

第 5 組 系統程式-期末報告

分工

1112914 宋冠穎 (25 %)

- Regular Expression -> Dotted Item -> NFA
- NFA -> DFA
- 使用 LaTeX 將 Dotted Item 手寫過程整理成數位檔案
- Demo
- 簡報製作

1112919 丁柏凱 (25 %)

- NFA -> DFA
- 使用 LaTeX 將 NFA -> DFA 手寫過程整理成數位檔案
- Python Program
- Demo
- inkscape 繪圖
- 簡報製作

1112929 詹佳瑜 (10 %)

- 想測資
- 檢查輸入 & 輸出

1112946 王嘉汶 (18 %)

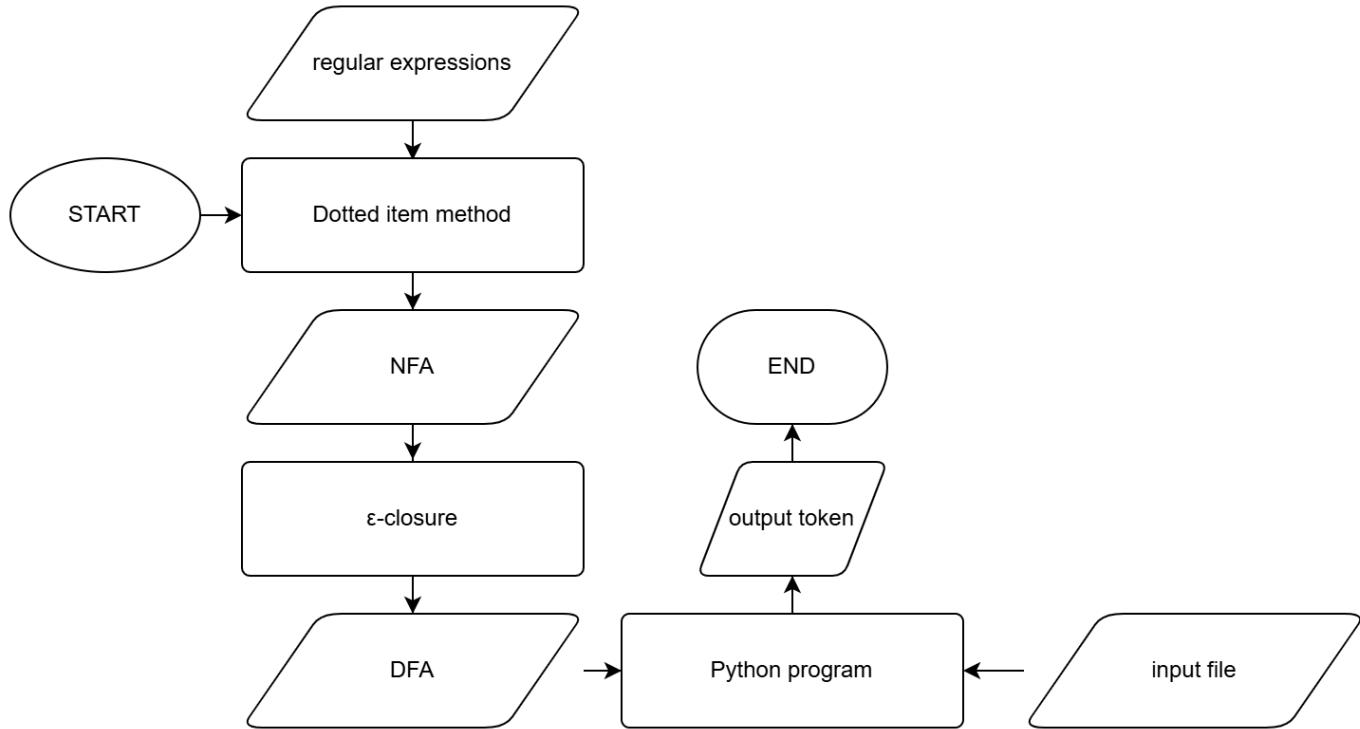
- inkscape 繪圖
- Demo
- 簡報製作

1112957 顏琦恩 (22 %)

- Regular Expression -> Dotted Item -> NFA
- NFA -> DFA

- 簡報製作
- 影片剪輯

題目 & 流程圖



1. 利用 dotted item method 把以下的 regular expressions 轉 NFA/DFA

underscore → _

letter → A | B ... | Z | a | b ... | z

digit → 0 | 1 | 2 ... | 9

underscore_tail → underscore (letter | digit)⁺

id → letter (letter | digit)* underscore_tail*

fraction → · digit⁺

optional_exponent → (E (+ | - | ε) digit⁺) | ε

integer → digit⁺

float → digit* fraction optional_exponent

2. 寫一個 lexical analyzer 來分析讀到的是哪個 token (Separators and Brackets, Binary Operators, Keywords, Comments, Identity, Integer and Float)

Separators and Brackets: the language uses ';;', '(, ')', '[,]', '{, }' as the separators or brackets.

Binary Operators: the language uses the following operators +, -, *, /, <=, >=, <>, <, >, =, ==.

Keywords: the language uses the following keywords int, float, bool, void, while, if, else, for, return.

Comments: the language supports C style comments. A comment starts with /* and end with */ with anything in between except a comment ending sequence (*/).

Layout: the language uses the following layout symbols: any number of blank, tabulation, carriage return and line feed.

注意: 不能用 lex 或 similar lexical analyzer

Regular Expression -> NFA

將 Regular Expression 利用 dotted item method 轉換成 NFA

NFA:

- 可以對同一個輸入字元有多個轉移方向
- 可以存在 ϵ -transition (空轉移，不需輸入符號即可移動)

dotted item method

- 定義
 - dotted item 是一個帶有「dot .」的產生式 (production)，標記目前已經分析到哪個位置。點號位於產生式右側的某個位置，表示目前已經「看到」這個位置前的所有符號，下一步將會分析點號後的符號。

- 範例
 - 一個正規表達式

$A \rightarrow XYZ$

- 會有以下 dotted items

- a. $A \rightarrow \cdot XYZ$: 表示尚未處理任何符號
- b. $A \rightarrow X \cdot YZ$: 表示已處理了 X，下一步要處理 Y
- c. $A \rightarrow XY \cdot Z$: 表示已處理了 X 和 Y
- d. $A \rightarrow XYZ \cdot$: 表示整個產生式已處理完

- Character moves

- 基本項目 (Basic item)

- 如果項目中的點 (dot) 位於基本項目前方，則在輸入字元與基本項目相符時，該項目會將點移動到下一個字元位置。

- 非基本項目 (Non-basic item)

- 項目中的點無法直接跨過一個字元，因為它沒有定義一組可比對的字元集合（無法判斷輸入字元是否符合）。
- 如 * 、 + 、 ? 、 |

dotted item 過程為簡化狀態機有進行部分簡化

- 一般 dotted item 中 Non-basic item $\varepsilon - move$ 規則

$$T \rightarrow \alpha \bullet (R)^* \beta \Rightarrow T \rightarrow \alpha (R)^* \bullet \beta \\ T \rightarrow \alpha (\bullet R)^* \beta$$

$$T \rightarrow \alpha (R \bullet)^* \beta \Rightarrow T \rightarrow \alpha (R)^* \bullet \beta \\ T \rightarrow \alpha (\bullet R)^* \beta$$

$$T \rightarrow \alpha \bullet (R)^+ \beta \Rightarrow T \rightarrow \alpha (\bullet R)^+ \beta$$

$$T \rightarrow \alpha (R \bullet)^+ \beta \Rightarrow T \rightarrow \alpha (R)^+ \bullet \beta \\ T \rightarrow \alpha (\bullet R)^+ \beta$$

- 簡化範例: (2) 讀完 letter 後，移到 (4) 表示括號內容讀取完畢，與移到 (1) 表示再次讀取括號內容。

$$\textcircled{1} \Rightarrow \textcircled{1} \text{ underscore_tail} \rightarrow \text{underscore} \bullet (letter \mid digit)^+$$

$$\textcircled{1} \Rightarrow \textcircled{2} \text{ underscore_tail} \rightarrow \text{underscore} (\bullet letter \mid digit)^+$$

$$\textcircled{1} \Rightarrow \textcircled{3} \text{ underscore_tail} \rightarrow \text{underscore} (letter \mid \bullet digit)^+$$

$$\textcircled{2} \Rightarrow \textcircled{4} \text{ underscore_tail} \rightarrow \text{underscore} (letter \mid digit)^+ \bullet <<< \text{recognized}$$

$$\textcircled{2} \Rightarrow \textcircled{1} \text{ underscore_tail} \rightarrow \text{underscore} \bullet (letter \mid digit)^+$$

1. underscore

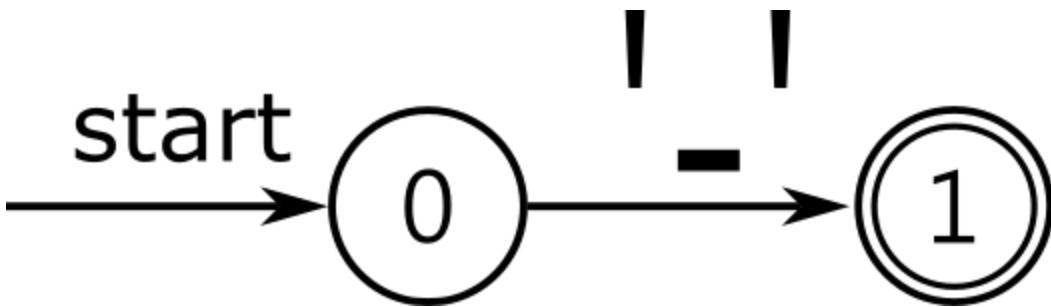
underscore有可能為

- ①:接下來要匹配'_'
- ②:已經完成匹配'_'

underscore → _

① *underscore* → •_

② *underscore* → _• <<<< recognized



2. letter

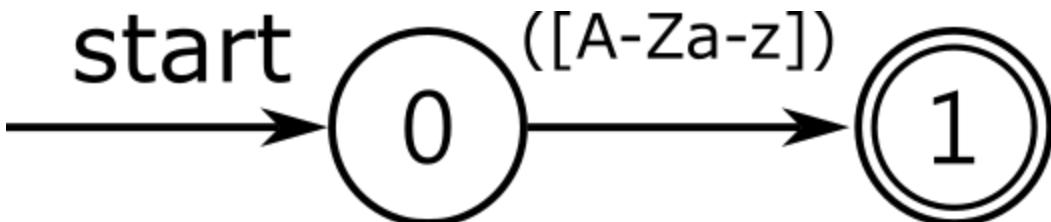
letter有可能為

- ①接下來要匹配[A-Za-z]
- ②已經完成匹配[A-Za-z]

letter → [A – Za – z]

① *letter* → •[A – Za – z]

② ⇒ ① *letter* → [A – Za – z]•

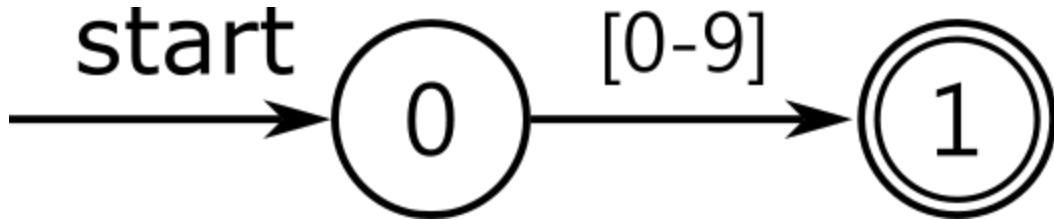


3. digit

$digit \rightarrow [0 - 9]$

① $digit \rightarrow \bullet[0 - 9]$

① \Rightarrow ① $digit \rightarrow [0 - 9]\bullet <<< \text{recognized}$



4. underscore_tail

$underscore_tail \rightarrow underscore(letter \mid digit)^+$

① $underscore_tail \rightarrow \bullet underscore(letter \mid digit)^+$

① \Rightarrow ① $underscore_tail \rightarrow underscore \bullet (letter \mid digit)^+$

② \Rightarrow ② $underscore_tail \rightarrow underscore(\bullet letter \mid digit)^+$

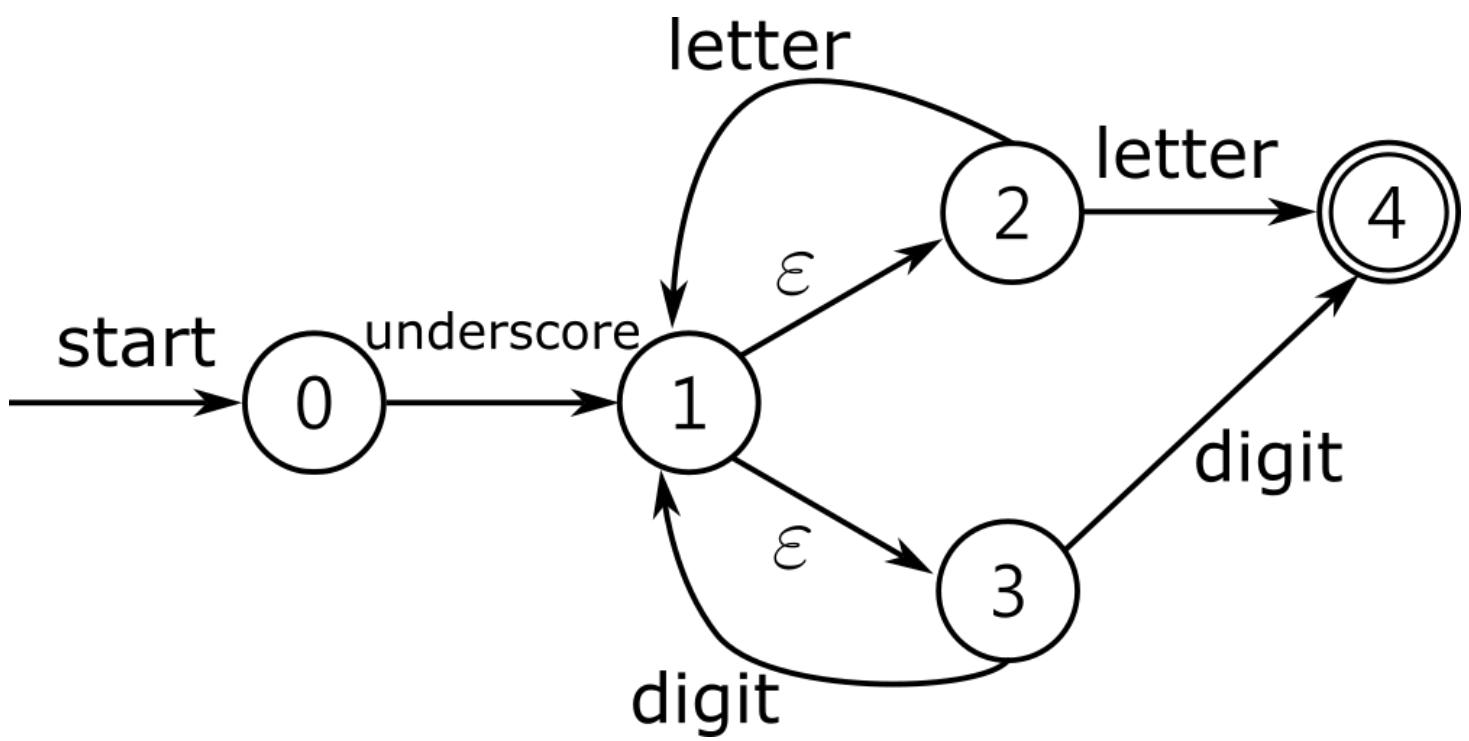
① \Rightarrow ③ $underscore_tail \rightarrow underscore(letter \mid \bullet digit)^+$

② \Rightarrow ④ $underscore_tail \rightarrow underscore(letter \mid digit)^+ \bullet <<< \text{recognized}$

② \Rightarrow ① $underscore_tail \rightarrow underscore \bullet (letter \mid digit)^+$

③ \Rightarrow ④ $underscore_tail \rightarrow underscore(letter \mid digit)^+ \bullet <<< \text{recognized}$

③ \Rightarrow ① $underscore_tail \rightarrow underscore \bullet (letter \mid digit)^+$



5. id

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

$$\textcircled{0} \quad id \rightarrow \bullet letter(letter \mid digit)^*(underscore_tail)$$

$$\textcircled{0} \Rightarrow \textcircled{1} \quad id \rightarrow letter \bullet (letter \mid digit)^*(underscore_tail)$$

$$\textcircled{1} \Rightarrow \textcircled{2} \quad id \rightarrow letter(\bullet letter \mid digit)^*(underscore_tail)$$

$$\textcircled{1} \Rightarrow \textcircled{3} \quad id \rightarrow letter(letter \mid \bullet digit)^*(underscore_tail)$$

$$\textcircled{1} \Rightarrow \textcircled{4} \quad id \rightarrow letter(letter \mid digit)^* \bullet (underscore_tail)$$

$$\textcircled{2} \Rightarrow \textcircled{4} \quad id \rightarrow letter(letter \mid digit)^* \bullet (underscore_tail)$$

$$\textcircled{2} \Rightarrow \textcircled{1} \quad id \rightarrow letter \bullet (letter \mid digit)^*(underscore_tail)$$

$$\textcircled{3} \Rightarrow \textcircled{4} \quad id \rightarrow letter(letter \mid digit)^* \bullet (underscore_tail)$$

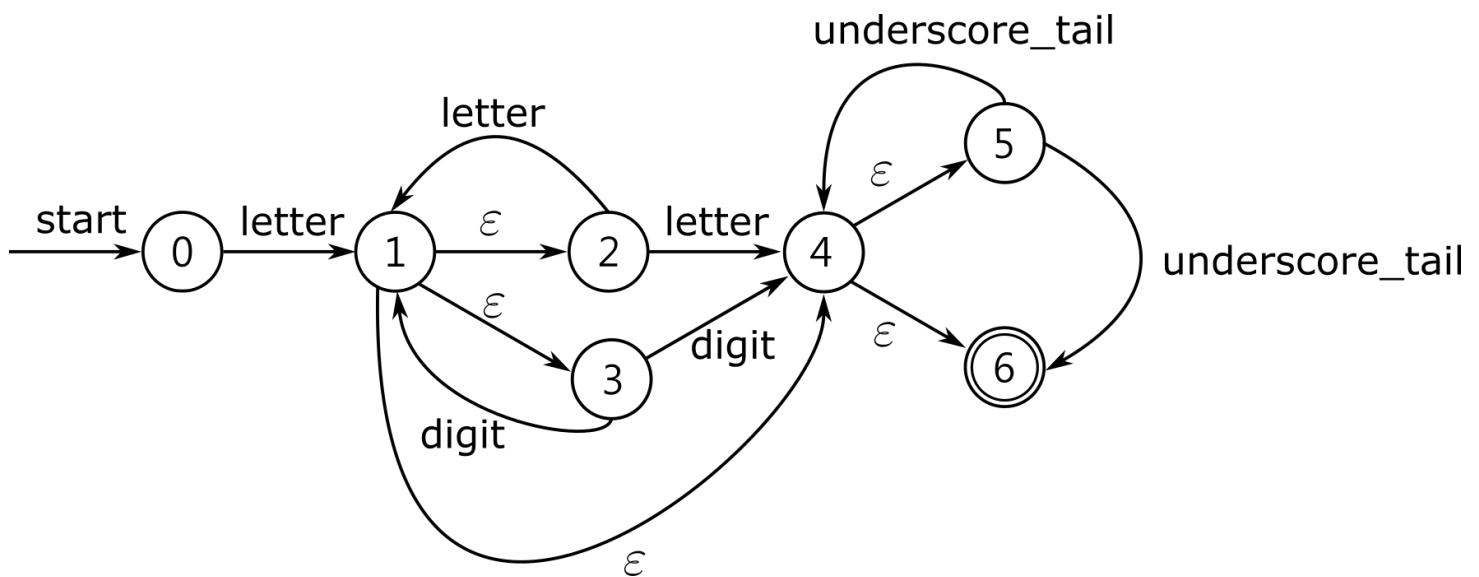
$$\textcircled{3} \Rightarrow \textcircled{1} \quad id \rightarrow letter \bullet (letter \mid digit)^*(underscore_tail)$$

$$\textcircled{4} \Rightarrow \textcircled{5} \quad id \rightarrow letter(letter \mid digit)^*(\bullet underscore_tail)$$

$$\textcircled{4} \Rightarrow \textcircled{6} \quad id \rightarrow letter(letter \mid digit)^*(underscore_tail) \bullet <<< \text{recognized}$$

$$\textcircled{5} \Rightarrow \textcircled{6} \quad id \rightarrow letter(letter \mid digit)^*(underscore_tail) \bullet <<< \text{recognized}$$

$$\textcircled{5} \Rightarrow \textcircled{4} \quad id \rightarrow letter(letter \mid digit)^* \bullet (underscore_tail)$$



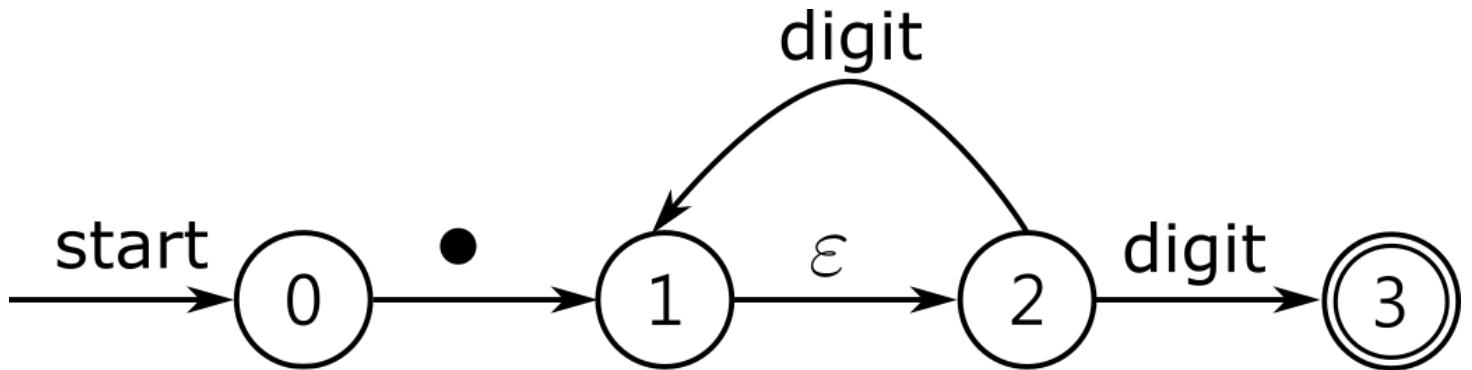
6. fraction

fraction \rightarrow $.(digit)^+$

① *fraction* \rightarrow $\bullet.(digit)^+$

① \Rightarrow ② *fraction* \rightarrow $. \bullet (digit)^+$

② \Rightarrow ③ *fraction* \rightarrow $.(digit)^+ \bullet <<< recognized$



7. optional_exponent

optional_exponent $\rightarrow (E(+ \mid - \mid \varepsilon)(digit)^+) \mid \varepsilon$

① *optional_exponent* $\rightarrow \bullet(E(+ \mid - \mid \varepsilon)(digit)^+) \mid \varepsilon$

① \Rightarrow ① *optional_exponent* $\rightarrow (\bullet E(+ \mid - \mid \varepsilon)(digit)^+) \mid \varepsilon$

① \Rightarrow ② *optional_exponent* $\rightarrow (E(+ \mid - \mid \varepsilon)(digit)^+) \mid \bullet \varepsilon$

② \Rightarrow ③ *optional_exponent* $\rightarrow (E(+ \mid - \mid \varepsilon)(digit)^+) \mid \varepsilon \bullet <<< \text{recognized}$

① \Rightarrow ④ *optional_exponent* $\rightarrow (E \bullet (+ \mid - \mid \varepsilon)(digit)^+) \mid \varepsilon$

④ \Rightarrow ⑤ *optional_exponent* $\rightarrow (E(\bullet + \mid - \mid \varepsilon)(digit)^+) \mid \varepsilon$

④ \Rightarrow ⑥ *optional_exponent* $\rightarrow (E(+ \mid \bullet - \mid \varepsilon)(digit)^+) \mid \varepsilon$

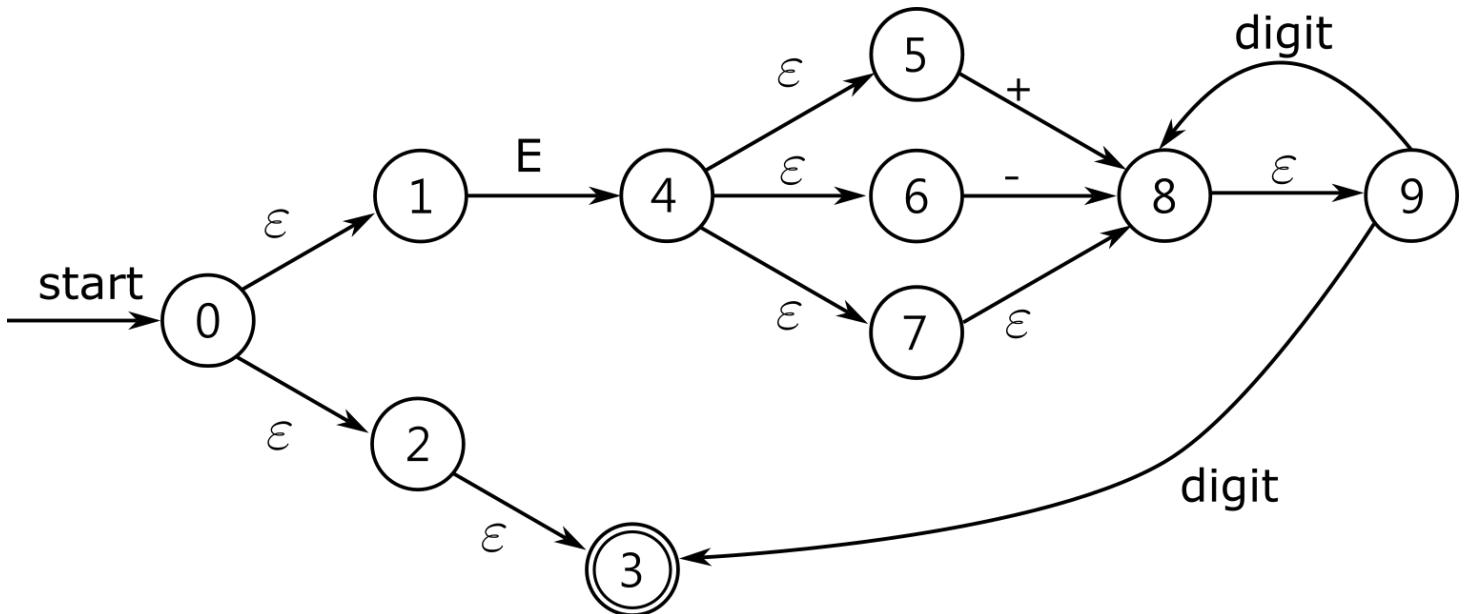
④ \Rightarrow ⑦ *optional_exponent* $\rightarrow (E(+ \mid - \mid \bullet \varepsilon)(digit)^+) \mid \varepsilon$

⑤⑥⑦ \Rightarrow ⑧ *optional_exponent* $\rightarrow (E(+ \mid - \mid \varepsilon) \bullet (digit)^+) \mid \varepsilon$

⑧ \Rightarrow ⑨ *optional_exponent* $\rightarrow (E(+ \mid - \mid \varepsilon)(\bullet digit)^+) \mid \varepsilon$

⑨ \Rightarrow ⑧ *optional_exponent* $\rightarrow (E(+ \mid - \mid \varepsilon) \bullet (digit)^+) \mid \varepsilon$

⑨ \Rightarrow ③ *optional_exponent* $\rightarrow (E(+ \mid - \mid \varepsilon)(digit)^+) \mid \varepsilon \bullet <<< \text{recognized}$



8. integer

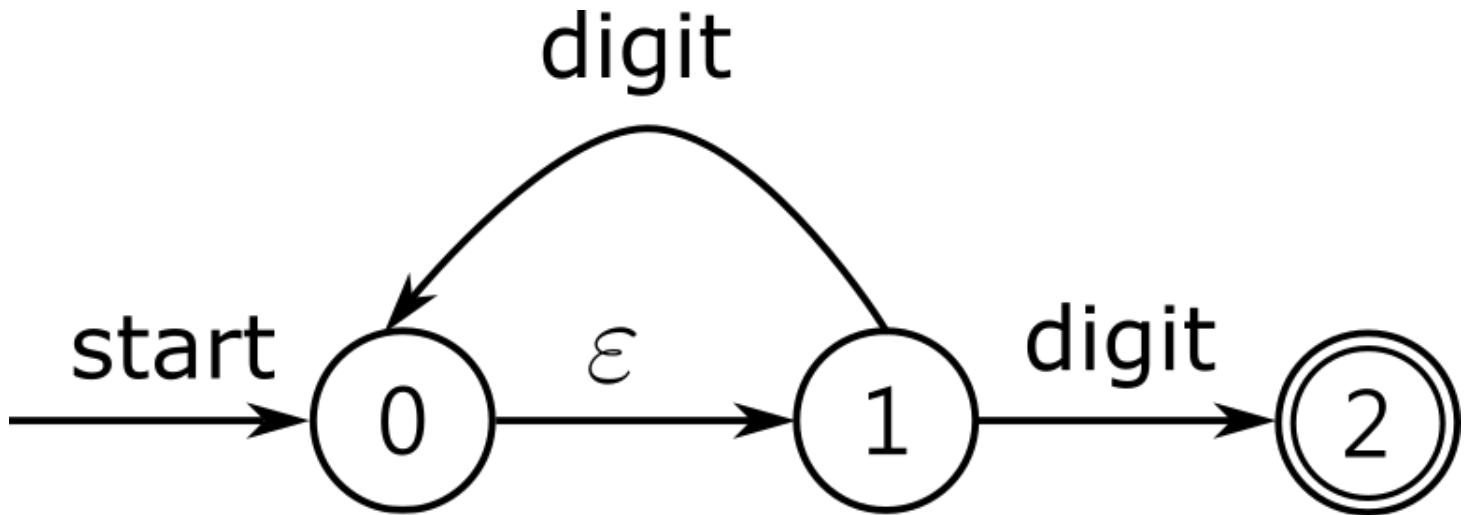
integer \rightarrow (*digit*)⁺

① *integer* \rightarrow •(*digit*)⁺

② \Rightarrow ① *integer* \rightarrow (•*digit*)⁺

① \Rightarrow ② *integer* \rightarrow •(*digit*)⁺

② \Rightarrow ③ *integer* \rightarrow (*digit*)⁺• <<< recognized



9. float

$float \rightarrow (digit)^* fraction optional_exponent$

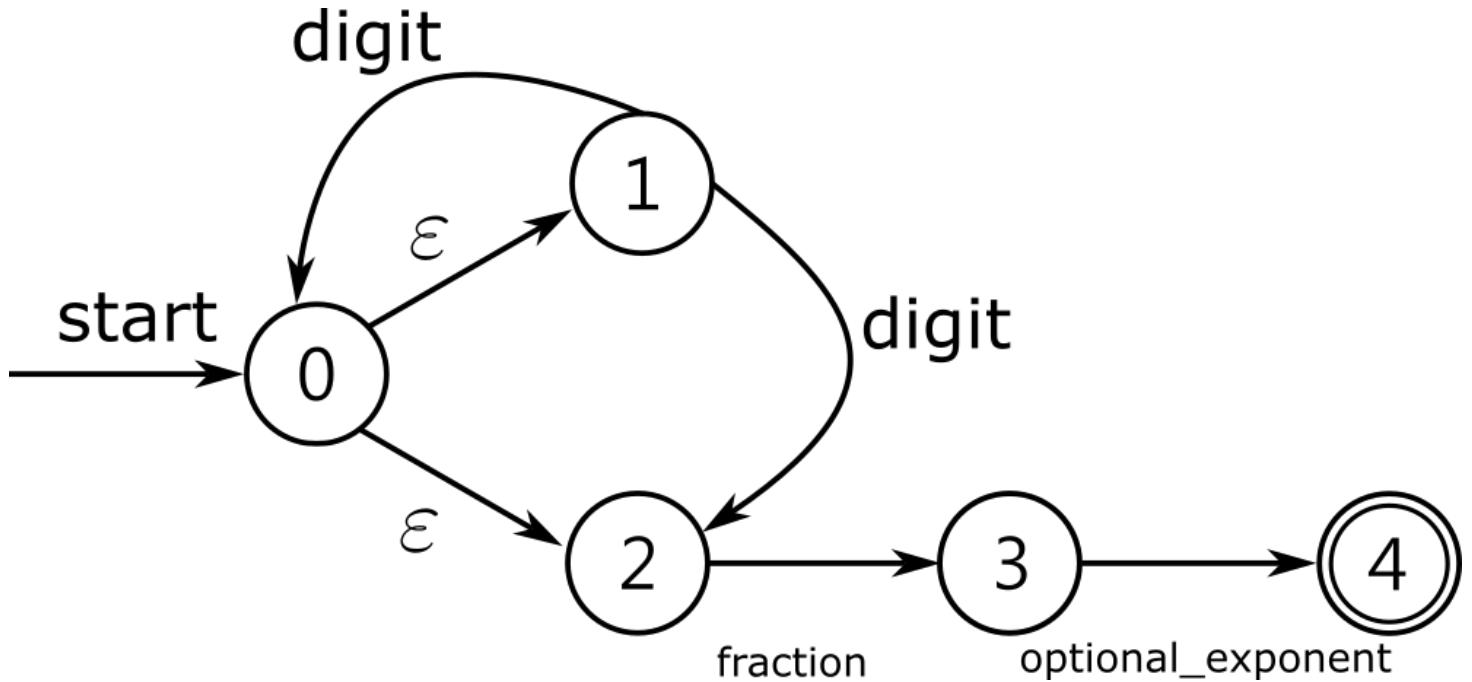
① $float \rightarrow \bullet (digit)^* fraction optional_exponent$

② $\Rightarrow ① float \rightarrow (\bullet digit)^* fraction optional_exponent$

③ $\Rightarrow ② float \rightarrow (digit)^* \bullet fraction optional_exponent$

④ $\Rightarrow ③ float \rightarrow (digit)^* fraction \bullet optional_exponent$

⑤ $\Rightarrow ④ float \rightarrow (digit)^* fraction optional_exponent \bullet <<< recognized$



NFA -> DFA

DFA

- 在任一狀態下，對於一個輸入字元，最多只能有一條對應的邊
- 不能對同一個輸入字元有兩種選擇

ϵ -closure

- ϵ -closure(S)：從S這個狀態出發，能沿著 ϵ 邊一直走下去的所有狀態就是 ϵ -closure(S)
- ϵ -closure(T)：T可以定義多個狀態來出發，就是可以不只closure一個狀態
- move(S, 字元)：從T中的所有狀態出發，根據輸入字元，所能轉移到的所有狀態

步驟

- 首先Initializations：一開始只放一個初始狀態，move後並設為 A

$$\epsilon - closure(\{0\}) = \{0\} \implies A$$

- 從初始狀態開始move(S, 字元)，表示從S中的所有狀態出發，根據輸入字元，所能轉移到的所有狀態
- 再 ϵ -closure(S)，也就是將 ϵ 消掉，如果有新集合，設定為新狀態(B、C、D...)，重複下去直到沒有新狀態。

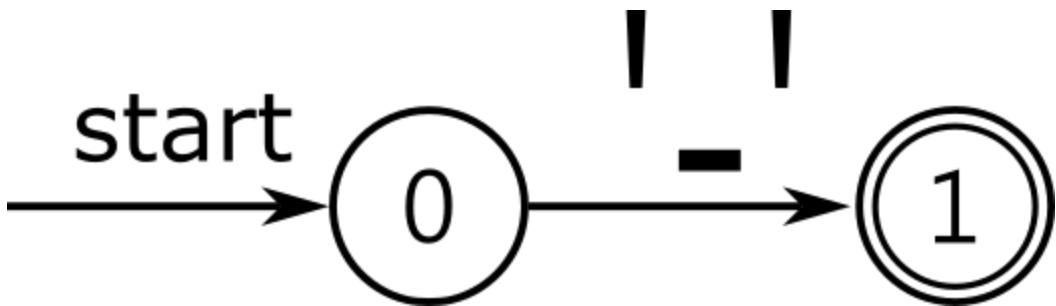
$$\begin{array}{ll} A = \{0, 1, 2, 4, 7\} & D = \{1, 2, 4, 5, 6, 7, 9\} \\ B = \{1, 2, 3, 4, 6, 7, 8\} & E = \{1, 2, 4, 5, 6, 7, 10\} \\ C = \{1, 2, 4, 5, 6, 7\} & \end{array}$$

- 再將每個 states 化成 transition diagram

STATE	INPUT SYMBOL	
	a	b
A	B	C
B	B	D
C	B	C
D	B	E
E	B	C

我們針對剛剛產生的 9 個 NFA 圖形，來產生 DFA：

1. underscore



首先他的初始狀態是0，對她closure後因為沒有 ϵ ，所以一樣是1的集合，然後把它設為A狀態

再來要對所有可能輸入的字元做move，那這邊所有可能讀到的只有底線，所以move(A,_)開始找可能轉移到的所有狀態，那他會到的只有1，

接下來對1， ϵ -closure，那他一樣因為沒有 ϵ 的邊，所以一樣只有1，且為新狀態所以設成B

再從狀態B重複以上步驟直到沒有新狀態，畫出轉移表

$$\epsilon - \text{closure}(\{0\}) = \{0\} \implies A$$

$$\text{move}(A,_) = \{1\}$$

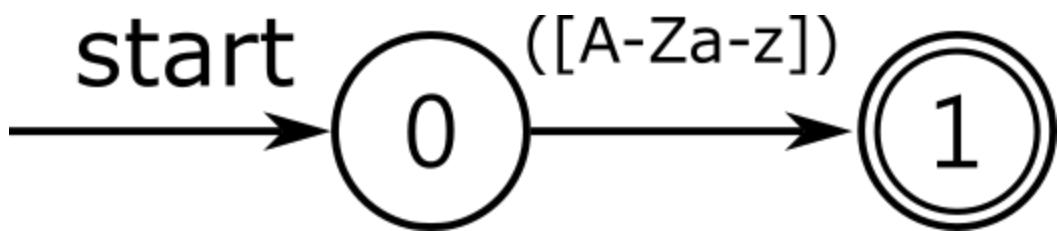
$$\epsilon - \text{closure}(\{1\}) = \{1\} \implies B$$

$$\text{move}(B,_) = \emptyset$$

$$\epsilon - \text{closure}(\emptyset) = \emptyset$$

No	State	_
0	A	B
1	B	-

2. letter



以此類推underscore

$$\varepsilon - \text{closure}(\{0\}) = \{0\} \implies A$$

$$\text{move}(A, [A - Za - z]) = \{1\}$$

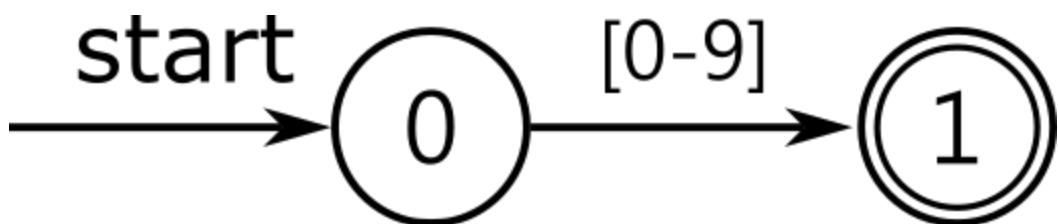
$$\varepsilon - \text{closure}(\{1\}) = \{1\} \implies B$$

$$\text{move}(B, [A - Za - z]) = \phi$$

$$\varepsilon - \text{closure}(\phi) = \phi$$

No	State	[A-Za-z]
0	A	B
1	B	-

3. digit



$$\varepsilon - closure(\{0\}) = \{0\} \implies A$$

$$move(A, [0 - 9]) = \{1\}$$

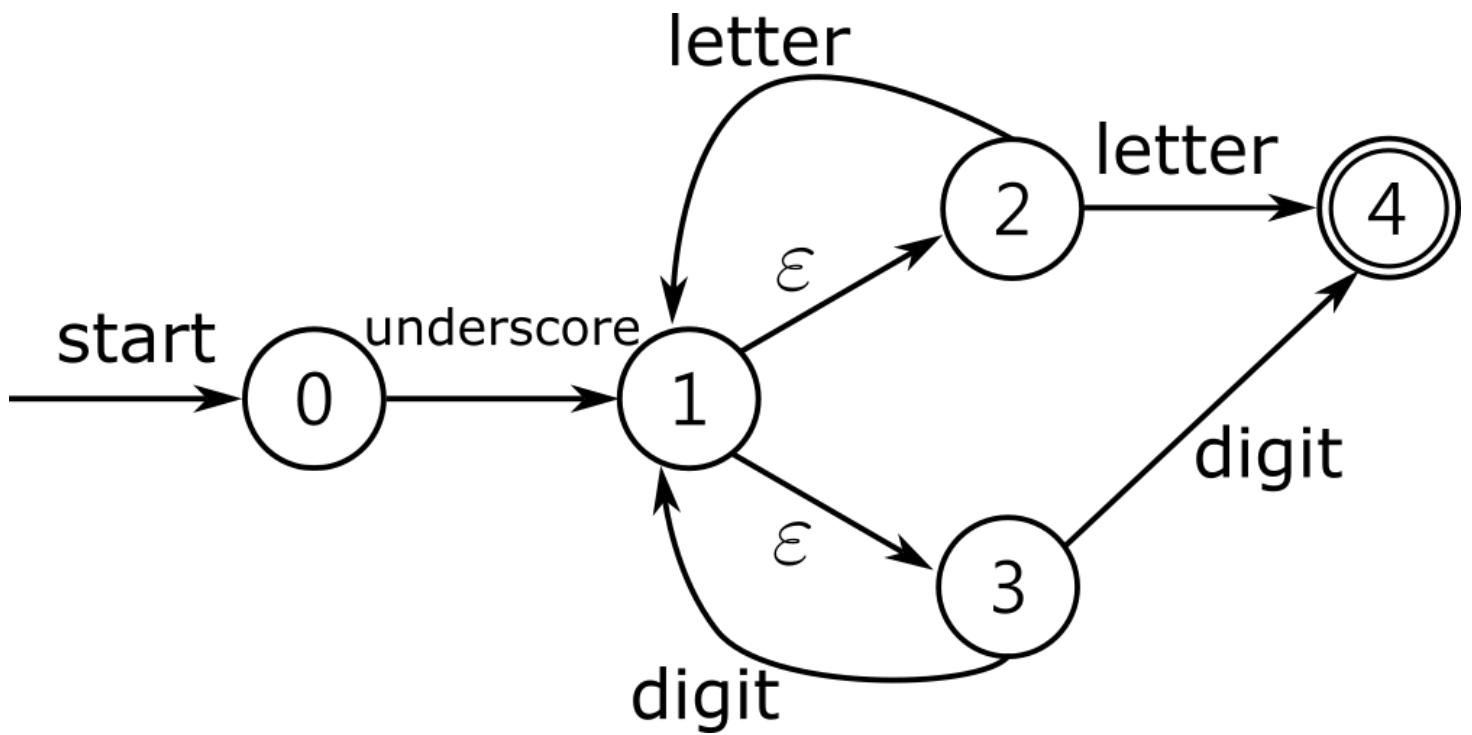
$$\varepsilon - closure(\{1\}) = \{1\} \implies B$$

$$move(B, [0 - 9]) = \phi$$

$$\varepsilon - closure(\phi) = \phi$$

No	State	[0-9]
0	A	B
1	B	-

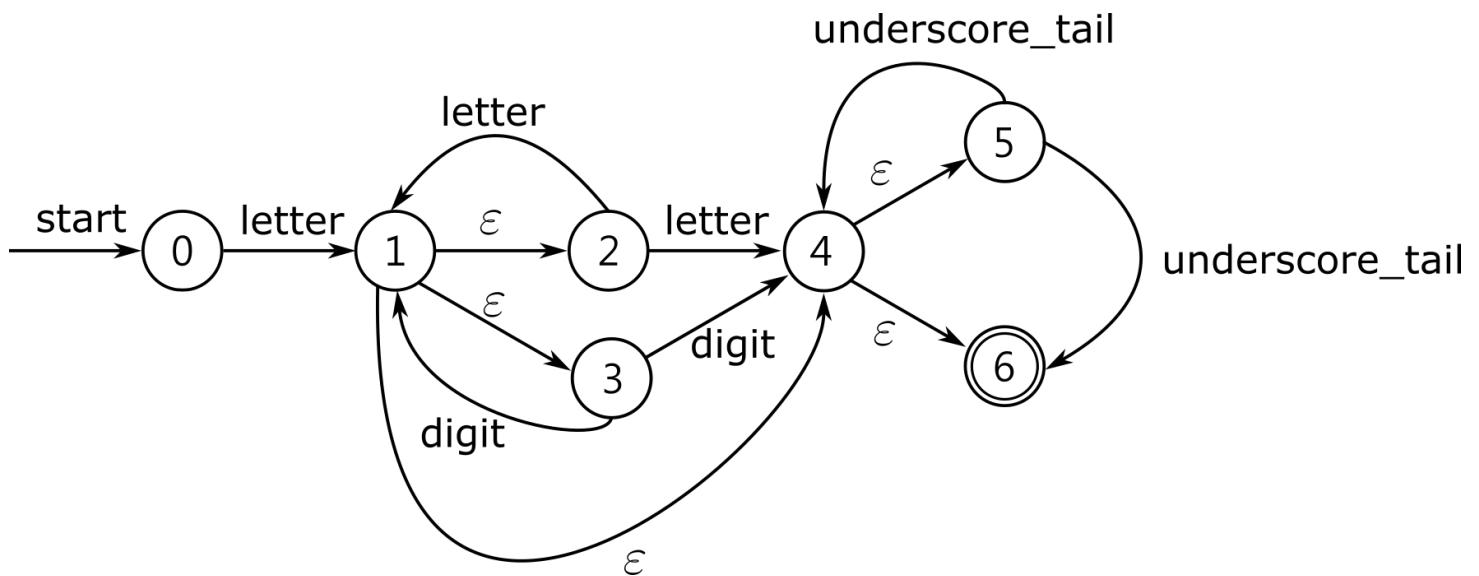
4. underscore_tail



$\varepsilon - closure(\{0\}) = \{0\} \implies A$ $move(A, underscore) = \{1\}$ $\varepsilon - closure(\{1\}) = \{1, 2, 3\} \implies B$ $move(A, letter) = \phi$ $\varepsilon - closure(\phi) = \phi$ $move(A, digit) = \phi$ $\varepsilon - closure(\phi) = \phi$ $move(B, underscore) = \phi$ $\varepsilon - closure(\phi) = \phi$ $move(B, letter) = \{1, 4\}$ $\varepsilon - closure(\{1, 4\}) = \{1, 2, 3, 4\} \implies C$ $move(B, digit) = \{1, 4\}$ $\varepsilon - closure(\{1, 4\}) = \{1, 2, 3, 4\} \implies C$ $move(C, underscore) = \phi$ $\varepsilon - closure(\phi) = \phi$ $move(C, letter) = \{1, 4\}$ $\varepsilon - closure(\{1, 4\}) = \{1, 2, 3, 4\} \implies C$ $move(C, digit) = \{1, 4\}$ $\varepsilon - closure(\{1, 4\}) = \{1, 2, 3, 4\} \implies C$

No	State	underscore	letter	digit
0	A	B	-	-
1	B	-	C	C
2	C	-	C	C

5. id



$$\varepsilon - \text{closure}(\{0\}) = \{0\} \implies A$$

$$\text{move}(A, \text{letter}) = \{1\}$$

$$\varepsilon - \text{closure}(\{1\}) = \{1, 2, 3, 4, 5, 6\} \implies B$$

$$\text{move}(A, \text{digit}) = \emptyset$$

$$\varepsilon - \text{closure}(\emptyset) = \emptyset$$

$$\text{move}(A, \text{underscore_tail}) = \emptyset$$

$$\varepsilon - \text{closure}(\emptyset) = \emptyset$$

$$move(B, letter) = \{1, 4\}$$

$$\varepsilon - closure(\{1, 4\}) = \{1, 2, 3, 4, 5, 6\} \implies B$$

$$move(B, digit) = \{1, 4\}$$

$$\varepsilon - closure(\{1, 4\}) = \{1, 2, 3, 4, 5, 6\} \implies B$$

$$move(B, underscore_tail) = \{4, 6\}$$

$$\varepsilon - closure(\{4, 6\}) = \{4, 5, 6\} \implies C$$

$$move(C, letter) = \phi$$

$$\varepsilon - closure(\phi) = \phi$$

$$move(C, digit) = \phi$$

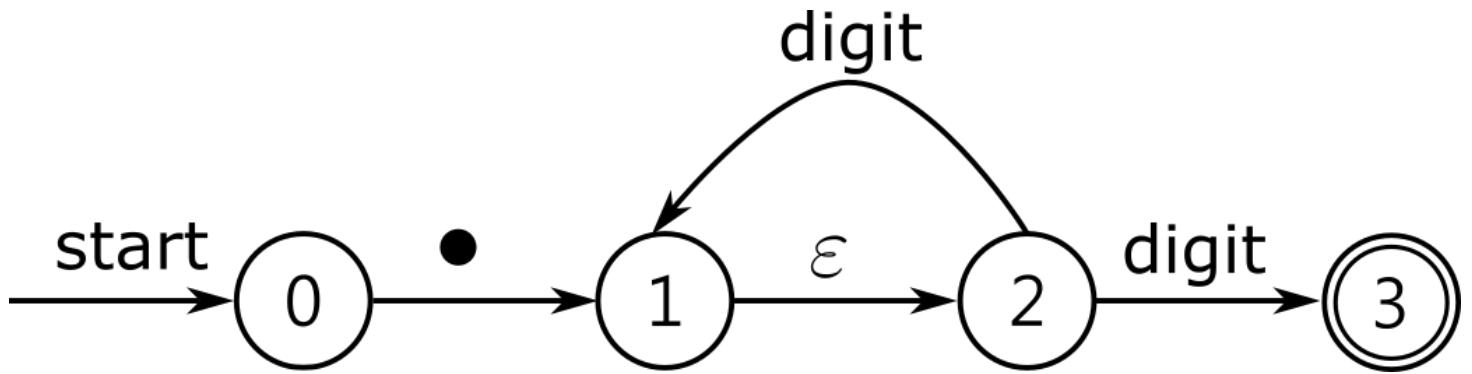
$$\varepsilon - closure(\phi) = \phi$$

$$move(C, underscore_tail) = \{4, 6\}$$

$$\varepsilon - closure(\{4, 6\}) = \{4, 5, 6\} \implies C$$

No	State	letter	digit	underscore_tail
0	A	B	-	-
1	B	B	B	C
2	C	-	-	C

6. fraction



$$\epsilon - \text{closure}(\{0\}) = \{0\} \implies A$$

$$\text{move}(A, \cdot) = \{1\}$$

$$\epsilon - \text{closure}(\{1\}) = \{1, 2\} \implies B$$

$$\text{move}(A, \text{digit}) = \phi$$

$$\epsilon - \text{closure}(\phi) = \phi$$

$$\text{move}(B, \cdot) = \phi$$

$$\epsilon - \text{closure}(\phi) = \phi$$

$$\text{move}(B, \text{digit}) = \{1, 3\}$$

$$\epsilon - \text{closure}(\{1, 3\}) = \{1, 2, 3\} \implies C$$

$$\text{move}(C, \cdot) = \phi$$

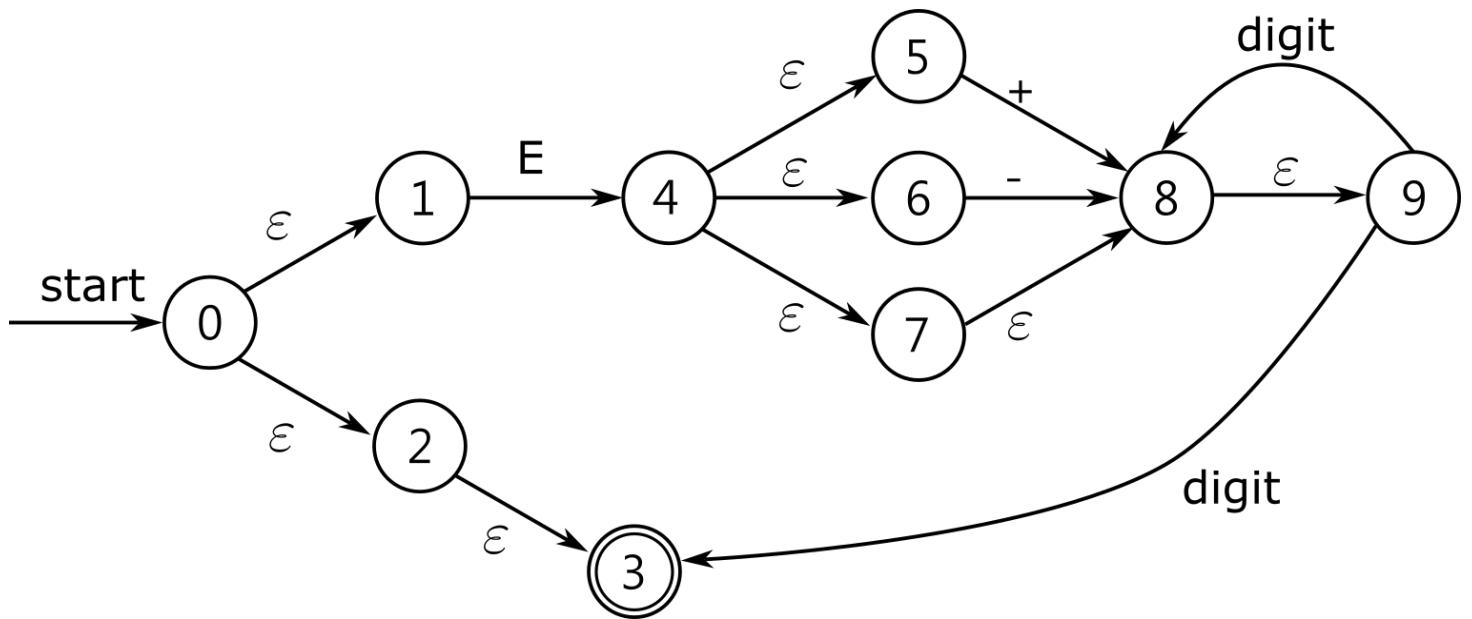
$$\epsilon - \text{closure}(\phi) = \phi$$

$$\text{move}(C, \text{digit}) = \{1, 3\}$$

$$\epsilon - \text{closure}(\{1, 3\}) = \{1, 2, 3\} \implies C$$

No	State	.	digit
0	A	B	-
1	B	-	C
2	C	-	C

7. optional_exponent



$\varepsilon - closure(\{0\}) = \{0, 1, 2, 3\} \implies A$ $move(A, E) = \{4\}$ $\varepsilon - closure(\{4\}) = \{4, 5, 6, 7, 8, 9\} \implies B$ $move(A, +) = \phi$ $\varepsilon - closure(\phi) = \phi$ $move(A, -) = \phi$ $\varepsilon - closure(\phi) = \phi$ $move(A, digit) = \phi$ $\varepsilon - closure(\phi) = \phi$ $move(B, E) = \phi$ $\varepsilon - closure(\phi) = \phi$ $move(B, +) = \{8\}$ $\varepsilon - closure(\{8\}) = \{8, 9\} \implies C$ $move(B, -) = \{8\}$ $\varepsilon - closure(\{8\}) = \{8, 9\} \implies C$ $move(B, digit)\{3, 8\}$ $\varepsilon - closure(\{3, 8\}) = \{3, 8, 9\} \implies D$

move(C, E) = ϕ

$\varepsilon - closure(\phi) = \phi$

move(C, +) = ϕ

$\varepsilon - closure(\phi) = \phi$

move(C, -) = ϕ

$\varepsilon - closure(\phi) = \phi$

move(C, digit){3, 8}

$\varepsilon - closure(\{3, 8\}) = \{3, 8, 9\} \implies D$

move(D, E) = ϕ

$\varepsilon - closure(\phi) = \phi$

move(D, +) = ϕ

$\varepsilon - closure(\phi) = \phi$

move(D, -) = ϕ

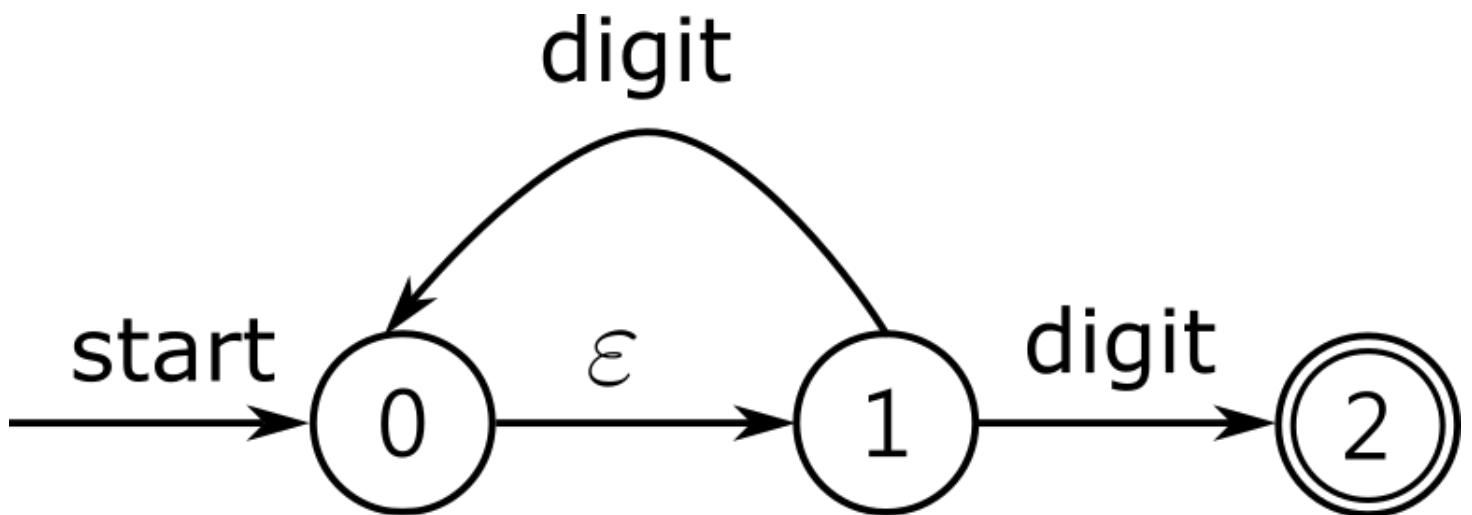
$\varepsilon - closure(\phi) = \phi$

move(D, digit){3, 8}

$\varepsilon - closure(\{3, 8\}) = \{3, 8, 9\} \implies D$

No	State	E	+	-	digit
0	A	B	-	-	-
1	B	-	C	C	D
2	C	-	-	-	D
3	D	-	-	-	D

8. integer



$$\epsilon - \text{closure}(\{0\}) = \{0, 1\} \implies A$$

$$\text{move}(A, \text{digit}) = \{0, 2\}$$

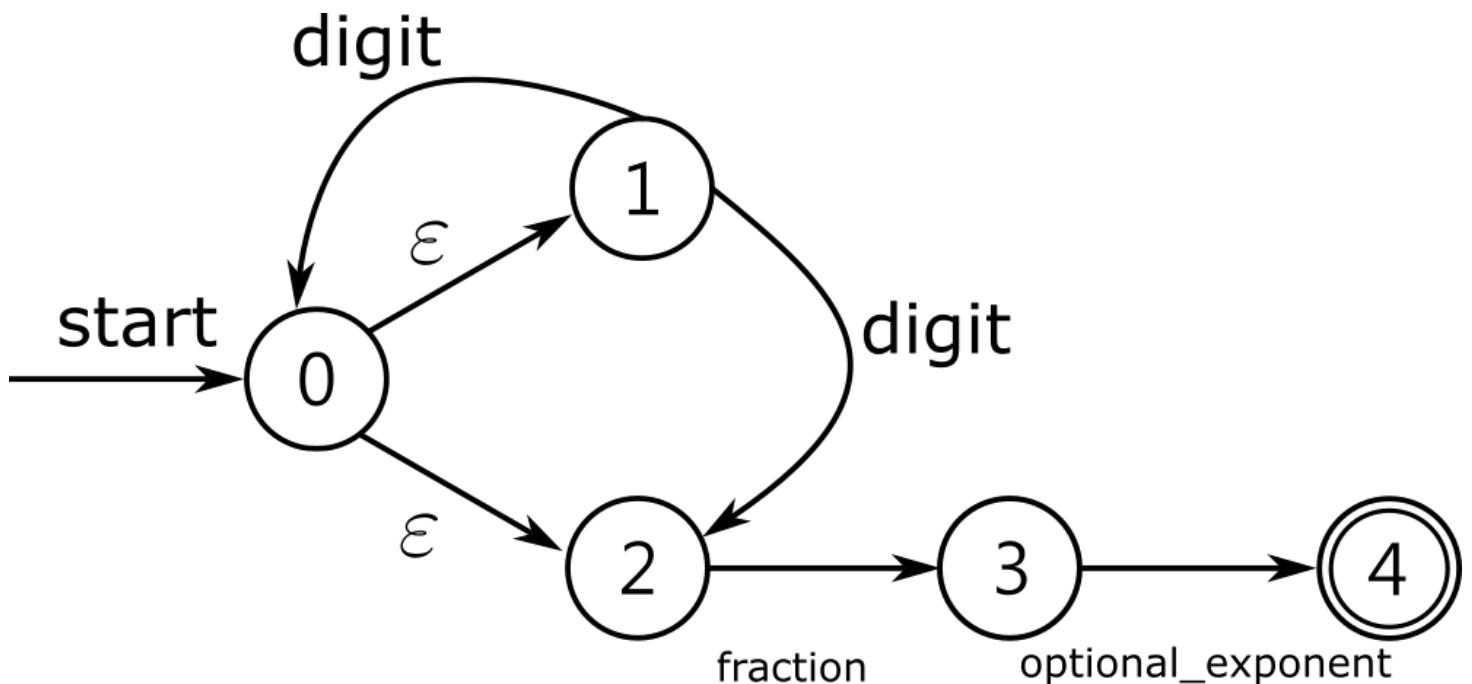
$$\epsilon - \text{closure}(\{0, 2\}) = \{0, 1, 2\} \implies B$$

$$\text{move}(B, \text{digit}) = \{0, 2\}$$

$$\epsilon - \text{closure}(\{0, 2\}) = \{0, 1, 2\} \implies B$$

No	State	digit
0	A	B
1	B	B

9. float



$$\varepsilon - \text{closure}(\{0\}) = \{0, 1, 2\} \implies A$$

$$\text{move}(A, \text{digit}) = \{0, 2\}$$

$$\varepsilon - \text{closure}(\{0, 2\}) = \{0, 1, 2\} \implies A$$

$$\text{move}(A, \text{fraction}) = \{3\}$$

$$\varepsilon - \text{closure}(\{3\}) = \{3\} \implies B$$

$$\text{move}(A, \text{optional_exponent}) = \phi$$

$$\varepsilon - \text{closure}(\phi) = \phi$$

$move(B, digit) = \phi$ $\varepsilon - closure(\phi) = \phi$ $move(B, fraction) = \phi$ $\varepsilon - closure(\phi) = \phi$ $move(B, optional_exponent) = \{4\}$ $\varepsilon - closure(\{4\}) = \{4\} \implies C$ $move(C, digit) = \phi$ $\varepsilon - closure(\phi) = \phi$ $move(C, fraction) = \phi$ $\varepsilon - closure(\phi) = \phi$ $move(C, optional_exponent) = \phi$ $\varepsilon - closure(\phi) = \phi$

No	State	digit	fraction	optional_exponent
0	A	A	B	-
1	B	-	-	C
2	C	-	-	-

DFA + DFA -> NFA -> DFA

接下來我們要處理一些複合型的 token

1. id

第一個是 id

- id

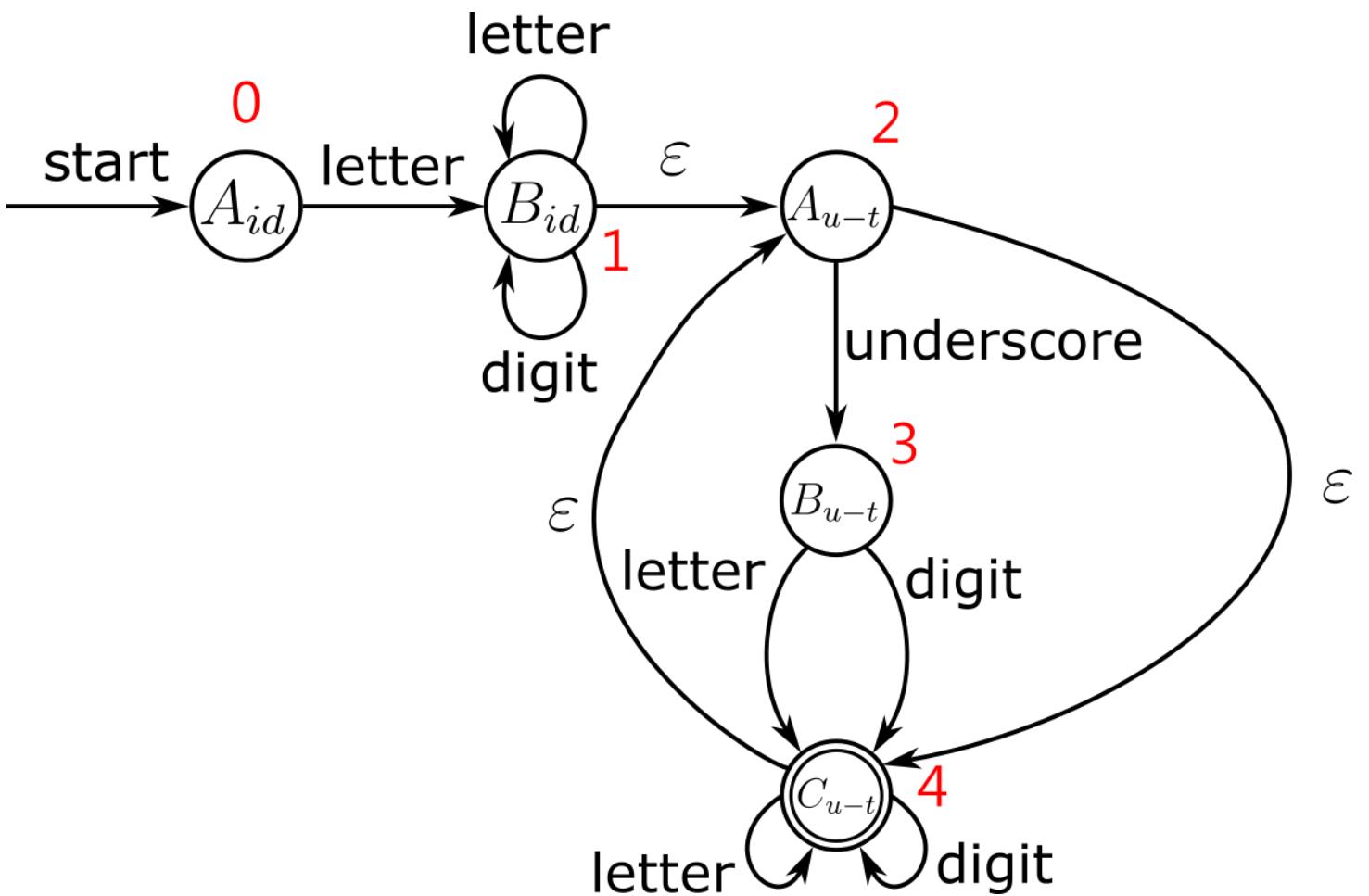
No	State	letter	digit	underscore_tail
0	A	B	-	-
1	B	B	B	C
2	C	-	-	C

可以看到，在 id 的 DFA transition table 裡面有一個叫做 underscore_tail 的東西，它還可以再拆解，我們前面的 DFA 也有提到

- underscore_tail

No	State	underscore	letter	digit
0	A	B	-	-
1	B	-	C	C
2	C	-	C	C

於是我們可以通過 id 和 underscore_tail 的 DFA Transition table，來重構出一張新的 NFA 圖形：



接著將這個 NFA 再轉換成 DFA，就像我們前面有做過的方式，這邊就不再重複。

$$\epsilon - closure(\{0\}) = \{0\} \implies A$$

$$move(A, letter) = \{1\}$$

$$\epsilon - closure(\{1\}) = \{1, 2, 4\} \implies B$$

$$move(A, digit) = \emptyset$$

$$\epsilon - closure(\emptyset) = \emptyset$$

$$move(A, underscore) = \emptyset$$

$$\epsilon - closure(\emptyset) = \emptyset$$

$move(B, letter) = \{1, 4\}$

$\varepsilon - closure(\{1, 4\}) = \{1, 2, 4\} \implies B$

$move(B, digit) = \{1, 4\}$

$\varepsilon - closure(\{1, 4\}) = \{1, 2, 4\} \implies B$

$move(B, underscore) = \{3\}$

$\varepsilon - closure(\{3\}) = \{3\} \implies C$

$move(C, letter) = \{4\}$

$\varepsilon - closure(\{4\}) = \{2, 4\} \implies D$

$move(C, digit) = \{4\}$

$\varepsilon - closure(\{4\}) = \{2, 4\} \implies D$

$move(C, underscore) = \phi$

$\varepsilon - closure(\phi) = \phi$

$move(D, letter) = \{4\}$

$\varepsilon - closure(\{4\}) = \{2, 4\} \implies D$

$move(D, digit) = \{4\}$

$\varepsilon - closure(\{4\}) = \{2, 4\} \implies D$

$move(D, underscore) = \{3\}$

$\varepsilon - closure(\{3\}) = \{3\} \implies C$

完成後，就可以得到新的 id 的 DFA Transition table:

No	State	letter	digit	underscore
0	A	B	-	-
1	B	B	B	C
2	C	D	D	-
3	D	D	D	C

2. float

接下來是 float 的部分，在 float 的 DFA transition table 中，可以看到也有一個 optional_exponent，也是可以再拆解

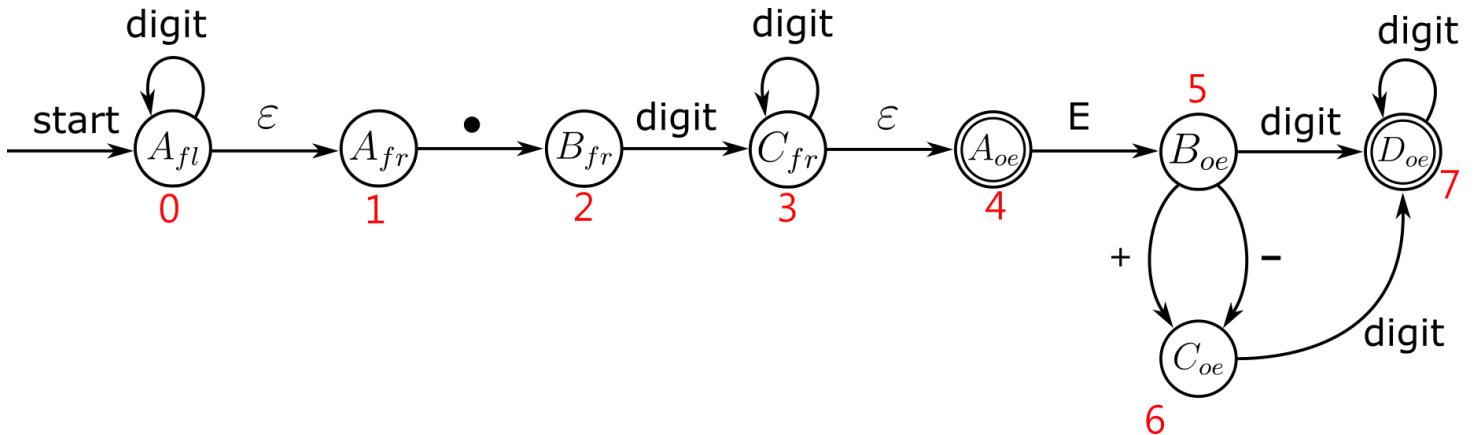
- float

No	State	digit	fraction	optional_exponent
0	A	A	B	-
1	B	-	-	C
2	C	-	-	-

- optional_exponent

No	State	E	+	-	digit
0	A	B	-	-	-
1	B	-	C	C	D
2	C	-	-	-	D
3	D	-	-	-	D

我們就一樣結合 float 和 optional_exponent 的 DFA transition table，產生新的 float NFA



再由這個新的 NFA 來產生 DFA

$$\varepsilon - closure(\{0\}) = \{0, 1\} \implies A$$

$$move(A, digit) = \{0\}$$

$$\varepsilon - closure(\{0\}) = \{0, 1\} \implies A$$

$$move(A, \cdot) = \{2\}$$

$$\varepsilon - closure(\{2\}) = \{2\} \implies B$$

$$move(A, E) = \phi$$

$$\varepsilon - closure(\phi) = \phi$$

$$move(A, +) = \phi$$

$$\varepsilon - closure(\phi) = \phi$$

$$move(A, -) = \phi$$

$$\varepsilon - closure(\phi) = \phi$$

$$move(B, digit) = \{3\}$$

$$\varepsilon - closure(\{3\}) = \{3, 4\} \implies C$$

$$move(B, \cdot) = \phi$$

$$\varepsilon - closure(\phi) = \phi$$

$$move(B, E) = \phi$$

$$\varepsilon - closure(\phi) = \phi$$

$$move(B, +) = \phi$$

$$\varepsilon - closure(\phi) = \phi$$

$$move(B, -) = \phi$$

$$\varepsilon - closure(\phi) = \phi$$

$$move(C, digit) = \{3\}$$

$$\varepsilon - closure(\{3\}) = \{3, 4\} \implies C$$

$$move(C, \cdot) = \phi$$

$$\varepsilon - closure(\phi) = \phi$$

$$move(C, E) = \{5\}$$

$$\varepsilon - closure(\{5\}) = \{5\} \implies D$$

$$move(C, +) = \phi$$

$$\varepsilon - closure(\phi) = \phi$$

$$move(C, -) = \phi$$

$$\varepsilon - closure(\phi) = \phi$$

$$move(D, digit) = \{7\}$$

$$\varepsilon - closure(\{7\}) = \{7\} \implies E$$

$$move(D, \cdot) = \phi$$

$$\varepsilon - closure(\phi) = \phi$$

$$move(D, E) = \phi$$

$$\varepsilon - closure(\phi) = \phi$$

$$move(D, +) = \{6\}$$

$$\varepsilon - closure(\{6\}) = \{6\} \implies F$$

$$move(D, -) = \{6\}$$

$$\varepsilon - closure(\{6\}) = \{6\} \implies F$$

$$move(E, digit) = \{7\}$$

$$\varepsilon - closure(\{7\}) = \{7\} \implies E$$

$$move(E, \cdot) = \phi$$

$$\varepsilon - closure(\phi) = \phi$$

$$move(E, E) = \phi$$

$$\varepsilon - closure(\phi) = \phi$$

$$move(E, +) = \phi$$

$$\varepsilon - closure(\phi) = \phi$$

$$move(E, -) = \phi$$

$$\varepsilon - closure(\phi) = \phi$$

$$move(F, digit) = \{7\}$$

$$\varepsilon - closure(\{7\}) = \{7\} \implies E$$

$$move(F, \cdot) = \phi$$

$$\varepsilon - closure(\phi) = \phi$$

$$move(F, E) = \phi$$

$$\varepsilon - closure(\phi) = \phi$$

$$move(F, +) = \phi$$

$$\varepsilon - closure(\phi) = \phi$$

$$move(F, -) = \phi$$

$$\varepsilon - closure(\phi) = \phi$$

No	State	.	digit	E	+	-
0	A	B	A	-	-	-
1	B	-	C	-	-	-
2	C	-	C	D	-	-
3	D	-	E	-	F	F
4	E	-	E	-	-	-
5	F	-	E	-	-	-

完成後就可以得到新的 DFA transition table

程式

我們的程式是用 python 來實作

讀取字元

我們首先要逐個字元來讀取，並判斷讀取到什麼類型的字元：



```
def char_type(ch):
    ascii_val = ord(ch) # 取得 ASCII 編碼

    # 判斷是否為字母 (A-Z, a-z)
    if (65 <= ascii_val <= 90) or (97 <= ascii_val <= 122):
        return 'letter'

    # 判斷是否為數字 (0-9)
    elif 48 <= ascii_val <= 57:
        return 'digit'

    # 判斷是否為底線 (_)
    elif ascii_val == 95:
        return 'underscore'

    # 判斷是否為空白字元 (空格、tab、換行、Enter)
    elif ascii_val in (9, 10, 13, 32):
        return 'layout'

    else:
        return ch # 其他字元
```

除了英文字母跟數字之外，還要判斷底線和空格等等。

如果讀取到的是空格的類別 layout，就直接跳過：



```
# 跳過 layout
if t == 'layout':
    i += 1
    continue
```

處理完空格之後，題目要求我們要可以辨識出：

(Separators and Brackets, Binary Operators,
Keywords, Comments, Identity, Integer, Float)

所以我們先定義一下每一個類別的 token

因為 Separators and Brackets, Binary Operators, Keywords 這幾種，裡面都是 or 的關係，像是
Keywords

Keywords: the language uses the following keywords int, float, bool, void, while, if, else,
for, return.

就會是有列舉出來的這幾種之一，所以簡單的用 list 的方式來處理



```
# 比較簡單的定義，只會是裡面的其中一個

# Separators
separators = {';'}

# Brackets
brackets = {'(', ')', '[', ']', '{', '}'}

# Binary Operators
binary_Operators = ['<=', '>=',
                    '<>', '==',
                    '+', '-',
                    '*', '/',
                    '<', '>',
                    '=']

# Keywords
keywords = {'int', 'float', 'bool', 'void',
            'while', 'if', 'else', 'for', 'return'}
```

那 Comments 的部分，因為需要識別開頭和結尾，所以額外寫他的規則:



```
if dfa is dfa_comment:
    if state == 0:
        if c == '/':
            state = 1
    else:
        break
elif state == 1:
    if c == '*':
        state = 2
    else:
        break
elif state == 2:
    if c == '*':
        state = 3
    else:
        # 只要不是 '*', 留在狀態 2
        state = 2
elif state == 3:
    if c == '/':
        state = 4
    elif c == '*':
        state = 3
    else:
        state = 2
else:
    break

if state in accept_states:
    last_accept = state
    last_accept_pos = i
i += 1
continue
```

讀取了開頭的 /* 後，判斷接下來是不是讀到 *，如果不是的話，就會一直停在狀態 2，這樣就可以判斷出，是不是完整的註解 /* */ 配對，也就是看註解符號有沒有閉合

id

接下來是 ID 的部分，這邊看到的是重構後的 transition table，我們需要把 DFA 的狀態轉移，放入到程式中

- DFA(重構後)

No	State	letter	digit	underscore
0	A	B	-	-
1	B	B	B	C
2	C	D	D	-
3	D	D	D	C

將原本的狀態 A~D 重新編號成 0~3 狀態，在 python 程式中，使用 dictionary 的方式，來進行對應



```
# id DFA
dfa_id = {
    0: {'letter': 1},
    1: {'letter': 1, 'digit': 1, 'underscore': 2},
    2: {'letter': 3, 'digit': 3},
    3: {'letter': 3, 'digit': 3, 'underscore': 2},
}
```

像這邊

從原本初始狀態 A，也就是 0，如果讀了 letter，就會跑到狀態 B，也就是狀態 1，那程式中就寫上，狀態 0，讀取 letter，就切換到狀態 1。

那在狀態 B，也就是 1，如果讀取了 letter 或是 digit，還是會留在狀態 B，讀取 underscore 會轉換到狀態 C，那程式中就寫上，狀態 1，讀取 letter，切換到狀態 1，讀取 digit，切換到狀態 1，讀取 underscore，切換到狀態 2。

後面依此類推，最後我們需要決定一下，在哪些狀態，是可以最為合法的結束狀態。

狀態 A 是初始狀態，當然不能當作是結束狀態。

而狀態 B，是從 A 讀了英文字母之後，或是 B 自己讀了字元或數字，也就是這串輸入的最後是英文字母或是數字，因此狀態 B 是可以作為結束狀態的。

狀態 C，是狀態 B 或 D 讀取了底線而來，也就是這串輸入的最後是底線，那根據規則：

underscore → _

letter → A|B...|Z|a|b...|z

digit → 0|1|2...|9

underscore_tail → underscore (*letter*|*digit*)⁺

id → *letter* (*letter*|*digit*)^{*} *underscore_tail*^{*}

fraction → . *digit*⁺

optional_exponent → (E (+| - |ε) *digit*⁺)|ε

integer → *digit*⁺

float → *digit*^{*} *fraction optional_exponent*

id 的結尾不可能是底線，所以狀態 C 不能當作結束狀態。

狀態 D，是由狀態 C 讀了英文字母或數字而來，也就是底線後面跟了英文字母或數字，或是從狀態 D 再讀英文字母或數字，因此這串輸入的最後面是英文字母或是數字，可以當作結束狀態。

id 這個 DFA，可以做為結束狀態的有 B 和 D，也就是狀態 1 和 3，因此我們把合法的結束狀態記錄下來：



`id_accept_states = {1, 3}`

另外的 *float* 和 *integer* 也是一樣的做法，把 DFA 的結果寫入程式中。

- *float*



```
# float DFA
dfa_float = {
    0: {'digit': 0, '.': 1},
    1: {'digit': 2},
    2: {'digit': 2, 'E': 3},
    3: {'digit': 5, '+': 4, '-': 4},
    4: {'digit': 5},
    5: {'digit': 5},
}
float_accept_states = {2, 5}
```

- integer



```
# integer DFA
dfa_int = {
    0: {'digit': 1},
    1: {'digit': 1}
}
int_accept_states = {1}
```

那程式中一些迭代的過程，我們就不去提細節。

完整程式碼

```
1 # 在同個目錄下，建立 input.txt 檔案
2 # 裡面放入要測試的輸入
3
4 def char_type(ch):
5     ascii_val = ord(ch) # 取得 ASCII 編碼
6
7     # 判斷是否為字母 (A-Z, a-z)
8     if (65 <= ascii_val <= 90) or (97 <= ascii_val <= 122):
9         return 'letter'
10
11    # 判斷是否為數字 (0-9)
12    elif 48 <= ascii_val <= 57:
13        return 'digit'
14
15    # 判斷是否為底線 (_)
16    elif ascii_val == 95:
17        return 'underscore'
18
19    # 判斷是否為空白字元 (空格、tab、換行、Enter)
20    elif ascii_val in (9, 10, 13, 32):
21        return 'layout'
22
23    else:
24        return ch # 其他字元
25
26
27    # 比較簡單的定義，只會是裡面的其中一個
28
29    # Separators
30    separators = {';'}
31
32    # Brackets
33    brackets = {'(', ')', '[', ']', '{', '}'}
34
35    # Binary Operators
36    binary_operators = ['<=', '>=', '<>', '==', '+', '-', '*', '/', '<', '>', '=']
37
38    # Keywords
39    keywords = {'int', 'float', 'bool', 'void',
40                'while', 'if', 'else', 'for', 'return'}
```

```

42 # Comments DFA
43 dfa_comments = {
44     0: {'/': 1},
45     1: {'*': 2},
46     2: {'*': 3, 'other': 2},
47     3: {'/': 4, '*': 3, 'other': 2}
48 }
49 comments_accept_states = {4}
50
51 # integer DFA
52 dfa_int = {
53     0: {'digit': 1},
54     1: {'digit': 1}
55 }
56 int_accept_states = {1}
57
58 # id DFA
59 dfa_id = {
60     0: {'letter': 1},
61     1: {'letter': 1, 'digit': 1, 'underscore': 2},
62     2: {'letter': 3, 'digit': 3},
63     3: {'letter': 3, 'digit': 3, 'underscore': 2},
64 }
65
66 id_accept_states = {1, 3}
67
68 # float DFA
69 dfa_float = {
70     0: {'digit': 0, '.': 1},
71     1: {'digit': 2},
72     2: {'digit': 2, 'E': 3},
73     3: {'digit': 5, '+': 4, '-': 4},
74     4: {'digit': 5},
75     5: {'digit': 5},
76 }
77 float_accept_states = {2, 5}
78
79
80 def run_dfa(dfa, accept_states, text, start):
81     state = 0
82     i = start
83     last_accept = -1
84     last_accept_pos = start - 1
85

```

```
86     length = len(text)
87
88     while i < length:
89         c = text[i]
90         t = char_type(c)
91
92         # Comments DFA 特殊狀態轉移
93         if dfa is dfa_comments:
94             if state == 0:
95                 if c == '/':
96                     state = 1
97                 else:
98                     break
99             elif state == 1:
100                 if c == '*':
101                     state = 2
102                 else:
103                     break
104             elif state == 2:
105                 if c == '*':
106                     state = 3
107                 else:
108                     # 只要不是 '*', 留在狀態 2
109                     state = 2
110             elif state == 3:
111                 if c == '/':
112                     state = 4
113                 elif c == '*':
114                     state = 3
115                 else:
116                     state = 2
117             else:
118                 break
119
120             if state in accept_states:
121                 last_accept = state
122                 last_accept_pos = i
123                 i += 1
124                 continue
125
126         # 一般 DFA 狀態轉移
127         transitions = dfa.get(state, {})
128
129
```

```
130     if t in transitions:
131         state = transitions[t]
132     elif c in transitions:
133         state = transitions[c]
134     else:
135         break
136
137     if state in accept_states:
138         last_accept = state
139         last_accept_pos = i
140     i += 1
141
142     if last_accept == -1:
143         return -1
144     else:
145         return last_accept_pos + 1
146
147
148 def lexer(text):
149     i = 0
150     length = len(text)
151     tokens = []
152
153     while i < length:
154         c = text[i]
155         t = char_type(c)
156
157         # 跳過 layout
158         if t == 'layout':
159             i += 1
160             continue
161
162         try:
163             # Comments DFA
164             if text.startswith('/*', i):
165                 end_pos = run_dfa(
166                     dfa_comments, comments_accept_states, text, i)
167                 if end_pos == -1:
168                     tokens.append((Binary Operators, '/'))
169                     i += 1
170                     continue
171                 else:
172                     tokens.append((Comments, text[i:end_pos]))
173
```

```
174     i = end_pos
175     continue
176
177     # id DFA + keyword 判斷
178     end_pos = run_dfa(dfa_id, id_accept_states, text, i)
179     if end_pos != -1:
180         word = text[i:end_pos]
181         if word in keywords:
182             tokens.append(('Keyword', word))
183         else:
184             tokens.append(('id', word))
185         i = end_pos
186         continue
187
188     # float DFA
189     end_pos = run_dfa(dfa_float, float_accept_states, text, i)
190     if end_pos != -1:
191         tokens.append(('float', text[i:end_pos]))
192         i = end_pos
193         continue
194
195     # integer DFA
196     end_pos = run_dfa(dfa_int, int_accept_states, text, i)
197     if end_pos != -1:
198         tokens.append(('integer', text[i:end_pos]))
199         i = end_pos
200         continue
201
202     # 運算子優先匹配（長字串優先）
203     matched_op = None
204     for op in sorted(binary_Operators, key=lambda x: -len(x)):
205         if text.startswith(op, i):
206             matched_op = op
207             break
208     if matched_op:
209         tokens.append(('Binary Operators', matched_op))
210         i += len(matched_op)
211         continue
212
213     if c in brackets:
214         tokens.append(('Brackets', c))
215         i += 1
216         continue
217
```

```
217
218     if c in separators:
219         tokens.append(( 'Separators' , c ))
220         i += 1
221         continue
222
223     # 錯誤處理
224     raise ValueError(f"Unknown token starting at position {i}: {c}")
225
226 except ValueError as e:
227     # 中止 lexer、印出目前 tokens
228     print_results(tokens)
229
230     raise e # 讓主程式捕捉錯誤
231
232
233     return tokens
234
235
236 def print_results(tokens):
237     for ttype, val in tokens:
238         print(f"{ttype:20}: {val}")
239
240
241 if __name__ == "__main__":
242     try:
243         with open("input.txt", "r", encoding="utf-8") as file:
244             code = file.read()
245
246             tokens = lexer(code)
247
248             print_results(tokens)
249
250     except ValueError as e:
251         # 有錯誤的話就輸出錯誤訊息
252         print(f"\nLexer Error: {e}")
```

Demo

我們用一個 input.txt 的檔案作為輸入



```
with open("input.txt", "r", encoding="utf-8") as file:  
    code = file.read()
```

正常情況

- 輸入



```
1 float x = 3.14E-2;  
2 if (x <= 10) {  
3 /* This is a comment */  
4 return x + 5;  
5 }
```

- 輸出

Final ➔ base

- > python index.py

```
Keyword : float
id      : x
Binary Operators : =
float   : 3.14E-2
Separators : ;
Keyword : if
Brackets : (
id      : x
Binary Operators : <=
integer : 10
Brackets : )
Brackets : {
Comment : /* This is a comment */
Keyword : return
id      : x
Binary Operators : +
integer : 5
Separators : ;
Brackets : }
```

特殊情況

最長字串匹配

輸入會從前往後，從最常匹配的字串來判斷，像是輸入三個等於，會變成是前面兩個等號一組，後面單獨一個等號一組

- 輸入



- 1 ===
- 2 >==
- 3 <=<
- 4 ==<

- 輸出

```
Final ➔ base
> python index.py
Binary Operators      : ==
Binary Operators      : =
Binary Operators      : >=
Binary Operators      : =
Binary Operators      : <=
↳ Binary Operators    : <
Binary Operators      : ==
Binary Operators      : <
```

這是因為他是用最常字串匹配的方式，他會先匹配最常符合的token類別
再來看第二個例子，他上會先匹配大於等於後再判斷另一個等於

錯誤輸出

像是我在原本正確的測試，多加了一個小老鼠 @ 符號，他就會顯示 Unknown token starting at position，但原本前面正確的 token 都還是會正確的輸出

- 輸入



```
1 if (x <= 10) {  
2 /* This is a comment */  
3 return x + 5; @  
4 }
```

- 輸出

```
Final ➔ base  
● > python index.py  
Keyword : if  
Brackets : (  
id : x  
Binary Operators : <=  
integer : 10  
Brackets : )  
Brackets : {  
Comment : /* This is a comment */  
Keyword : return  
id : x  
Binary Operators : +  
integer : 5  
Separators : ;  
  
Lexer Error: Unknown token starting at position 61: @  
Final ➔ base
```

Comments

Comments 不會去管裡面寫些什麼，只要當作是 Comments 就會完整印出內容

- 輸入



```
1 /* comments 一段奇怪的符號 @$$@%^@#$!# */
```

- 輸出

```
Final ➔ base
python index.py
Comments : /* comments 一段奇怪的符號 @$$@%^@#$!# */
Final ➔ base
o > [REDACTED]
```

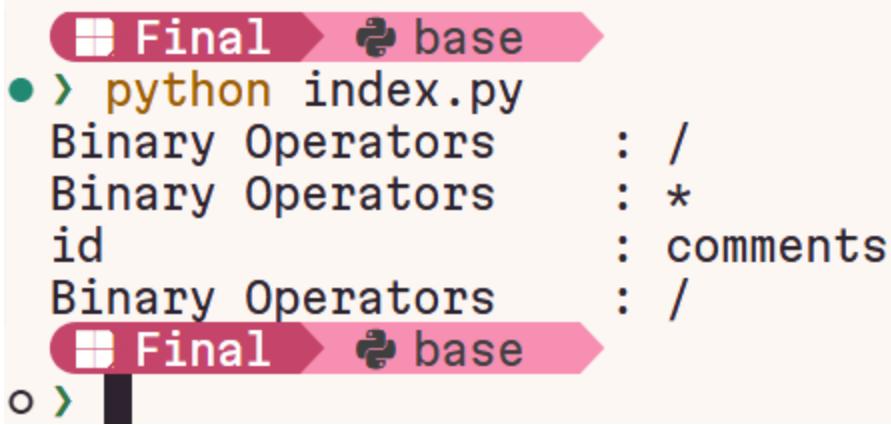
如果是沒有閉合的註解，那 / 和 * 就會被判斷為 Binary Operator，而 comments 也就會被判斷成 id

- 輸入



```
1 /* comments /
```

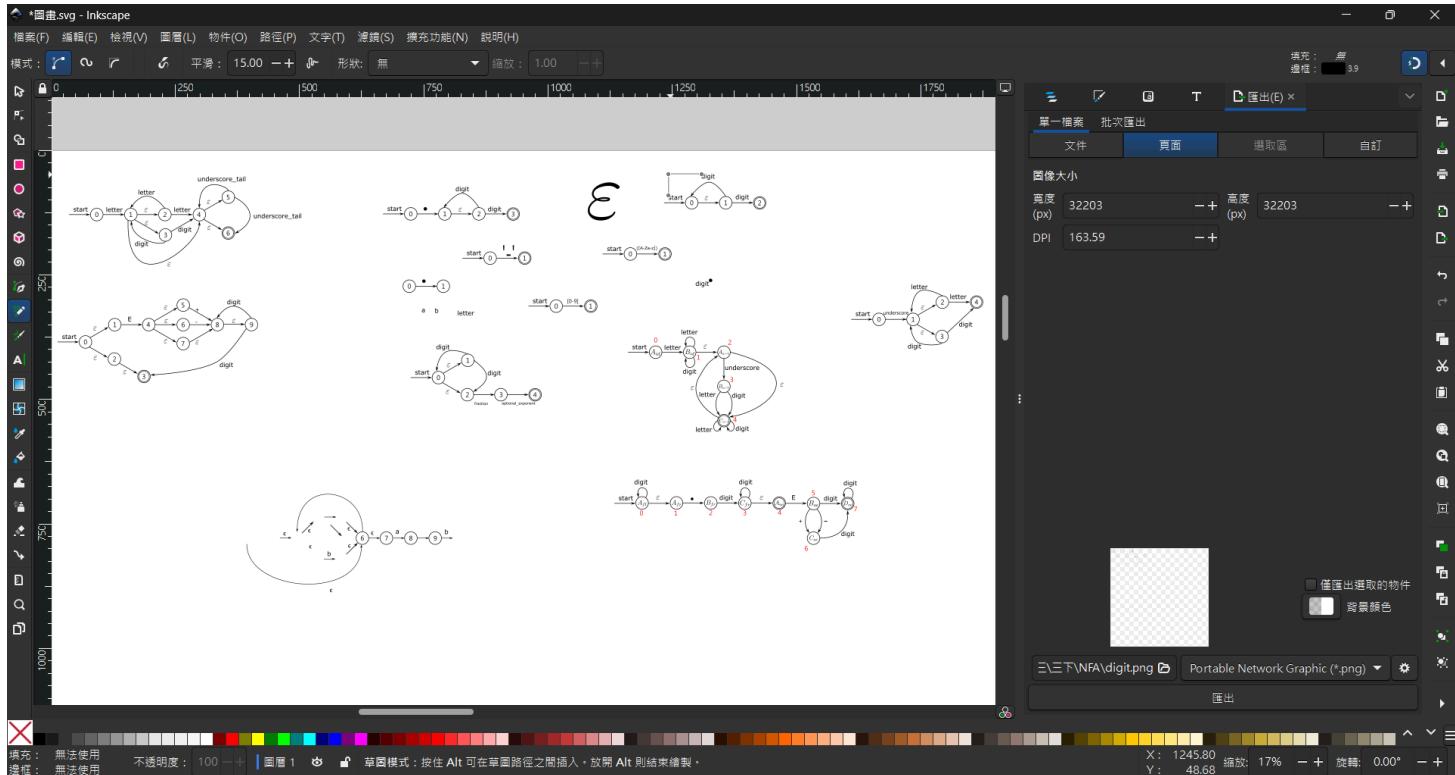
- 輸出



結論

我們用手繪的方式，將 Regular Expression 透過 Dotted Item 的方法轉換成 NFA，並畫出 NFA 圖形

我們使用的工具是 Inkscape 將手繪圖形繪製成電子檔案



再用 NFA 圖形，畫出 DFA 的 transition table，接著將 DFA 的過程加入程式中，再不使用 lex 或 similar lexical analyzer的情況下，完成 lexical analysis 的實作去分析讀到哪些Token

參考資料

- 老師上課 PPT