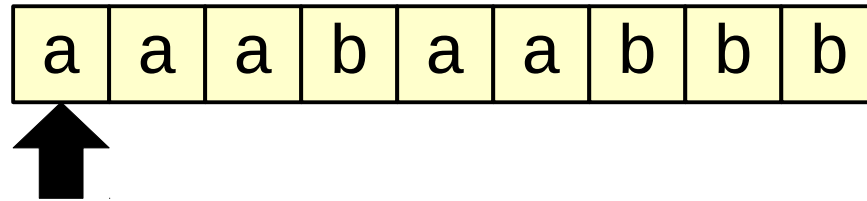
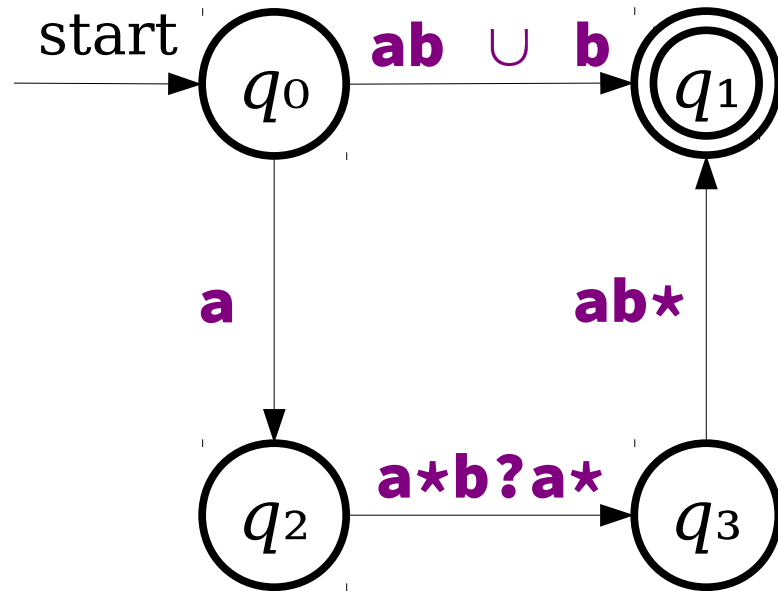
Note: Actual NFAs aren't allowed to have transitions like these. This is just a thought experiment.

# Generalizing NFAs



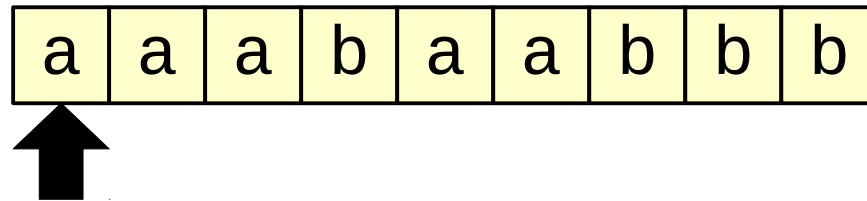
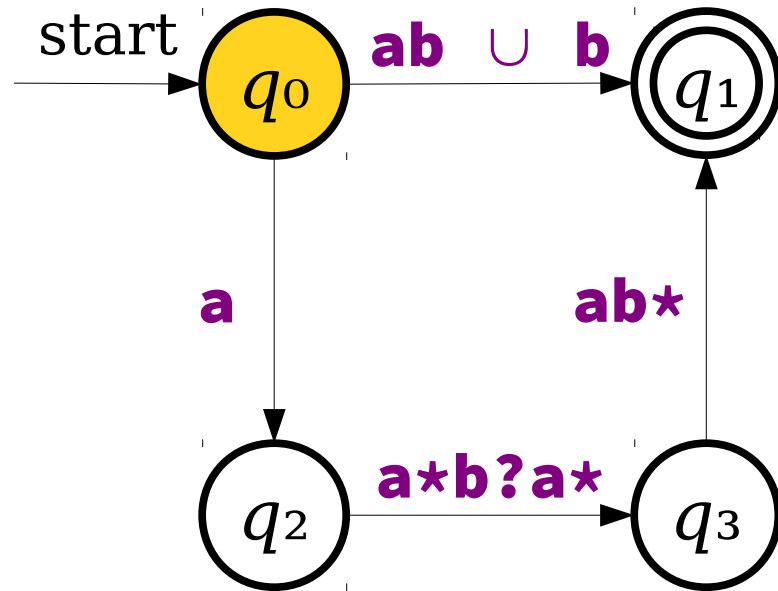
a	a	a	b	a	a	b	b	b
---	---	---	---	---	---	---	---	---

# Generalizing NFAs

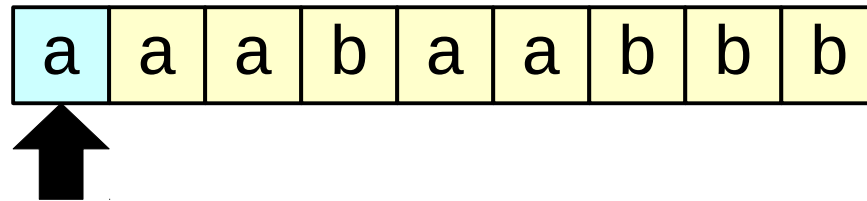
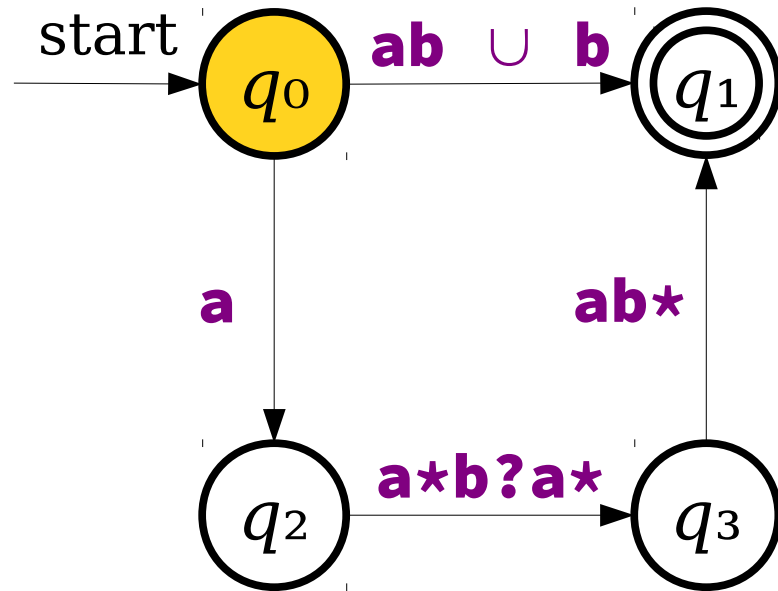




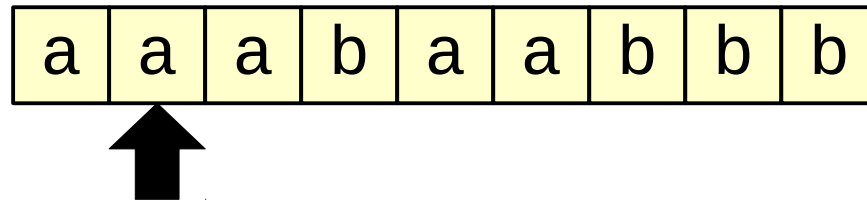
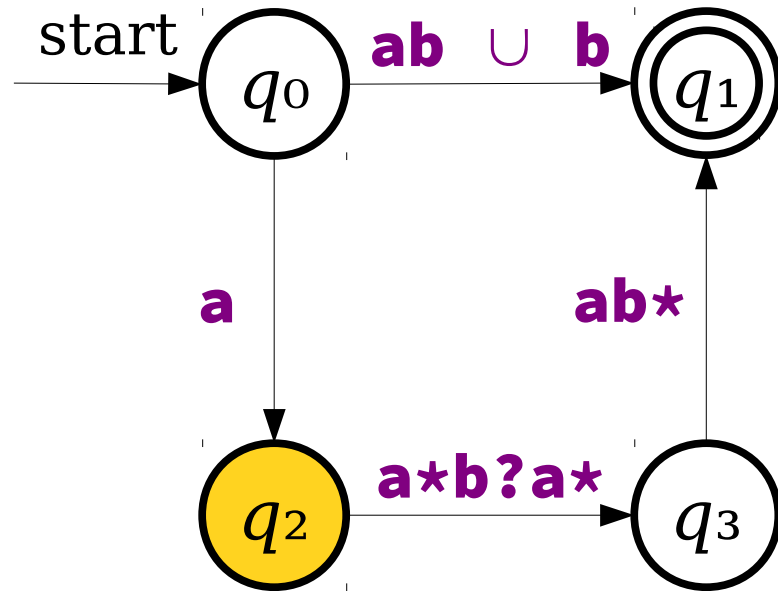
# Generalizing NFAs



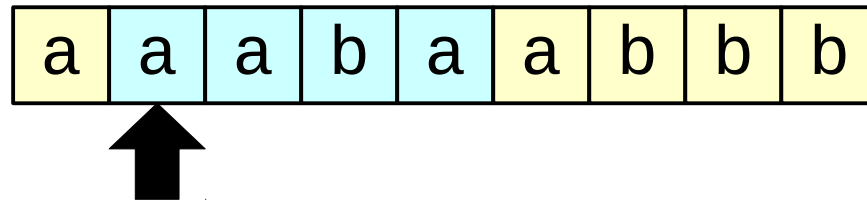
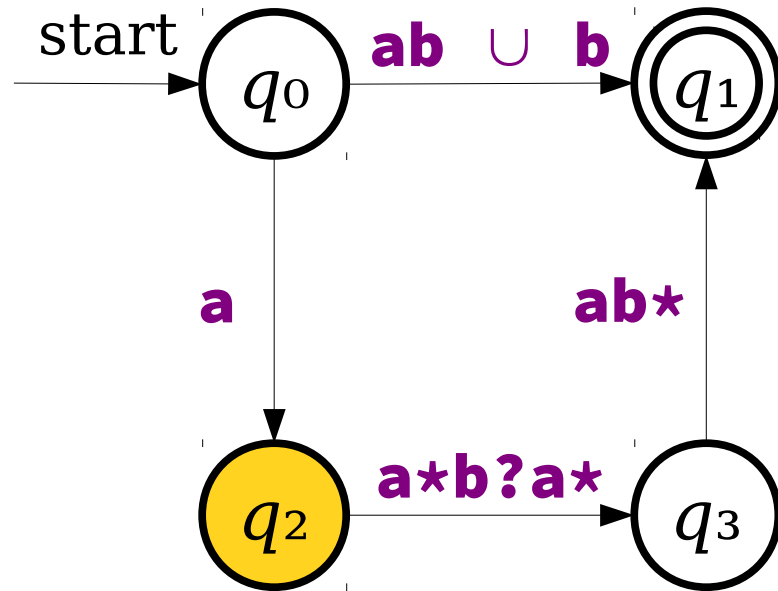
# Generalizing NFAs



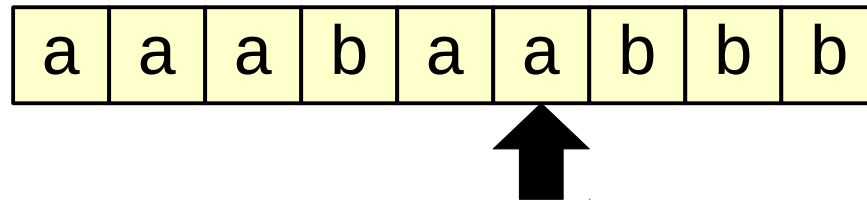
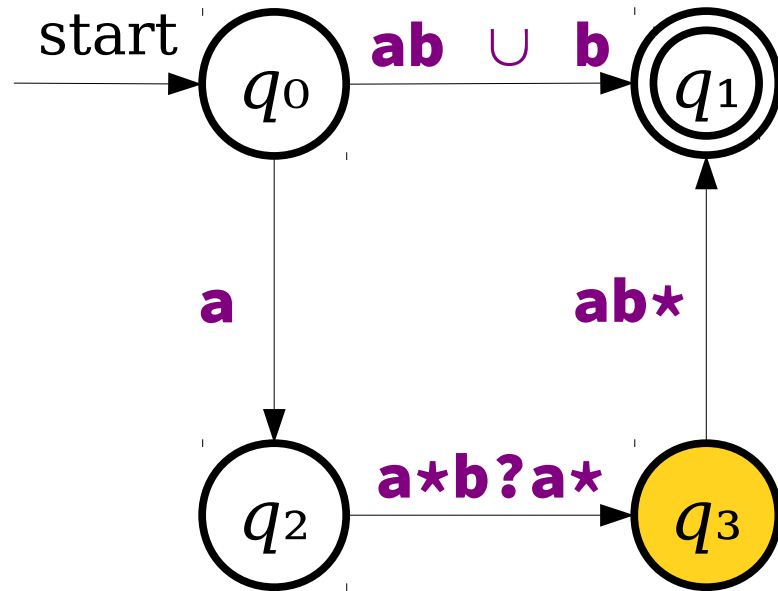
# Generalizing NFAs



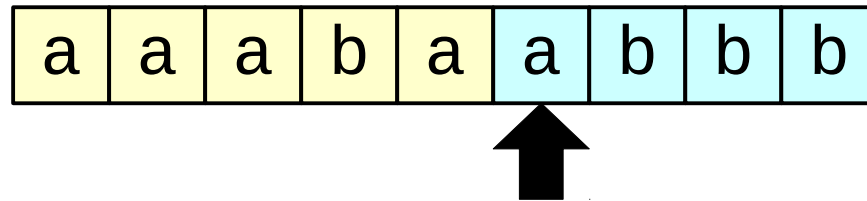
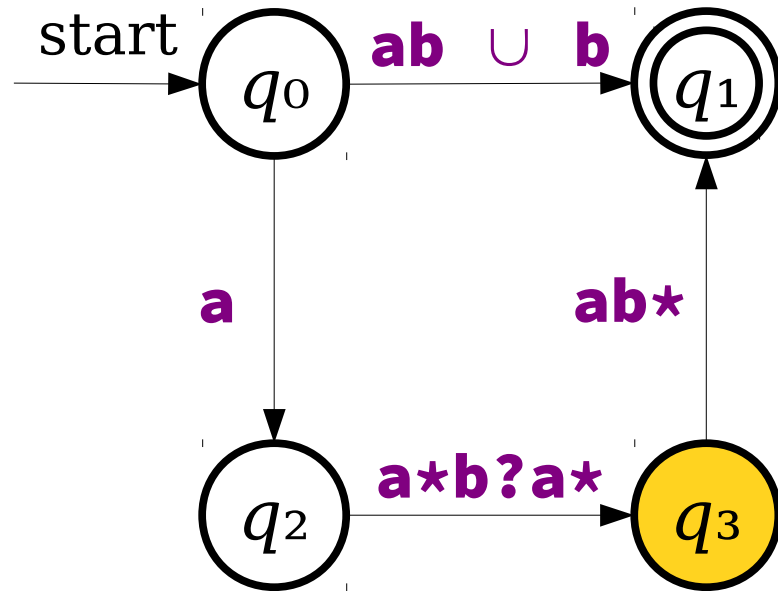
# Generalizing NFAs



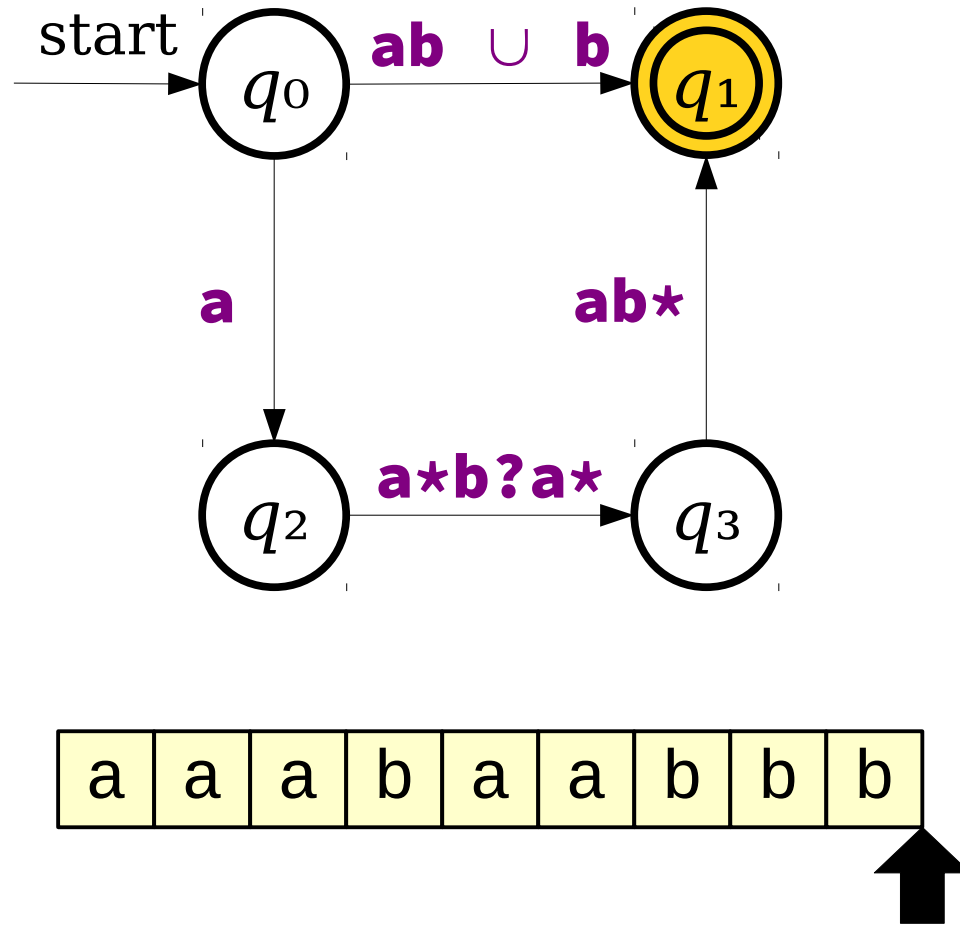
# Generalizing NFAs



# Generalizing NFAs



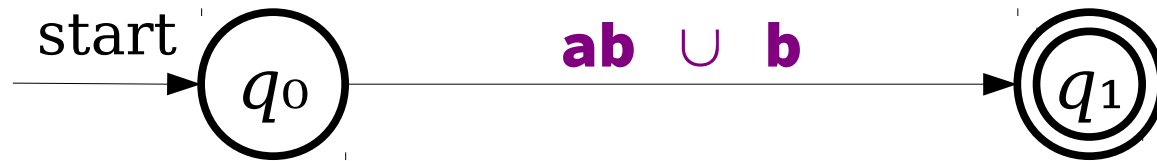
# Generalizing NFAs



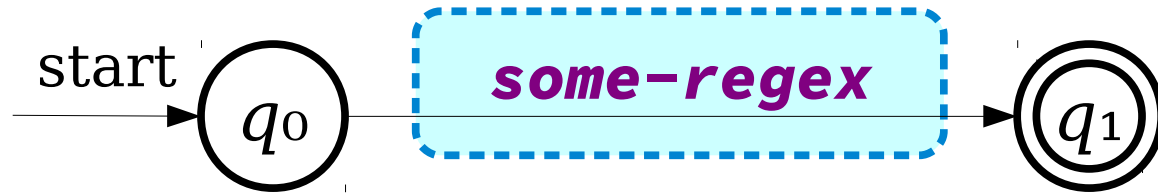
***Key Idea 1:*** Imagine that we can label transitions in an NFA with arbitrary regular expressions.



# Generalizing NFAs

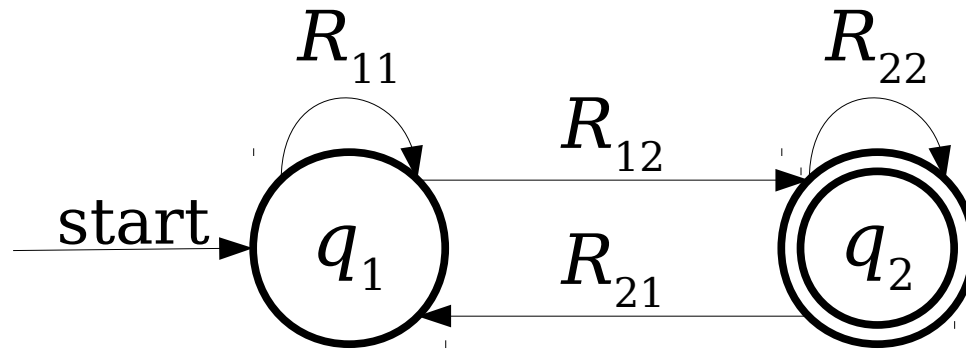


**Key Idea 2:** If we can convert an NFA into a generalized NFA that looks like this...

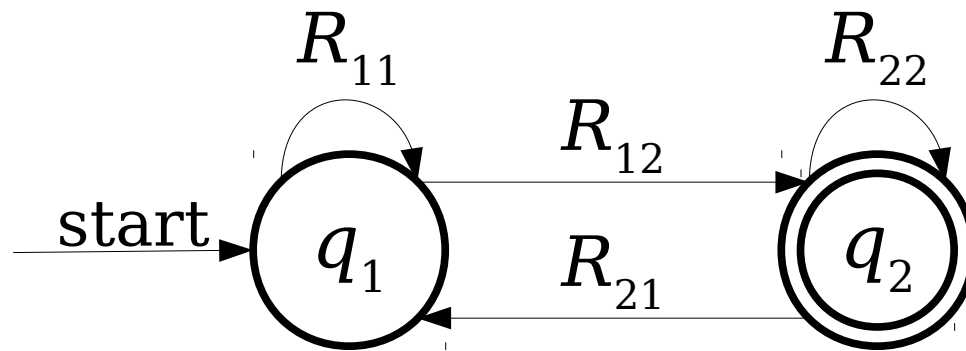


...then we can easily read off a regular expression for the original NFA.

# From NFAs to Regular Expressions

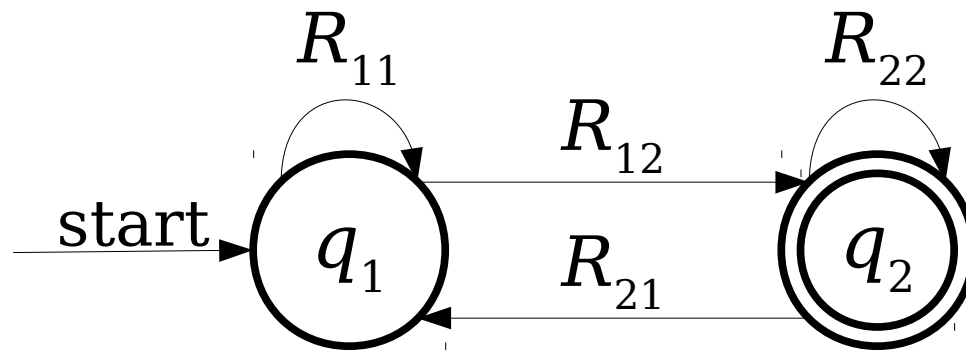


# From NFAs to Regular Expressions



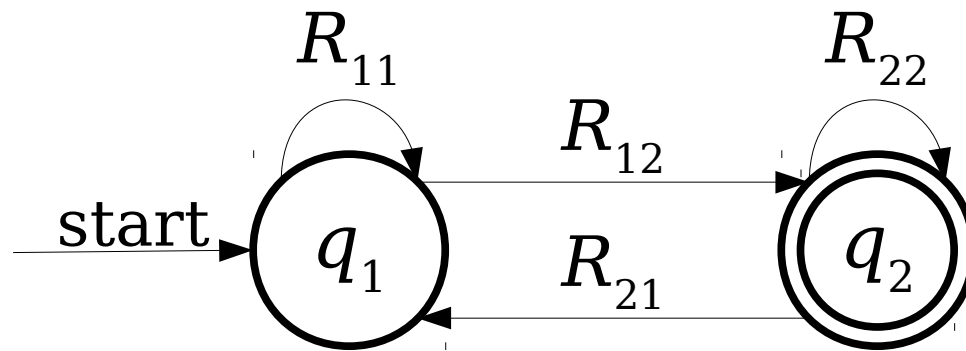
Here,  $R_{11}$ ,  $R_{12}$ ,  $R_{21}$ , and  $R_{22}$  are arbitrary regular expressions.

# From NFAs to Regular Expressions

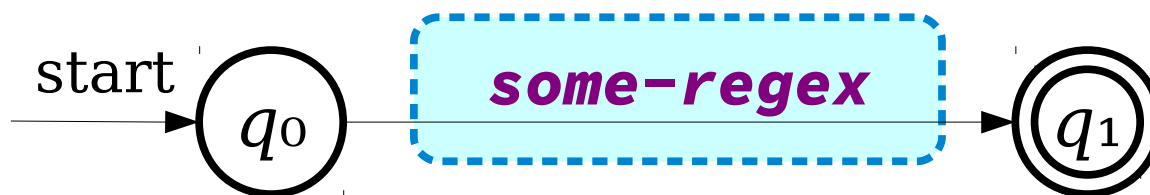


Question: Can we get a clean regular expression from this NFA?

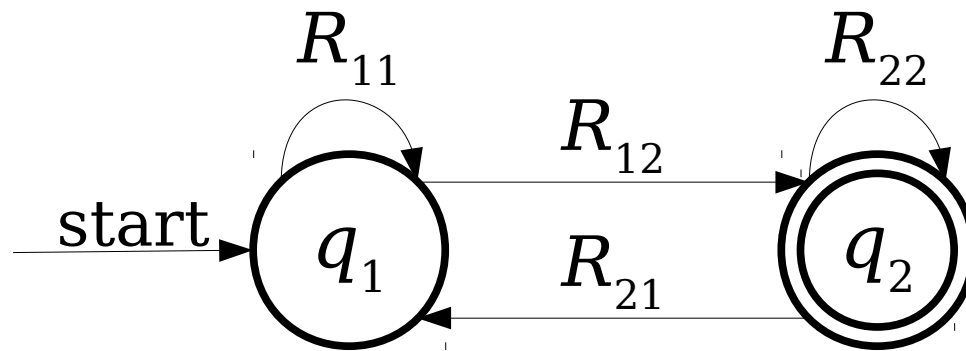
# From NFAs to Regular Expressions



**Key Idea 3:** Somehow transform this NFA so that it looks like this:

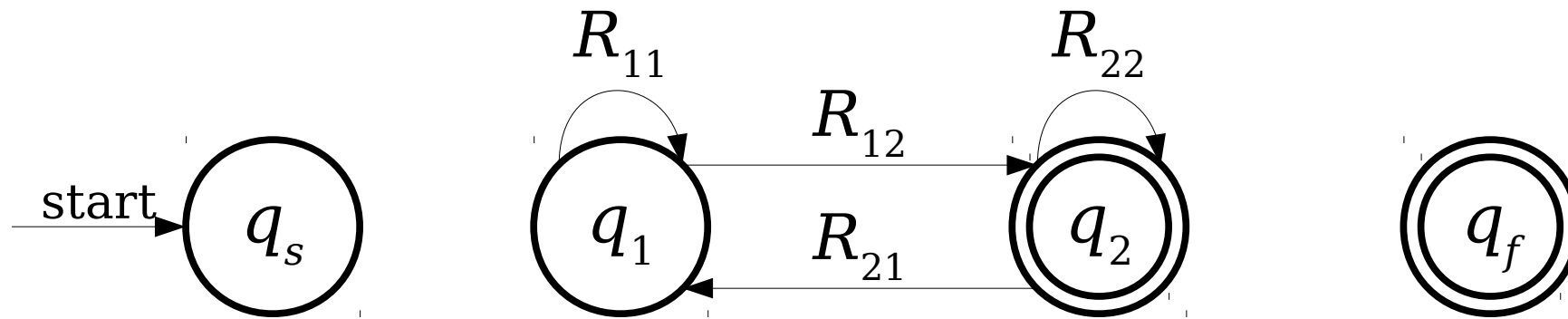


# From NFAs to Regular Expressions



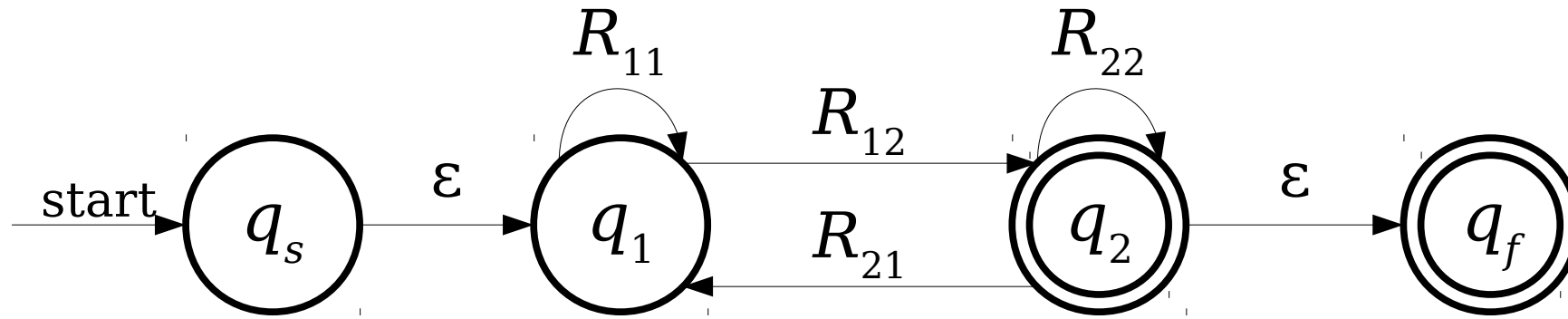
The first step is going to be a  
bit weird...

# From NFAs to Regular Expressions

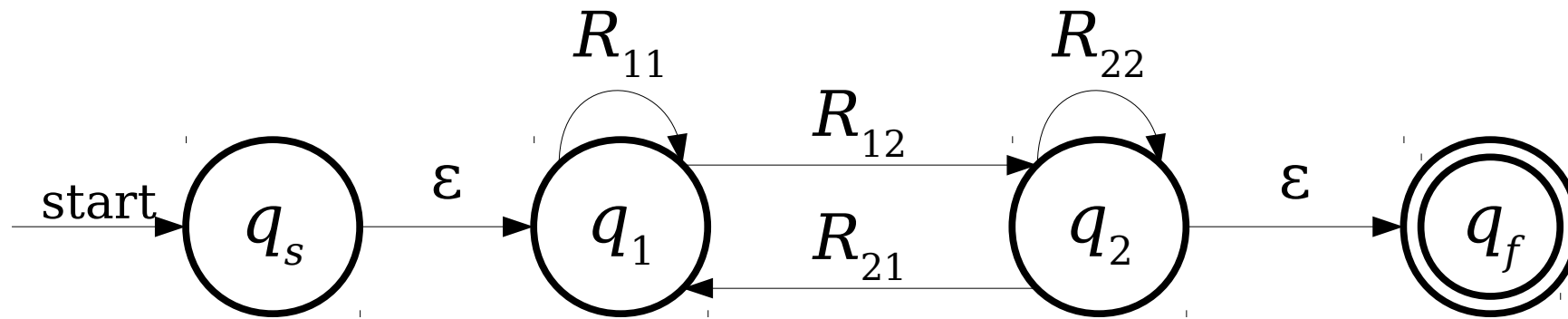




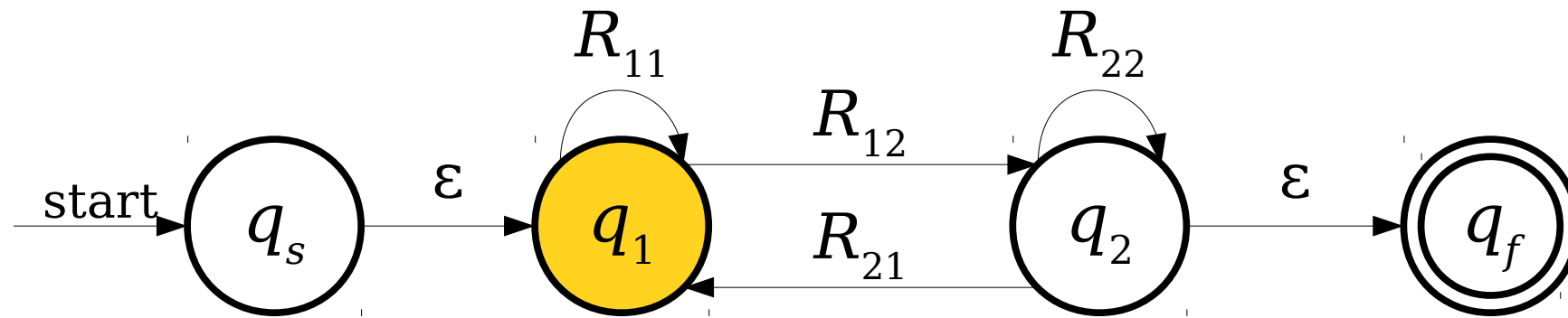
# From NFAs to Regular Expressions



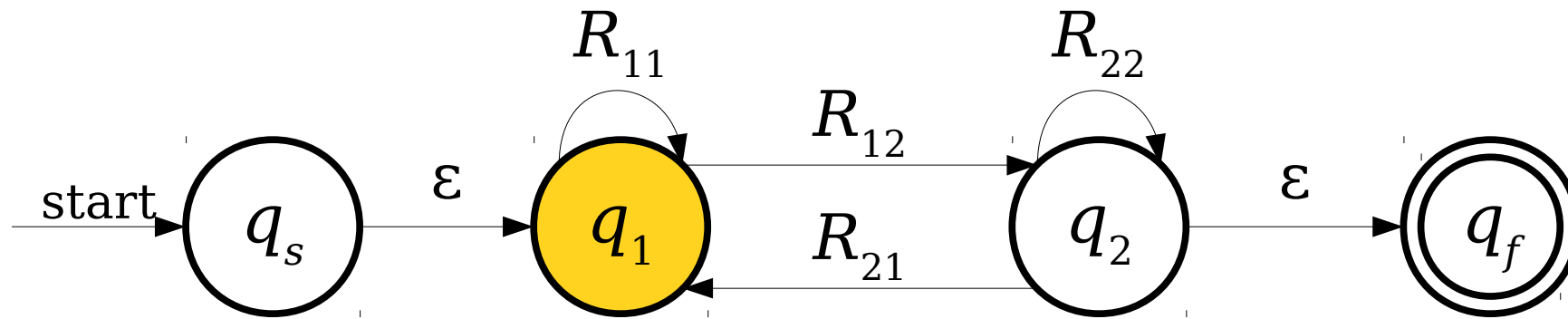
# From NFAs to Regular Expressions



# From NFAs to Regular Expressions

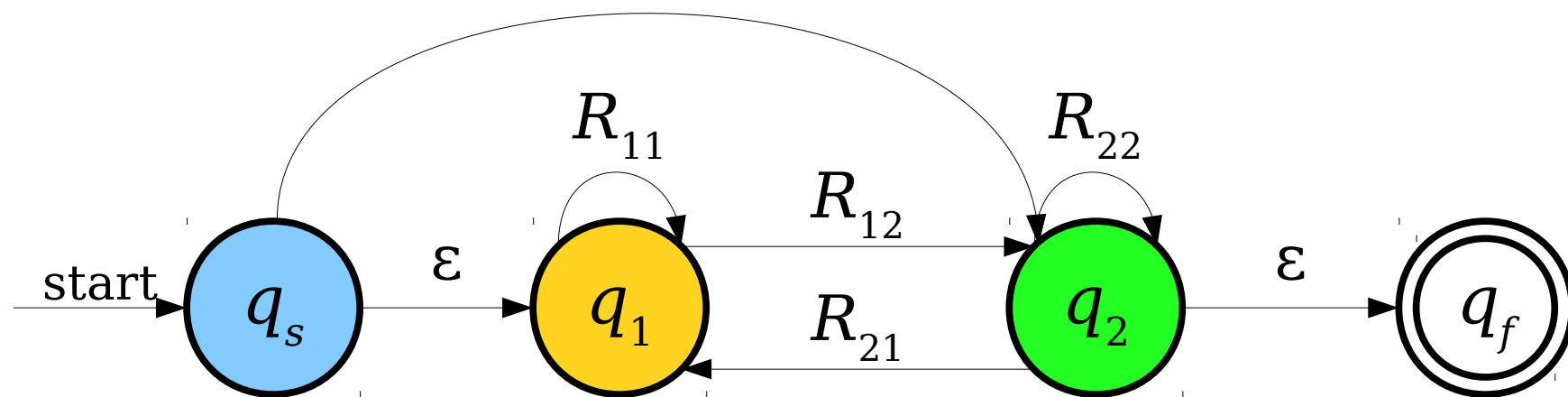


# From NFAs to Regular Expressions

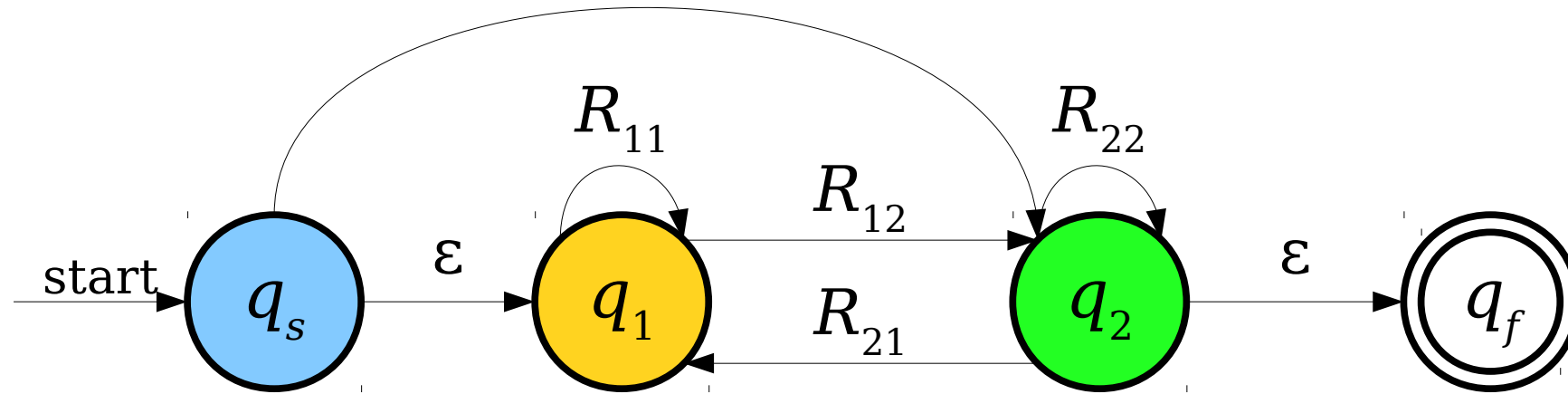


Could we eliminate  
this state from the  
NFA?

# From NFAs to Regular Expressions

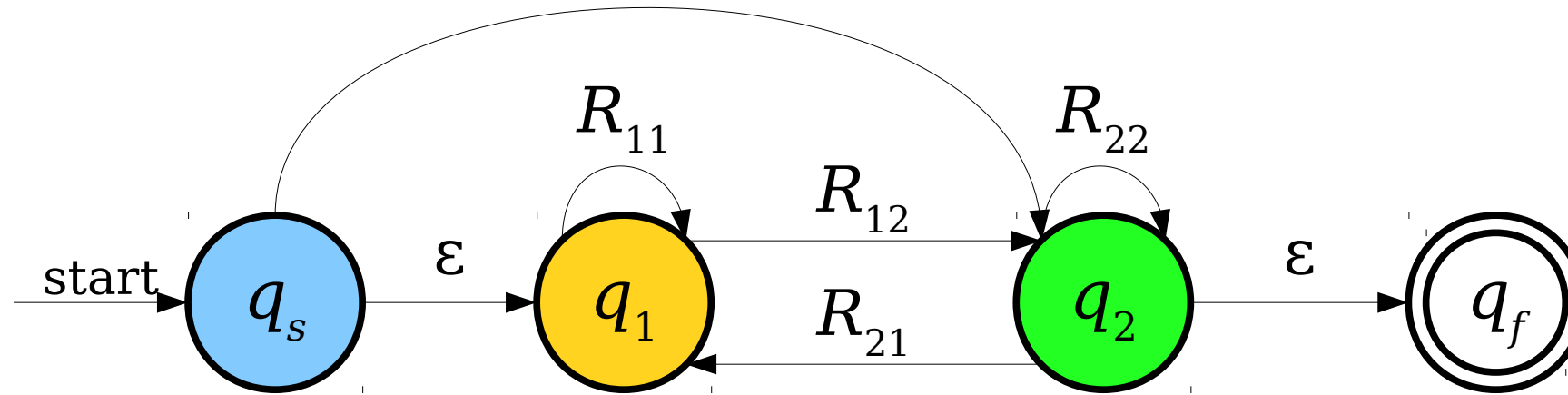


# From NFAs to Regular Expressions



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 quick check:

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

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

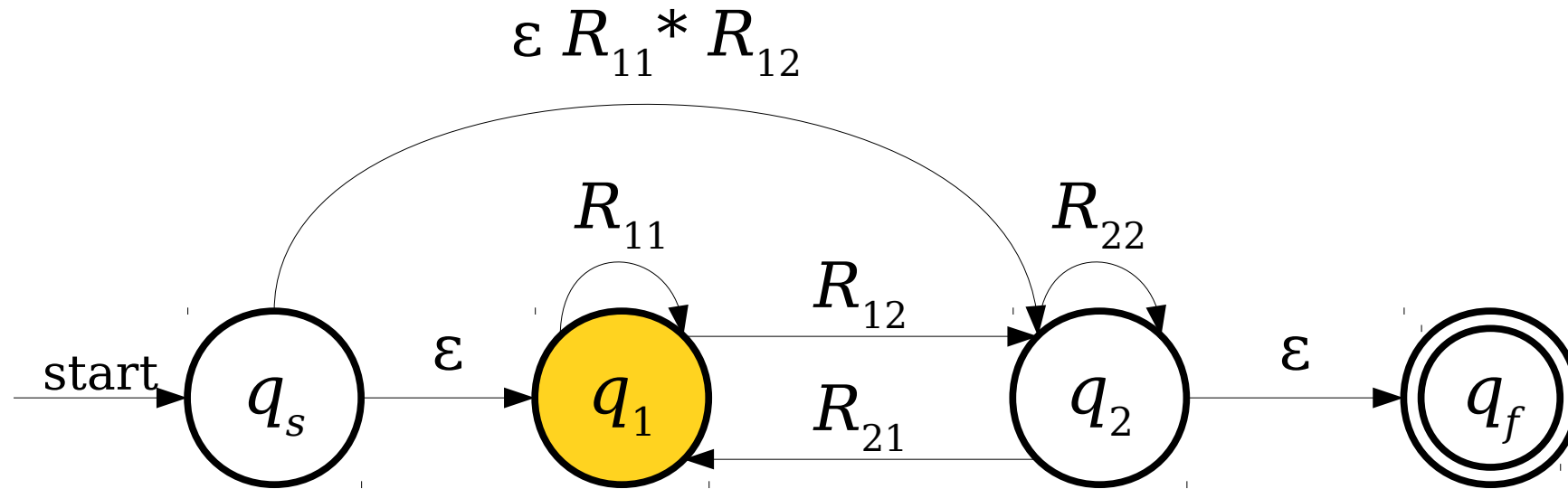
$\epsilon R_{11} R_{11} R_{12}$

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

$\epsilon R_{11} R_{12} R_{21}$

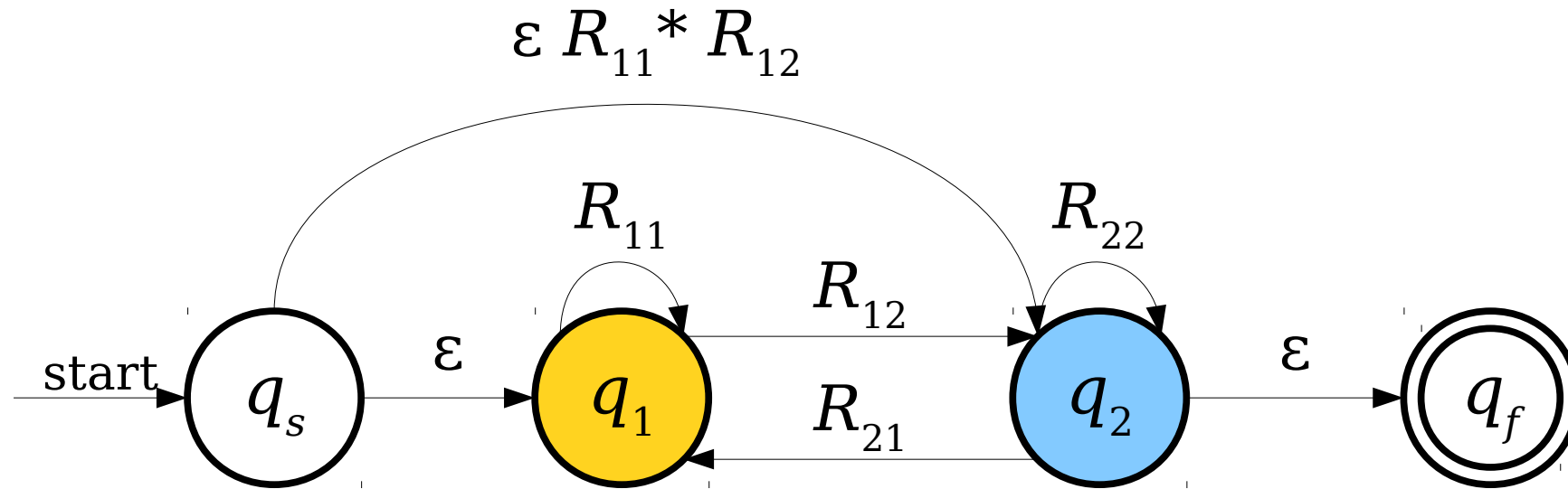
$\epsilon R_{11} R_{12} R_{21} R_{11} R_{12}$

# From NFAs to Regular Expressions

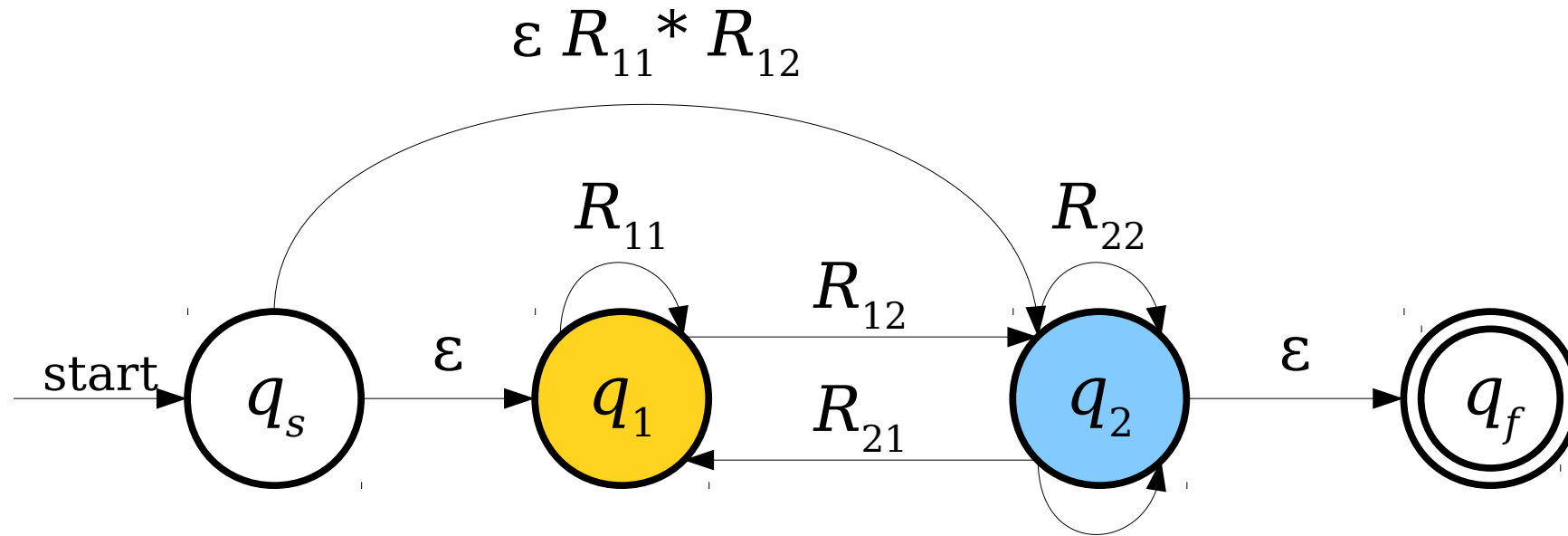




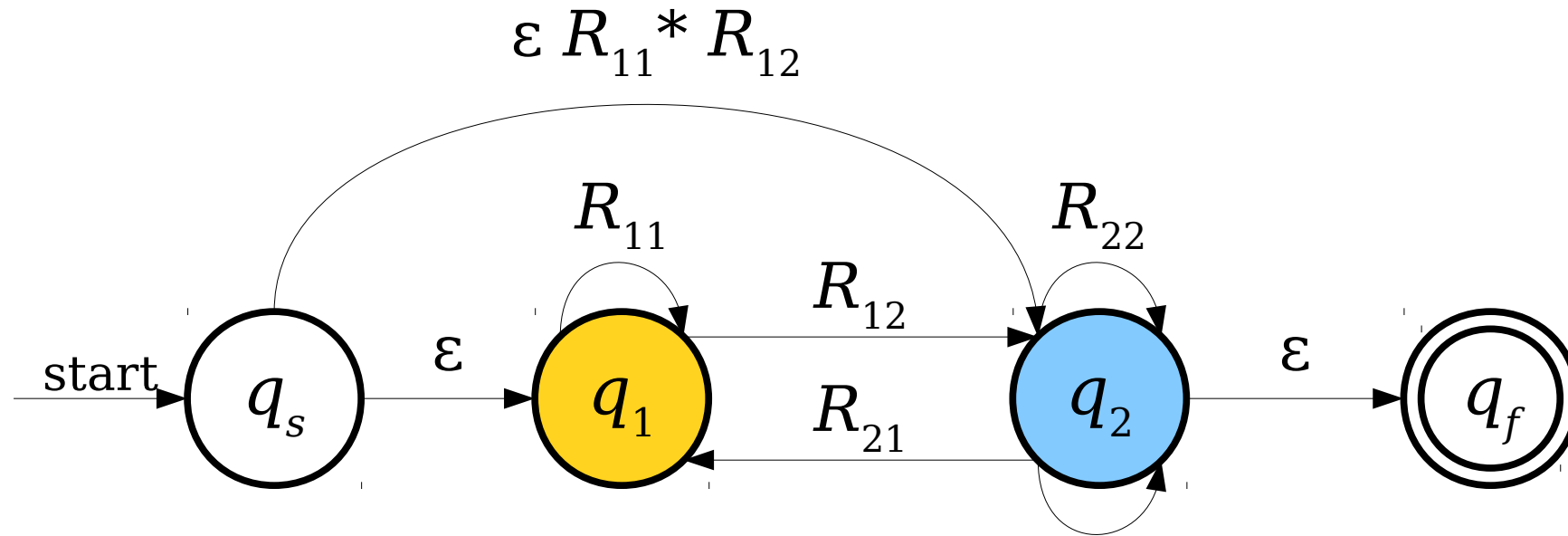
# From NFAs to Regular Expressions



# From NFAs to Regular Expressions

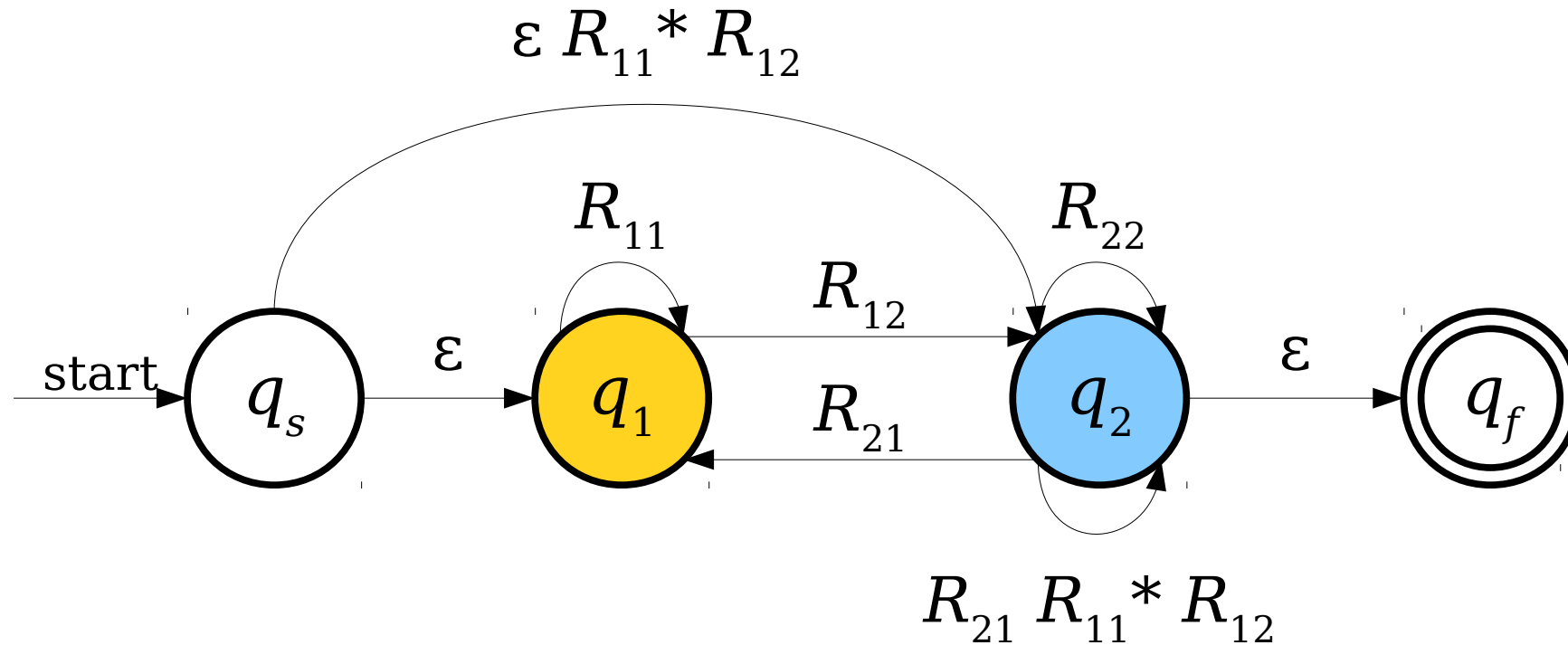


# From NFAs to Regular Expressions



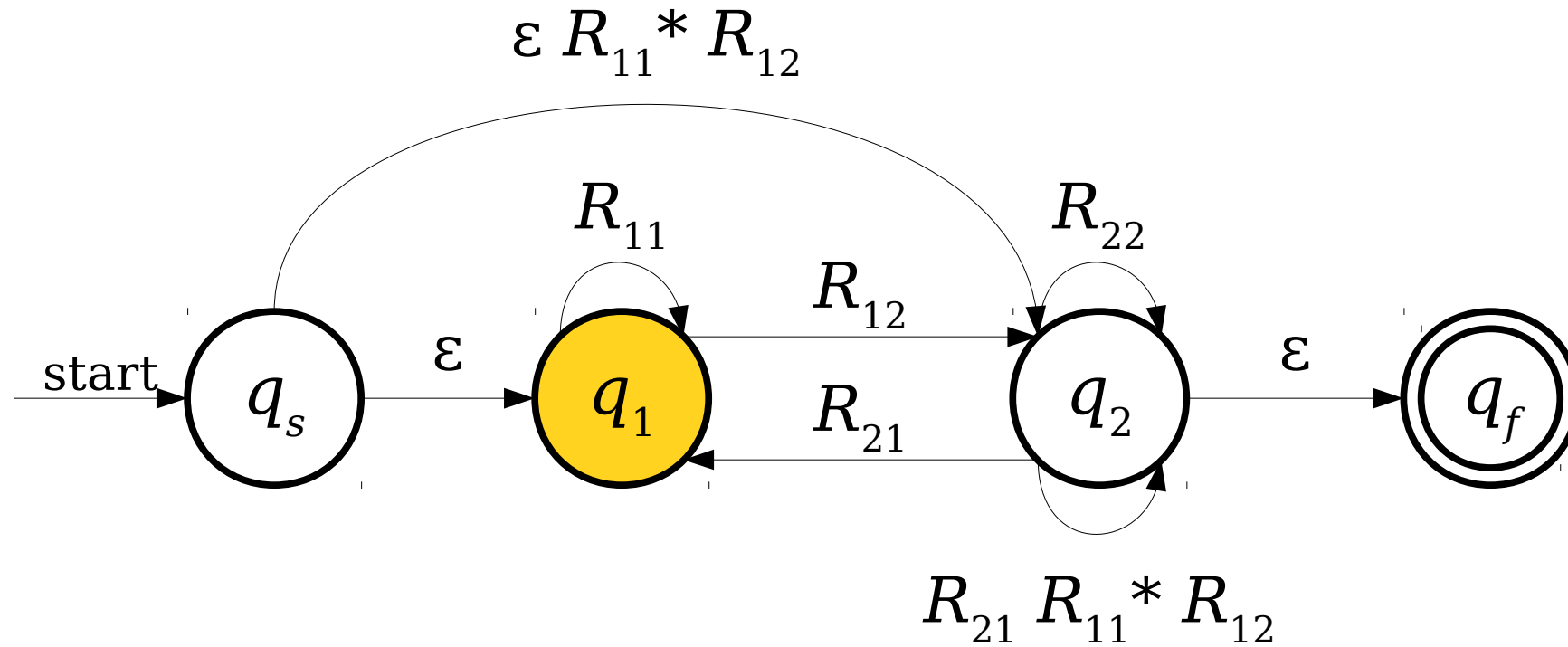
Here is a pattern that we might process  
when going from  $q_2$  to  $q_2$ :  $R_{21} R_{11} R_{12}$

# From NFAs to Regular Expressions

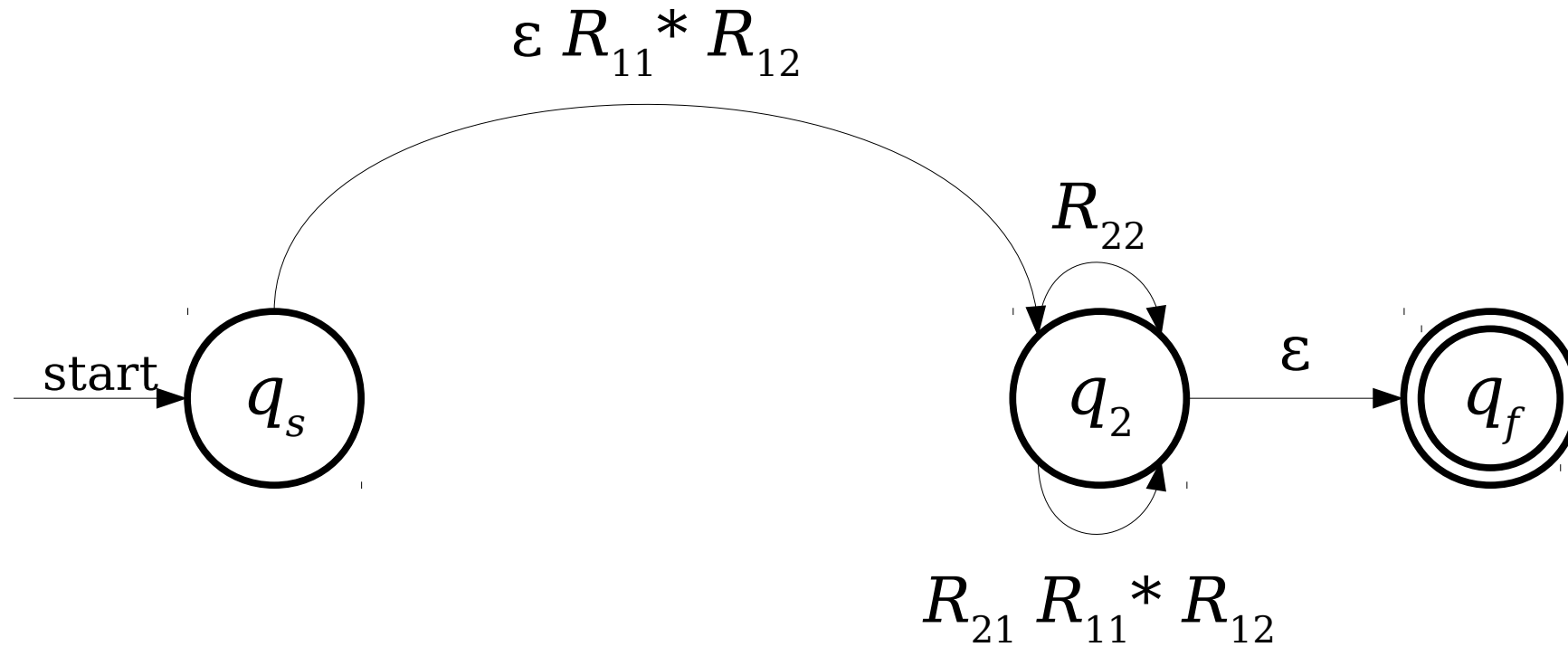


Here is a pattern that we might process  
when going from  $q_2$  to  $q_2$ :  $R_{21} R_{11} R_{12}$

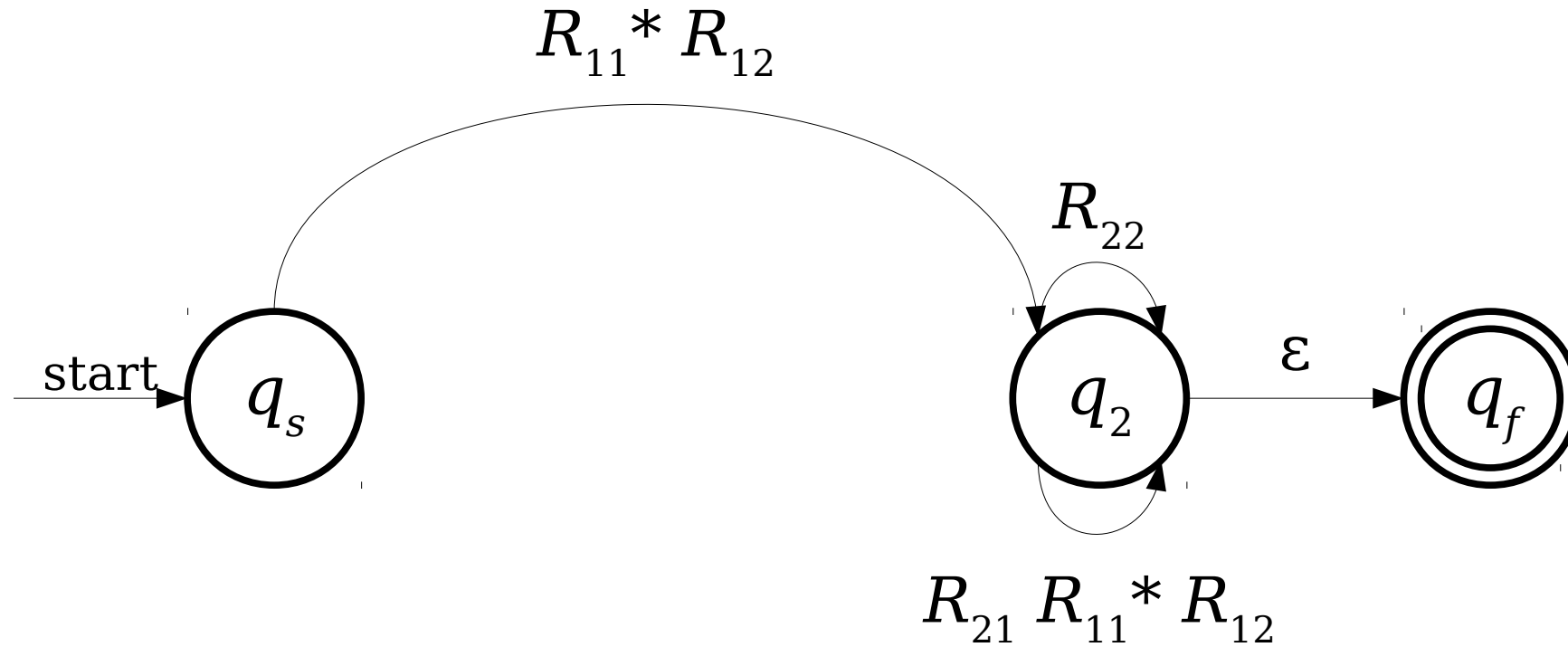
# From NFAs to Regular Expressions



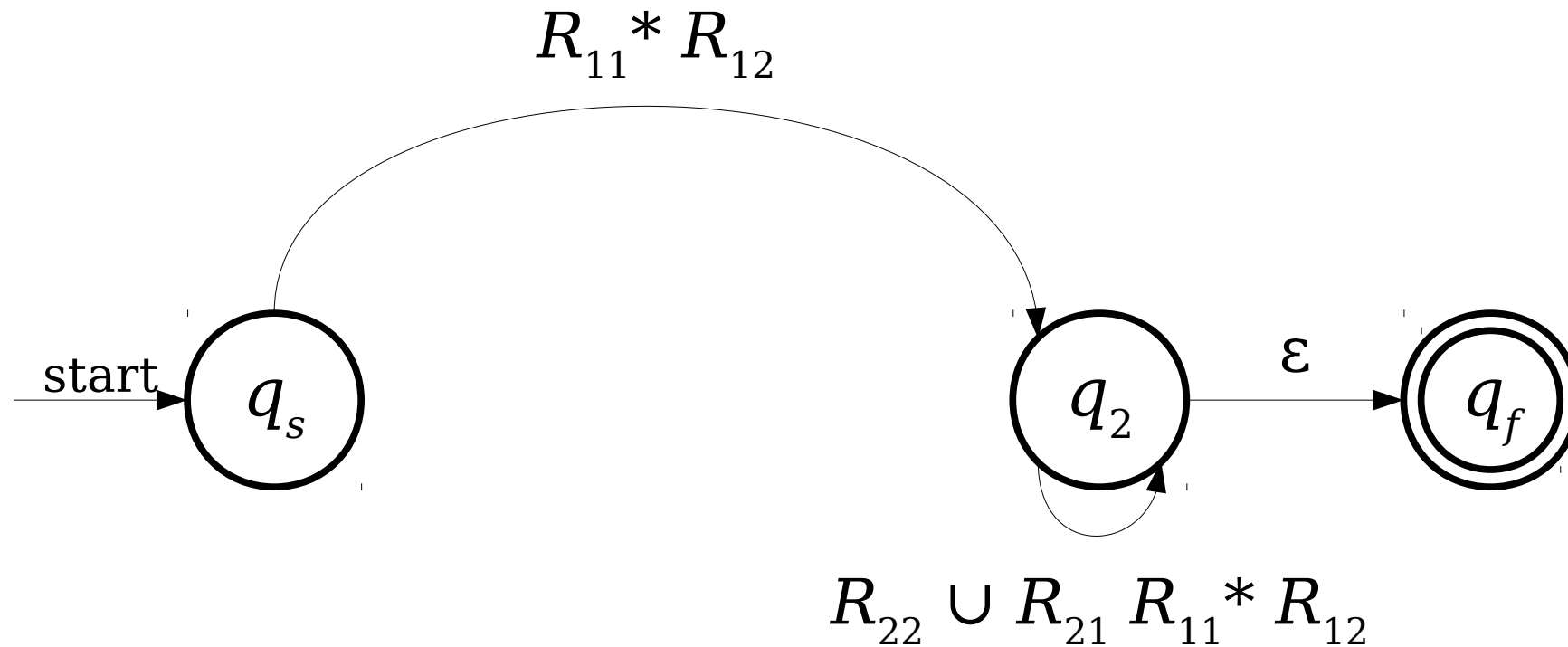
# From NFAs to Regular Expressions



# From NFAs to Regular Expressions



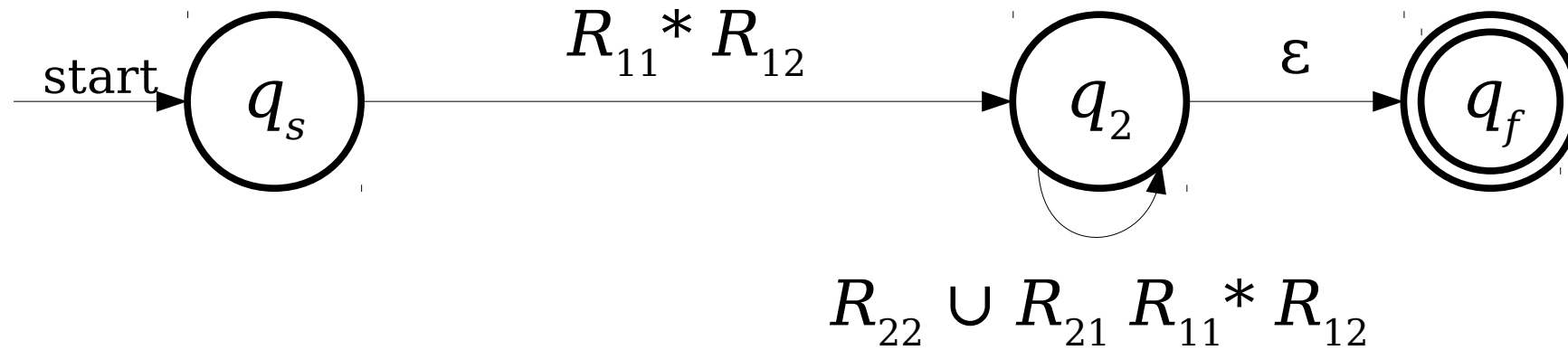
# From NFAs to Regular Expressions



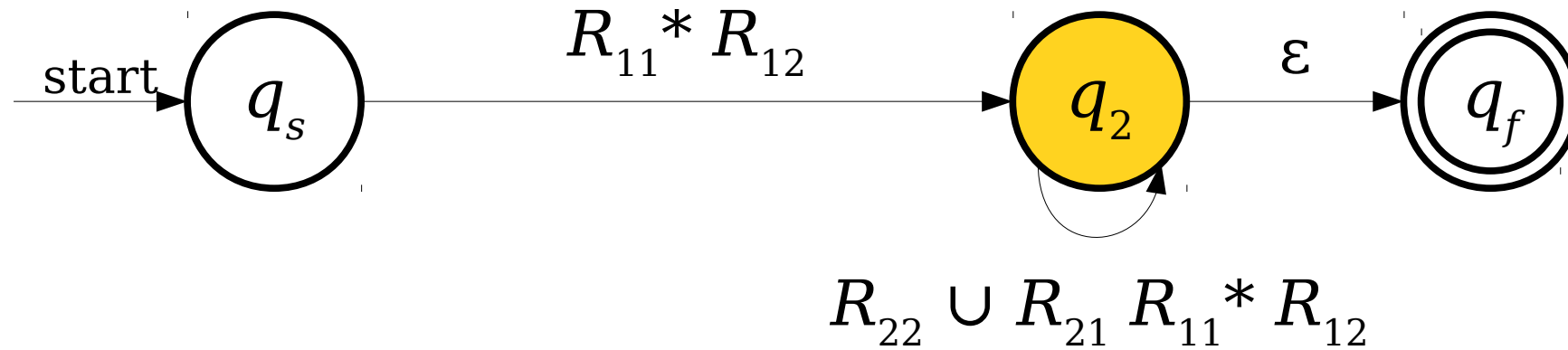
Note: We're using **union** to combine these transitions together.



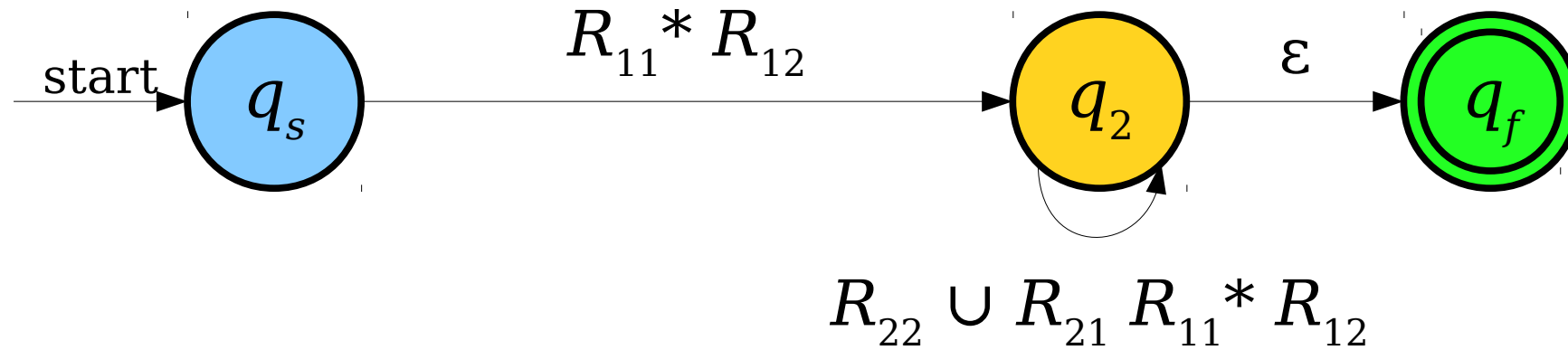
# From NFAs to Regular Expressions



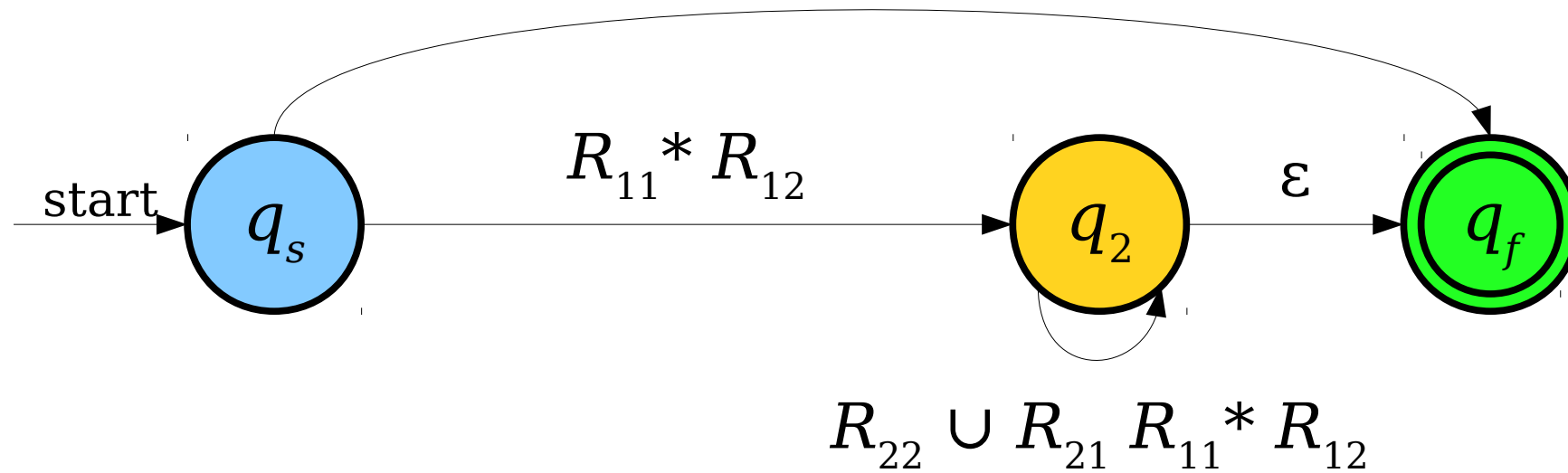
# From NFAs to Regular Expressions



# From NFAs to Regular Expressions

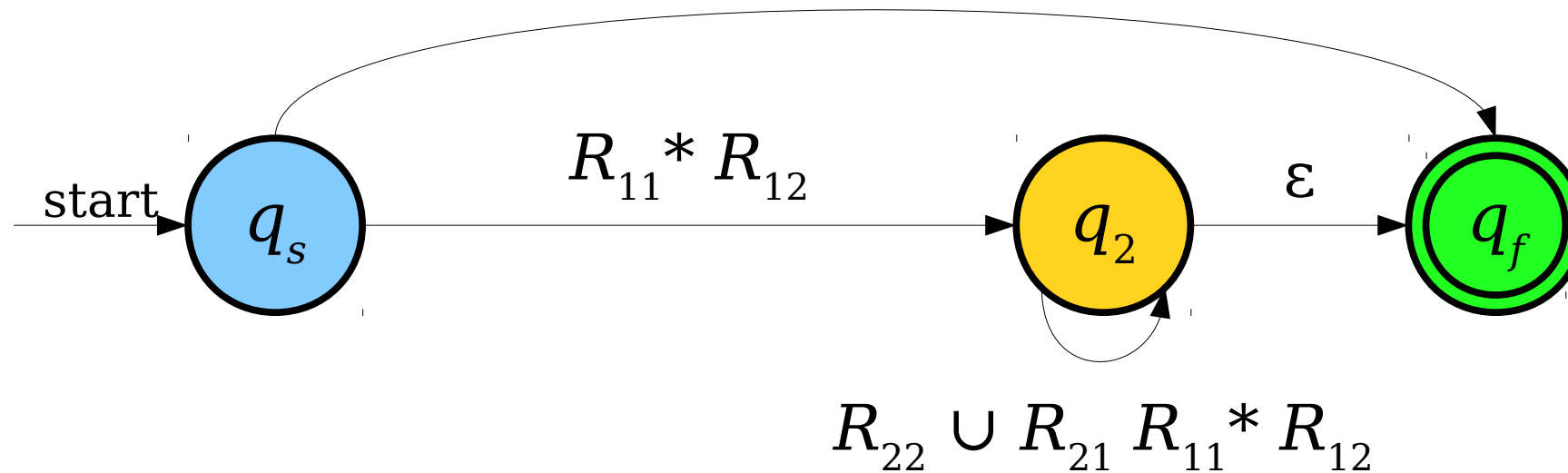


# From NFAs to Regular Expressions

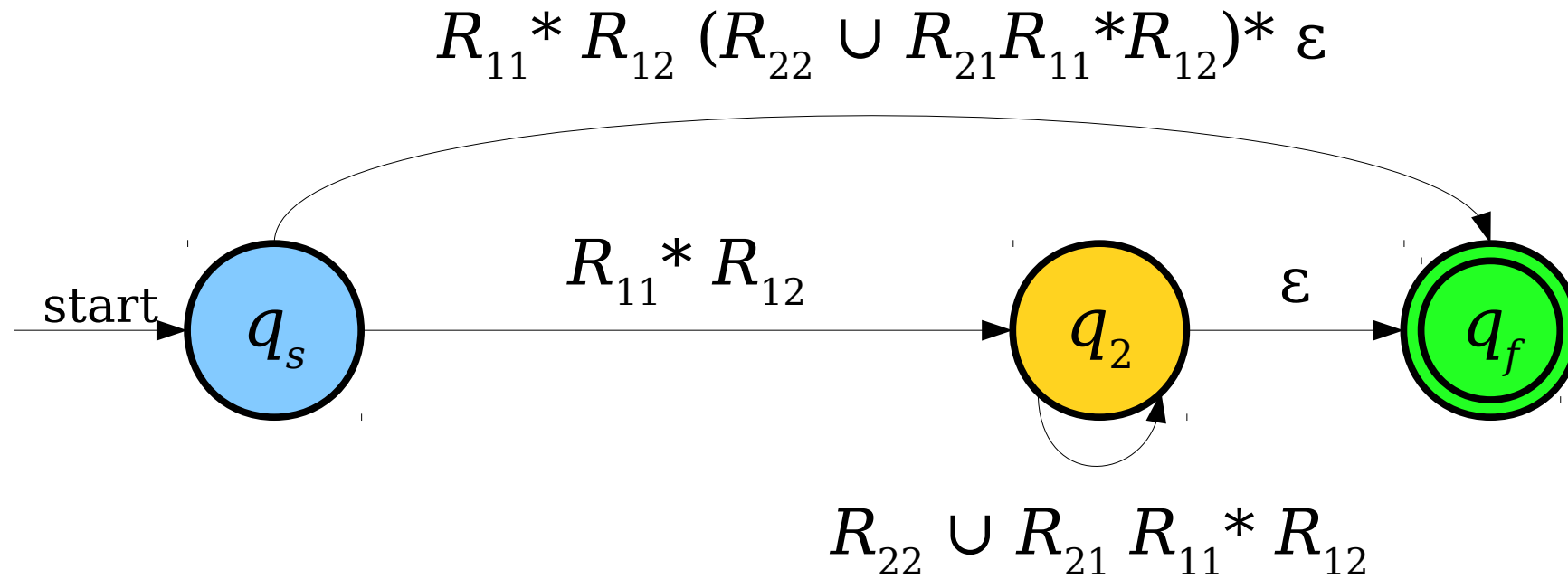


# From NFAs to Regular Expressions

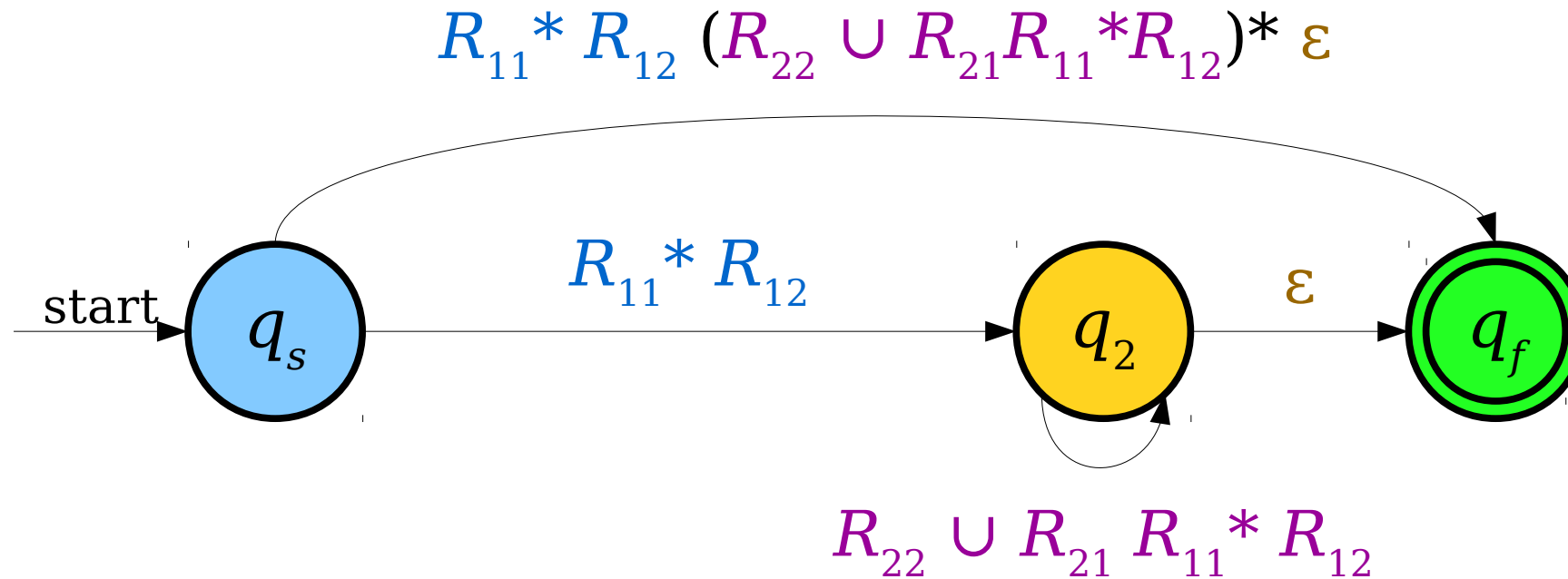
**Quick check:** what goes on this transition?



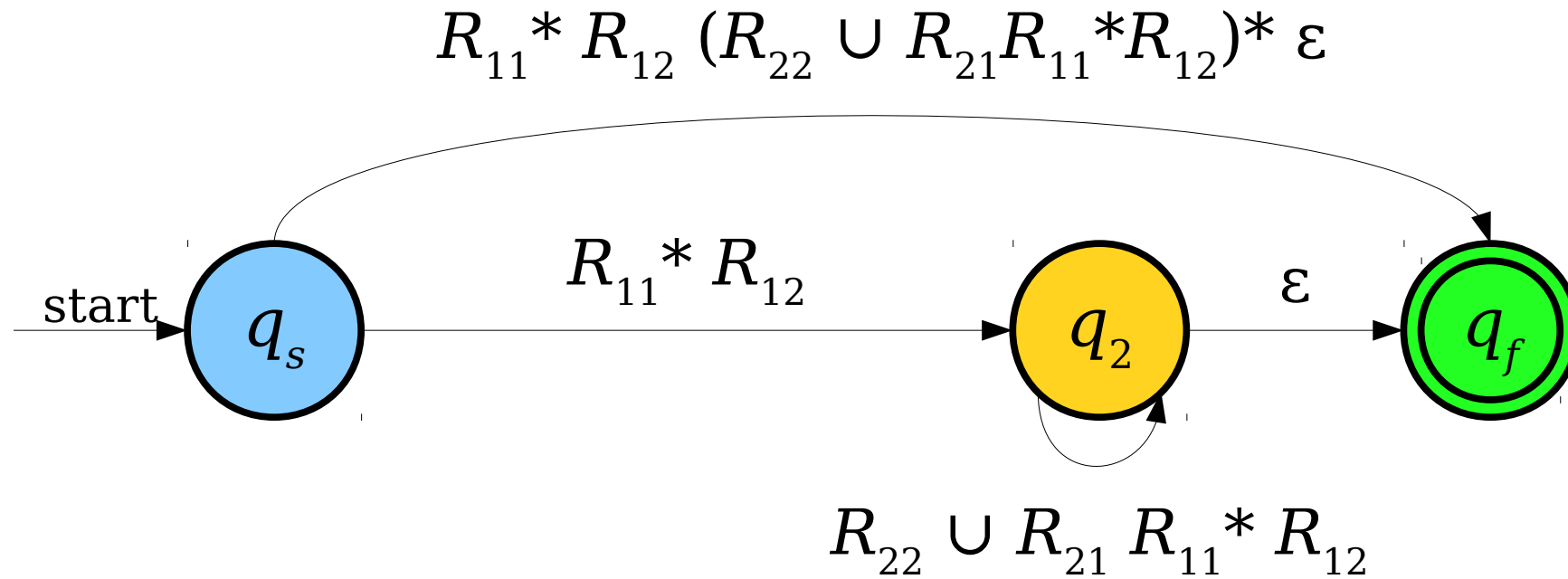
# From NFAs to Regular Expressions



# From NFAs to Regular Expressions

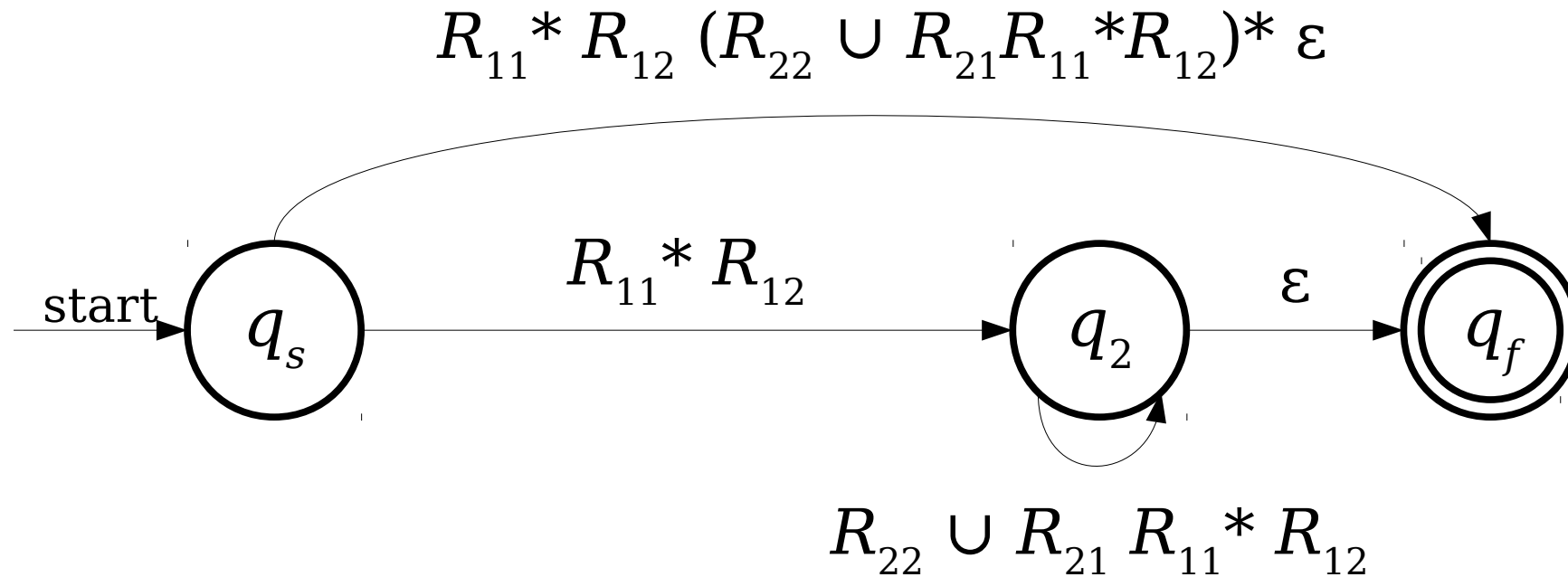


# From NFAs to Regular Expressions

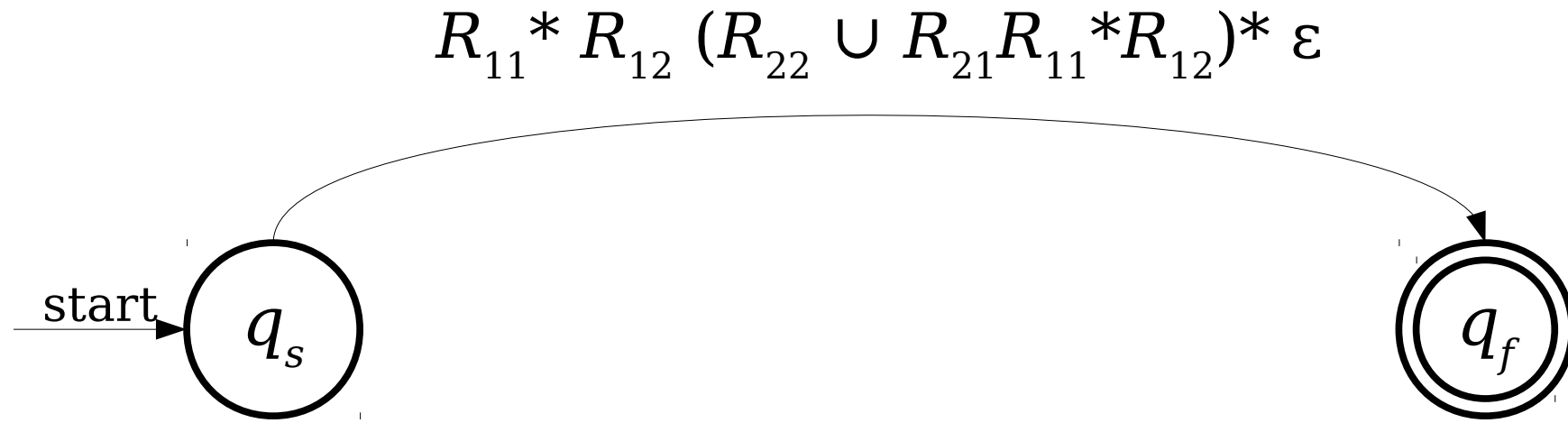




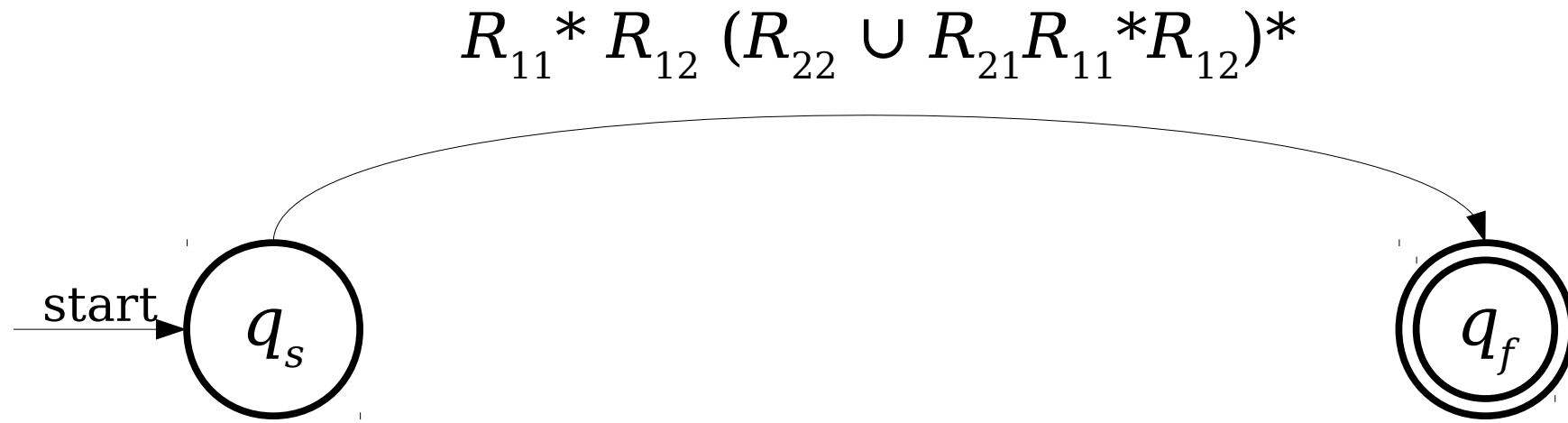
# From NFAs to Regular Expressions



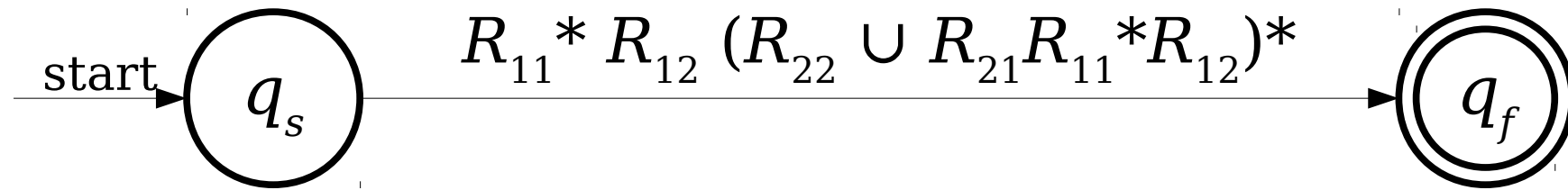
# From NFAs to Regular Expressions



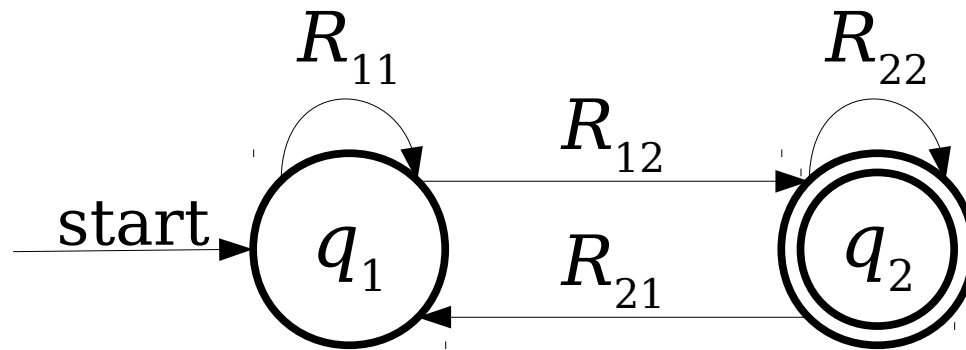
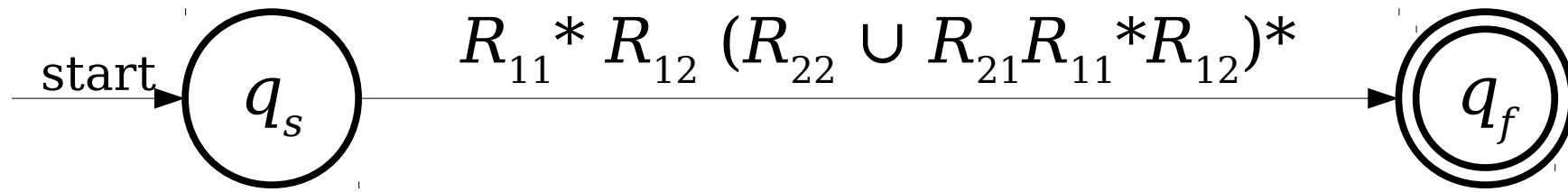
# From NFAs to Regular Expressions



# From NFAs to Regular Expressions



# From NFAs to Regular Expressions



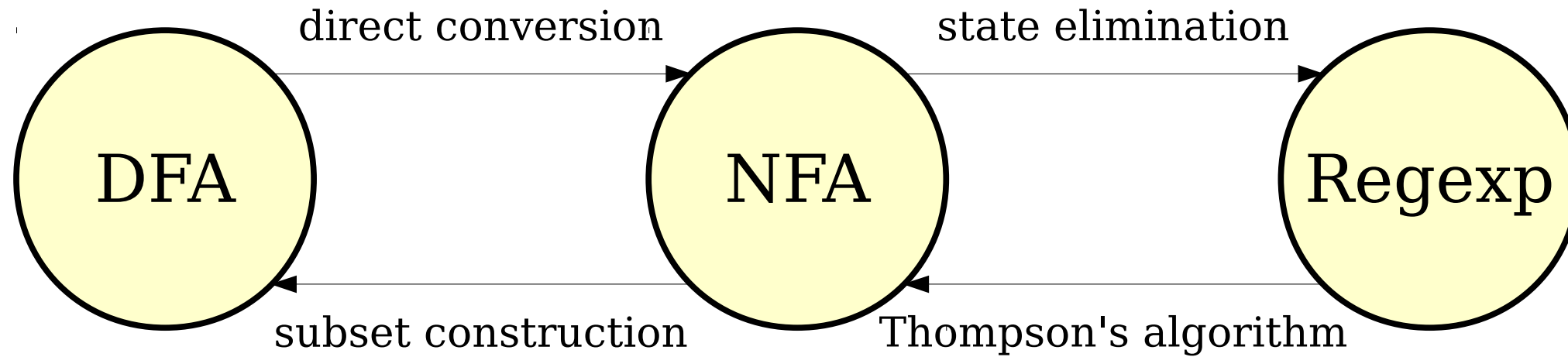
# The State-Elimination Algorithm

- 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  $R_1 \cup R_2 \cup \dots \cup R_k$ .

# 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  $\mathcal{L}(R) = L$ .

# Why This Matters

- The equivalence of regular expressions and finite automata has practical relevance.
  - Regular expression matchers have all the power available to them of DFAs and NFAs.
- This also is hugely theoretically significant: the regular languages can be assembled “from scratch” using a small number of operations!

# Next Time

- ***Applications of Regular Languages***
  - Answering “so what?”
- ***Intuiting Regular Languages***
  - What makes a language regular?
- ***The Myhill-Nerode Theorem***
  - The limits of regular languages.