تاج الدين سر الختم عبدالواحد عثمان 2023170793 Level 3 2 Section

احمد سفيان احمد حسن 2023170789 Level 3 1 Section

احمد ياسر محمد الفكي 2023170811 Level 3 1 Section

ايمن محجوب عبدالله محمد 2023170121 Level 3 1 Section

إسماعيل إبراهيم بخيت أحمد مسعود 2023170092 Level 3 1 Section

علا عبدالرحيم عبدالعزيز أحمد 2023170780 Level 3 4 Section

Table 1: TinyProgLang Token List & Regular Expressions

| Token Category | Token Type or "Lexeme" | Regular Expression or "RE" | Additional Notes |
|---|---|---|---|
| Literals | T_NUMBER | [0-9]+(\.[0-9]+)? | |
| | T_STRING | "[^"]*" | |
| Identifier | T_IDENTIFIER | ([a-z]\|[A-Z])([a-z]\|[A-Z]\|[0-9])* | Check against keyword list |
| Keywords | T_INT | int | These are literal strings. |
| | T_FLOAT | float | |
| | T_STRING | string | |
| | T_READ | read | |
| | T_WRITE | write | |
| | T_REPEAT | repeat | |
| | T_UNTIL | until | |
| | T_IF | if | |
| | T_ELSEIF | elseif | |
| | T_ELSE | else | |
| | T_THEN | then | |
| | T_RETURN | return | |
| | T_ENDL | endl | |
| | T_MAIN | main | |
| Double Char Operators | T_ASSIGN | := | |
| | T_NOTEQUAL | <> | |
| | T_AND | && | |
| | T_OR | \|\| | |
| Single Char Operators | T_PLUS | \+ | |
| | T_MINUS | - | |
| | T_TIMES | \* | |
| | T_DIVIDE | / | |
| | T_LESSTHAN | < | |
| | T_GREATERTHAN | > | |
| | T_EQUAL | = | |
| Symbols | T_LPAREN | \( | Left Parenthesis |
| | T_RPAREN | \) | Right Parenthesis |
| | T_LCURLY | { | Left Curly Brace |
| | T_RCURLY | } | Right Curly Brace |
| | T_COMMA | , | |
| | T_SEMICOLON | ; | |
| Ignored | COMMENT | /\*.*?\*/ | Discard this token. |
| | WHITESPACE | [ \t\n\r]+ *or* \s+ | Discard this token. |
| Special | T_EOF | N/A (Generated by scanner) | Virtualized; Signals the En |

# 1 Alphabet of the TinyProgLang

The alphabet of the TINY language, denoted by $\Sigma$, is the set of all valid input symbols from which strings in the language are formed.

$$\Sigma = \{\texttt{a-z}, \texttt{A-Z}, \texttt{0-9}, +, -, *, /, :, =, <, >, \&, |, ., ", (, ), \{, \}, ,, ;, \texttt{space}, \texttt{\textbackslash t}, \texttt{\textbackslash n}, \texttt{\textbackslash r}\}$$

The words in the language are composed of sequences of strings over $\Sigma$ . And the RE for the tokens in our language are built from this alphabet. The main character classes are enumerated as follows: a) Letters; defined by the RE [a-zA-Z]. and b) Digits, defined by the RE [0-9]; and d) Symbols, where each character is a literal (e.g., \+, ;); and e) Whitespace, defined by the RE [ \t\n\r]. The Language allows Comments matched with the RE (/* .*?  */). The Table 1 shows the full description of tokens TingProgLang does allow.

# Token Priority

When constructing a scanner, the order in which token rules are defined is critical. The scanner follows two fundamental principles to resolve ambiguity,

1. **The Longest Match Rule:** The scanner will always attempt to match the longest possible string of characters from the input. For example, given the input <>, the scanner will match the single token T_NOTEQUAL rather than two separate tokens, T_LESSTHAN < and T_GREATERTHAN >.

2. **The Priority Rule (First Match Wins):** Whenever two or more rules match the *same* longest string, the scanner will use the rule that was defined *first* in the priority list.

We apply the second rule to TinyProgLang by defining our tokens from most-specific to least-specific. Prioritizing keywords over identifiers and distinguishing both of these entities is one classical example ( i.e. with the input if which matches both the T_IF and the T_IDENTIFIER rules (RE [a-zA-Z][a-zA-Z0-9]*) with the same length.)

## Token Priority Order for TinyProgLang

We list the relevant tokens in order of priority, from highest to lowest. The scanner consults this priority list to decide which regular expression to use when it is taking input symbols.

## Token Priority Rationale

1. **Reserved Keywords:** These are the most specific literal strings. They must be listed before the general T_IDENTIFIER rule.

2. **Multi-Character Operators:** These must be listed before their single-character sub-parts to satisfy the Longest Match Rule.

3. **General Pattern Tokens:** These compose the match-all patterns in our TinyProgLang, and they must come after all keywords.

4. **Single-Character Operators & Symbols:** The relative order of these unit tokens does not matter. "Unit" here refers to the fact that they are composed of complete, separate and unambigious patterns.

5. **Ignored Patterns:** These are matched last and discarded Or by the scanner.

## Token Priority List And Their RE Mapping

From above rules we declare the enumerated list, and sign them into four priority classes (from P1 to P5, with P1 the highest) as follows:
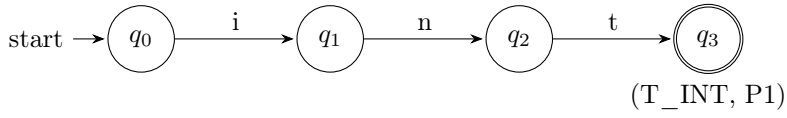{T_INT, T_FLOAT, T_STRING_KEY, T_READ, T_WRITE, T_REPEAT, T_UNTIL, T_IF, T_ELSEIF, T_ELSE, T_THEN, T_RETURN, T_ENDL, T_MAIN } { T_ASSIGN, T_NOTEQUAL, T_AND, T_OR} { T_NUMBER, T_IDENTIFIER, T_STRING } { T_PLUS, T_MINUS, T_TIMES, T_DIVIDE, T_LESSTHAN, T_GREATERTHAN, T_EQUAL, T_LPAREN, T_RPAREN, T_LCURLY, T_RCURLY, T_COMMA, T_SEMICOLON } {COMMENT, WHITESPACE}

Mapped to RE list:
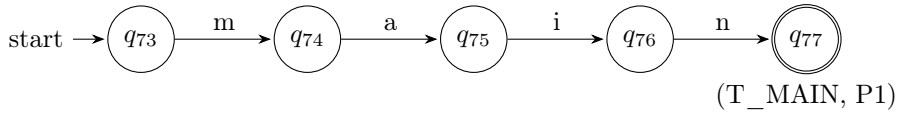{ (int), (float), (string), (read), (write), (repeat), (until), (if), (elseif), (else), (then), (return), (endl), (main) } { (:=), (<>), (&&), (||) } { ([0-9]+(\.[0-9]+)?), (([a-z]|[A-Z])([a-z]|[A-Z]|[0-9])*), ("[^"]*") } { (\+), (-), (\*), (/), (<), (>), (=), (\(), (\)), ({), (}), (,), (;) } { (/\*.*?\*/), ([ \t\n]+) }

From those two lists we attempt to Build our NFAs using Thompson Construction Method, which is used to build the subsequent DFA via the Subset Construction Method. The prioritized classes would then help us resolve conflicts label the end state for each sub-finite automata.

# Sub-NFAs for Keywords; Priority 1:

start $\rightarrow$ $q_0$ $\xrightarrow{\text{i}}$ $q_1$ $\xrightarrow{\text{n}}$ $q_2$ $\xrightarrow{\text{t}}$ (($q_3$))

(T_INT, P1)

NFA$_1$: RE = int

start $\rightarrow$ $q_4$ $\xrightarrow{\text{f}}$ $q_5$ $\xrightarrow{\text{l}}$ $q_6$ $\xrightarrow{\text{o}}$ $q_7$ $\xrightarrow{\text{a}}$ $q_8$ $\xrightarrow{\text{t}}$ (($q_9$))

(T_FLOAT, P1)

start $\rightarrow$ $q_{10}$ $\xrightarrow{\text{s}}$ $q_{11}$ $\xrightarrow{\text{t}}$ $q_{12}$ $\xrightarrow{\text{r}}$ $q_{13}$ $\xrightarrow{\text{i}}$ $q_{14}$ $\xrightarrow{\text{n}}$ $q_{15}$ $\xrightarrow{\text{g}}$ (($q_{16}$))

(T_STRING_KEY, P1)

start $\rightarrow$ $q_{17}$ $\xrightarrow{\text{r}}$ $q_{18}$ $\xrightarrow{\text{e}}$ $q_{19}$ $\xrightarrow{\text{a}}$ $q_{20}$ $\xrightarrow{\text{d}}$ (($q_{21}$))

(T_READ, P1)

start $\rightarrow$ $q_{22}$ $\xrightarrow{\text{w}}$ $q_{23}$ $\xrightarrow{\text{r}}$ $q_{24}$ $\xrightarrow{\text{i}}$ $q_{25}$ $\xrightarrow{\text{t}}$ $q_{26}$ $\xrightarrow{\text{e}}$ (($q_{27}$))

(T_WRITE, P1)

start $\rightarrow$ $q_{28}$ $\xrightarrow{\text{r}}$ $q_{29}$ $\xrightarrow{\text{e}}$ $q_{30}$ $\xrightarrow{\text{p}}$ $q_{31}$ $\xrightarrow{\text{e}}$ $q_{32}$ $\xrightarrow{\text{a}}$ $q_{33}$ $\xrightarrow{\text{t}}$ (($q_{34}$))

(T_REPEAT, P1)

start $\rightarrow$ $q_{35}$ $\xrightarrow{\text{u}}$ $q_{36}$ $\xrightarrow{\text{n}}$ $q_{37}$ $\xrightarrow{\text{t}}$ $q_{38}$ $\xrightarrow{\text{i}}$ $q_{39}$ $\xrightarrow{\text{l}}$ (($q_{40}$))

(T_UNTIL, P1)

start $\rightarrow$ $q_{41}$ $\xrightarrow{\text{i}}$ $q_{42}$ $\xrightarrow{\text{f}}$ (($q_{43}$))

(T_IF, P1)

start $\rightarrow$ $q_{44}$ $\xrightarrow{\text{e}}$ $q_{45}$ $\xrightarrow{\text{l}}$ $q_{46}$ $\xrightarrow{\text{s}}$ $q_{47}$ $\xrightarrow{\text{e}}$ $q_{48}$ $\xrightarrow{\text{i}}$ $q_{49}$ $\xrightarrow{\text{f}}$ (($q_{50}$))

(T_ELSEIF, P1)

start $\rightarrow$ $q_{51}$ $\xrightarrow{\text{e}}$ $q_{52}$ $\xrightarrow{\text{l}}$ $q_{53}$ $\xrightarrow{\text{s}}$ $q_{54}$ $\xrightarrow{\text{e}}$ (($q_{55}$))

(T_ELSE, P1)

start $\rightarrow$ $q_{56}$ $\xrightarrow{\text{t}}$ $q_{57}$ $\xrightarrow{\text{h}}$ $q_{58}$ $\xrightarrow{\text{e}}$ $q_{59}$ $\xrightarrow{\text{n}}$ (($q_{60}$))

(T_THEN, P1)

start $\rightarrow$ $q_{61}$ $\xrightarrow{\text{r}}$ $q_{62}$ $\xrightarrow{\text{e}}$ $q_{63}$ $\xrightarrow{\text{t}}$ $q_{64}$ $\xrightarrow{\text{u}}$ $q_{65}$ $\xrightarrow{\text{r}}$ $q_{66}$ $\xrightarrow{\text{n}}$ (($q_{67}$))

(T_RETURN, P1)

start $\rightarrow$ $q_{68}$ $\xrightarrow{\text{e}}$ $q_{69}$ $\xrightarrow{\text{n}}$ $q_{70}$ $\xrightarrow{\text{d}}$ $q_{71}$ $\xrightarrow{\text{l}}$ (($q_{72}$))

(T_ENDL, P1)

start $\rightarrow$ $q_{73}$ $\xrightarrow{\text{m}}$ $q_{74}$ $\xrightarrow{\text{a}}$ $q_{75}$ $\xrightarrow{\text{i}}$ $q_{76}$ $\xrightarrow{\text{n}}$ (($q_{77}$))

(T_MAIN, P1)

# Sub-NFAs for Double Char Operators; Priority 2:

start $\rightarrow$ $q_{78}$ $\xrightarrow{:}$ $q_{79}$ $\xrightarrow{=}$ $q_{80}$

(T_ASSIGN, P2)

start $\rightarrow$ $q_{81}$ $\xrightarrow{<}$ $q_{82}$ $\xrightarrow{>}$ $q_{83}$

(T_NOTEQUAL, P2)

start $\rightarrow$ $q_{84}$ $\xrightarrow{\&}$ $q_{85}$ $\xrightarrow{\&}$ $q_{86}$

(T_AND, P2)

start $\rightarrow$ $q_{87}$ $\xrightarrow{|}$ $q_{88}$ $\xrightarrow{|}$ $q_{89}$

(T_OR, P2)

# Sub-NFAs for Literals and Identifiers; Priority 3:

start $\rightarrow$ $q_{90}$ $\xrightarrow{\epsilon}$ $q_{91}$ $\xrightarrow{[0-9]}$ $q_{92}$ $\xrightarrow{\epsilon}$ $q_{93}$ $\xrightarrow{\epsilon}$ $q_{99}$

(T_NUMBER, P3)

$q_{94}$ $\xrightarrow{"."}$ $q_{95}$ $\xrightarrow{\epsilon}$ $q_{96}$ $\xrightarrow{[0-9]}$ $q_{97}$ $\xrightarrow{\epsilon}$ $q_{98}$

start $\rightarrow$ $q_{100}$

$q_{101}$ $\xrightarrow{[a-z]}$ $q_{102}$

$q_{103}$ $\xrightarrow{[A-Z]}$ $q_{104}$

$q_{105}$ $\xrightarrow{\epsilon}$ $q_{106}$ $\xrightarrow{\epsilon}$ $q_{107}$

$q_{108}$ $\xrightarrow{[a-z]}$ $q_{109}$

$q_{110}$ $\xrightarrow{[A-Z]}$ $q_{111}$

$q_{112}$ $\xrightarrow{[0-9]}$ $q_{113}$

$q_{114}$ $\xrightarrow{\epsilon}$ $q_{115}$

(T_IDENTIFIER, P3)

start $\rightarrow$ $q_{116}$ $\xrightarrow{"}$ $q_{117}$ $\xrightarrow{\epsilon}$ $q_{118}$ $\xrightarrow{[\char94 "]}$ $q_{119}$ $\xrightarrow{\epsilon}$ $q_{120}$ $\xrightarrow{"}$ $q_{121}$

(T_STRING, P3)

# Sub-NFAs for Single Char Operators; Priority 4:

start $\to$ $q_{121}$ $\xrightarrow{+}$ $q_{122}$

(T_PLUS, P4)

start $\to$ $q_{123}$ $\xrightarrow{-}$ $q_{124}$

(T_MINUS, P4)

start $\to$ $q_{125}$ $\xrightarrow{*}$ $q_{126}$

(T_TIMES, P4)

start $\to$ $q_{127}$ $\xrightarrow{/}$ $q_{128}$

(T_DIVIDE, P4)

start $\to$ $q_{129}$ $\xrightarrow{<}$ $q_{130}$

(T_LESSTHAN, P4)

start $\to$ $q_{131}$ $\xrightarrow{>}$ $q_{132}$

(T_GREATERTHAN, P4)

start $\to$ $q_{133}$ $\xrightarrow{=}$ $q_{134}$
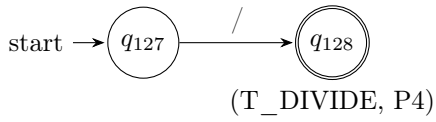
(T_EQUAL, P4)

start $\to$ $q_{135}$ $\xrightarrow{(}$ $q_{136}$

(T_LPAREN, P4)

start $\to$ $q_{137}$ $\xrightarrow{)}$ $q_{138}$

(T_RPAREN, P4)

start $\to$ $q_{139}$ $\xrightarrow{\{}$ $q_{140}$

(T_LCURLY, P4)

start $\to$ $q_{141}$ $\xrightarrow{\}}$ $q_{142}$

(T_RCURLY, P4)

start $\to$ $q_{143}$ $\xrightarrow{,}$ $q_{144}$

(T_COMMA, P4)

start $\to$ $q_{145}$ $\xrightarrow{;}$ $q_{146}$

(T_SEMICOLON, P4)

# Ignored Tokens; Priority 5:

start $\rightarrow$ $q_{147}$ $\xrightarrow{/}$ $q_{148}$ $\xrightarrow{*}$ $q_{149}$ $\xrightarrow{\epsilon}$ $q_{152}$ $\xrightarrow{*}$ $q_{153}$ $\xrightarrow{/}$ $q_{154}$

(COMMENT, P5)

$q_{149}$ $\xrightarrow{\epsilon}$ $q_{150}$ $\xrightarrow{.}$ $q_{151}$ $\xrightarrow{\epsilon}$ $q_{152}$

$q_{151}$ $\xrightarrow{\epsilon}$ $q_{150}$

start $\rightarrow$ $q_{155}$ $\xrightarrow{\epsilon}$ $q_{156}$ $\xrightarrow{[ \backslash t\backslash n\backslash r]}$ $q_{157}$ $\xrightarrow{\epsilon}$ $q_{158}$

(WHITESPACE, P5)

$q_{157}$ $\xrightarrow{\epsilon}$ $q_{156}$

[ \t\n\r]

Any Char

start $\rightarrow$ A $\xrightarrow{\text{Other Chars}}$ B

(WHITESPACE, P5)    TRAP

# Constructed and Minimized DFAs

## Keywords; Priority 1

(T_INT, P1)

start → (A) --i--> (B) --n--> (C) --t--> ((D))

A --n,t--> E
B --i,t--> E
C --i,n--> E
D --i,n,t--> E

E
TRAP
i,n,t

(T_FLOAT, P1)

start → (A) --f--> (B) --l--> (C) --o--> (D) --a--> (E) --t--> ((F))

A --l,o,a,t--> G
B --f,o,a,t--> G
C --f,l,a,t--> G
D --f,l,o,t--> G
E --f,l,o,a--> G
F --f,l,o,a,t--> G

G
TRAP
f,l,o,a,t

(T_STRING_KEY, P1)

start → (A) --s--> (B) --t--> (C) --r--> (D) --i--> (E) --n--> (F) --g--> ((G))

A --t,r,i,n,g--> H
B --s,r,i,n,g--> H
C --s,t,i,n,g--> H
D --s,t,r,n,g--> H
E --s,t,r,i,g--> H
F --s,t,r,i,n--> H
G --s,t,r,i,n,g--> H

H
TRAP
s,t,r,i,n,g

(T_READ, P1)

start → (A) --r--> (B) --e--> (C) --a--> (D) --d--> ((E))

A --e,a,d--> F
B --r,a,d--> F
C --r,e,d--> F
D --r,e,a--> F
E --r,e,a,d--> F

F
TRAP
r,e,a,d

(T_WRITE, P1)

start → A --w--> B --r--> C --i--> D --t--> E --e--> F

w,i,t,e  w,r,t,e  w,r,i,e  w,r,i,t  w,r,i,t,e

r,i,t,e

G
TRAP
w,r,i,t,e

(T_REPEAT, P1)

start → A --r--> B --e--> C --p--> D --e--> E --a--> F --t--> G

r,e,a,t  r,p,a,t  r,e,p,t  r,e,p,a  r,e,p,a,t

r,p,a,t

e,p,a,t

H
TRAP
r,e,p,a,t

(T_UNTIL, P1)

start → A --u--> B --n--> C --t--> D --i--> E --l--> F

u,t,i,l  u,n,i,l  u,n,t,i  u,n,t,i  u,n,t,i,l

n,t,i,l

G
TRAP
u,n,t,i,l

start → A --i--> B --f--> C  (T_IF, P1)

f  i  i,f

D
TRAP
i,f

(T_ELSEIF, P1)

start → A --e--> B --l--> C --s--> D --e--> E --i--> F --f--> G

e,l,i,f  e,l,s,i,f  e,l,s,f  e,l,s,e,i  e,l,s,e,i,f

e,s,i,f

l,s,i,f

I
TRAP
e,l,s,i,f

9

**(T_ELSE, P1)**

start → A —e→ B —l→ C —s→ D

A —l,s→ E
B —e,s→ E
C —e,l→ E
D —e,l,s→ E
E —e,l,s→ E (self loop)

E TRAP

---

**(T_THEN, P1)**

start → A —t→ B —h→ C —e→ D —n→ E

A —h,e,n→ F
B —t,e,n→ F
C —t,h,n→ F
D —t,h,e→ F
E —t,h,e,n→ F
F —T,h,e,n→ F (self loop)

F TRAP

---

**(T_RETURN, P1)**

start → A —r→ B —e→ C —t→ D —u→ E —r→ F

A —e,t,u,r,n→ G
B —r,t,u,r,n→ G
C —r,e,u,r,n→ G
D —r,e,t,r,n→ G
E —r,e,t,u,n→ G
F —r,e,t,u,r,n→ G
G —r,e,t,u,r,n→ G (self loop)

G TRAP

---

**(T_ENDL, P1)**

start → A —e→ B —n→ C —d→ D —l→ E

A —n,d,l→ F
B —e,d,l→ F
C —e,n,l→ F
D —e,n,d→ F
E —e,n,d,l→ F
F —e,n,d,l→ F (self loop)

F TRAP

---

**(T_MAIN, P1)**

start → A —m→ B —a→ C —i→ D —n→ E

A —a,i,n→ F
B —m,i,n→ F
C —m,a,n→ F
D —m,a,i→ F
E —m,a,i,n→ F
F —m,a,i,n→ F (self loop)

F TRAP

10

## Double Char Operators;Priority 2

(T_ASSIGN, P2)

start → (A) —:→ (B) —=→ ((C))

A →=, <>, &, | → (D)

B →:, <>, &, | → (D)

C →:, =, <>, &, | → (D)

(D) TRAP
:, =, <>, &, |

(T_NOTEQUAL, P2)

start → (A) —<→ (B) —>→ ((C))

A →>, :, =, &, | → (D)

B →<, :, =, &, | → (D)

C →<, >, :, =, &, | → (D)

(D) TRAP
<, >, :, =, &, |

(T_AND, P2)

start → (A) —&→ (B) —&→ ((C))

A →<, >, :, =, | → (D)

B →<, >, :, =, | → (D)

C →&, <, >, :, =, | → (D)

(D) TRAP
&, <, >, :, =, |

(T_OR, P2)

start → (A) —|→ (B) —|→ ((C))

A →<, >, :, =, & → (D)

B →<, >, :, =, & → (D)

C →|, <, >, :, =, & → (D)

(D) TRAP
|, <, >, :, =, &

en

11

## Literals And Identifiers ; Priority 3



## String Recognition; Priority 3

## Single Char Operators ;Priority 4

start → ( A )  — + →  (( B ))  — Other Chars →  ( C ) ↺ Any Char

(T_PLUS, P4)     TRAP

start → ( A )  — - →  (( B ))  — Other Chars →  ( C ) ↺ Any Char

(T_MINUS, P4)     TRAP

start → ( A )  — * →  (( B ))  — Other Chars →  ( C ) ↺ Any Char

(T_TIMES, P4)     TRAP

start → ( A )  — / →  (( B ))  — Other Chars →  ( C ) ↺ Any Char

(T_DIVIDE, P4)     TRAP

start → ( A )  — < →  (( B ))  — Other Chars →  ( C ) ↺ Any Char

(T_LESSTHAN, P4)     TRAP

start → ( A )  — > →  (( B ))  — Other Chars →  ( C ) ↺ Any Char

(T_GREATERTHAN, P4)     TRAP

start → ( A )  — = →  (( B ))  — Other Chars →  ( C ) ↺ Any Char

(T_EQUAL, P4)     TRAP

start → ( A )  — ( →  (( B ))  — Other Chars →  ( C ) ↺ Any Char

(T_LPAREN, P4)     TRAP

start → ( A )  — ) →  (( B ))  — Other Chars →  ( C ) ↺ Any Char

(T_RPAREN, P4)     TRAP

Any Char

start → ( A )  --{-->  (( B ))  --Other Chars-->  ( C )  ⟲ Any Char

(T_LCURLY, P4)    TRAP

Any Char

start → ( A )  --}-->  (( B ))  --Other Chars-->  ( C )  ⟲ Any Char

(T_RCURLY, P4)    TRAP

Any Char

start → ( A )  --,-->  (( B ))  --Other Chars-->  ( C )  ⟲ Any Char

(T_COMMA, P4)    TRAP

Any Char

start → ( A )  --;-->  (( B ))  --Other Chars-->  ( C )  ⟲ Any Char

(T_SEMICOLON, P4)    TRAP

---

[1]Two labels on the transition state, namely "Any Char" and "Other char" are not formal REs in TinyProgLang but are shorthand conventions used to satisfy requirement of the DFAs:

† **Any Char:** Shorthand for $\Sigma$.

‡ **Other Chars:** Shorthand for $\Sigma$ minus the single successful symbol expected from input (e.g., $\Sigma \setminus \{+\}$ at start state of `T_PLUS`).

# Ignored Tokens; Priority 5

(COMMENT, P5)

start → (A) —/→ (B) —*→ (C) —*→ (D) —/→ ((F))

C —*, .→ G

G TRAP

/, *, .

icture

[ \t\n\r]

start → ((A)) —Other Chars→ (B)

(WHITESPACE, P5)

TRAP

Any Char