

CS 314 Principles of Programming Languages

Lecture 2: Syntax Analysis (Scanning)

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Prof. Zheng Zhang



Rutgers University

September 6, 2018

Announcement

- **First recitation starts this coming Wednesday**
- **Homework 1 will be released after lecture 3.**
- **My office hour:**
Thursday 2pm - 3pm at CoRE 315
- **TA office hours will be announced soon.**

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Last Class

- Overview of compilation
- Syntax and semantics
- Formal language definition
- A rule-based rewriting system
- Introduction to regular expression

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Compilation and Interpretation

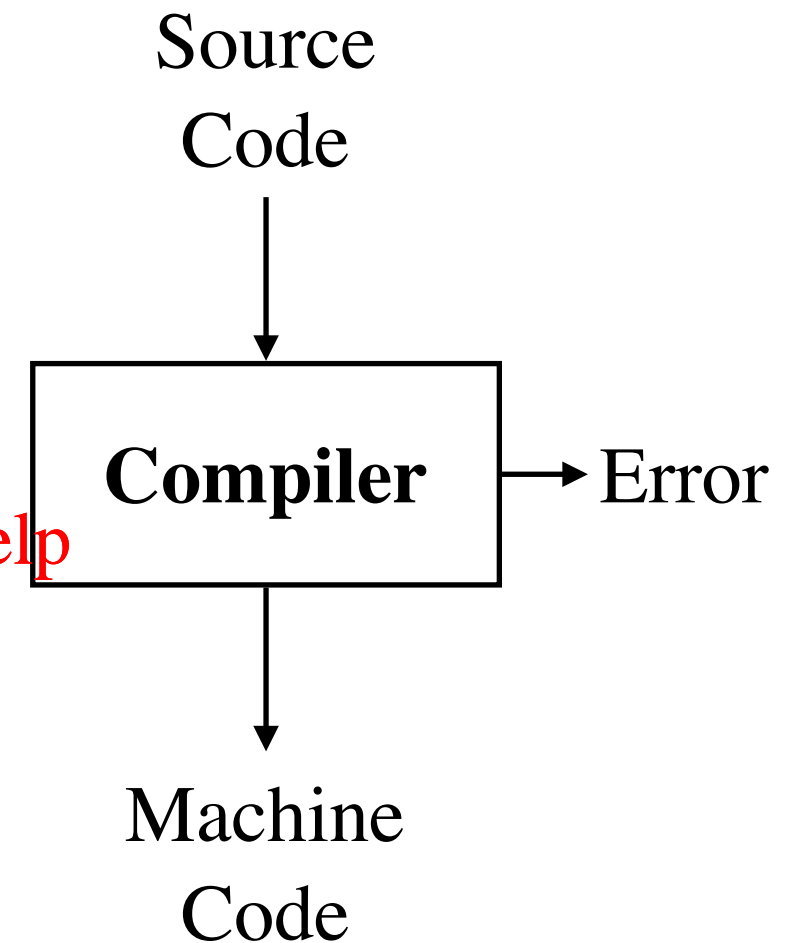
Compiler

- Recognize legal (and illegal) programs
- Generate correct code
- Manage storage of all variables and code
- Need format for object (or assembly) code

Assignment Project Exam Help

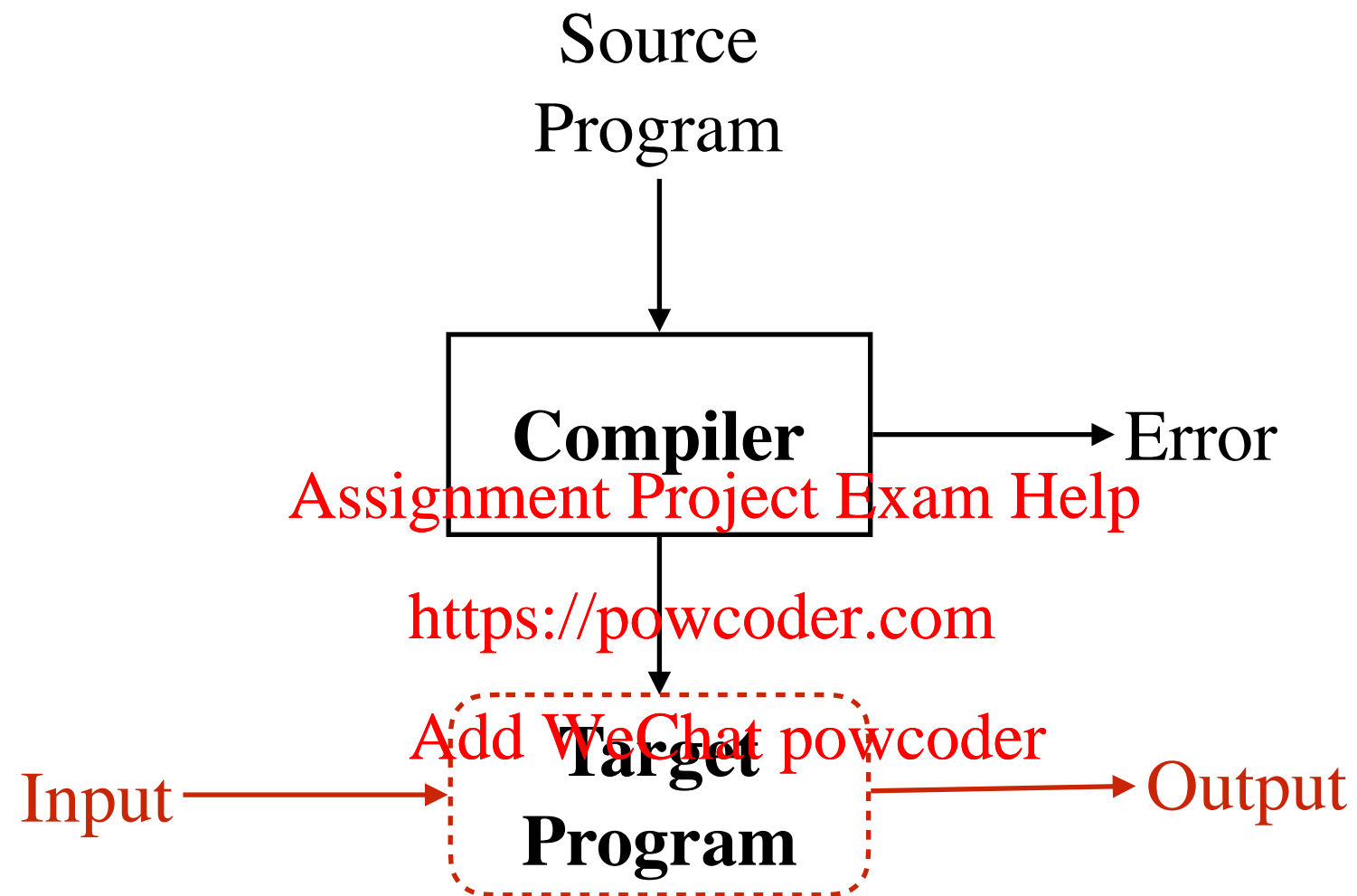
<https://powcoder.com>

Add WeChat powcoder



Big step up from assembler to higher level notations

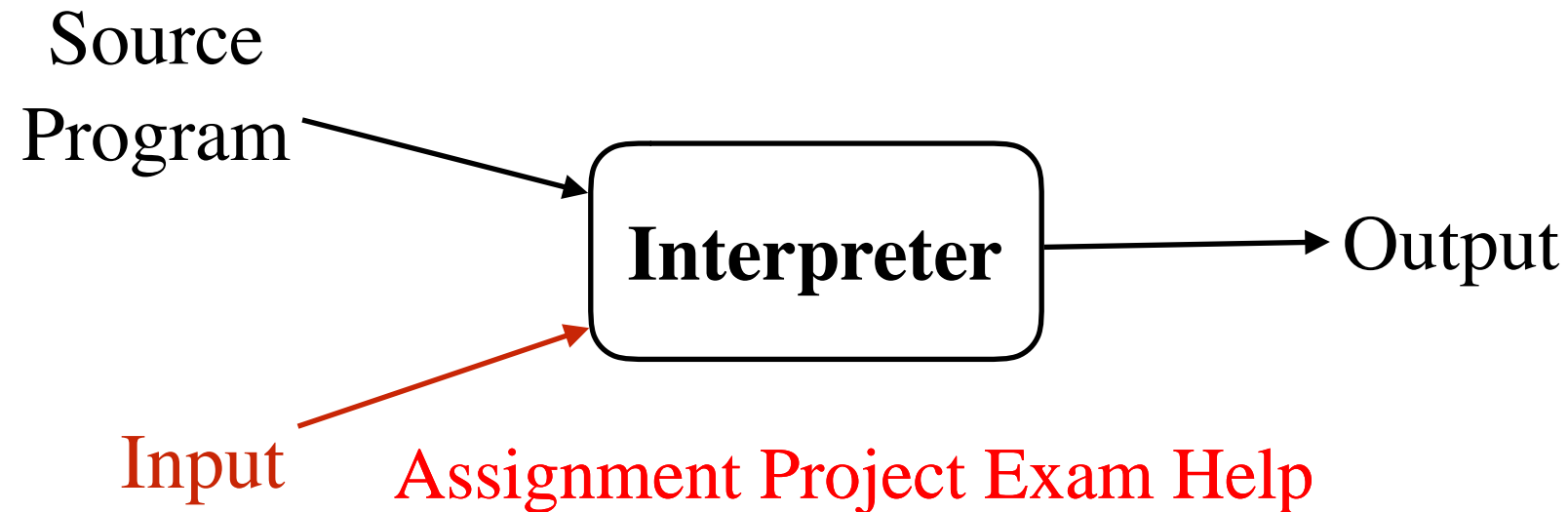
Compilation and Interpretation



Pure Compilation

- Mainly refers to translation
- Take a program in source language, output a program in target language (usually machine code)

Compilation and Interpretation



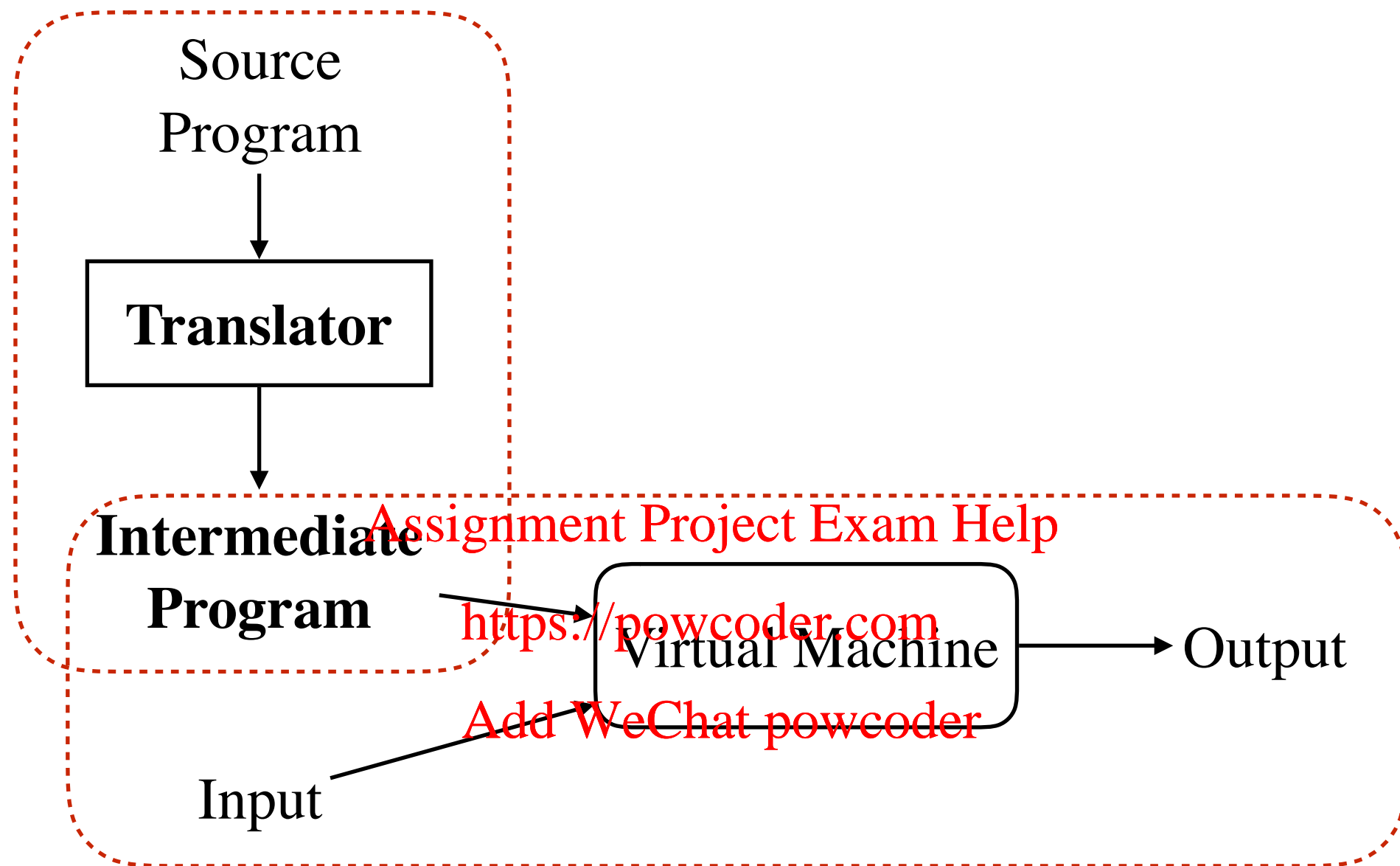
<https://powcoder.com>

Add WeChat powcoder

Interpretation

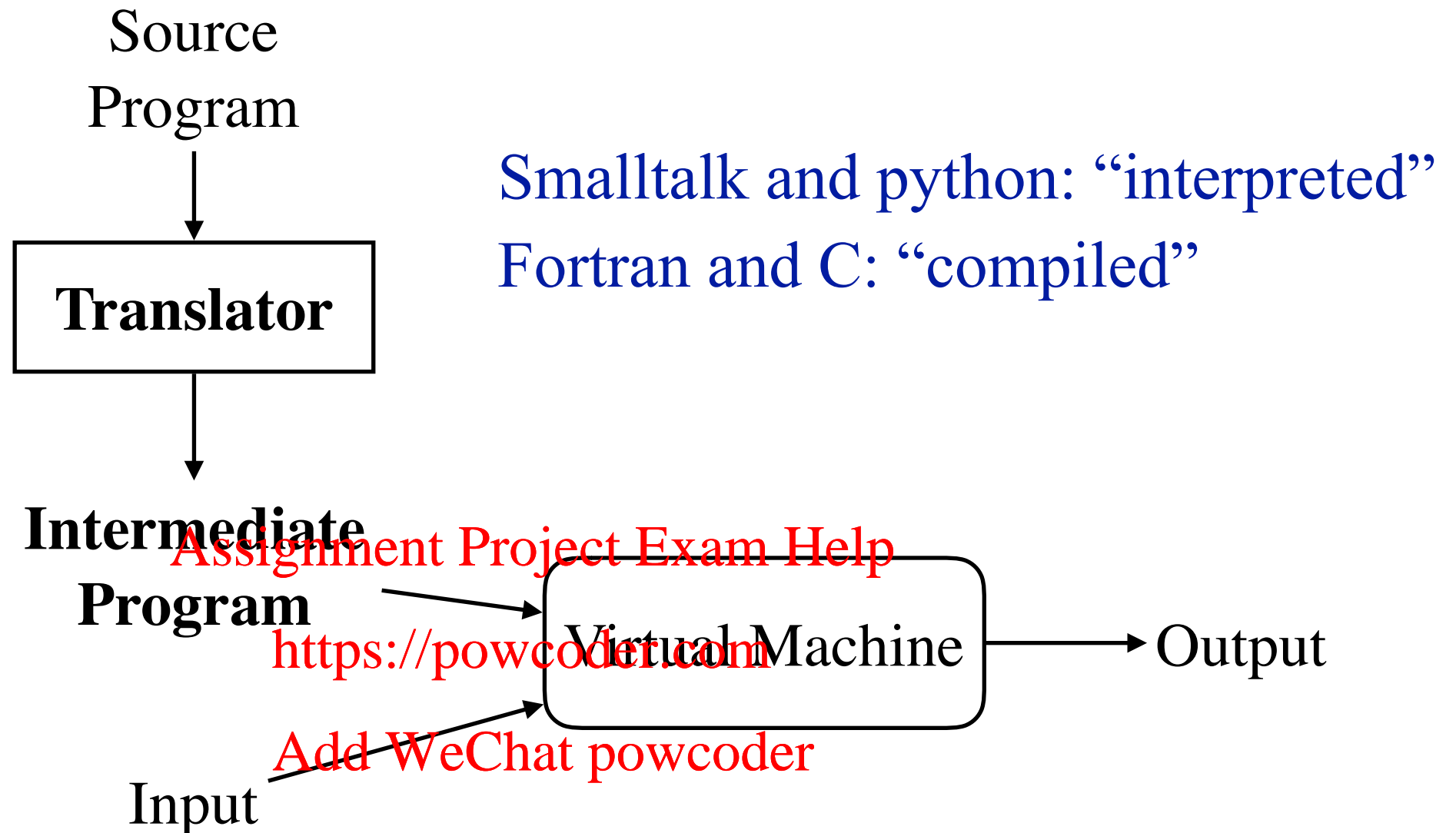
- Interpreter stays around for the execution of the program
- Interpreter is the locus of control during execution

Compilation and Interpretation



- Most language implementations include a mixture of both compilation and interpretation.
- Common case is compilation or simple pre-processing, followed by interpretation.

Compiled V.S. Interpreted Languages



- We generally say that:
A language is “interpreted” if the initial translator is “simple”,
or “compiled” if the initial translator is “complicated”
- Very subjective, but a language is still “compiled” if the translator has thorough analysis and non-trivial transformation.

Syntax and Semantics of Programming Languages

Syntax:

Describes what a legal program looks like

Semantics:

Describes what a correct (legal) program means

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Syntax of Programming Languages

The syntax of programming languages is often defined in two layers:

tokens and *sentences*.

- *tokens* - legal combination of characters in the language

Question: How to spell a token (word)?

Assignment Project Exam Help

Answer: Regular expressions

<https://powcoder.com>

- *sentences* - legal combinations of tokens in the language

Question: How to build correct sentences with tokens?

Answer: (Context - free) grammars (CFG)

Lexical Analysis (Scott 2.1, 2.2)

Character Sequence:

i f a < = b t h e n c : = d

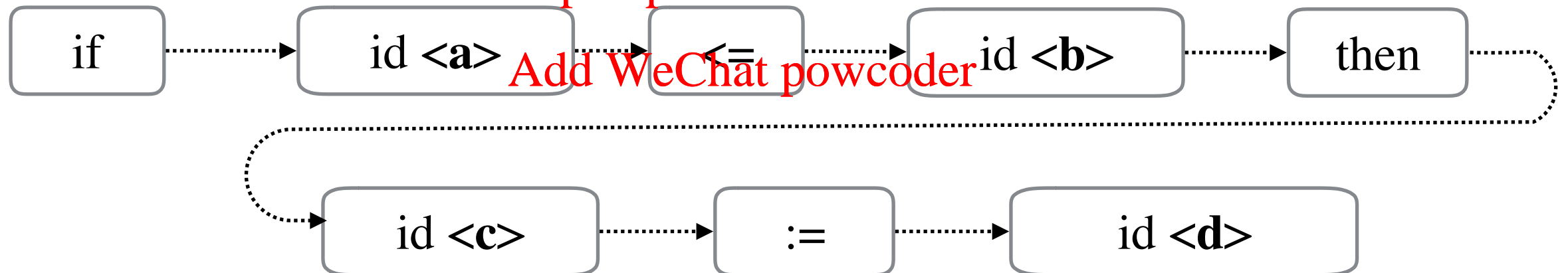
scanner

Token Sequence:

Assignment Project Exam Help

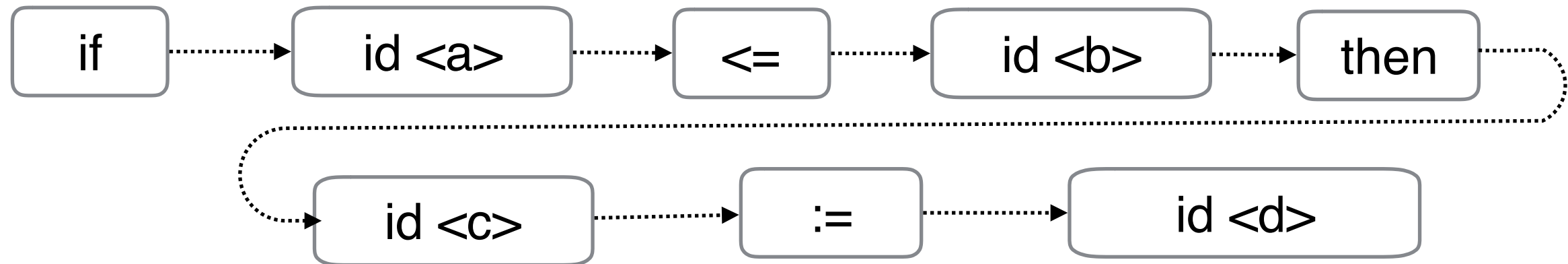
<https://powcoder.com>

Add WeChat powcoder



Syntax Analysis (Scott, Chapter 2.3)

Token Sequence:

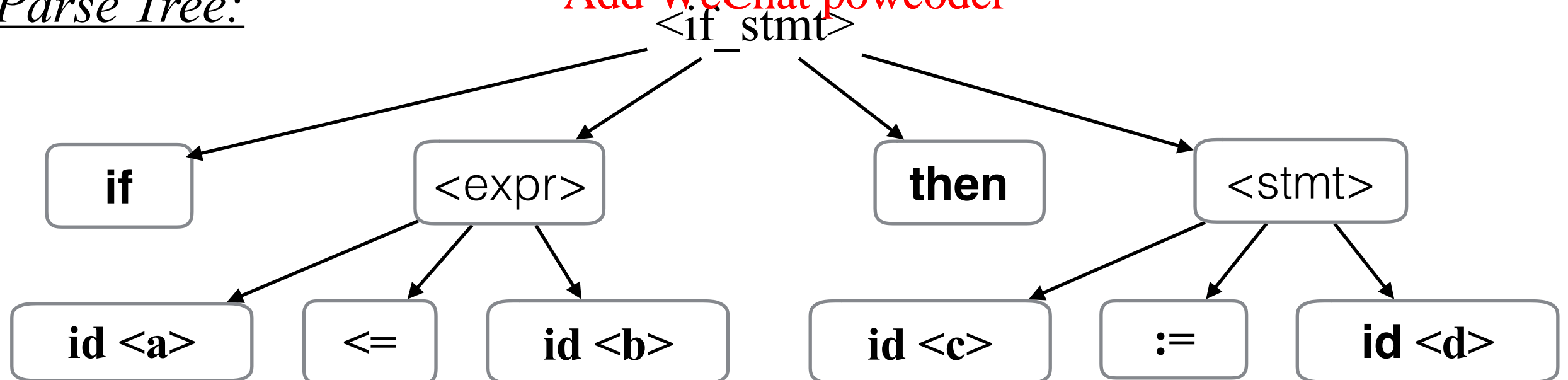


Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Parse Tree:



Tokens (Scott 2.1, 2.2)

Tokens (Analogous to *Words of Language*)

- Smallest “atomic” units of further syntax analysis
- Used to build all the other constructs
- Example, in C:

Keywords: **for if goto volatile...**

= * / - < > == <= >= <> () [] ; := . , ...

Number: (Example: 3.14 28 ...)

Identifier: (Example: b square addEntry ...)

Formalisms for Lexical and Syntactic Analysis

Two issues in *Formal Languages*:

- Language Specification → formalism to describe what a valid program (word/sentence) looks like.
- Language Recognition → formalism to describe a machine and an algorithm that can verify that a program is valid or not.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Formalisms for Lexical and Syntactic Analysis

Two issues in *Formal Languages*:

- Language Specification → formalism to describe what a valid program (word/sentence) looks like.
- Language Recognition → formalism to describe a machine and an algorithm that can verify that a program is valid or not.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

We use **regular expression** to specify tokens (words)

Regular Expressions

A syntax (notation) to specify regular languages.

<i>RE</i> p	Language $L(p)$
---------------	-----------------

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

$L(p)$: the set of strings that can be represented using the expression p

Regular Expressions

A syntax (notation) to specify regular languages.

<i>RE</i> p	Language $L(p)$
---------------	-----------------

$r \mid s$	$L(r) \cup L(s)$
------------	------------------

Either r or s is a regular expression, i.e. $0|11$

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

$L(p)$: the set of strings that can be represented using the expression p

Regular Expressions

A syntax (notation) to specify regular languages.

<i>RE p</i>	Language $L(p)$
-------------	-----------------

$r \mid s$	$L(r) \cup L(s)$
------------	------------------

rs	$\{RS \mid R \in L(r), S \in L(s)\}$
------	--------------------------------------

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

$L(p)$: the set of strings that can be represented using the expression p

Regular Expressions

A syntax (notation) to specify regular languages.

<i>RE</i> p	Language $L(p)$
---------------	-----------------

$r \mid s$	$L(r) \cup L(s)$
------------	------------------

rs	$\{RS \mid R \in L(r), S \in L(s)\}$
------	--------------------------------------

r^+	$L(r) \cup L(rr) \cup L(rrr) \cup \dots$
-------	--

<https://powcoder.com>

Add WeChat powcoder

$L(p)$: the set of strings that can be represented using the expression p

Regular Expressions

A syntax (notation) to specify regular languages.

<i>RE p</i>	Language $L(p)$
---------------------------------	-----------------------------------

$r \mid s$	$L(r) \cup L(s)$
------------	------------------

rs	$\{RS \mid R \in L(r), S \in L(s)\}$
------	--------------------------------------

r^+	$L(r) \cup L(rr) \cup L(rrr) \cup \dots$
-------	--

Assignment Project Exam Help

<https://powcoder.com>

$r^* (r^* = r^+ \mid \epsilon)$

Add WeChat powcoder
 $\{\epsilon\} \cup L(r) \cup L(rr) \cup \dots$



Any number of r 's concatenated.

(s)	$L(s)$
-------	--------

Regular Expressions

A syntax (notation) to specify regular languages.

<i>RE</i> p	Language $L(p)$
---------------	-----------------

$r \mid s$	$L(r) \cup L(s)$
------------	------------------

rs	$\{RS \mid R \in L(r), S \in L(s)\}$
------	--------------------------------------

r^+	$L(r) \cup L(rr) \cup L(rrr) \cup \dots$ https://powcoder.com
-------	---

Assignment Project Exam Help

$L(r) \cup L(rr) \cup L(rrr) \cup \dots$
<https://powcoder.com>

r^* ($r^* = r^+ \mid \epsilon$)	$\{\epsilon\} \cup L(r) \cup L(rr) \cup \dots$
-------------------------------------	--

Add WeChat powcoder

Any number of r 's concatenated.

(s)	$L(s)$
-------	--------

a	$\{\mathbf{a}\}$
ϵ	$\{\epsilon\}$

A RE can simply be a letter from the alphabet Σ or an empty string ϵ

Examples of Expressions— Solution

RE

Language

$a|bc$

$\{a, bc\}$

$(b|c)a$

$\{ba, ca\}$

$a \in$

$\{a\}$

Assignment Project Exam Help

$a^*|b$

<https://powcoder.com>

ab^*

Add WeChat powcoder

$ab^*|c^+$

$(a|b)^*$

$(0|1)^*0$

Examples of Expressions— Solution

RE

Language

$a|bc$

$\{a, bc\}$

$(b|c)a$

$\{ba, ca\}$

$a \in$

$\{a\}$

Assignment Project Exam Help

$a^*|b$

$\{a, aa, aaa, aaaa, \dots\} \cup \{b\}$

Add WeChat powcoder

ab^*

$ab^*|c^+$

$(a|b)^*$

$(0|1)^*0$

Examples of Expressions— Solution

RE

Language

$a|bc$

$\{a, bc\}$

$(b|c)a$

$\{ba, ca\}$

$a \in$

$\{a\}$

Assignment Project Exam Help

$a^*|b$

$\{a, aa, aaa, aaaa, \dots\} \cup \{b\}$

Add WeChat powcoder

ab^*

$\{a, ab, abb, abbb, abbbb, \dots\}$

$ab^*|c^+$

$(a|b)^*$

$(0|1)^*0$

Examples of Expressions— Solution

RE	Language
$a bc$	$\{a, bc\}$
$(b c)a$	$\{ba, ca\}$
$a \epsilon$	$\{a\}$
$a^* b$	$\{a, aa, aaa, aaaa, \dots\} \cup \{b\}$
ab^*	$\{a, ab, abb, abbb, abbbb, \dots\}$
$ab^* c^+$	$\{a, ab, abb, abbb, abbbb, \dots\} \cup \{c, cc, ccc, \dots\}$
$(a b)^*$	
$(0 1)^*0$	

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Examples of Expressions— Solution

RE	Language
$a bc$	$\{a, bc\}$
$(b c)a$	$\{ba, ca\}$
$a \epsilon$	$\{a\}$
$a^* b$	$\{\epsilon, a, aa, aaa, aaaa, \dots\} \cup \{b\}$
ab^*	$\{a, ab, abb, abbb, abbbb, \dots\}$
$ab^* c^+$	$\{a, ab, abb, abbb, abbbb, \dots\} \cup \{c, cc, ccc, \dots\}$
$(a b)^*$	$\{\epsilon, a, b, aa, ab, ba, bb, aaa, aab, \dots\}$
$(0 1)^*0$	

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Examples of Expressions— Solution

RE

Language

a|bc

{a, bc}



*Concatenation has
higher precedence
over alternation |.*

(b|c)a

{ba, ca}

a €

{a}

Assignment Project Exam Help

a*|b

{€ , a, aa, aaa, aaaa, ...} ∪ {b}

Add WeChat powcoder

ab*

{a, ab, abb, abbb, abbbb, ...}

ab*|c⁺

{a, ab, abb, abbb, abbbb, ...} ∪ {c, cc, ccc, ...}

(a|b)*

{€, a, b, aa, ab, ba, bb, aaa, aab, ...}

(0|1)*0

binary numbers ending in 0

Regular Expressions for Programming Languages

Let *letter* stand for $A \mid B \mid C \mid \dots \mid Z$

Let *digit* stand for $0 \mid 1 \mid 2 \mid \dots \mid 9$

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Regular Expressions for Programming Languages

Let *letter* stand for $A \mid B \mid C \mid \dots \mid Z$

Let *digit* stand for $0 \mid 1 \mid 2 \mid \dots \mid 9$

integer constant: digit^+

identifier: $\text{letter}(\text{letter} \mid \text{digit})^*$

real constant: $\text{digit}^*.\text{digit}^+$

Formalisms for Lexical and Syntactic Analysis

Two issues in *Formal Languages*:

- Language Specification → formalism to describe what a valid program (word/sentence) looks like.
- Language Recognition → formalism to describe a machine and an algorithm that can verify that a program is valid or not.

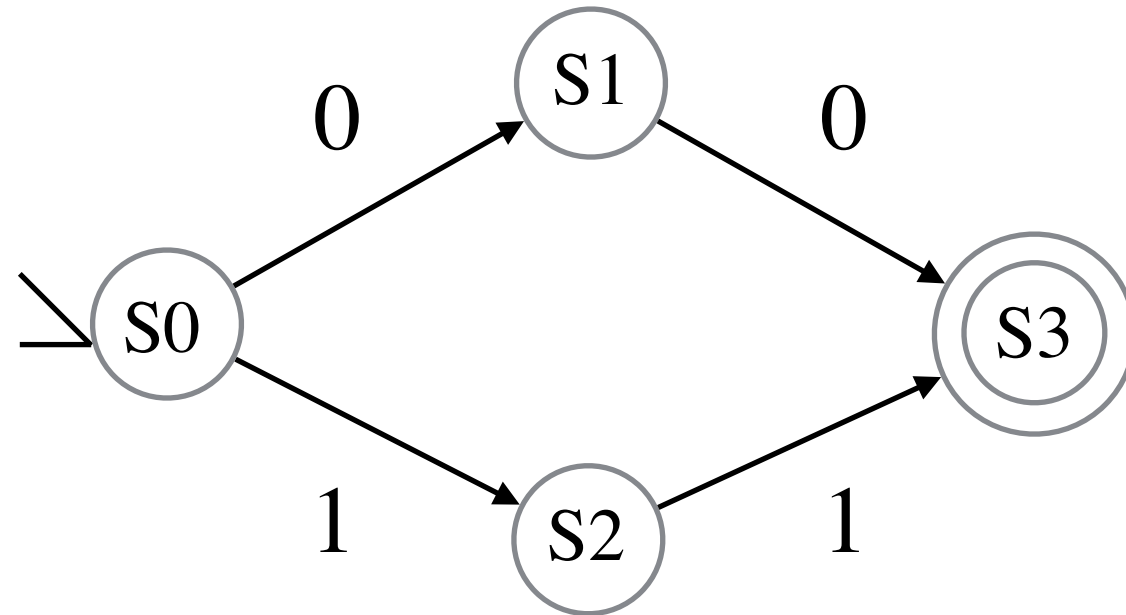
Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

We use **finite state automata** to recognize regular language

Finite State Automata



Assignment Project Exam Help

A Finite-State Automaton is a quadruple: $\langle S, s, F, T \rangle$

<https://powcoder.com>

- S is a set of states, e.g., $\{S0, S1, S2, S3\}$
- s is the start state, e.g., $S0$
- F is a set of final states, e.g., $\{S3\}$
- T is a set of labeled transitions, of the form
(state, input) \rightarrow state

formally,

$$S \times \Sigma \rightarrow S$$

Regular Expressions for Programming Languages

Let *letter* stand for $A \mid B \mid C \mid \dots \mid Z$

Let *digit* stand for $0 \mid 1 \mid 2 \mid \dots \mid 9$

integer constant: digit^+

identifier: $\text{letter}(\text{letter} \mid \text{digit})^*$

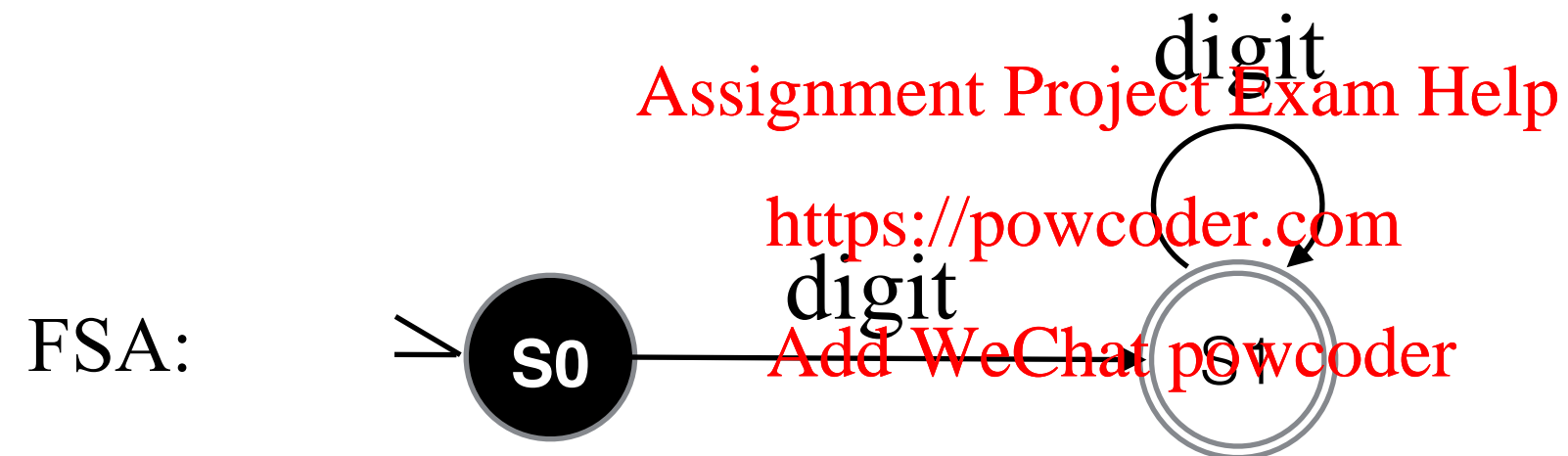
real constant: $\text{digit}^*.\text{digit}^+$

Recognizers for Regular Expressions

Example 1:

Integer Constant

RE: digit^+



Recognizers for Regular Expressions(Cont.)

Example 2:

Identifier

RE: letter(letter | digit)*

Assignment Project Exam Help

<https://powcoder.com>

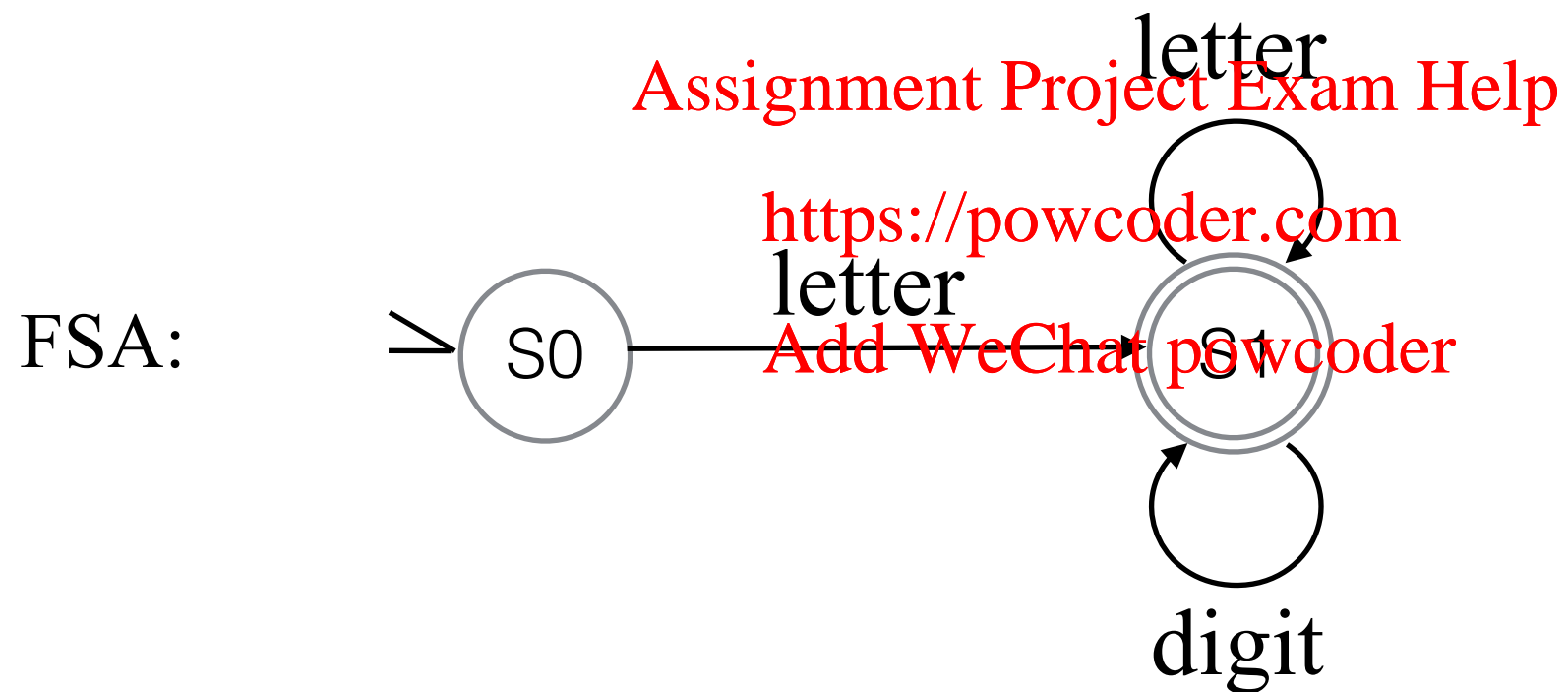
Add WeChat powcoder

Recognizers for Regular Expressions(Cont.)

Example 2:

Identifier

RE: letter(letter | digit)*



Recognizers for Regular Expressions(Cont.)

Example 3:

Real constant

RE: $\text{digit}^*.\text{digit}^+$

Assignment Project Exam Help

<https://powcoder.com>

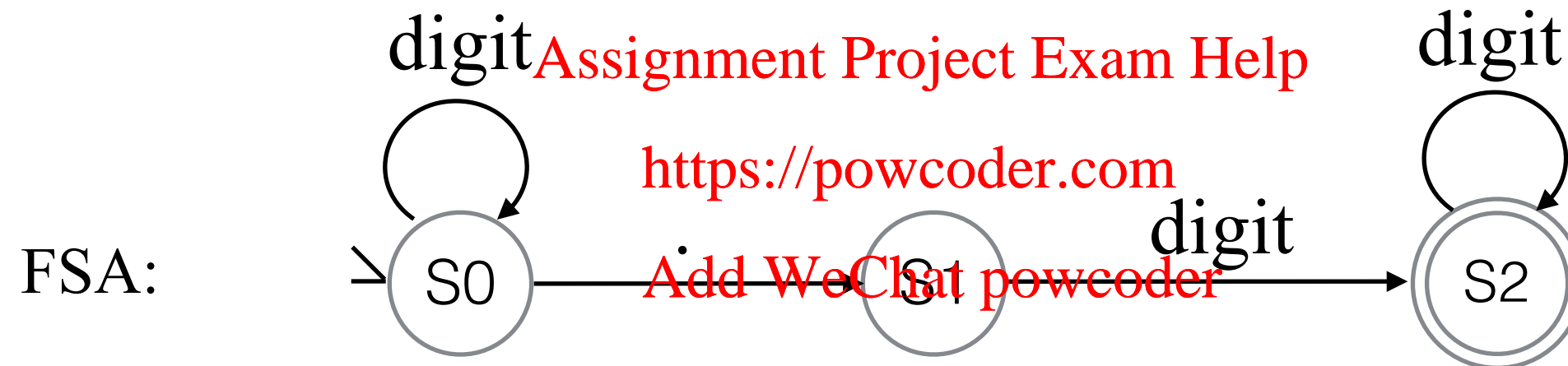
Add WeChat powcoder

Recognizers for Regular Expressions(Cont.)

Example 3:

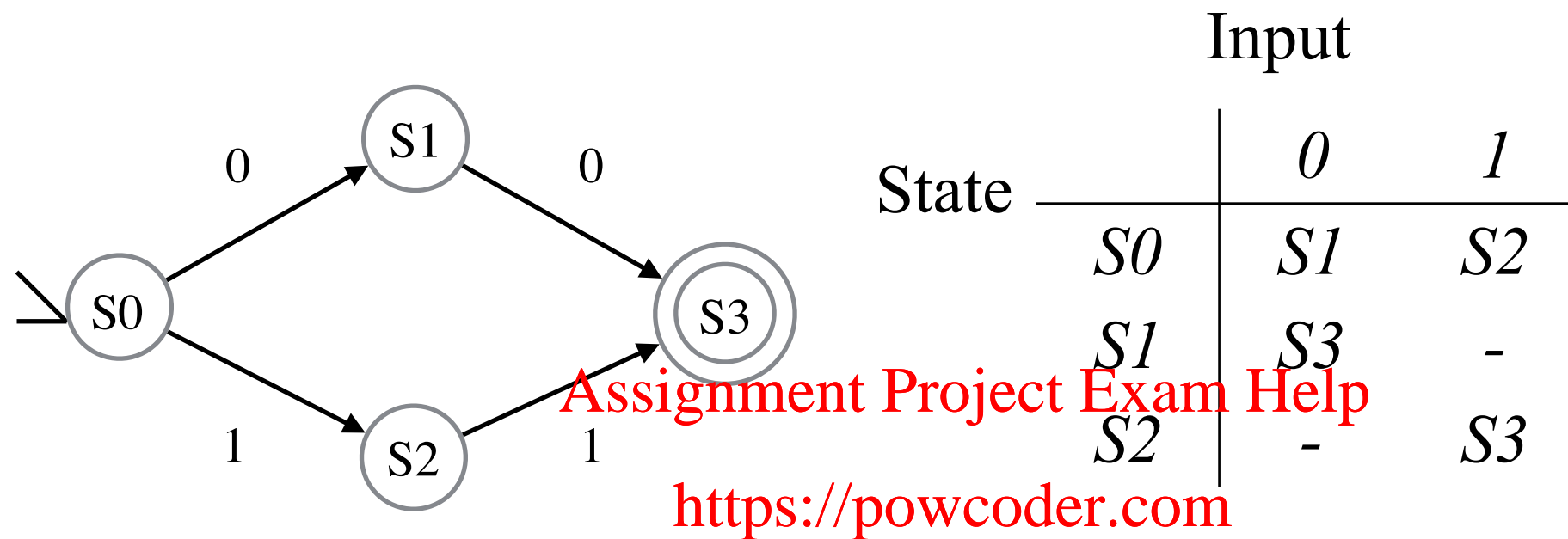
Real constant

RE: $\text{digit}^*.\text{digit}^+$



Finite State Automata

Transitions can be represented using a transition table:

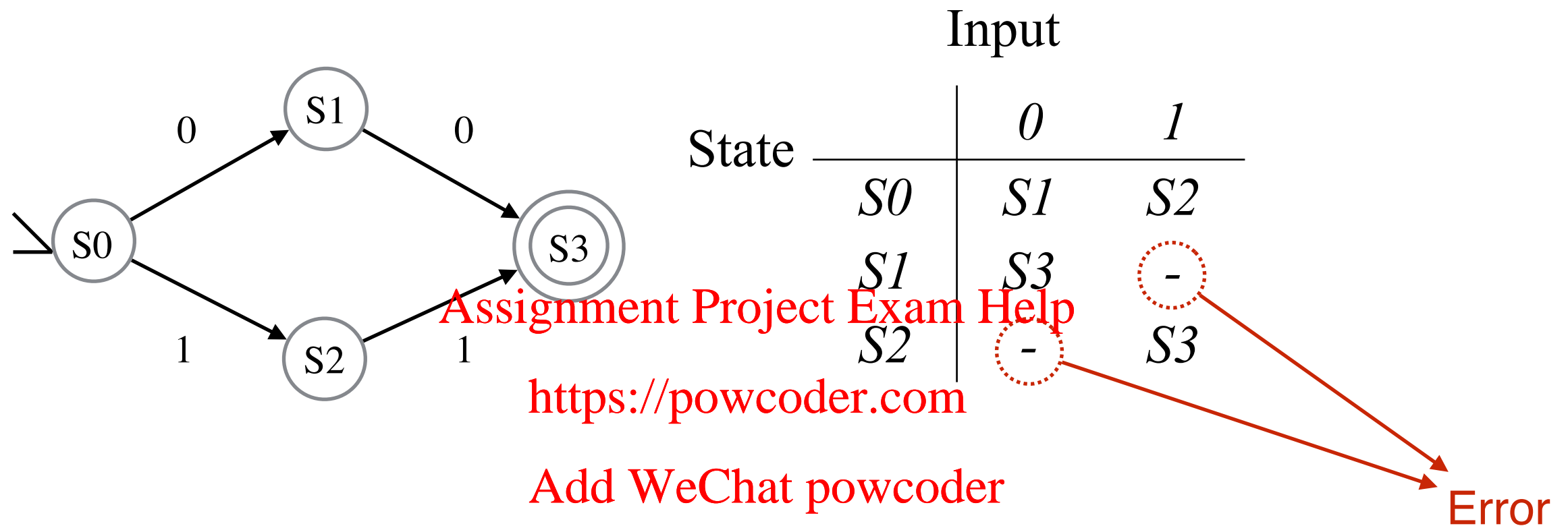


Add WeChat powcoder

An FSA *accepts* or *recognizes* an input string N **iff** there is some path from start state to a final state such that the labels on the path spell N.

Finite State Automata

Transitions can be represented using a transition table:



An FSA *accepts* or *recognizes* an input string **N** **iff** there is some path from start state to a final state such that the labels on the path spell N.

Lack of entry in the table (or no arc for a given character) indicates an *error—reject*.

Practical Recognizers

- Recognizer should be a deterministic finite automaton (DFA)
- Read until the end of a token
- Report errors (error recovery)

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Practical Recognizers

“identifier” regular expression:

$letter \rightarrow (a \mid b \mid c \mid \dots \mid z \mid A \mid B \mid C \mid \dots \mid Z)$

$digit \rightarrow (0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9)$

$id \rightarrow letter (letter \mid digit)^*$

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Practical Recognizers

“identifier” regular expression:

$letter \rightarrow (a \mid b \mid c \mid \dots \mid z \mid A \mid B \mid C \mid \dots \mid Z)$

$digit \rightarrow (0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9)$

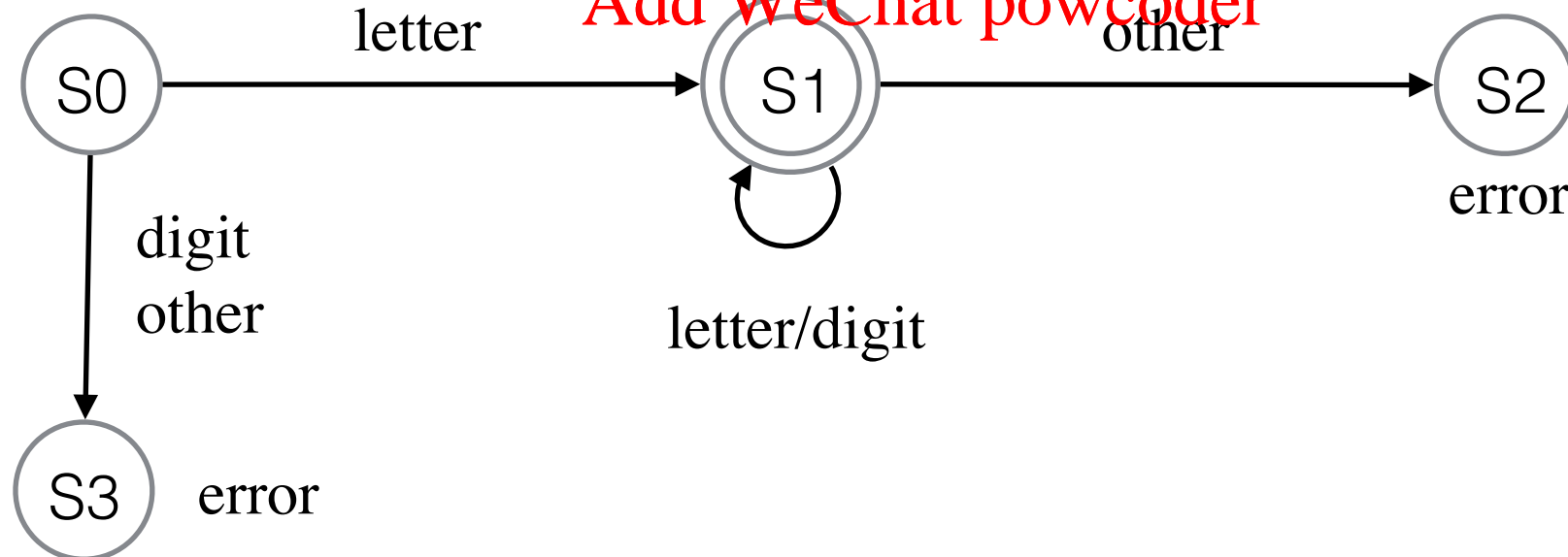
$id \rightarrow letter (letter \mid digit)^*$

Recognizer for “identifier”:

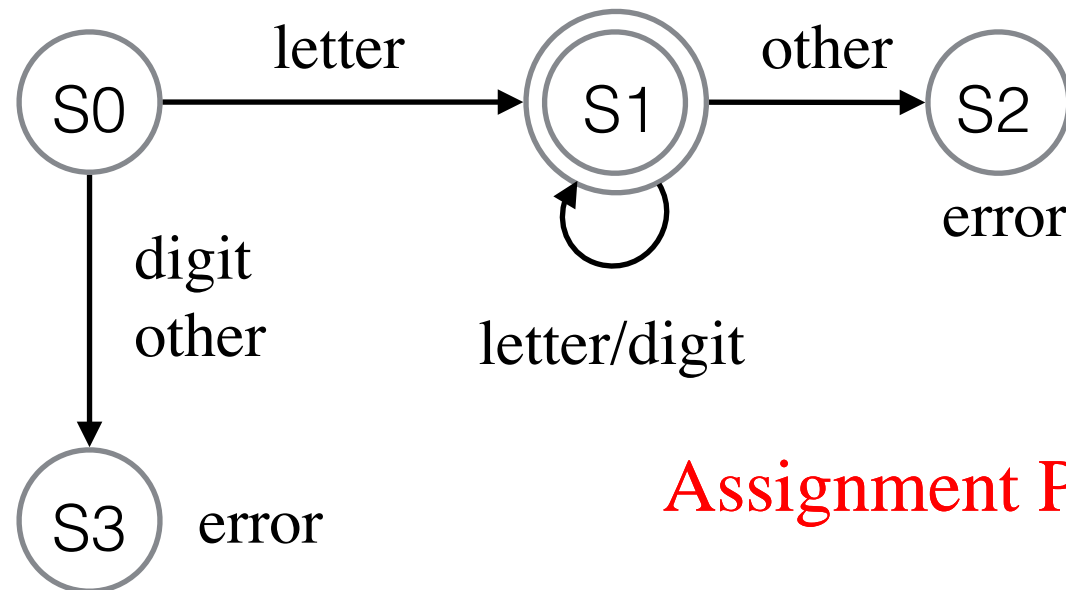
Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



Implementation: Code for the recognizer



```

char ← next_char();
state ← S0;
done ← false;
while( not done ) {
    class ← char_class[char];
    state ← next_state[class,state];
    switch(state) {
        case S1:
            /* building an id */
            token_value ← token_value + char;
            char ← next_char();
            if (char == DELIMITER)
                done = true;
            break;
        case S2: /* error state */
        case S3: /* error state */
            token type = error;
            done = true;
            break;
    }
}
return token_type;
  
```

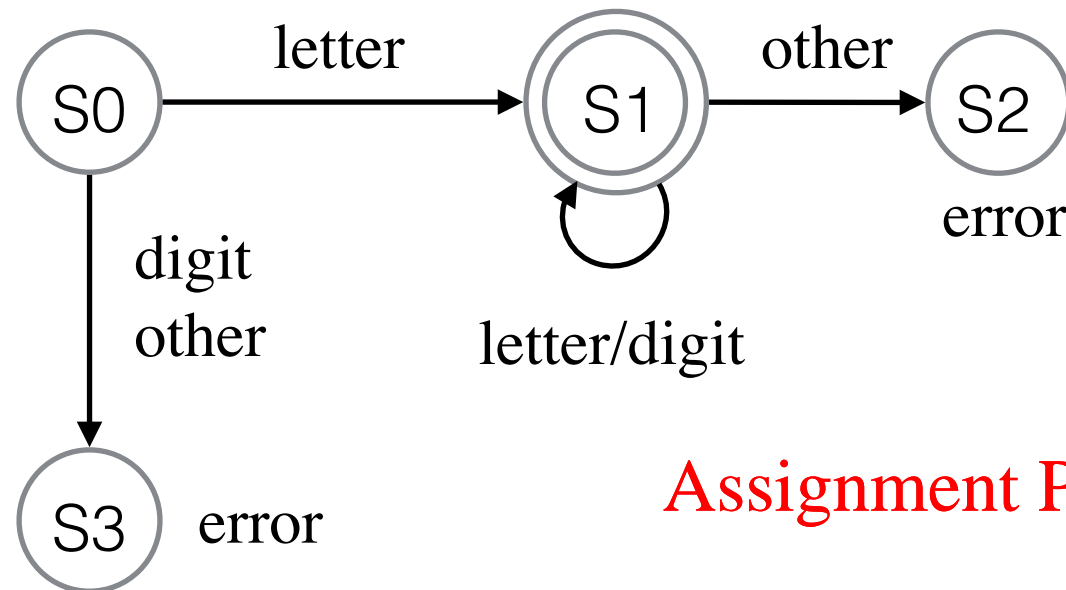
Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

<i>class</i>	<i>S0</i>	<i>S1</i>	<i>S2</i>	<i>S3</i>
<i>letter</i>	<i>S1</i>	<i>S1</i>	—	—
<i>digit</i>	<i>S3</i>	<i>S1</i>	—	—
<i>other</i>	<i>S3</i>	<i>S2</i>	—	—

Implementation: Code for the recognizer



```
char ← next_char();
```

```
state ← S0;
```

```
done ← false;
```

```
while( not done ) {
```

```
    class ← char_class[char];
```

```
    state ← next_state[class,state];
```

```
    switch(state) {
```

```
        case S1:
```

```
            /* building an id */
```

```
            token_value ← token_value + char;
```

```
            char ← next_char();
```

```
            if (char == DELIMITER)
```

```
                done = true;
```

```
            break;
```

```
        case S2: /* error state */
```

```
        case S3: /* error state */
```

```
            token type = error;
```

```
            done = true;
```

```
            break;
```

```
    }
```

```
}
```

```
return token_type;
```

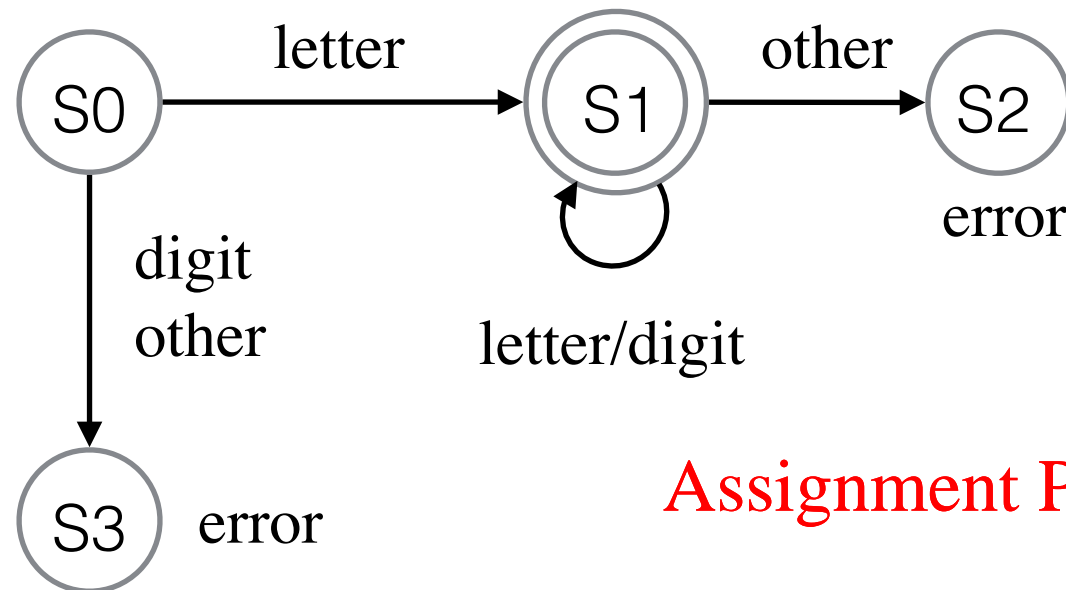
Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

<i>class</i>	<i>S0</i>	<i>S1</i>	<i>S2</i>	<i>S3</i>
<i>letter</i>	<i>S1</i>	<i>S1</i>	—	—
<i>digit</i>	<i>S3</i>	<i>S1</i>	—	—
<i>other</i>	<i>S3</i>	<i>S2</i>	—	—

Implementation: Code for the recognizer



```

char ← next_char();
state ← S0;
done ← false;
while( not done ) {
    class ← char_class[char];
    state ← next_state[class, state];
    switch(state) {

```

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat

powcoder

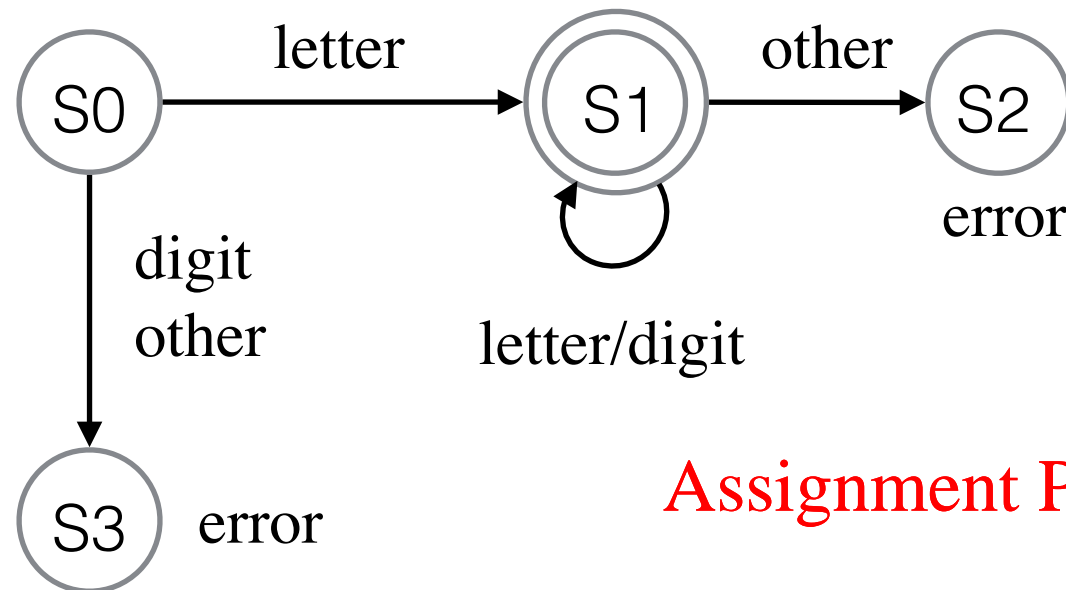
<i>class</i>	<i>S0</i>	<i>S1</i>	<i>S2</i>	<i>S3</i>
<i>letter</i>	<i>S1</i>	<i>S1</i>	—	—
<i>digit</i>	<i>S3</i>	<i>S1</i>	—	—
<i>other</i>	<i>S3</i>	<i>S2</i>	—	—

```

        case S1: /* building an id */
            token_value ← token_value + char;
            char ← next_char();
            if (char == DELIMITER)
                done = true;
            break;
        case S2: /* error state */
        case S3: /* error state */
            token type = error;
            done = true;
            break;
    }
}
return token_type;

```

Implementation: Code for the recognizer



```

char ← next_char();
state ← S0;
done ← false;
while( not done ) {
    class ← char_class[char];
    state ← next_state[class, state];
    switch(state) {
        case S1:
            /* building an id */
            token_value ← token_value + char;
            char ← next_char();
            if (char == DELIMITER)
                done = true;
            break;
        case S2: /* error state */
        case S3: /* error state */
            token type = error;
            done = true;
            break;
    }
}
return token_type;
  
```

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

<i>class</i>	<i>S0</i>	<i>S1</i>	<i>S2</i>	<i>S3</i>
<i>letter</i>	<i>S1</i>	<i>S1</i>	—	—
<i>digit</i>	<i>S3</i>	<i>S1</i>	—	—
<i>other</i>	<i>S3</i>	<i>S2</i>	—	—

Next Lecture

Things to do:

- Read Scott, Chapters 2.3 - 2.5

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder