

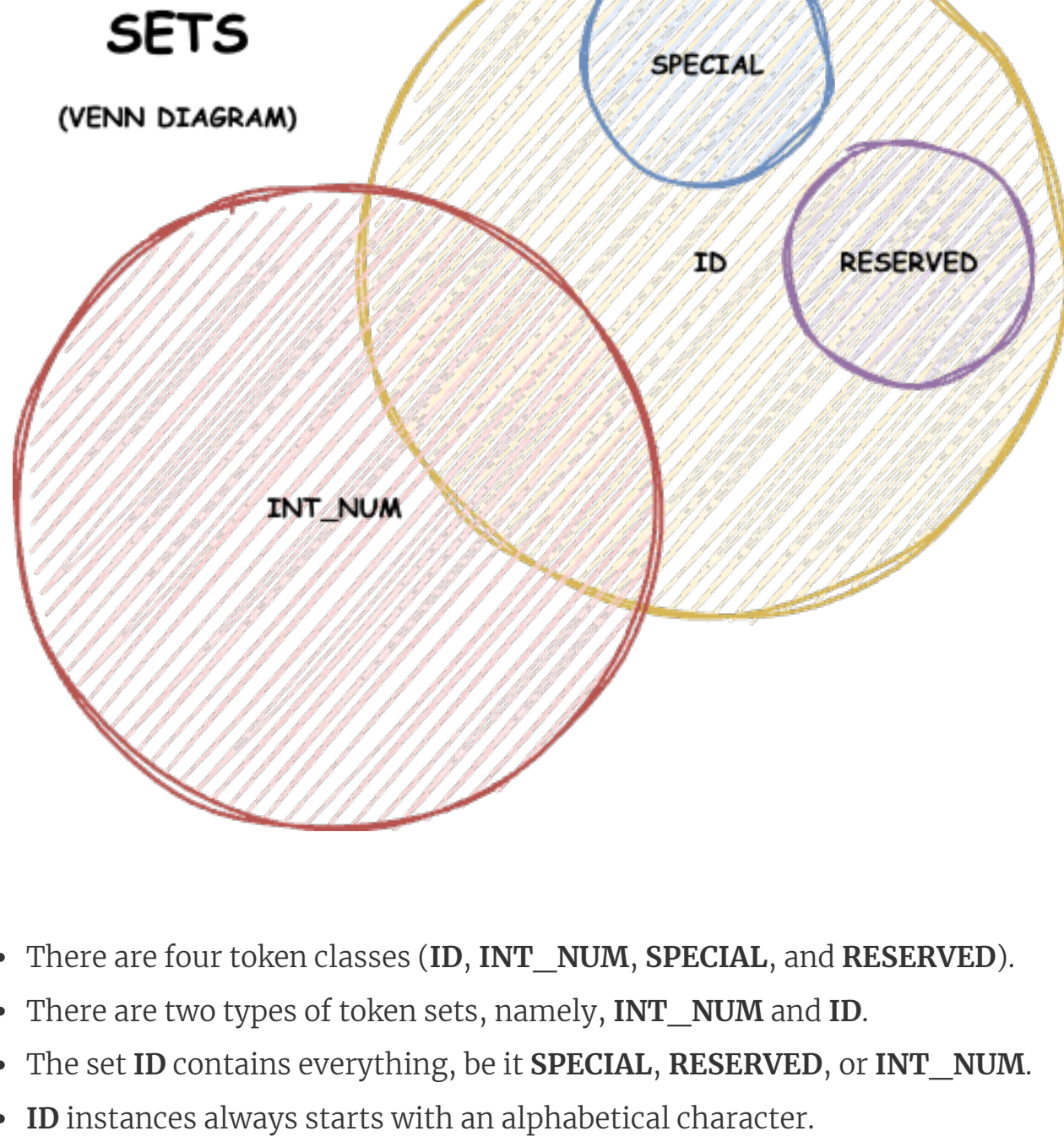
Scanner

Jose Andreas - 119010501

BASIC IDEAS

TOKEN SETS

(VENN DIAGRAM)



- There are four token classes (**ID**, **INT_NUM**, **SPECIAL**, and **RESERVED**).
- There are two types of token sets, namely, **INT_NUM** and **ID**.
- The set **ID** contains everything, be it **SPECIAL**, **RESERVED**, or **INT_NUM**.
- **ID** instances always starts with an alphabetical character.
- **ID** instances cannot contain any **SPECIAL** characters.
- `Scanner::scan()` only scans for **ID** after it exhausts every other possibility.

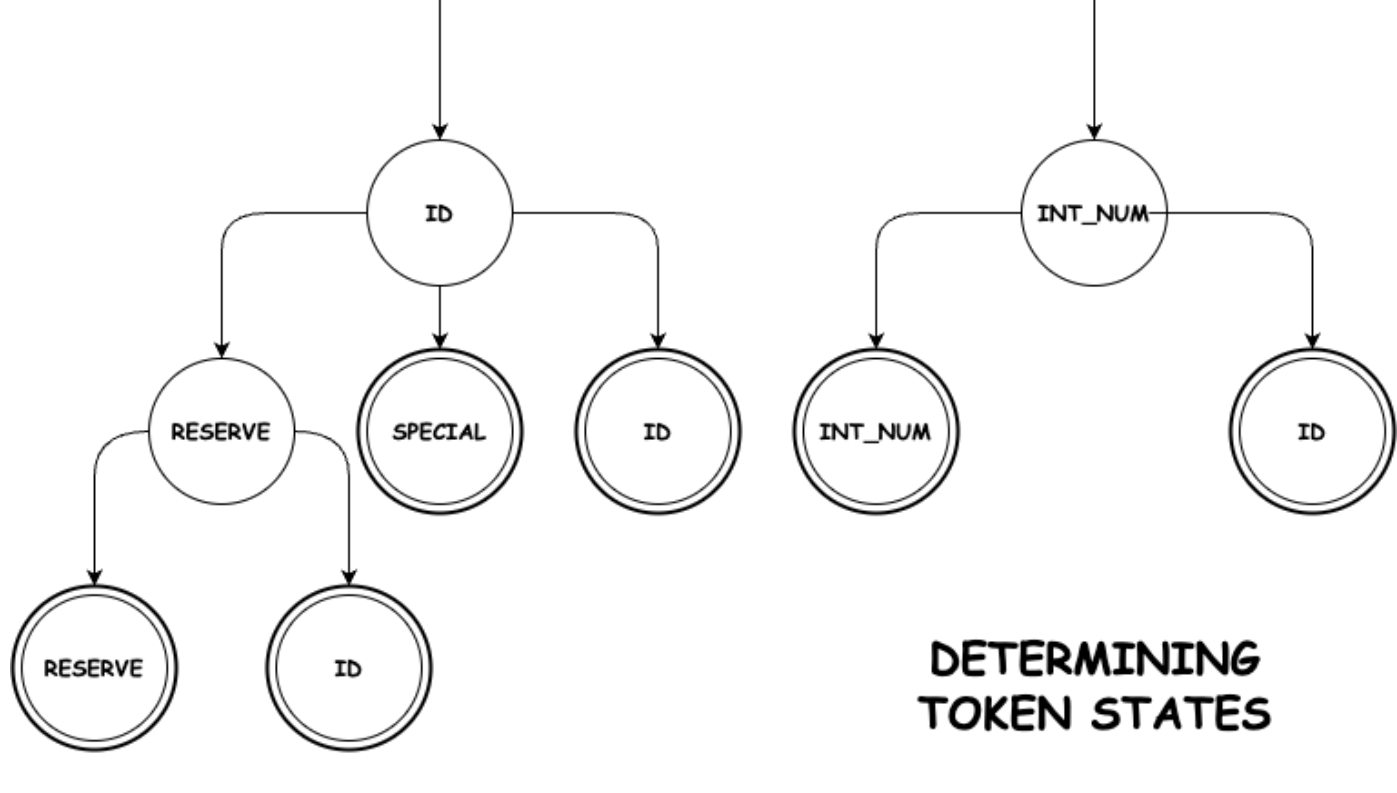
SCANNER CODE LOGIC

GOBBLE-GOBBLE



The picture above explains it all. The code traverses through the whole file character by character (albeit it also peeks at the next character for each iteration per se) and tries to matches it with the set of **SPECIAL** and **RESERVED** keywords. If it doesn't match, it stays in the same state (either **ID** or **INT_NUM**), eats up the current character, and continues to the next character.

TOP-LEVEL MATCHING RULES



HANDLING SPECIAL KEYWORDS

Note that **SPECIAL** entries may contain more than one character. Hence, there need to be rules for dealing with such cases. Fortunately, the following bit of code does just that.

```
1 #define str(character) std::string(1, character)
2 if (is_special(str(ch))) {
3     // s_char -> s_char_tail.
4     if (is_special(str(ch) + tail)) {
5         // terminating state -> id.
6         save();
7
8         // store the first special char.
9         buffer += ch;
10        break;
11    };
12
13    // buffered_s_char -> s_char.
14    if (is_special(buffer + ch)) {
15        buffer += ch;
16        // terminating state -> two length special.
17        save();
18        break;
19    };
20
21    // terminating state -> id | int_num on nonempty buffer.
22    save();
23
24    // terminating state -> single length special.
25    buffer = ch;
26    save();
27    break;
28 };
```

The code essentially enforces the following rules:

1. **IF** the next input character is also **SPECIAL**, **IF** the combination matches a different **SPECIAL** keyword, then remember and continue to next character.
2. **IF** the next input character is also **SPECIAL**, yet, the combination doesn't form another **SPECIAL** keyword, return **SPECIAL**.
3. **IF** the next input character is not **SPECIAL**, return **SPECIAL**.

Or, in grammar form:

```
1 //
2 // WARNING:
3 // the following snippet is only for
4 // explanatory purposes, not much else.
5 //
6
7 expr:
8     SPECIAL tail;
9
10 tail:
11     SPECIAL {
12         if (matches) {
13             add_to_buffer(ch);
14             return;
15         };
16         create_token("special", buffer);
17         clear_buffer();
18         return SPECIAL;
19     };
20     | ID {
21         create_token("special", ch);
22         clear_buffer();
23         return SPECIAL;
24     };
25     | INT_NUM {
26         create_token("special", ch);
27         clear_buffer();
28         return SPECIAL;
29     }
30 };
```

So, the bottom line is, the scanner checks if the current input and characters that follow can merge into a **SPECIAL** character. If they can't, the scanner will just follow the top-level matching rules (in other words, return **SPECIAL**).

HANDLING RESERVED WORDS

Like its needy token class counterpart, **RESERVED** words also needs to be handled with extra care. It is because an **ID** may or may not contain **RESERVED** words. Fret not, however, as the solution is really simple; Simply check if the merger between the current **RESERVED** word and the next character is yet another **RESERVED** word.

```
1 //
2 // context:
3 // buffer + ch forms the current word.
4 // tail refers to the next character.
5 // EOF denotes the end of file.
6 //
7
8 // terminating state -> reserved word.
9 if (is_reserved(buffer + ch)
10     && (tail == ' ' || tail == '\n' || tail == EOF)) {
11     buffer += ch;
12     save();
13     break;
14 };
```

HANDLING INT_NUM

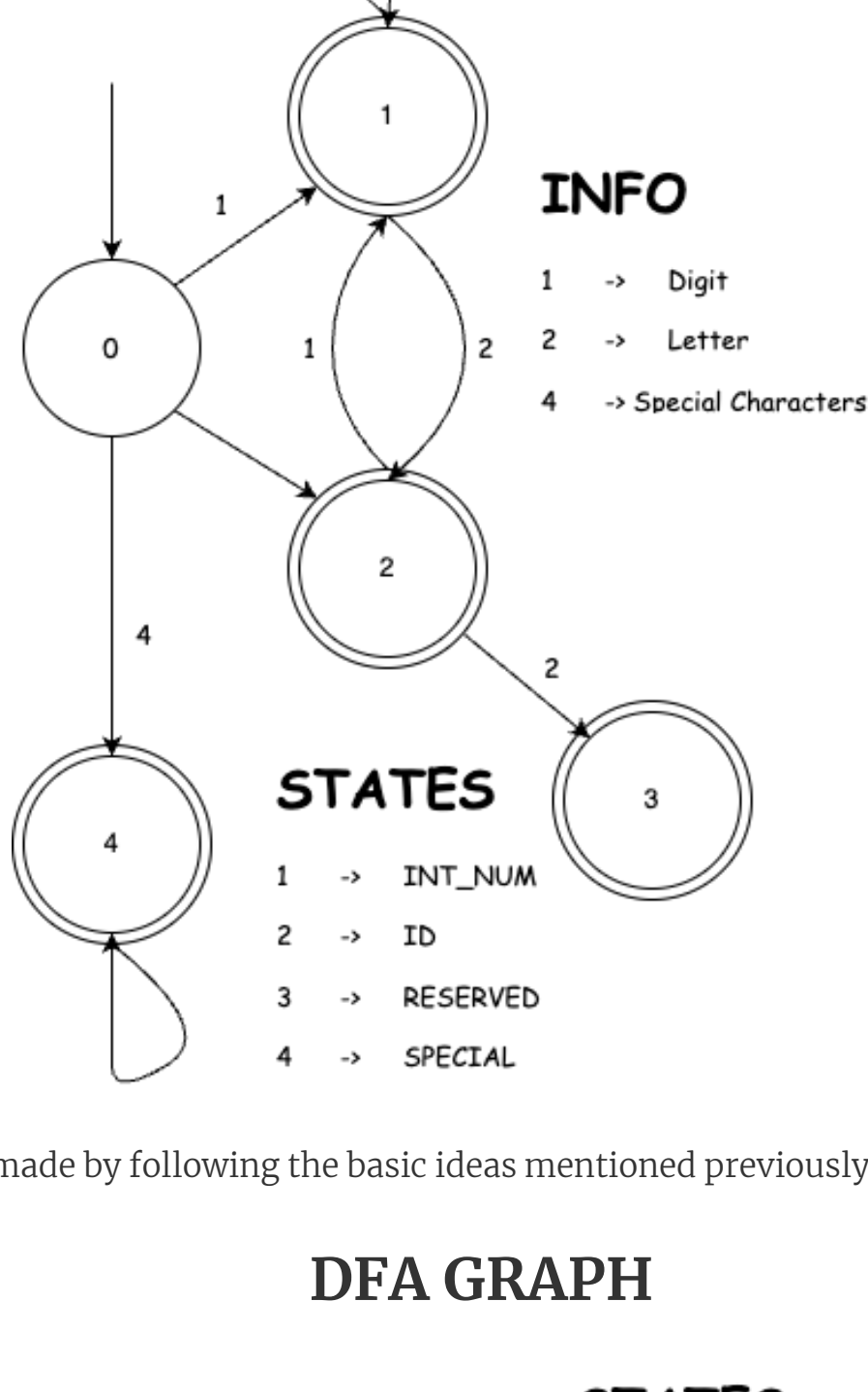
Because of the way **INT_NUM** behaves, it is the most clear-cut token class to manage. It's always through multiple copies of itself (a chain of **INT_NUM**s) or a single copy of itself. So, when the next character is no longer an **INT_NUM**, the scanner will only return **INT_NUM** if, and **ONLY** if, the every character is an **INT_NUM**. Below is a snippet that handles **INT_NUM** situated at the edge.

```
1 // edge case handling for when int_num is at the edge.
2
3 // breaking point for integers.
4 // edge case number -> space or end of line.
5 if (tail == ' ' || tail == '\n' || tail == EOF) {
6     save();
7 }
8
9 break;
```

HANDLING ID

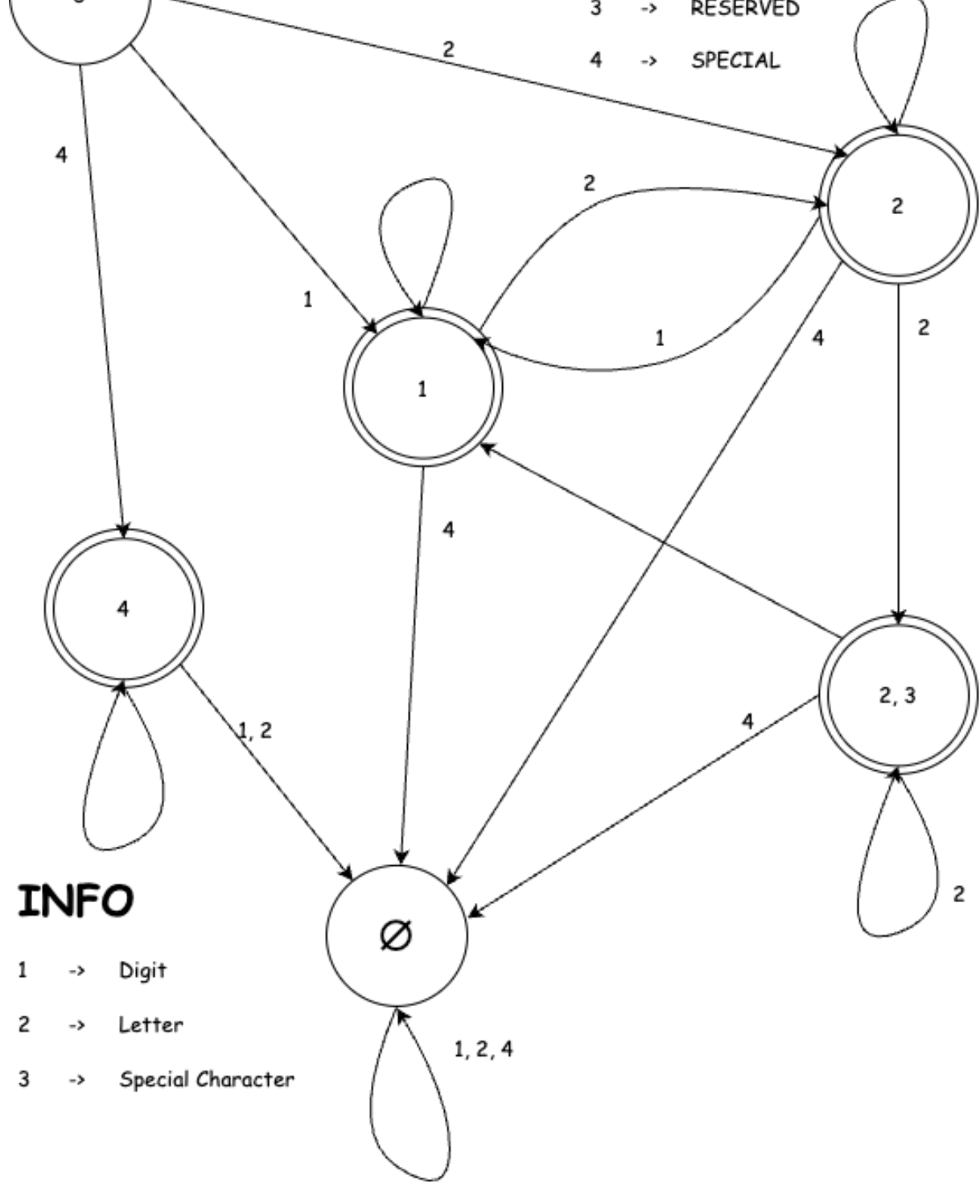
EVERYTHING that is not **INT_NUM**, **SPECIAL**, or **RESERVED**, Quite simple, really.

NFA GRAPH



The NFA is made by following the basic ideas mentioned previously.

DFA GRAPH



The DFA is constructed by adding some extra conditions and a terminating state. Do note that the code implementation differs from the NFA/DFA because the code implementation only categorizes the input types into **ID** and **INT_NUM** (the third, **SPECIAL** is expressed through several if-statements and acts as an indicator for breaking points).