# 6.1400

## Lecture 3:
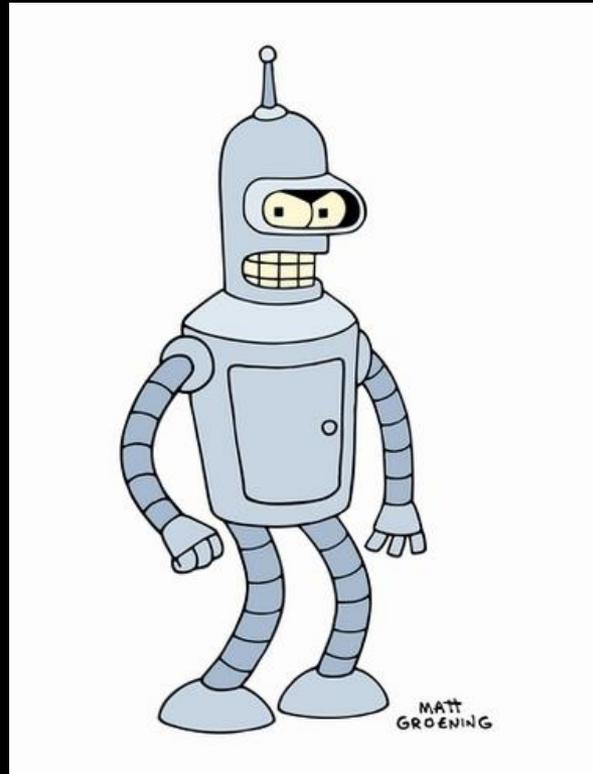## Nondeterminism
## and Regular Expressions

**These slides are on the webpage!**

# 6.1400

## Announcements:
- Pset 0 is out, due tomorrow 11:59pm
- Latex source of hw on Piazza
- Pset 1 coming out Thursday
- No class next Tuesday *(...because next week Monday classes will be on Tuesday)*

# **D**eterministic **F**inite **A**utomata

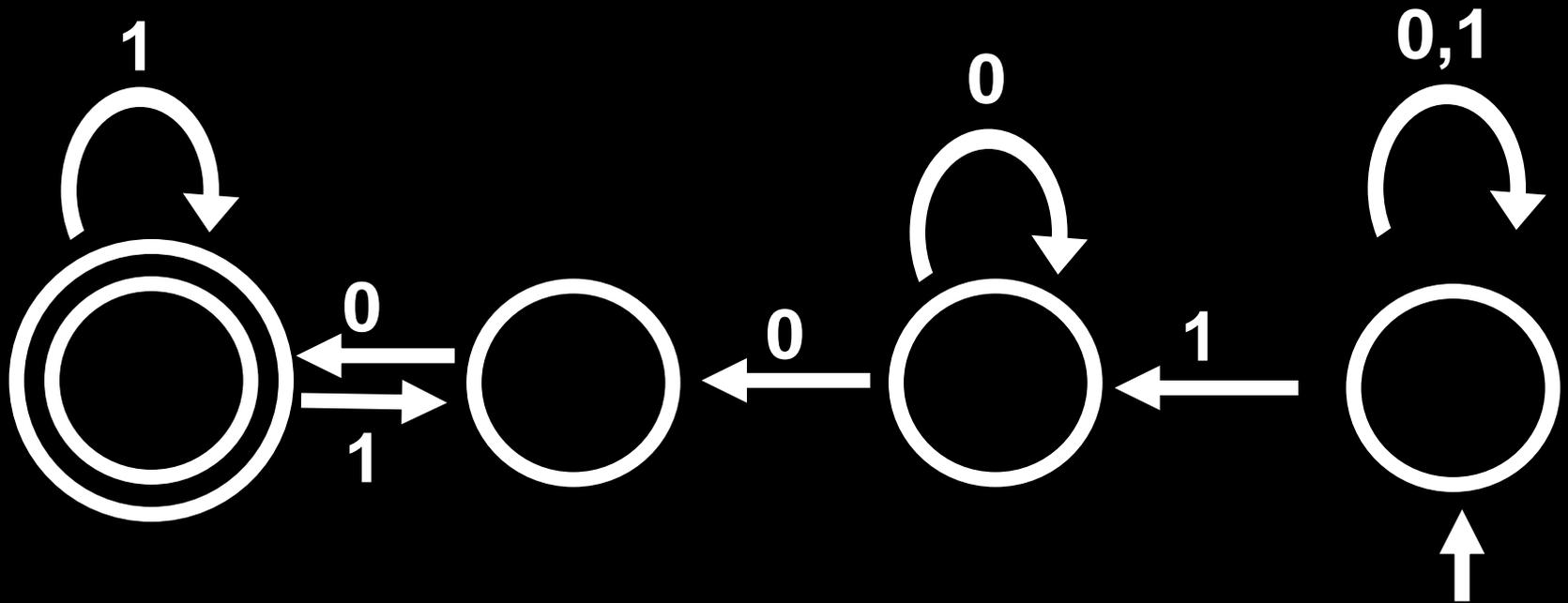

## **Computation with finite memory**

**Regular Languages = Sets of strings accepted by DFAs**

# Non-Deterministic Finite Automata



**Computation with finite memory**

*and magical guessing*

# Non-deterministic Finite Automata (NFA)

This NFA recognizes: {w | w contains 100}

An NFA **accepts string x**
if *there is some path reading in x* that
reaches *some accept state* from *some start state*

# Every NFA can be perfectly simulated by some DFA! 😲

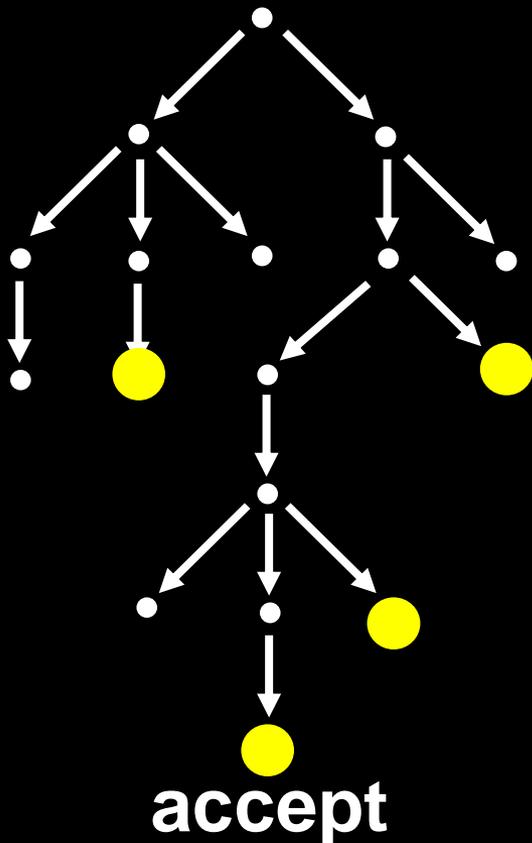**Theorem:** For every NFA N, there is a DFA M such that L(M) = L(N)

**Corollary:** A language A is regular if and only if A is recognized by an NFA

**Corollary:** A is regular iff $A^R$ is regular

left-to-right DFAs $\equiv$ right-to-left DFAs

# From NFAs to DFAs

**Input: NFA $N$ = $(Q, \Sigma, \delta, Q_0, F)$**

**Output: DFA $M$ = $(Q', \Sigma, \delta', q_0', F')$**

**To learn if NFA N accepts, we could do the computation of N *in parallel*, maintaining the set of *all* possible states that can be reached**

**accept**

**Idea:**

**Set $Q' = 2^Q$**

# From NFAs to DFAs: Subset Construction

Input: NFA $N = (Q, \Sigma, \delta, Q_0, F)$

Output: DFA $M = (Q', \Sigma, \delta', q_0', F')$

$$Q' = 2^Q$$

$$\delta' : Q' \times \Sigma \to Q'$$

For $S \in Q'$, $\sigma \in \Sigma$: $\quad \delta'(S,\sigma) = \bigcup_{q \in S} \varepsilon( \delta(q,\sigma) )$ *

$$q_0' = \varepsilon(Q_0)$$

$$F' = \{ S \in Q' \mid f \in S \text{ for } \textit{some } f \in F \}$$

* For $S \subseteq Q$, the **ε-closure of S is**
$\varepsilon(S) = \{r \in Q \text{ reachable from some } q \in S$
by taking zero or more **ε**-transitions$\}$

# Reverse Theorem for Regular Languages

**The reverse of a regular language
is also a regular language**

***If*** **a language can be recognized by a DFA that
reads strings from *right* to *left*,
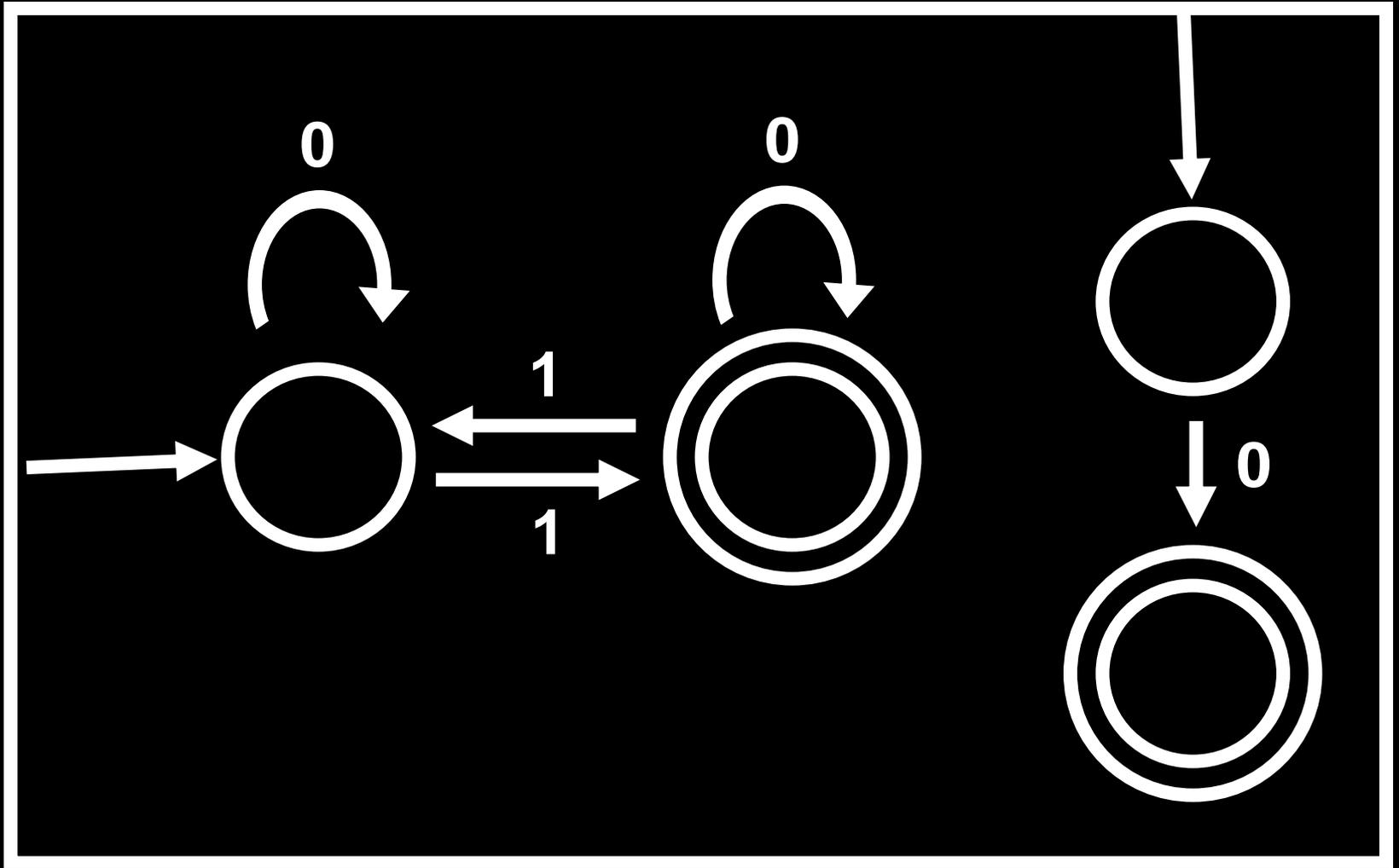*then* there is an "normal" DFA that accepts the
same language**

## Proof Sketch?

*Given a DFA for a language L, "reverse" its arrows,
and flip its start and accept states, getting an NFA.
Convert that NFA back to a DFA!*

# Using NFAs instead of DFAs can make proofs about regular languages *much* easier!

Remember this on homework/exams!

# Union Theorem using NFAs?

# Some Operations on Languages

➡ **Union:** $A \cup B = \{ w \mid w \in A \text{ or } w \in B \}$

➡ **Intersection:** $A \cap B = \{ w \mid w \in A \text{ and } w \in B \}$

➡ **Complement:** $\neg A = \{ w \in \Sigma^* \mid w \notin A \}$

➡ **Reverse:** $A^R = \{ w_1 \ldots w_k \mid w_k \ldots w_1 \in A, w_i \in \Sigma \}$

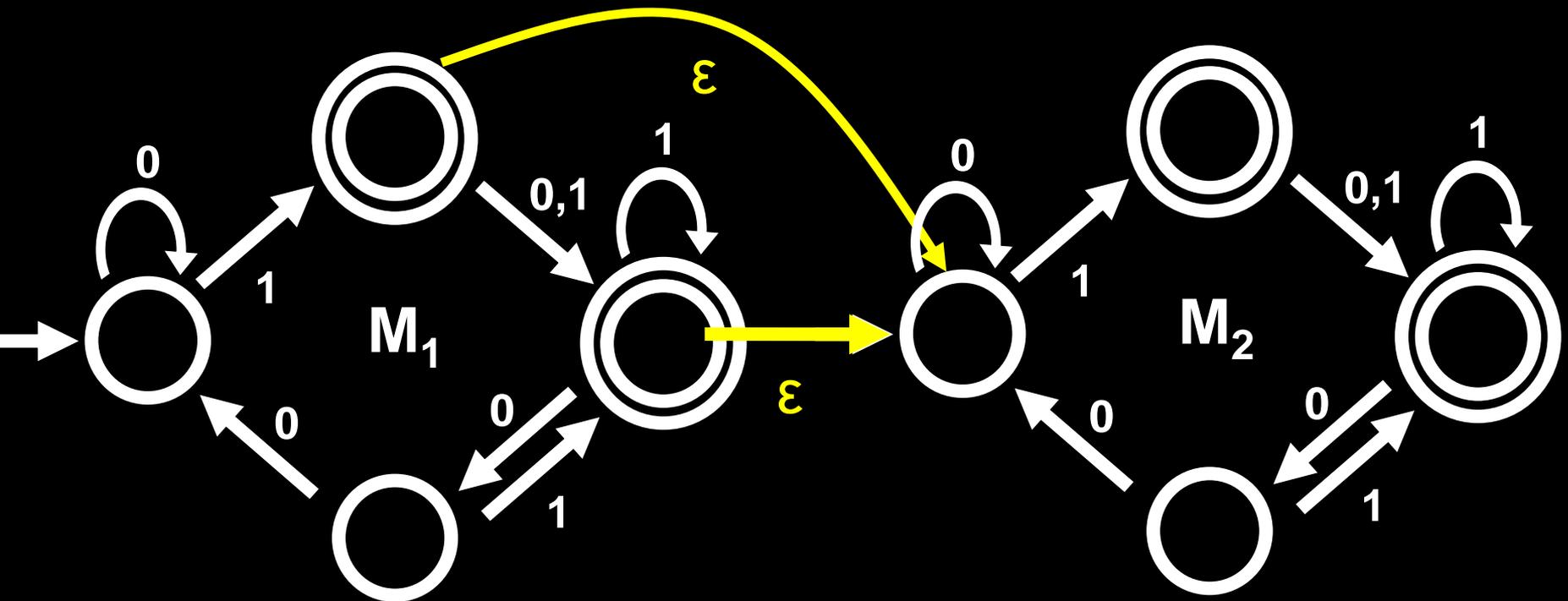**Concatenation:** $A \cdot B = \{ vw \mid v \in A \text{ and } w \in B \}$

**Star:** $A^* = \{ s_1 \ldots s_k \mid k \geq 0 \text{ and each } s_i \in A \}$

$A^* =$ set of all strings over alphabet A

# Regular Languages are closed under concatenation

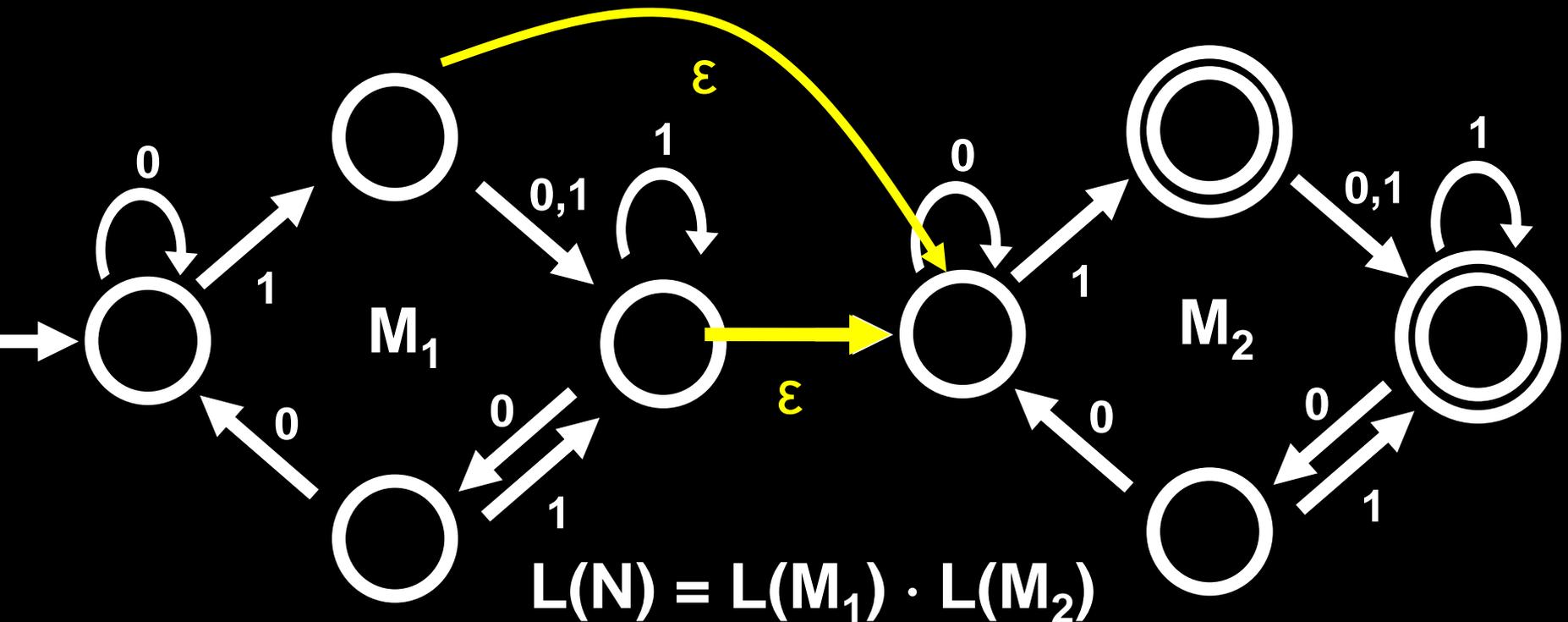**Concatenation:** $A \cdot B = \{ vw \mid v \in A \text{ and } w \in B \}$

**Given DFAs $M_1$ for A and $M_2$ for B, connect**

**the accept states of $M_1$ to the start state of $M_2$**
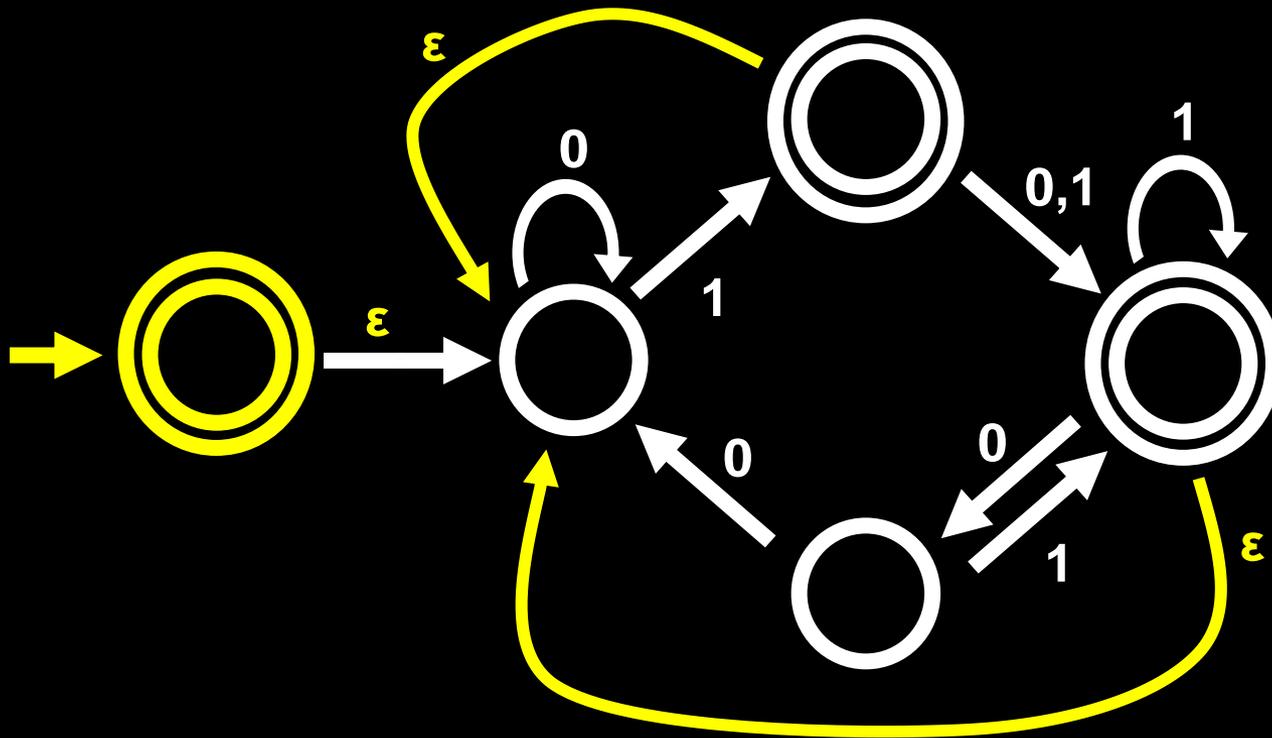


$$L(N) = L(M_1) \cdot L(M_2)$$

# Regular Languages are closed under star

$$A^* = \{ s_1 \dots s_k \mid k \geq 0 \text{ and each } s_i \in A \}$$

## Let M be a DFA

## We construct an NFA N that recognizes L(M)*

**Formally, the construction is:**

Input: DFA $M = (Q, \Sigma, \delta, q_1, F)$

Output: NFA $N = (Q', \Sigma, \delta', \{q_0\}, F')$

$Q' = Q \cup \{q_0\}$

$F' = F \cup \{q_0\}$

$$\delta'(q,a) = \begin{cases} \{\delta(q,a)\} & \text{if } q \in Q \text{ and } a \neq \varepsilon \\ \{q_1\} & \text{if } q \in F \text{ and } a = \varepsilon \\ \{q_1\} & \text{if } q = q_0 \text{ and } a = \varepsilon \\ \varnothing & \text{if } q = q_0 \text{ and } a \neq \varepsilon \\ \varnothing & \text{else} \end{cases}$$

# Regular Languages are closed under star

How would we *prove* that the NFA construction works? 🤔

Want to show:   L(N) = L(M)*

1. L(N) ⊇ L(M)*

2.  L(N) ⊆ L(M)*

# 1. L(N) ⊇ L(M)*

Let $w = w_1 \cdots w_k$ be in **L(M)*** where $w_1, \ldots, w_k \in$ **L(M)**

We show: **N** accepts **w** by induction on k

**Base Cases:**

✓ **k = 0**     **(w = ε)**
✓ **k = 1**     **(w ∈ L(M) and L(M) ⊆ L(N))**

**Inductive Step:** Let k ≥ 1 be an integer

**I.H. N** accepts all strings $v = v_1 \cdots v_k \in$ **L(M)***, $v_i \in$ **L(M)**

Let $u = u_1 \cdots u_k u_{k+1} \in$ **L(M)*** , $u_i \in$ **L(M)**

**N** accepts $u_1 \cdots u_k$ (by I.H.) and **M** accepts $u_{k+1}$

imply that **N** also accepts **u**

since N has $\varepsilon$-transitions from final states to start state of M

# 2. L(N) ⊆ L(M)*

**Let w be accepted by N; we want to show w ∈ L(M)***

**If w = ε, then w ∈ L(M)*** ✓

**I.H. If N accepts u and takes k ε-transitions, then u ∈ L(M)***

**Let w be accepted by N with k+1 ε-transitions.**

**Write w as w=uv, where v is the substring read after the *last* ε-transition**

ε

u

**final**

ε

**start**

v

**accept**

**N accepts u, so**

**u ∈ L(M)***

**By I.H.**

**w = uv ∈ L(M)***

**v ∈ L(M)**

# Regular Languages are closed under all of the following operations:

**Union:** $A \cup B = \{ w \mid w \in A \text{ or } w \in B \}$

**Intersection:** $A \cap B = \{ w \mid w \in A \text{ and } w \in B \}$

**Complement:** $\neg A = \{ w \in \Sigma^* \mid w \notin A \}$

**Reverse:** $A^R = \{ w_1 \ldots w_k \mid w_k \ldots w_1 \in A, w_i \in \Sigma \}$

**Concatenation:** $A \cdot B = \{ vw \mid v \in A \text{ and } w \in B \}$

**Star:** $A^* = \{ s_1 \ldots s_k \mid k \geq 0 \text{ and each } s_i \in A \}$

# Regular Expressions: Computation as Description

**A different way of thinking about computation:**
*What is the complexity of describing
the strings in the language?*

*DFAs find "patterns" in strings;
regular expressions give us a way to describe them precisely*

# During the "nerve net" hype in the 1950s…

U. S. AIR FORCE
## PROJECT RAND

## RESEARCH MEMORANDUM

REPRESENTATION OF EVENTS IN NERVE NETS AND
FINITE AUTOMATA

S. C. Kleene

RM-704

15 December 1951

# Inductive Definition of Regexp

**Let Σ be an alphabet. We define the regular expressions over Σ inductively:**

For all $\sigma \in \Sigma$, $\sigma$ is a regexp

ε is a regexp

$\varnothing$ is a regexp

If $R_1$ and $R_2$ are both regexps, then

$(R_1 R_2)$, $(R_1 + R_2)$, and $(R_1)^*$ are regexps

Examples: ε, 0, (1)*, (0+1)*, ((((0)*1)*1) + (10))

**Precedence Order:** $*$

**then** $\cdot$

**then** $+$

**Example:** $R_1 * R_2 + R_3 = ((R_1*) \cdot R_2) + R_3$

# Definition: Regexps Represent Languages

The regexp $\sigma \in \Sigma$ *represents* the language $\{\sigma\}$

The regexp $\varepsilon$ represents $\{\varepsilon\}$

The regexp $\varnothing$ represents $\varnothing$

If $R_1$ and $R_2$ are regular expressions representing $L_1$ and $L_2$ then:

$(R_1 R_2)$ represents $L_1 \cdot L_2$

$(R_1 + R_2)$ represents $L_1 \cup L_2$

$(R_1)^*$ represents $L_1^*$

Example: $(10 + 0^*1)$ represents $\{10\} \cup \{0^k 1 \mid k \geq 0\}$

# Regexps Represent Languages

**For every regexp R,
define L(R) to be the language that R represents**

**A string w ∈ Σ\* is *accepted by R*
(or, *w matches R)* if w ∈ L(R)**

**Examples: 0, 010, and 01010 match (01)\*0
1101011101011100 matches (0+1)\*0**

**L((0+1)\*0) = {w ∈ {0,1}\* | w ends in a 0}**

# Give a regular expression for …

**{ w | w has exactly a single 1 }**

**0*10***

**{ w | w contains 001 }**

**(0+1)*001(0+1)***

$\Sigma = \{0,1\}$

**What language does
the regexp $\varnothing^*$ represent?**

**{ε}**

# Give a regular expression for …

**{ w | w has length ≥ 3 and its 3rd symbol is 0 }**

**(0+1)(0+1)0(0+1)***

$$\Sigma = \{0,1\}$$

# Give a regular expression for …

{ w | w = ε  or every odd position in w is a 1 }

$$(1(0 + 1))*(1 + ε)$$

## How expressive are regular expressions?

# DFAs ≡ NFAs ≡ Regular Expressions!

**L can be represented by some regexp**

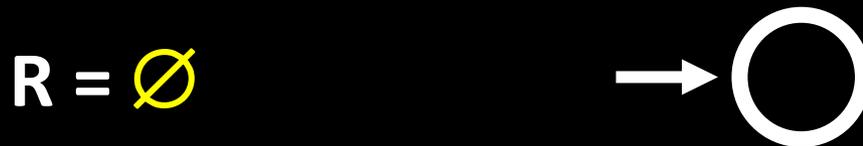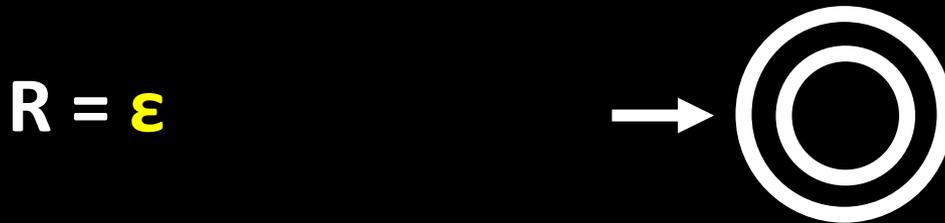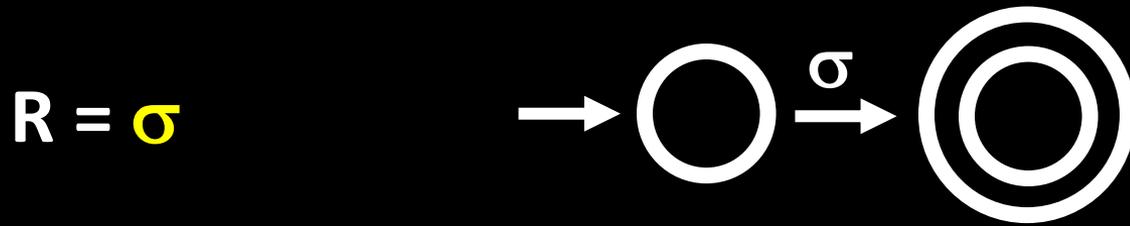$$\Leftrightarrow \quad \text{L is regular}$$

# L can be represented by some regexp
## $\Rightarrow$ L is regular

**Given any regexp R, we will construct an NFA N such that R represents L(N)**

**Proof by induction on the *length* of the regexp R:**

**Base Cases (R has length 1):**

**Induction Step:**   Suppose every regexp of length $< k$
represents some regular language.

**Consider a regexp R of length $k > 1$**

**Three possibilities for R:**

$R = R_1 + R_2$

$R = R_1 R_2$

$R = (R_1)^*$

**Induction Step:** Suppose every regexp of length **< k** represents some regular language.

**Consider a regexp R of length k > 1**

**Three possibilities for R:**

$R = R_1 + R_2$

$R = R_1 R_2$

$R = (R_1)^*$

**Induction Step:** Suppose every regexp of length $< k$ represents some regular language.

Consider a regexp R of length $k > 1$

Three possibilities for R:

$R = R_1 + R_2$      By induction, $R_1$ and $R_2$ represent some regular languages, $L_1$ and $L_2$

$R = R_1 R_2$           But $L(R) = L(R_1 + R_2) = L_1 \cup L_2$

$R = (R_1)^*$      so $L(R)$ is regular, by the union theorem!

**Induction Step:** Suppose every regexp of length **< k** represents some regular language.

Consider a regexp R of length **k > 1**

Three possibilities for R:

$R = R_1 + R_2$     By induction, $R_1$ and $R_2$ represent some regular languages, $L_1$ and $L_2$

$R = R_1 R_2$     But $L(R) = L(R_1 \cdot R_2) = L_1 \cdot L_2$

$R = (R_1)^*$     Thus L(R) is regular because regular languages are closed under concatenation

**Induction Step:** Suppose every regexp of length $< k$ represents some regular language.

Consider a regexp R of length $k > 1$

Three possibilities for R:

$R = R_1 + R_2$

$R = R_1 R_2$

$R = (R_1)^*$

By induction, $R_1$ represents a regular language $L_1$

But $L(R) = L(R_1^*) = L_1^*$

Thus $L(R)$ is regular because regular languages are closed under star

**Induction Step:**   Suppose every regexp of length $< k$
represents some regular language.

**Consider a regexp R of length $k > 1$**

**Three possibilities for R:**

$R = R_1 + R_2$      By induction, $R_1$ represents
a regular language $L_1$

$R = R_1 R_2$

But $L(R) = L(R_1*) = L_1*$
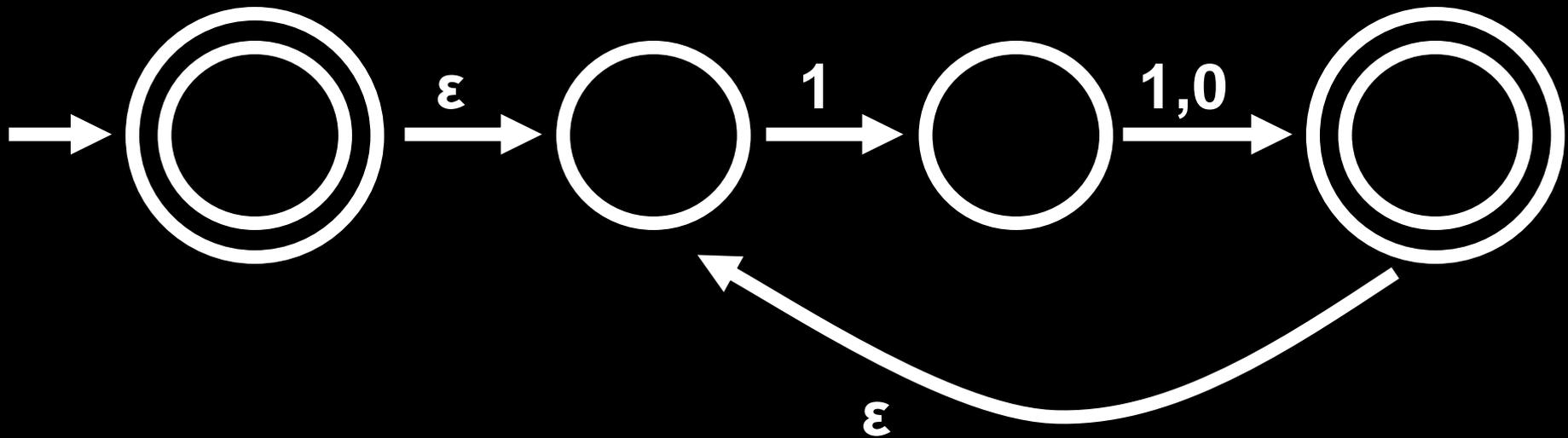
$R = (R_1)*$      Thus L(R) is regular because regular
languages are closed under star

**Therefore:**  If L is represented by a regexp,
then L is regular

# An Algorithm for Converting Regexps to NFAs!

**Give an NFA that accepts the language represented by (1(0 + 1))\***



**Regular expression:  (1(0+1))\***

# Generalized NFAs (GNFA)

**L can be represented by a regexp**

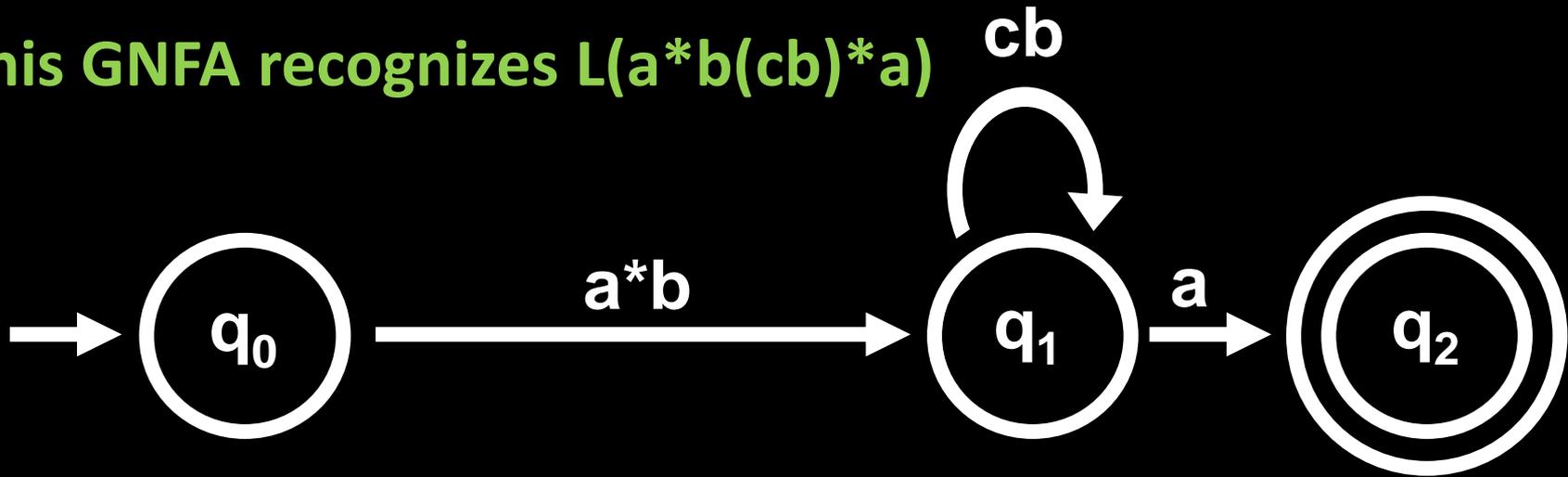$\Longleftarrow$

**L is a regular language**

**Idea:** Transform a DFA for L into a regular expression by *removing states* and re-labeling the arcs with *regular expressions*

Rather than reading in just 0 or 1 letters from the string on an arc, we can read in *entire substrings*

# Generalized NFA (GNFA)

**This GNFA recognizes L(a*b(cb)*a)**



**Accept string $x \Leftrightarrow$ there is *some path* of regexps $R_1, \ldots, R_k$ from start state to final state such that $x$ matches $R_1 \cdots R_k$**

Is aaabcbcba accepted or rejected?   **accepted**

Is bba accepted or rejected?   **rejected**

Is bcba accepted or rejected?   **accepted**

# Generalized NFA (GNFA)

This GNFA recognizes L(a*b(cb)*a)



Accept string $x \Leftrightarrow$ there is *some path* of regexps $R_1, \ldots, R_k$ from start state to final state such that $x$ matches $R_1 \cdots R_k$

Every NFA is also a GNFA.
Every regexp can be converted into
a GNFA with just two states!