# Pushdown automata and context-free languages

**Sides based on material by Kani, Erickson, Chekuri, et. al.**

**All mistakes are my own! - Ivan Abraham (Fall 2024)**

Image by ChatGPT (probably collaborated with DALL-E)

# Introduction

Will this code execute successfully?    **YES**

# Introduction

Will this code execute successfully?     **NO**
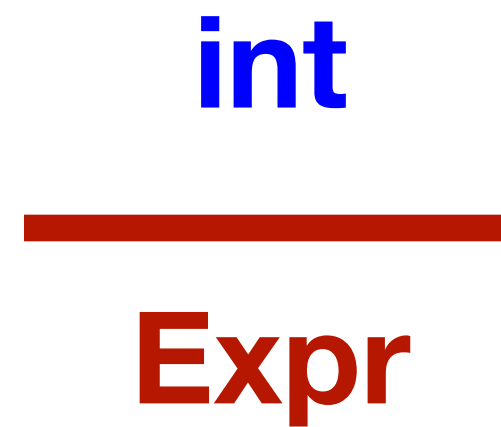
```
main.c                    [ ]  (C)  < Share    Run

1  // Online C compiler to run C program online
2  #include <stdio.h>
3
4▾ int main() {
5      int a = 2;
6      int b =  ;
7  }
```

```
Output                              Clear

/tmp/Zqon05DwrB.c: In function 'main':
ERROR!
/tmp/Zqon05DwrB.c:6:14: error: expected expression
    before ';' token
    6 |      int b =  ;
      |               ^



=== Code Exited With Errors ===
```
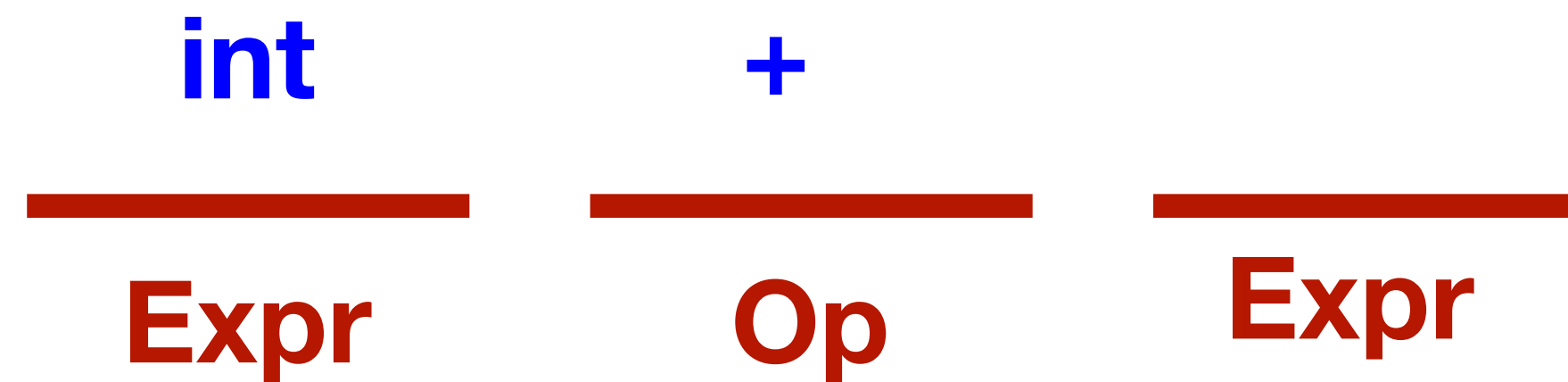
What was the compiler expecting?

# What can an arithmetic expression be?

- **int** - A single number.

<div align="center">

**int**
———
**Expr**

</div>

- **Expr Op Expr** - Two expressions joined by an operator.

<div align="center">

**int**     **+**

**Expr**    **Op**    **Expr**

</div>

- **Recursive Expression**

<div align="center">

**int**     **+**

**Expr**   **Op**   **Expr**   **Op**   **Expr**

</div>

# Arithmetic expressions

- Here is one way to express these rules

- **Expr** → **int**

- **Expr** → **Expr Op Expr** ────→ This is called a ***production rule***. It says "if you see **Expr**, you can replace it with **Expr Op Expr**."

- **Expr** → **( Expr )**

- **Op** → **+** | **-** | **×** | **/** ────→ This one says "if you see **Op**, you can replace it with **+** or **-** or **×** or **/**
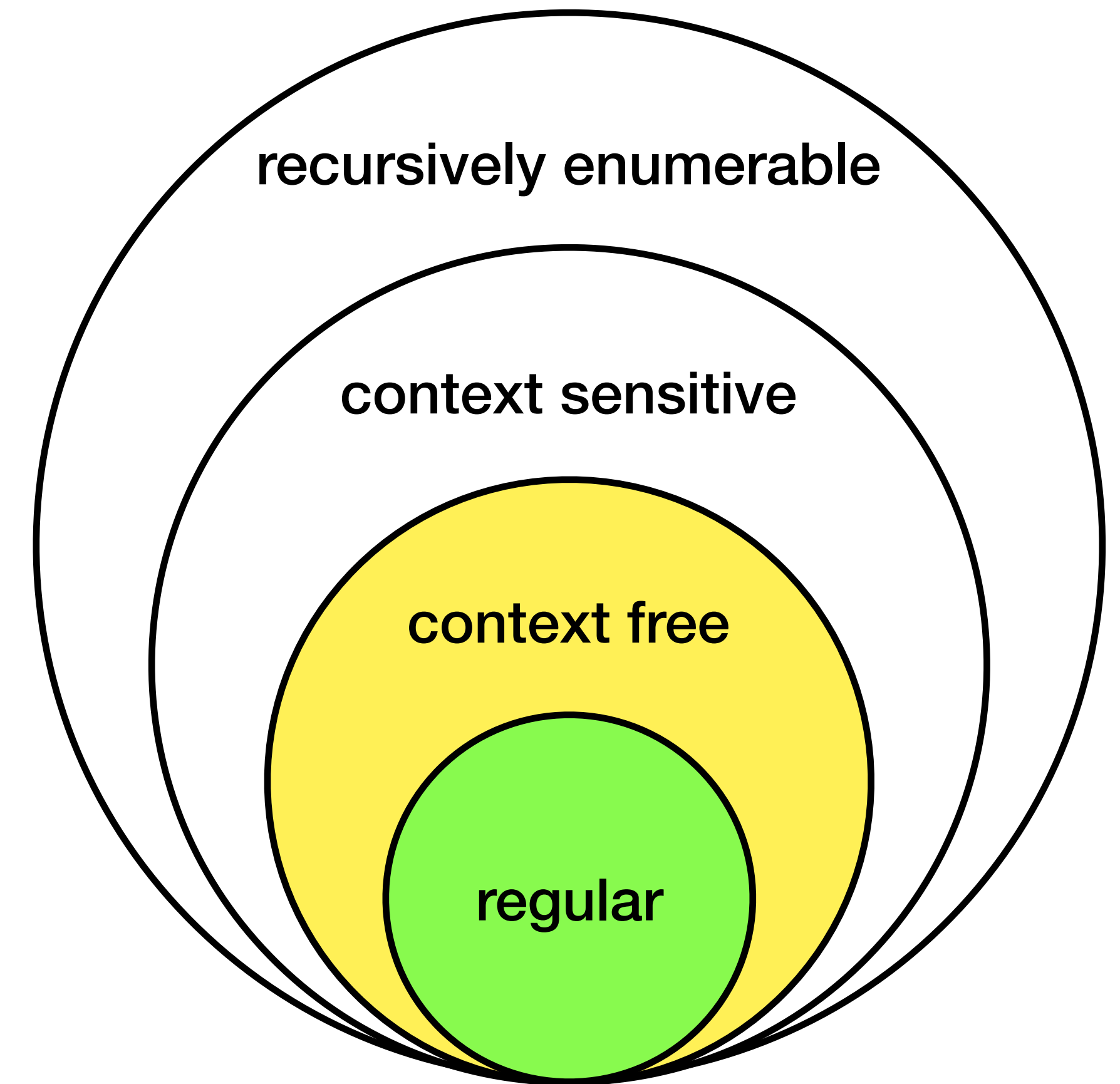
# Grammar - rules for a language

A ***context-free grammar*** (or ***CFG***) is a recursive set of rules that define a language.

**Definition:**

A **CFG** is a quadruple $G = (V, T, P, S)$ where

$G = (Variables, Terminals, Productions, Start\ Var)$

# Context Free Grammar

**Expr → int**

**Expr → Expr Op Expr**

**Expr → ( Expr )**

**Op → + | - | × | /**

**Definition:** A **CFG** is a quadruple $G = (V, T, P, S)$

- $V$ is a finite set of *non-terminal* (*variable*) *symbols*

- $T$ is a finite set of *terminal symbols* (*alphabet*)

- $P$ is a finite set of *productions*, each of the form $A \rightarrow \alpha$ where $A \in V$ and $\alpha$ is a string in $(V \cup T)*$ . Formally, $P \subset V \times (V \cup T)^*$

- $S \in V$ is a *start symbol* .

# Context Free Grammar
## Example

- $V = \{S\}$

- $T = \{0,1\}$

- $P = \{S \rightarrow \varepsilon \mid 0S0 \mid 1S1\}$ *( abbrev. for $S \rightarrow \epsilon$, $S \rightarrow 0S0$, $S \rightarrow 1S1$ )*

- $S = S$

$$S \longrightarrow 0S0 \longrightarrow 01S10 \longrightarrow 011S110 \longrightarrow 011\epsilon110 \longrightarrow 011110$$

$$P = \{S \to \varepsilon \mid 0S0 \mid 1S1\}$$

# *Derives* relation
## Formalism for how strings are derived/generated

**Definition:** Let $G = (V, T, P, S)$ be a CFG. For strings, $\alpha_1, \alpha_2 \in (V \cup T)*$ we say $\alpha_2$ *derives from* $\alpha_1$, denoted by $\alpha_1 \rightsquigarrow \alpha_2$, if there exist strings $\beta, \gamma, \delta$ in $(V \cup T)*$ such that

- $\alpha_1 = \beta A \delta$

- $\alpha_2 = \beta \gamma \delta$

- $A \to \gamma \in P$

**Examples:** $S \rightsquigarrow \epsilon$, $S \rightsquigarrow 0S0$, $0S1 \rightsquigarrow 01S11$, $0S1 \rightsquigarrow 01$ .

# *Derives* relation
## Formalism for how strings are derived/generated

**Definition:** For integers $k \geq 0$, define $\alpha_1 \overset{k}{\rightsquigarrow} \alpha_2$ inductively as follows:

- $\alpha_1 \overset{0}{\rightsquigarrow} \alpha_2$ if $\alpha_1 = \alpha_2$

- $\alpha_1 \overset{k}{\rightsquigarrow} \alpha_2$ if $\alpha_1 \rightsquigarrow \beta_1$ and $\beta_1 \overset{k-1}{\rightsquigarrow} \alpha_2$

- Alternatively, $\alpha_1 \overset{k}{\rightsquigarrow} \alpha_2$ if $\alpha_1 \overset{k-1}{\rightsquigarrow} \beta_1$ and $\beta_1 \rightsquigarrow \alpha_2$

Finally, we use the notation $\alpha_1 \overset{*}{\rightsquigarrow} \alpha_2$ to mean that $\alpha_2$ can be derived from $\alpha_1$. In other words,

$$\alpha_1 \overset{*}{\rightsquigarrow} \alpha_2 \text{ if } \alpha_1 \overset{k}{\rightsquigarrow} \alpha_2 \text{ for some } k$$

# Context Free Languages

**Definition:** Let $G = (V, T, P, S)$ be a CFG. Then the *language generated* by G, denoted by $L(G)$ is the set

$$L(G) := \left\{ w \in T^* \; \mid \; S \overset{*}{\rightsquigarrow} w \right\}.$$

Thus, a language $L$ is context-free (called a context-free language or CFL) if it is generated by a context-free grammar.

Alternatively, a language $L$ is said to be a CFL, if there exists a CFG $G$ such that $L = L(G)$.
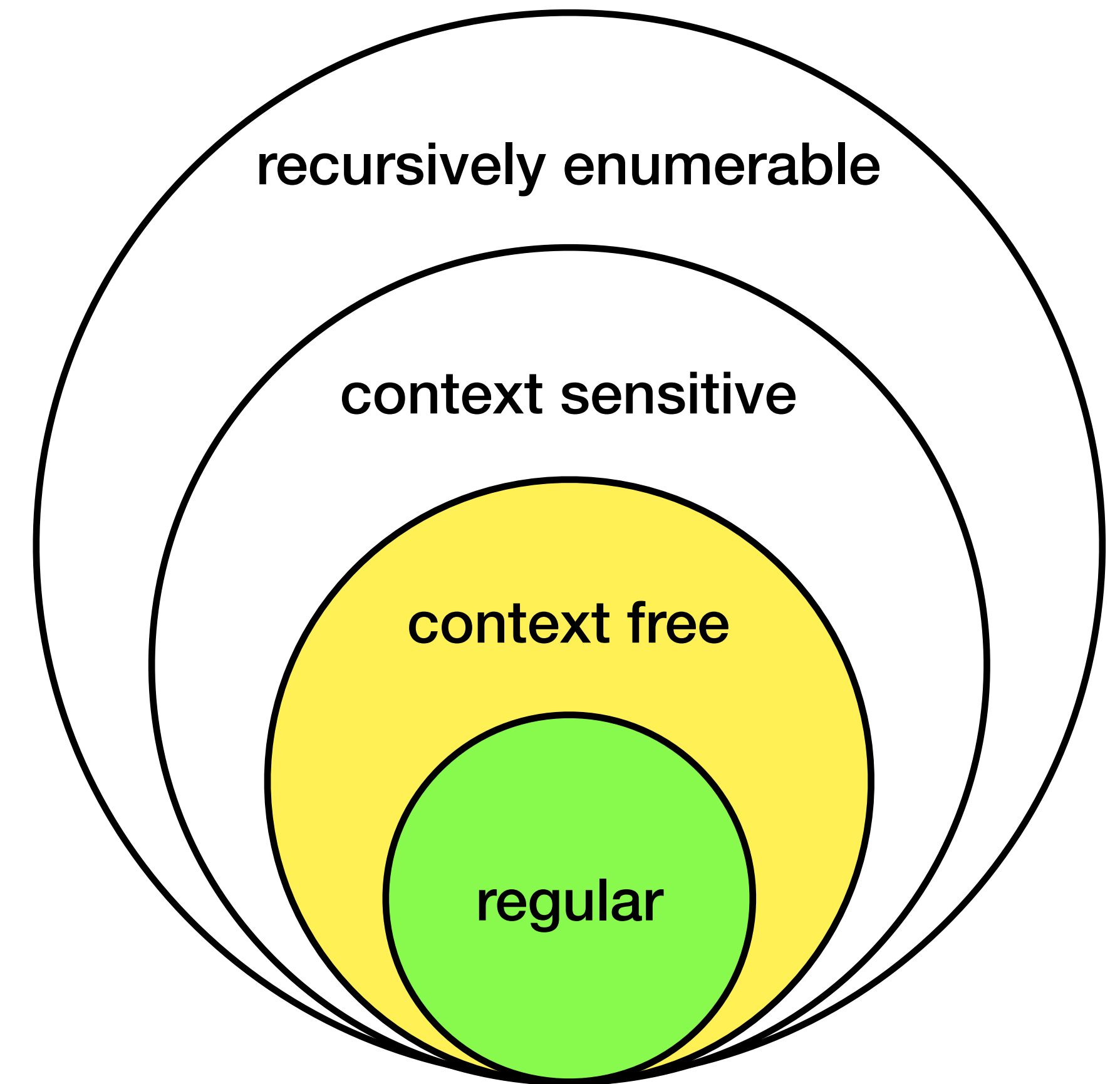
# Context Free Languages
## Production rule examples

- $L = \{0^n 1^n \mid n \geq 0\}$

- $L = \{0^n 1^m \mid m \geq n\}$

- $L = \{0^n 1^m \mid m, n \geq 0\}$

# CFL/CFGs and regular languages
## Recall Chomsky Heirarchy

- The picture depicts regular languages as a proper subset of context-free languages.

- Thus, all regular languages are also CFLs.

  - What was the grammar that generated a regular language?

    - We can start with the DFA recognizing a regular language.

    - Then, extend the algebraic method.



recursively enumerable

context sensitive

context free

regular

# Converting DFAs into CFL



$$G = \left( \{A, B, C, D, E\}, \{a, b\}, \left\{ \begin{array}{r} A \to aA, A \to bA, A \to aB \\ B \to bC \\ C \to aD, \\ D \to bE, \\ E \to aE, E \to bE, E \to \varepsilon \end{array} \right\}, F \right)$$

# Converting regular languages into CFL



$M = (Q, \Sigma, \delta, q_0, F)$: DFA for regular language $L$

$$G = \left( \overbrace{Q}^{\text{Variables}} , \overbrace{\Sigma}^{\text{Terminals}} , \overbrace{\{q \to a\delta(q, a) \mid q \in Q, a \in \Sigma\} \bigcup_{q \in F} \{q \to \varepsilon\}}^{\text{Productions}} , \overbrace{q_0}^{\text{Start var}} \right)$$

15

# Converting regular languages into CFL

$$G = \left( \{A, B, C, D, E\}, \{a, b\} \left\{ \begin{array}{r} A \rightarrow aA, A \rightarrow bA, A \rightarrow aB \\ B \rightarrow bC \\ C \rightarrow aD, \\ D \rightarrow bE, \\ E \rightarrow aE, E \rightarrow bE, E \rightarrow \varepsilon \end{array} \right\}, A \right)$$

In regular languages:

- Terminals can only appear on one side of the production string

- Only one variable allowed in the production result

# Converting regular languages into CFL

$$G = \left( \{A, B, C, D, E\}, \{a, b\} \left\{ \begin{array}{r} A \to aA, A \to bA, A \to aB \\ B \to bC \\ C \to aD, \\ D \to bE, \\ E \to aE, E \to bE, E \to \varepsilon \end{array} \right\}, A \right)$$

In regular languages:

- Terminals can only appear on one side of the production string

- Only one variable allowed in the production result

# Closure Properties of CFL

Let $G_1 = (V_1, T, P_1, S_1)$ and $G_2 = (V_2, T, P_2, S_2)$ be CFGs for $L_1 = L(G_1)$ and $L_2 = L(G_2)$

**Simplifying assumption:** $V_1 \cap V_2 = \varnothing$, that is, non-terminals are not shared

- CFLs are closed under union: $L_1 \cup L_2$ is a CFL.

- CFLs are closed under concatenation: $L_1 \cdot L_2$ is a CFL.

- CFLs are closed under Kleene star: $L_k$ CFL implies $L_k^*$ is a CFL.

# Pushdown automata

# Pushdown automata
## The machine that recognizes CFGs

We established that $\{0^n 1^n | n \geq 0\}$ is a CFL but not a regular language.

We have NFAs from regular languages. What can we add to enable them to recognize CFLs?

The key idea is that CFGs allow recursive definitions.

(ECE 220) What enables recursion in programming?

We need a stack!
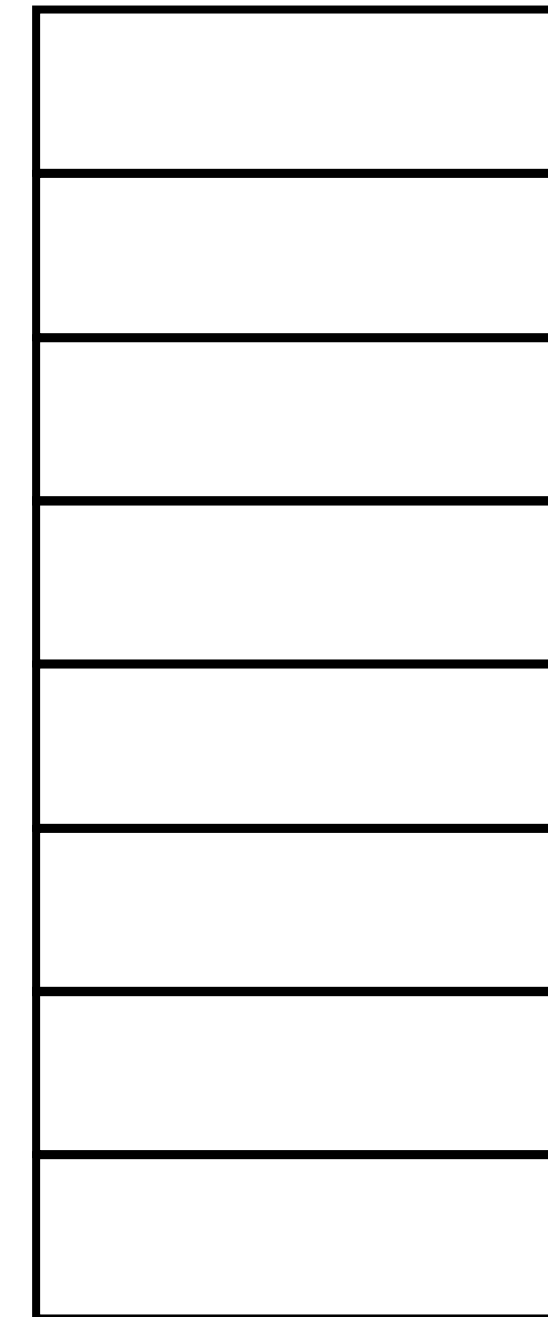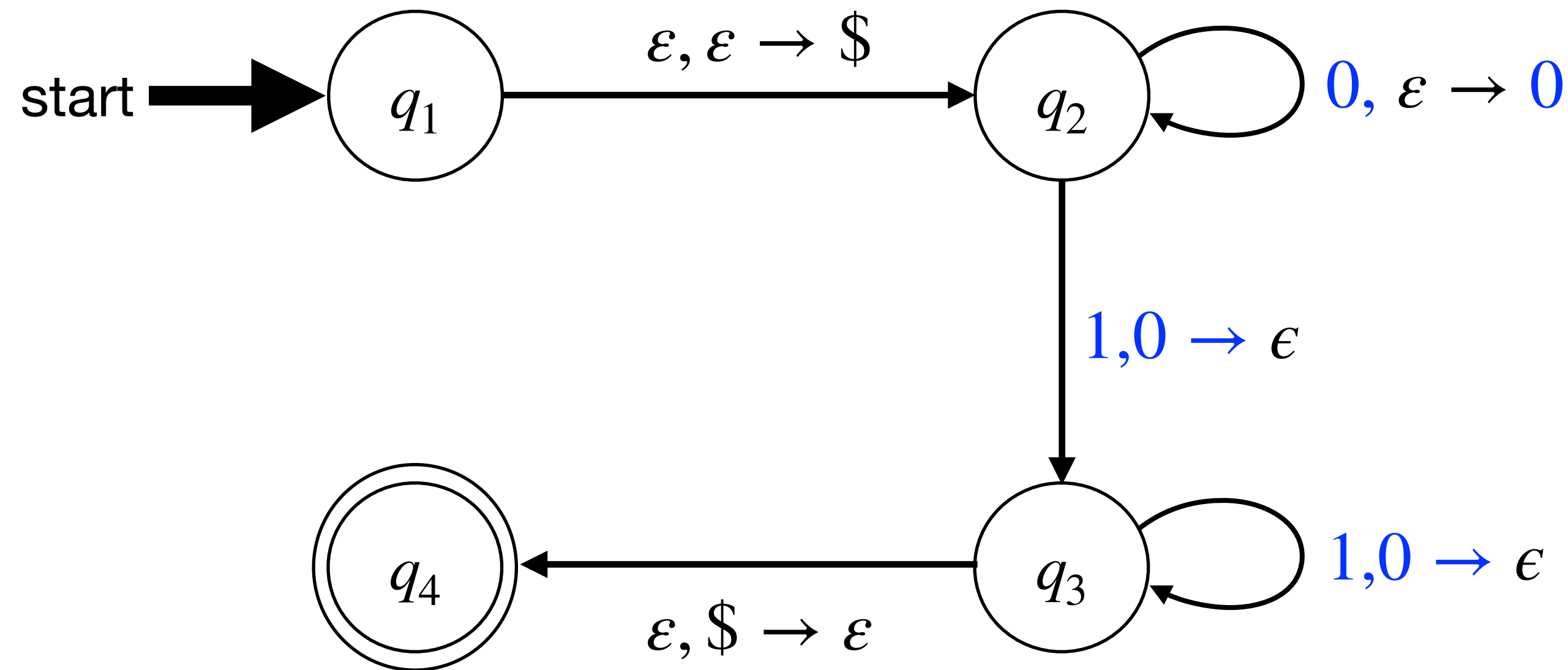
# Push-down Automata
## The machine that generates CFGs



Each transition is formatted as:

$$\text{<token read>, <stack pop>} \rightarrow \text{<stack push>}$$
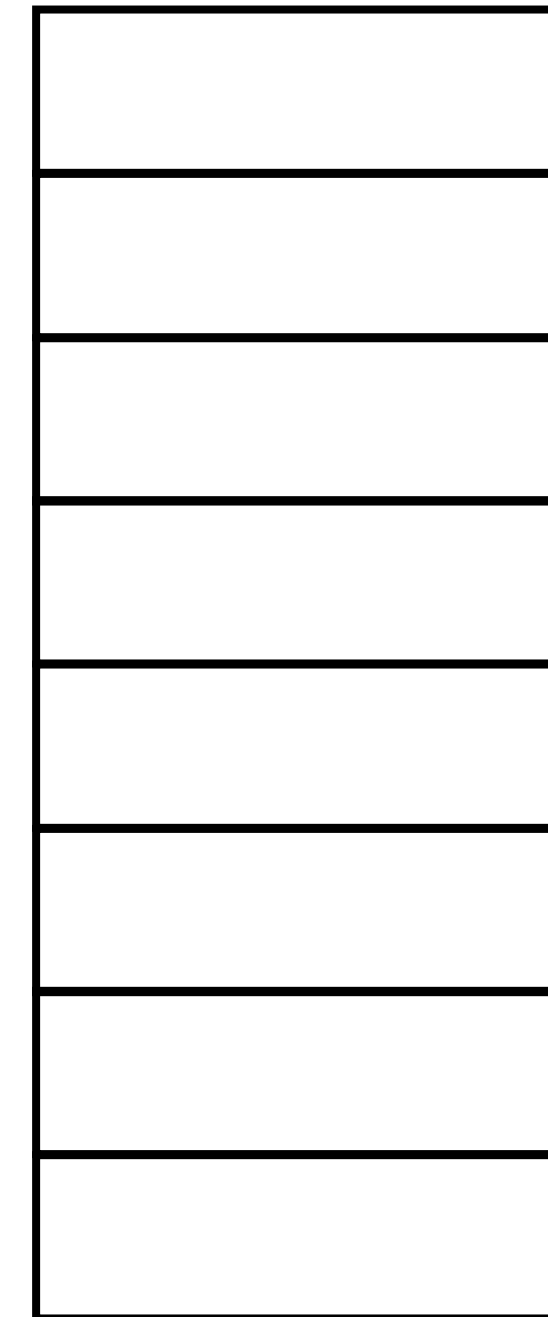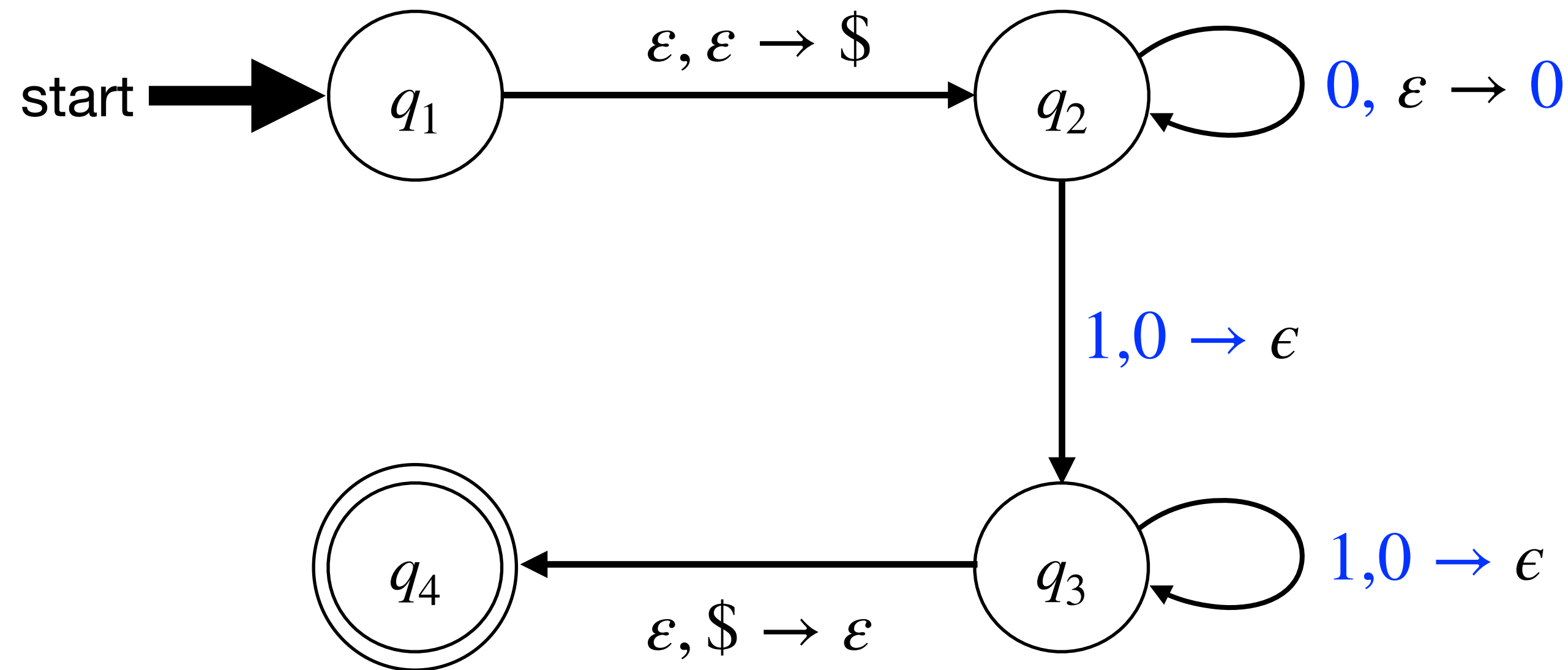
# Push-down Automata
## The machine that generates CFGs



Does this machine recognize $0011$?

# Push-down Automata
## The machine that generates CFGs



Does this machine recognize $0101$?

# Formal Tuple Notation

**Definition:** A non-deterministic push-down automaton $P = \big( Q, \Sigma, \Gamma, \delta, s, A \big)$ is a 6-tuple where

- $Q$ is a finite set whose elements are called states,

- $\Sigma$ is a finite set called the input alphabet,

- $\Gamma$ is a finite set called the stack alphabet,

- $\delta : Q \times \Sigma \cup \{\varepsilon\} \times \Gamma \cup \{\varepsilon\} \to \mathscr{P}\big( Q \times (\Gamma \cup \{\varepsilon\}) \big)$ is the transition function

- $s$ is the start state

- $A$ is the set of accepting states

# CFGs and PDAs
## Convert a CFG to a PDA

Consider,

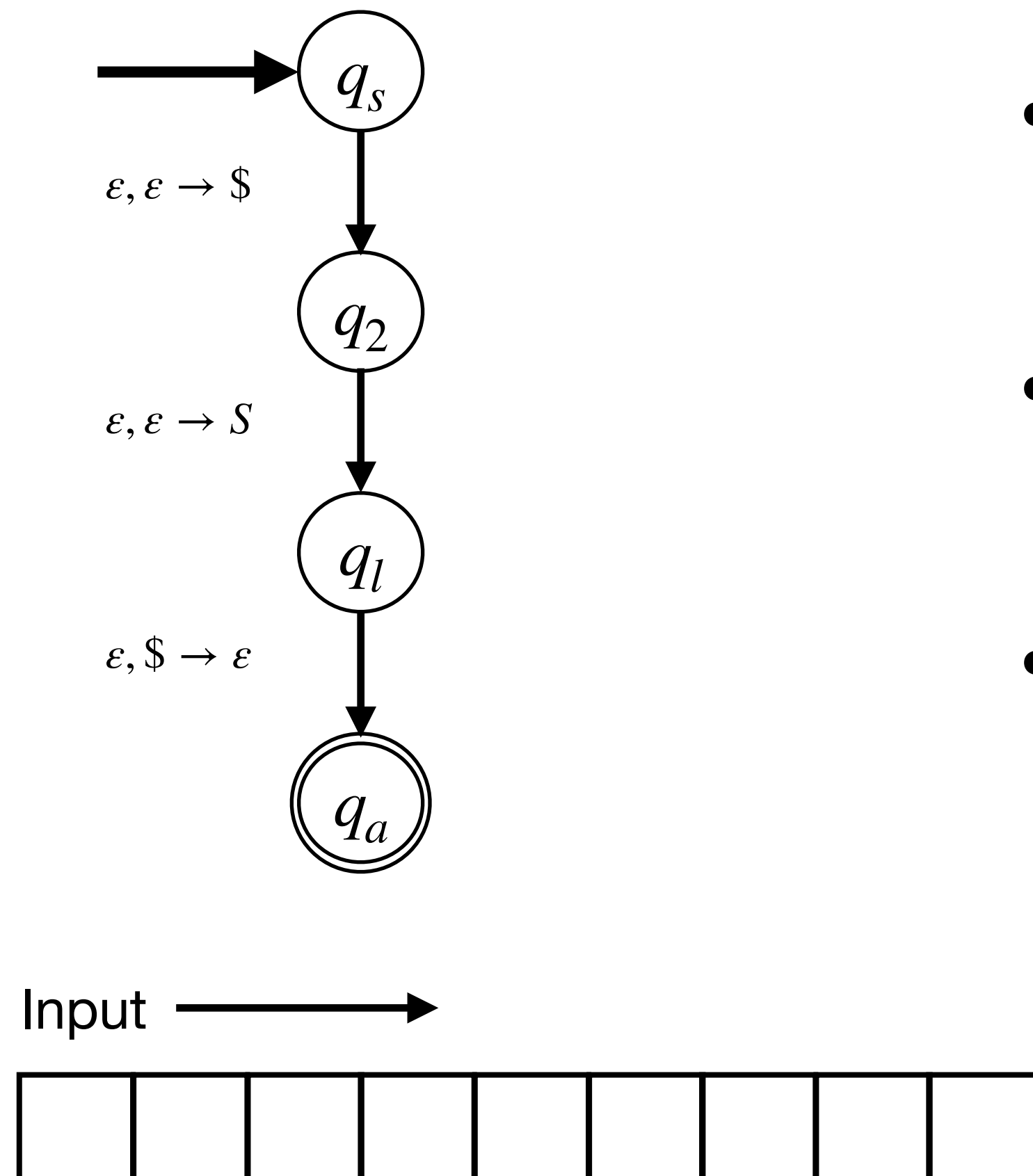$$S \rightarrow 0S \mid 1 \mid \epsilon$$

What is a PDA for this?

Key idea: Recreate the string on the stack

- Every time we see a non-terminal, we replace it with one of the replacement rules.

- Every time we see a terminal symbol, we take that symbol from the input.

- If we reach a point where the stack and input are empty, then we accept the string.

# CFGs and PDAs
## Convert a CFG to a PDA

$$S \to 0S \ | \ 1 \ | \ \epsilon$$

- First let's put in a $\$$ to mark the end of the string

- Also let's put in the start symbol on the stack.

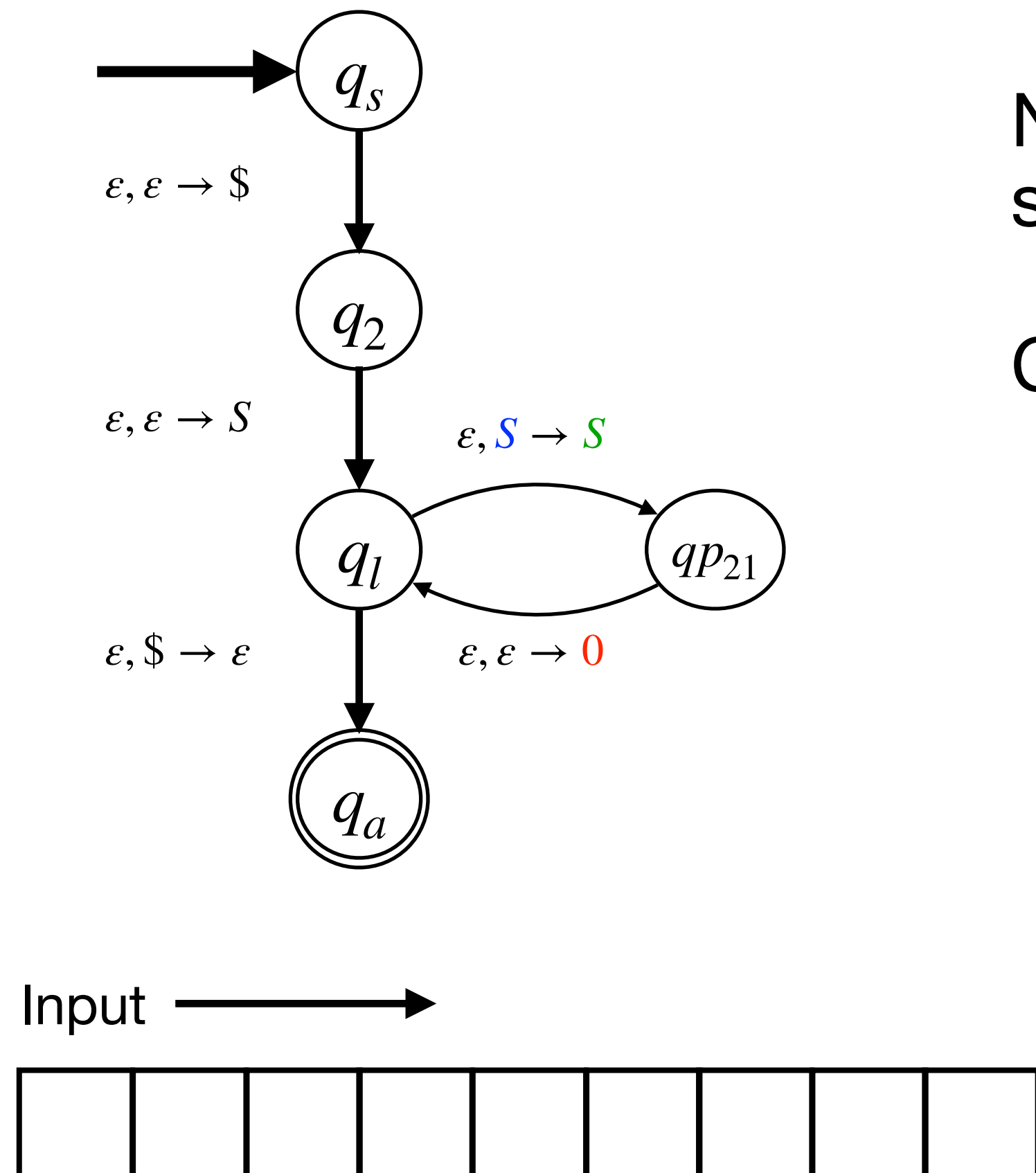- We can accept if nothing left to read and stack is empty.

$q_s$

$\varepsilon, \varepsilon \to \$$

$q_2$

$\varepsilon, \varepsilon \to S$

$q_l$

$\varepsilon, \$ \to \varepsilon$

$q_a$

Input $\longrightarrow$

Stack $\longrightarrow$

# CFGs and PDAs
## Convert a CFG to a PDA

$$S \to 0S \ \mid \ 1 \ \mid \ \epsilon$$



Next we want to add a loop for every non-terminal symbol that replaces that non-terminal with the result.
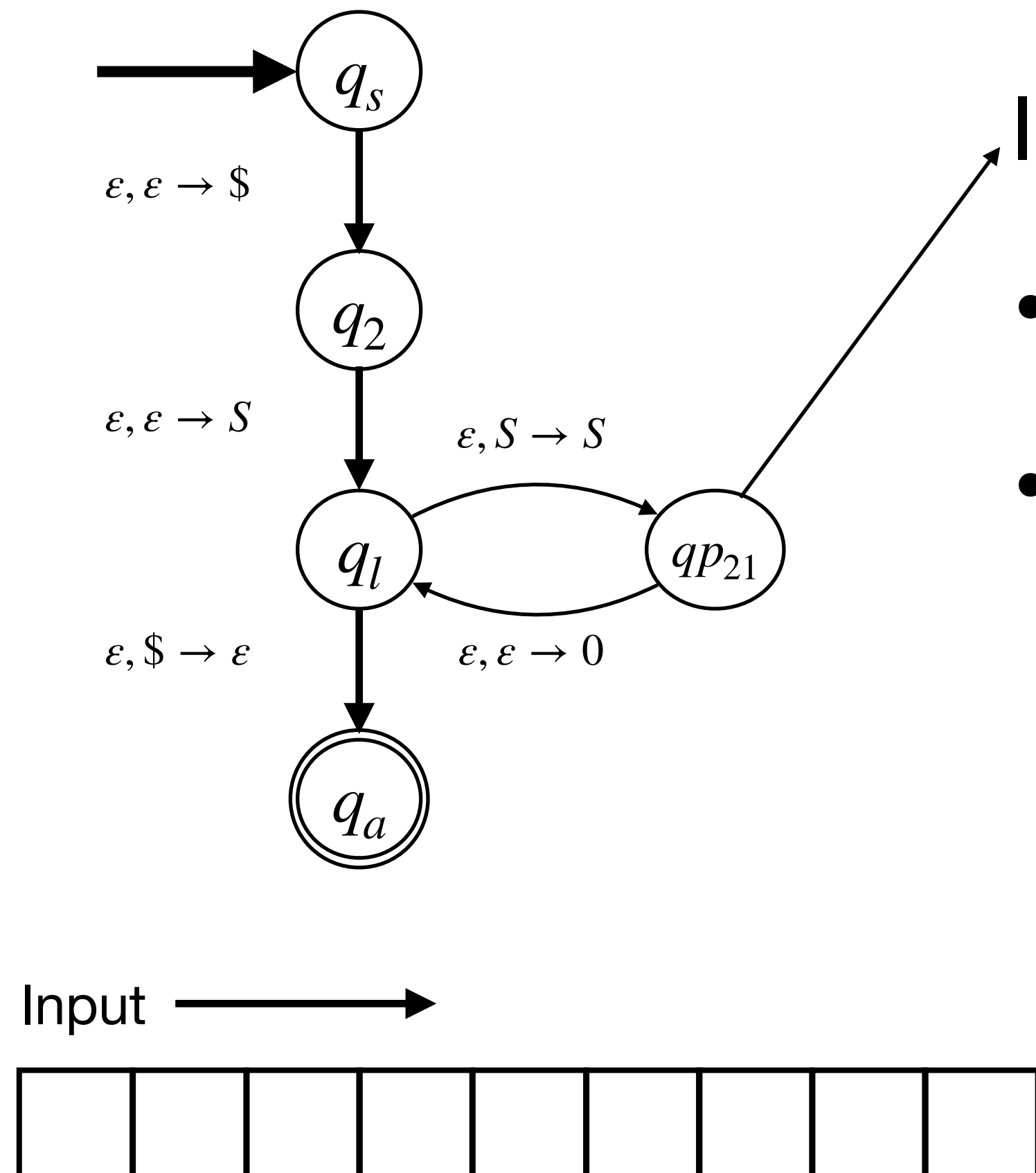
Consider the rule: $S \to 0S$

- So we got to pop $S$ the non-terminal and …

- Add a non-terminal $S$ to the stack.

- And add a terminal $0$ to the stack.
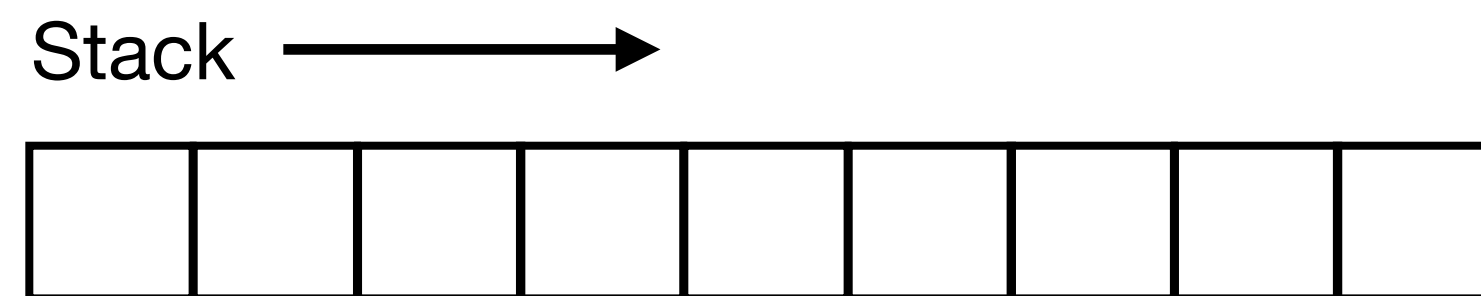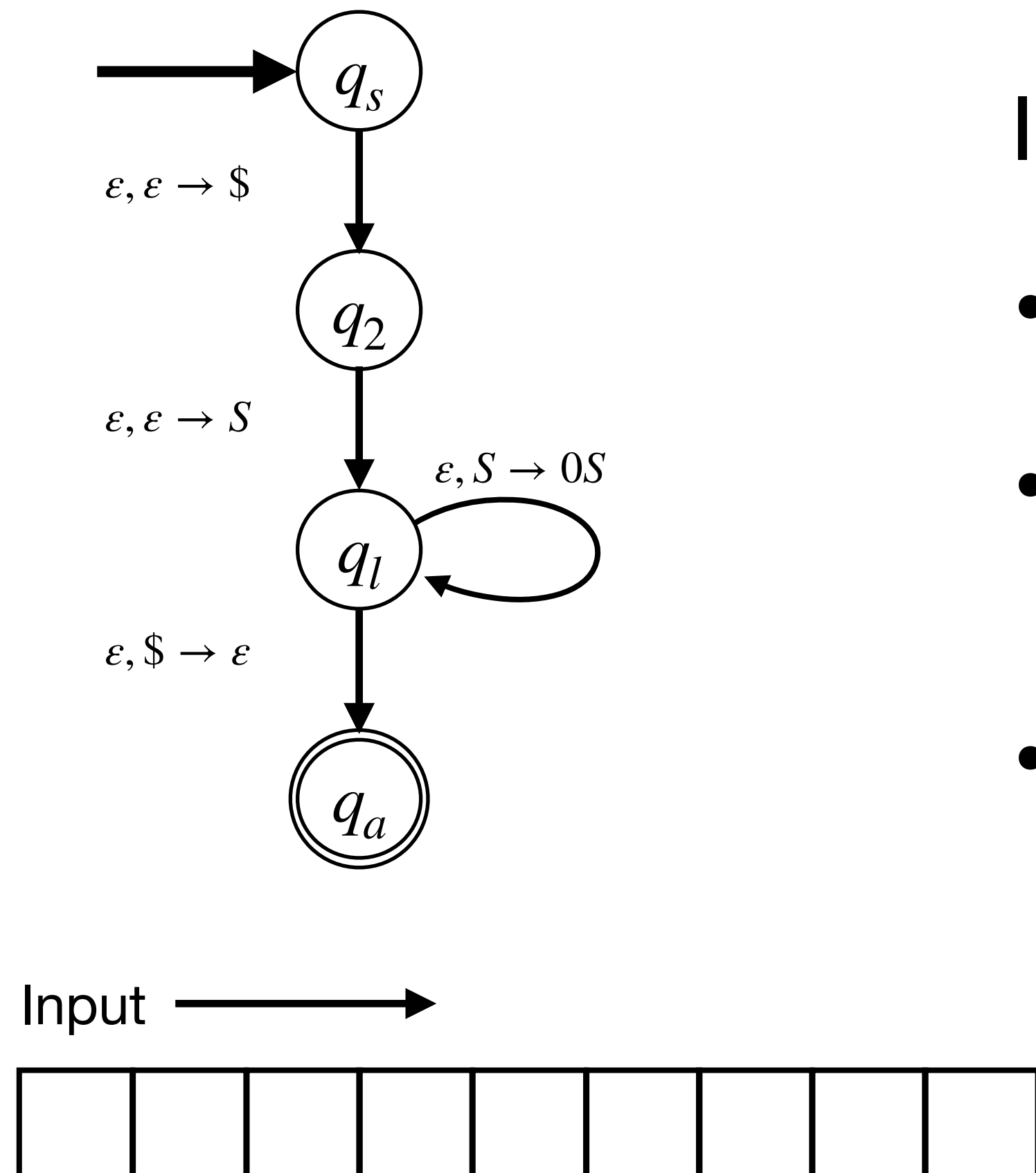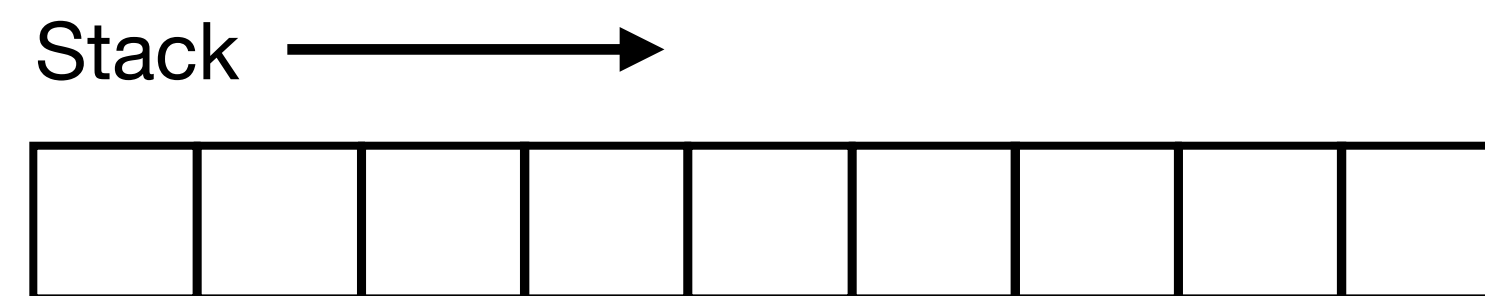
Input ⟶

Stack ⟶

27

# CFGs and PDAs
## Convert a CFG to a PDA

$$S \rightarrow 0S \mid 1 \mid \epsilon$$



Is this state necessary?

- Recall generalized NFAs?

- Can follow same route to allow entire strings to be pushed onto stack.

States and transitions shown:
- $q_s$ (start state)
- $\varepsilon, \varepsilon \rightarrow \$$ from $q_s$ to $q_2$
- $q_2$
- $\varepsilon, \varepsilon \rightarrow S$ from $q_2$ to $q_l$
- $q_l$
- $\varepsilon, S \rightarrow S$ from $q_l$ to $qp_{21}$
- $qp_{21}$
- $\varepsilon, \varepsilon \rightarrow 0$ from $qp_{21}$ to $q_l$
- $\varepsilon, \$ \rightarrow \varepsilon$ from $q_l$ to $q_a$
- $q_a$ (accept state)

Input →

Stack →

# CFGs and PDAs
## Convert a CFG to a PDA

$$S \rightarrow 0S \;\; | \;\; 1 \;\; | \;\; \epsilon$$



$\varepsilon, \varepsilon \rightarrow \$$

$\varepsilon, \varepsilon \rightarrow S$

$\varepsilon, S \rightarrow 0S$

$\varepsilon, \$ \rightarrow \varepsilon$

Is this state necessary?

- Recall generalized NFAs?

- Can follow same route to allow entire strings to be pushed onto stack.
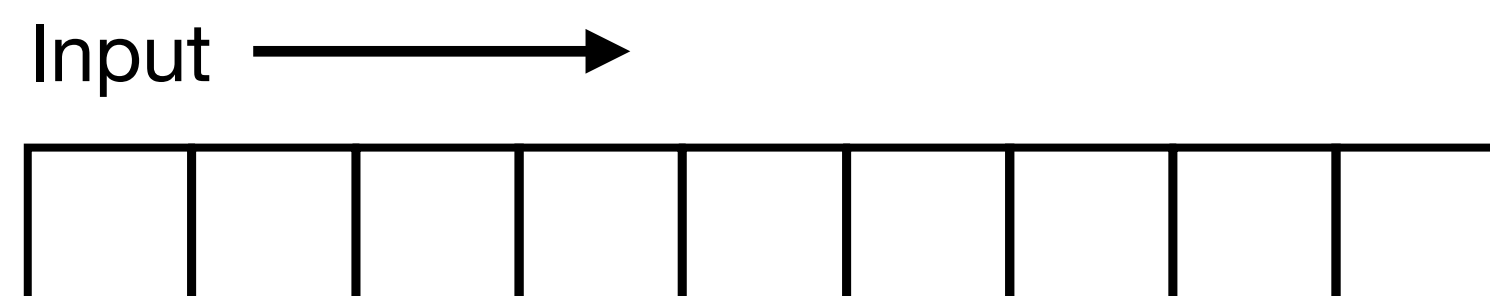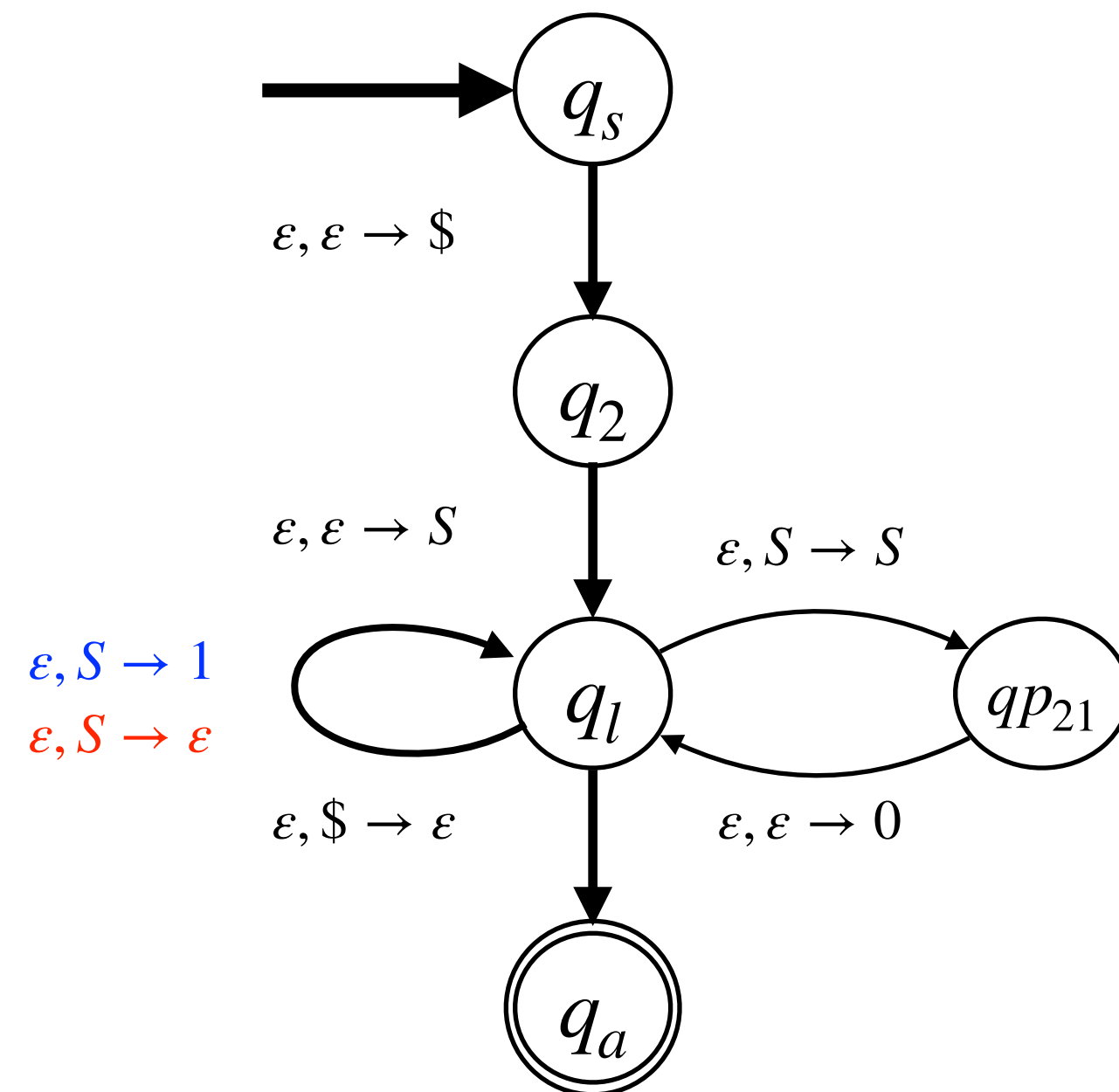
- But we are going to stick with PDAs.

Input $\longrightarrow$
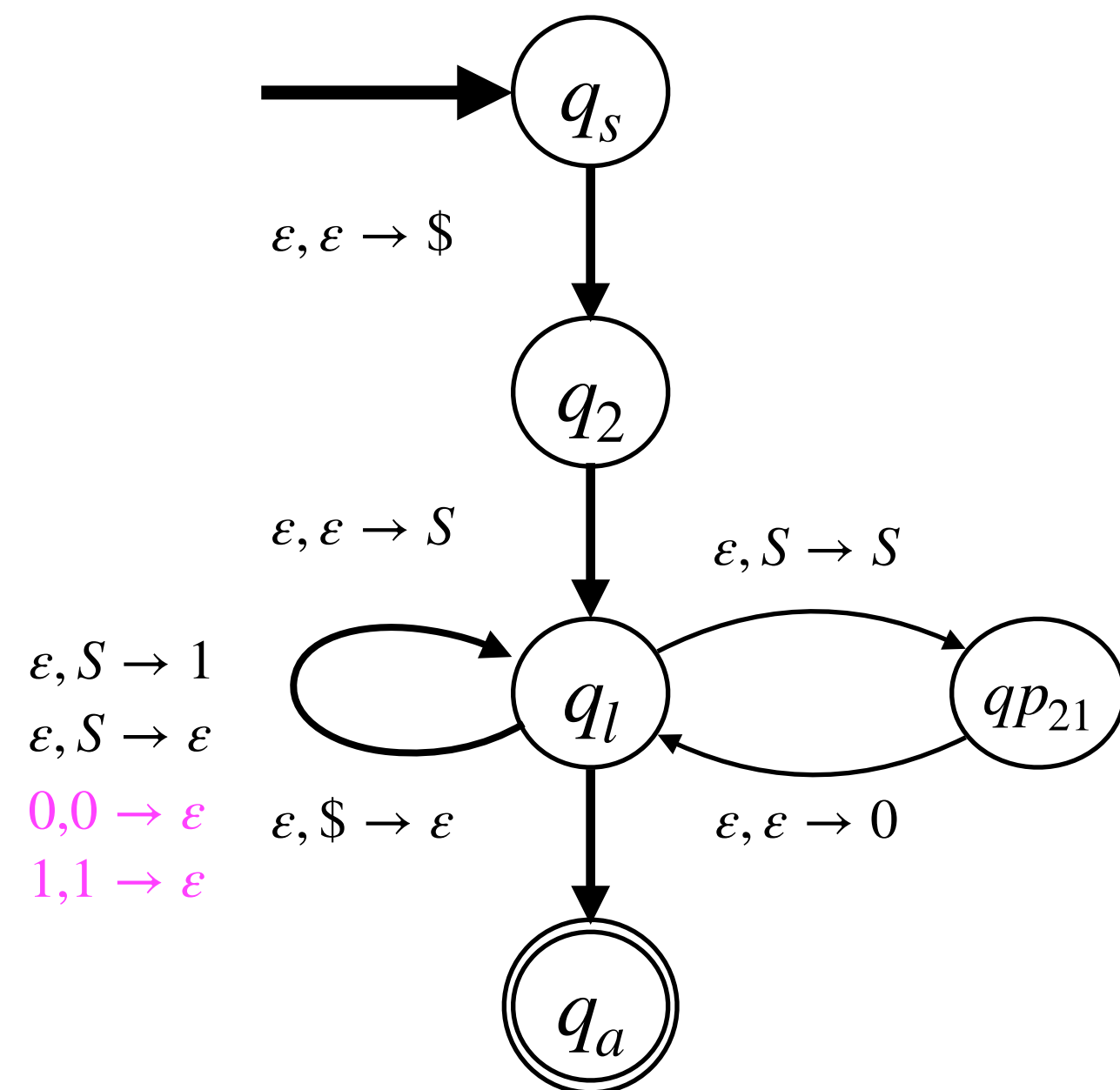
Stack $\longrightarrow$

# CFGs and PDAs
## Convert a CFG to a PDA



$$S \rightarrow 0S \mid 1 \mid \varepsilon$$

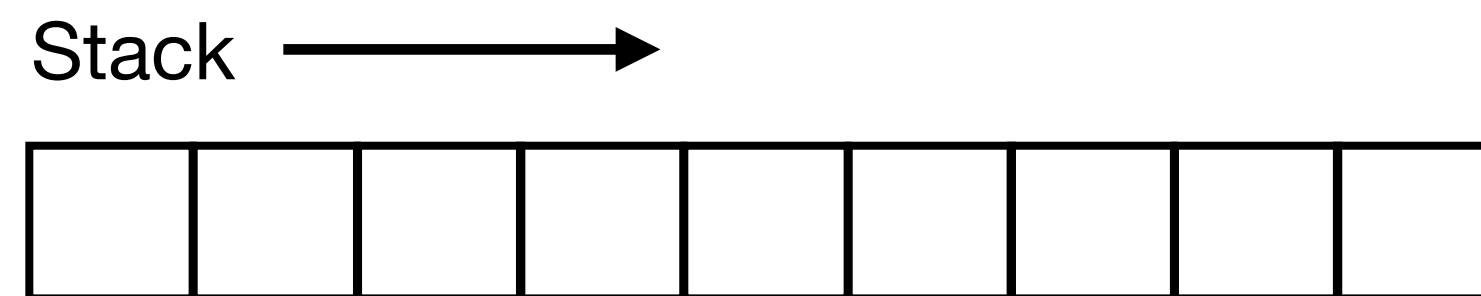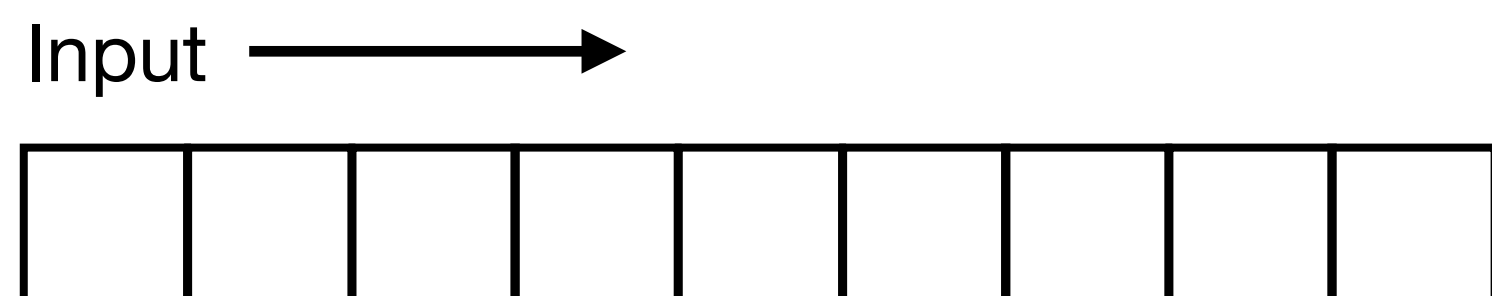- Do the same thing for $S \rightarrow 1$ and $S \rightarrow \varepsilon$

The diagram shows states $q_s$, $q_2$, $q_l$, $qp_{21}$, and $q_a$ with transitions:
- $\varepsilon, \varepsilon \rightarrow \$$
- $\varepsilon, \varepsilon \rightarrow S$
- $\varepsilon, S \rightarrow S$
- $\varepsilon, S \rightarrow 1$
- $\varepsilon, S \rightarrow \varepsilon$
- $\varepsilon, \$ \rightarrow \varepsilon$
- $\varepsilon, \varepsilon \rightarrow 0$

Input →

Stack →

# alf
## Convert a CFG to a PDA

$$S \to 0S \ | \ 1 \ | \ \varepsilon$$

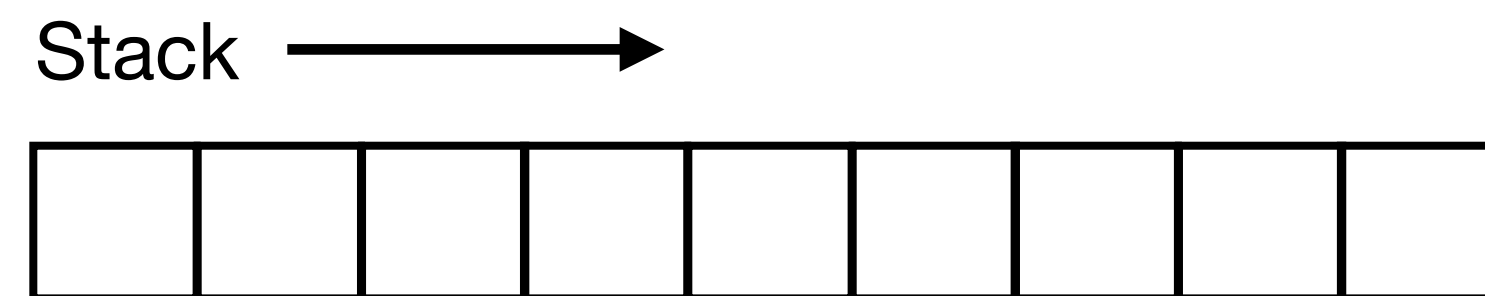- If we see a non-terminal symbol on the stack, then we can cross that symbol from the input.
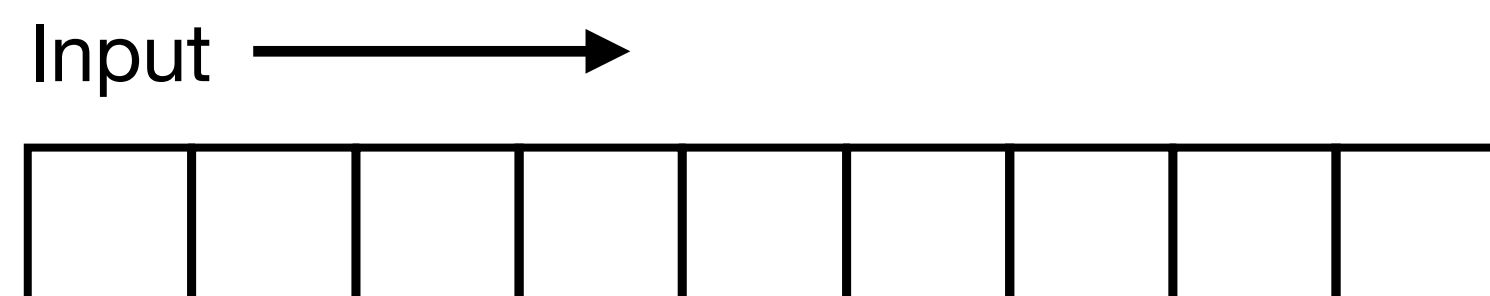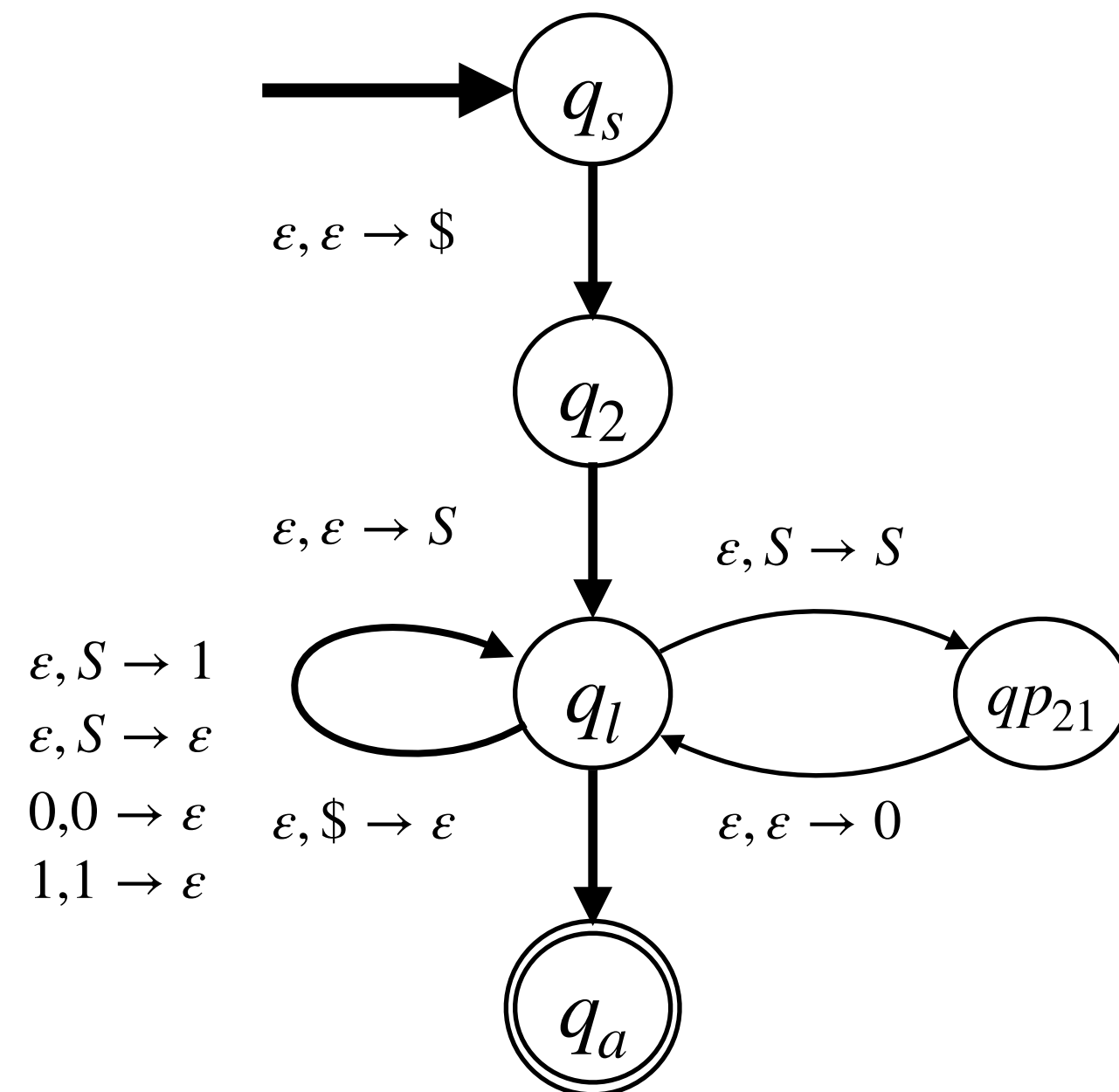


Input ⟶

Stack ⟶

# CFGs and PDAs
## Convert a CFG to a PDA

$$S \to 0S \ | \ 1 \ | \ \varepsilon$$

Study the automata to verify:

- Does this automata accept $001$?

- Does this automata accept $010$?



$\varepsilon, \varepsilon \to \$$

$q_s$

$q_2$

$\varepsilon, \varepsilon \to S$

$\varepsilon, S \to S$

$\varepsilon, S \to 1$
$\varepsilon, S \to \varepsilon$
$0,0 \to \varepsilon$
$1,1 \to \varepsilon$

$\varepsilon, \$ \to \varepsilon$

$q_l$

$qp_{21}$

$\varepsilon, \varepsilon \to 0$

$q_a$

Input ⟶

Stack ⟶

# Convert a CFG to a PDA
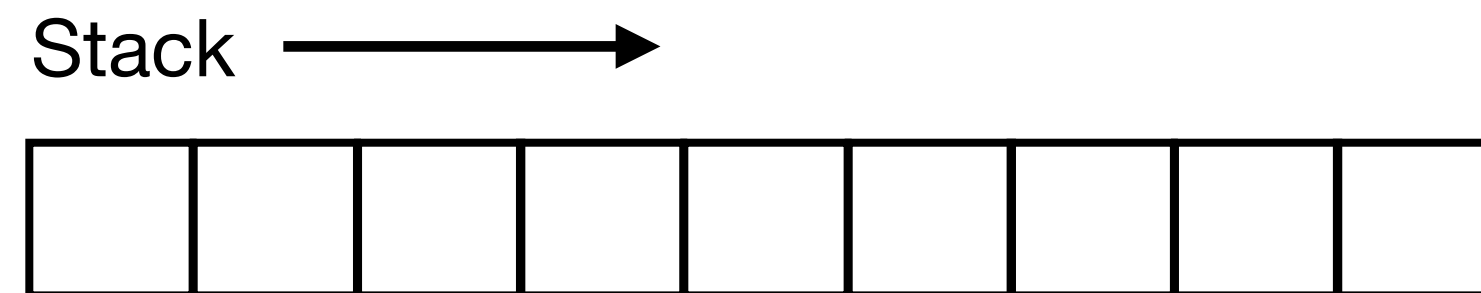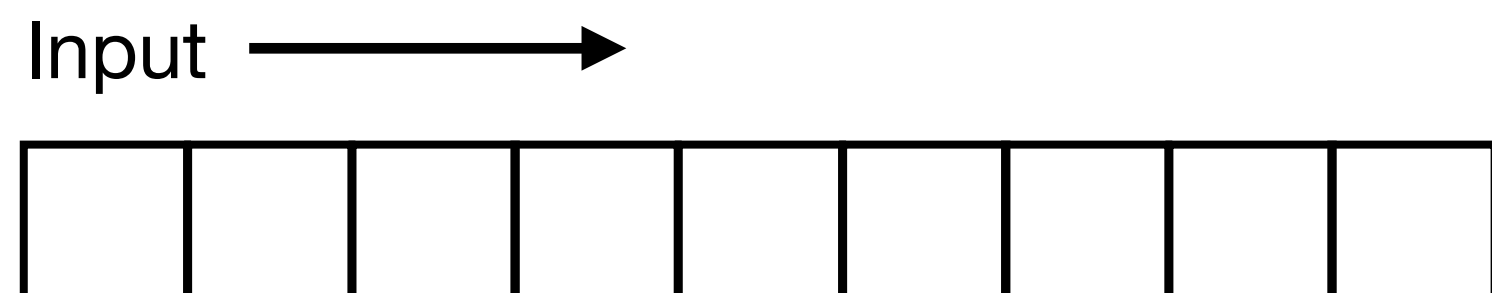
## Another example



$S \to 0T1 \mid 1$

$T \to T0 \mid \varepsilon$

- Insert transitions for initialization, start symbol & accept state

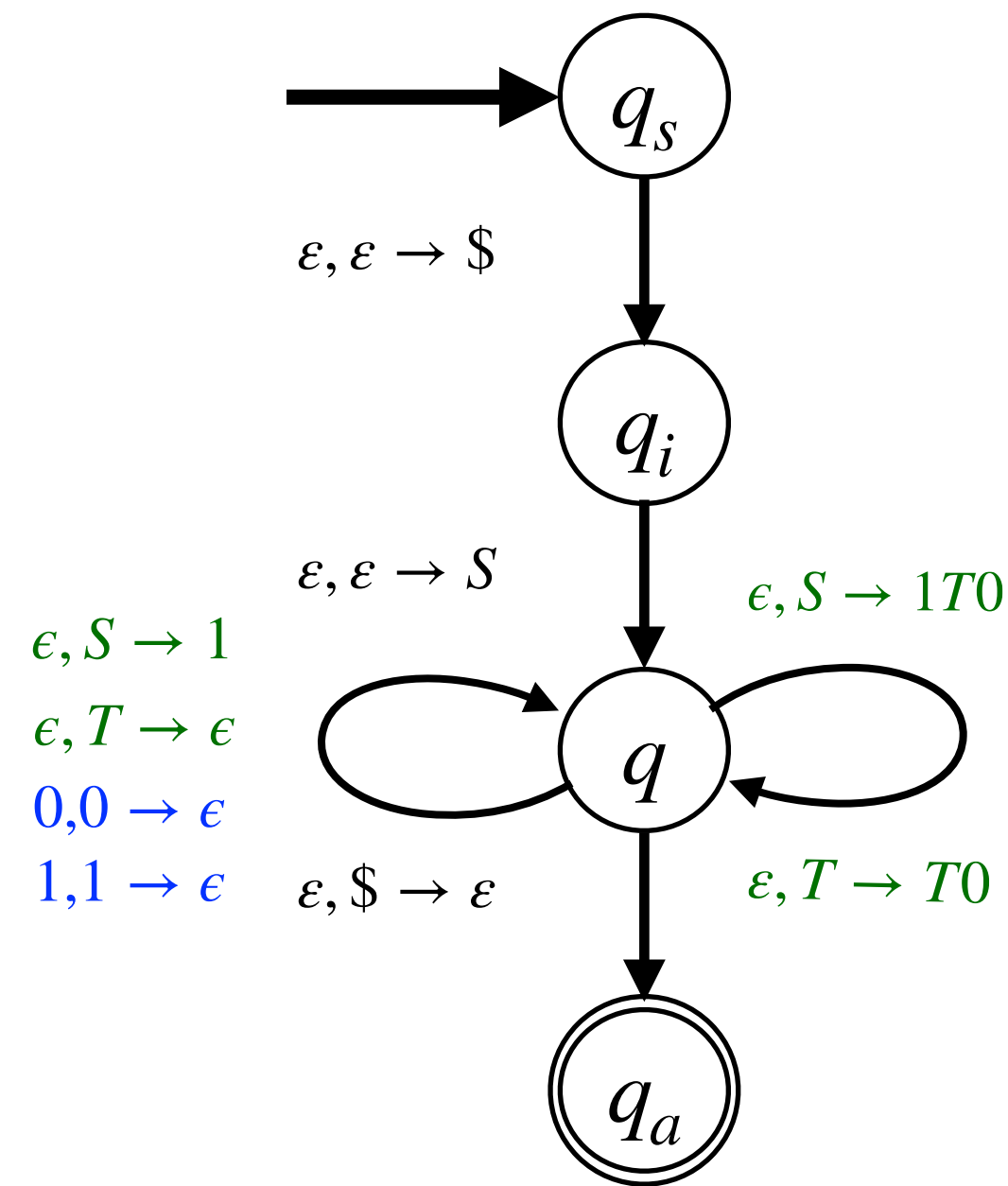- Add all production rules
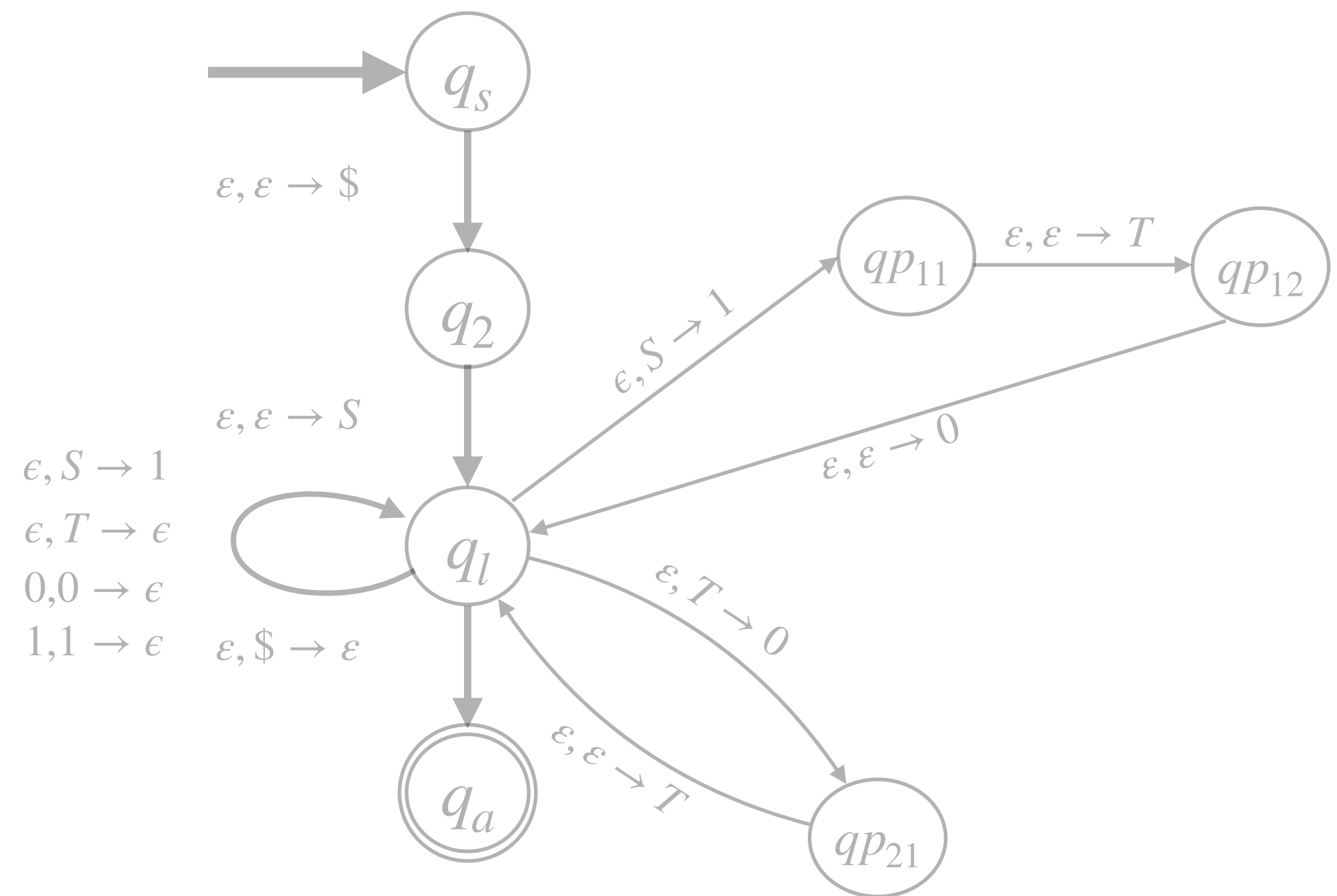
- Take care of terminals

# Convert a CFG to a PDA
## With generalized PDAs

- Start with the grammar $G = (V, T, P, S)$ and consider the PDA

$$M = \left( \{q_s, q_l, q, q_a\}, T, V \cup T, \delta, q_s, \{q_a\} \right)$$

- Define $\delta$ as follows:

  - Insert transitions for initialization, start symbol & accept state.

  - For every production rule $A \to \beta$ in $P$, add a transition from $q$ to $q$, consuming $\varepsilon$, popping $A$ and pushing $\beta$.

  - For every terminal $t \in T$, add a transition from $q$ to $q$, consuming $t$, popping $t$ and pushing $\varepsilon$.

# Next class of languages
**Canonical non-CFL**

- $L = \{a^n b^n c^n \mid n \geq 0\}$

  - Intuition why a PDA cannot recognize this language.

  - This is in fact what we call a context-sensitive language.

    - Corresponding automaton is called *Linear Bounded Automaton* (LBA)

      - We will not discuss LBAs

- Next class: **Turing Machines**