# CS4100 - Project 1 Report
Jansen Craft & Kaleb Demaline

## What regular expressions/tokens did you identify, and how does your program react to each one.

### Macros

| | |
|---|---|
| F | [a-zA-Z0-9_] |
| D | [0-9] |
| INT | "int"\|"short"\|"unsigned int"\|"unsigned short"\|"long"\|"unsigned long"\|"long long"\|"unsigned long long" |
| FLOAT | "float"\|"double"\|"long double" |

### Regular Expressions and Reactions with Token names

| Regular Expression | Reaction | Token Returned (If applicable) |
|---|---|---|
| \/\/.*\n | Increment linecount | no token – single-line comment |
| "/*"([^*]*\|[*][^/])*"*/" | Add number of newlines to linecount | no token – multi-line comment |
| #include\ *<[a-z.]*> | Do nothing | no token – include statement |
| #include\ *\"[a-z.]*\" | Do nothing | no token – include statement |
| \( | Return token for open parenthesis | TK_OPEN_P |
| \) | Return token for closed parenthesis | TK_CLOSED_P |
| \{ | Return token for open curly brace | TK_OPEN_CB |
| \} | Return token for closed curly brace | TK_CLOSED_CB |
| \[ | Return token for open square bracket | TK_OPEN_SB |
| \] | Return token for closed square bracket | TK_CLOSED_SB |
| \; | Return token for semicolon | TK_SEMICOLON |
| \, | Return token for comma | TK_COMMA |
| \"[^\"]*\" | Add number of newlines to linecount. Return token for string literal | TK_STRING_LITERAL |
| \'[^\']{1,2}\' | Return token for character literal | TK_CHAR_LITERAL |
| ({D}*\.{D}*)\|{D}+ | Return token for number literal | TK_NUMBER_LITERAL |
| \+ | Return token for plus | TK_PLUS |
| \- | Return token for minus | TK_MINUS |
| \! | Return token for exclamation | TK_EXCLAMATION |
| \~ | Return token for tilde | TK_TILDA |
| \& | Return token for ampersand | TK_AMPERSAND |
| \* | Return token for asterisk | TK_ASTERISK |
| \= | Return token for equal | TK_EQUAL |

| Pattern | Action | Token |
|---|---|---|
| \/ | Return token for forward slash | TK_FORWARD_SLASH |
| \% | Return token for percent | TK_PERCENT |
| \^ | Return token for caret | TK_CARET |
| \> | Return token for greater | TK_GREATER |
| \< | Return token for less | TK_LESS |
| \? | Return token for question | TK_QUESTION |
| \| | Return token for pipe | TK_PIPE |
| {INT} | Return token for int | TK_INT |
| {FLOAT} | Return token for float | TK_FLOAT |
| "void" | Return token for void | TK_VOID |
| "char" | Return token for char | TK_CHAR |
| "while" | Return token for while | TK_WHILE |
| "for" | Return token for for | TK_FOR |
| "do" | Return token for do | TK_DO |
| "if" | Return token for if | TK_IF |
| "else" | Return token for else | TK_ELSE |
| "switch" | Return token for switch | TK_SWITCH |
| "case" | Return token for case | TK_CASE |
| "default" | Return token for default | TK_DEFAULT |
| "break" | Return token for break | TK_BREAK |
| "continue" | Return token for continue | TK_CONTINUE |
| "return" | Return token for return | TK_RETURN |
| "goto" | Return token for goto | TK_GOTO |
| {F}+ | Return token for identifier | TK_IDENTIFIER |
| [ \t\n\f] | Add number of newlines to linecount | no token - whitespace |
| . | Do nothing | no token - unmatched |

In the main of the lex program, we print out the number associated with the token (1-42) with width 2. We prepend the number with zeros if it is < 10.

**A brief description in your own words of how you implemented the Winnowing algorithm**

```cpp
int main(){
    int k = 20, w = 20;
    hash<string> hasher;
    vector<string> filenames;
    vector<vector<size_t>> fingerprints;

    string filename, tokens;
```

```cpp
while(cin >> filename && getline(cin, tokens)){
    filenames.push_back(filename);
    auto end = remove_if(tokens.begin(),tokens.end(),[](char c){return c == ' ';});
    tokens.erase(end,tokens.end());


    vector<string> k_grams;
    for (int i = 0; i <= tokens.size()-k; i++)
        k_grams.push_back(tokens.substr(i, k));


    auto fp = winnow(w,k_grams,hasher);
    fingerprints.push_back(fp);
}
```

In our main, we read in all the tokens of a file as one space-delimited string, and then remove all the whitespace. After this, we add the k_grams (substrings of this mega-string) to a vector. We have a hasher object from the stl library which maps strings to SIZE_T. We then pass the window size, the vector of k_grams, and the hasher object to out winnowing function.

```cpp
vector<size_t> winnow (int w, const vector<string>& k_grams, hash<string>& hasher){
    deque<size_t> buffer(w, SIZE_T_MAX);
    vector<size_t> fingerprints;
    int min_hash_index = 0;


    // Load initial w-1 k grams
    for (int k_idx = 0; k_idx < w-1; k_idx++){
        buffer.push_back(hasher(k_grams[k_idx]));
        buffer.pop_front();
    }


    for (int k_idx = w; k_idx < k_grams.size(); k_idx++){
        buffer.push_back(hasher(k_grams[k_idx]));
        buffer.pop_front();
        min_hash_index--;


        if (min_hash_index == -1){
            min_hash_index = w-1;
            for(int i = w-1; i >= 0; i--)
                if (buffer[i] < buffer[min_hash_index]) min_hash_index = i;
            fingerprints.push_back(buffer[min_hash_index]);
        } else {
```

```cpp
            if (buffer.back() <= buffer[min_hash_index]){
                min_hash_index = w-1;
                fingerprints.push_back(buffer[min_hash_index]);
            }
        }
    }
    return fingerprints;
}
```

In the winnowing function, we initialize a deque to w nodes with SIZE_T_MAX. We then load in the first w-1 hashed k_grams into the deque by adding to the back and popping off the front. After this, starting at the wth hash, we slide add the next hash, remove the last one, increment our min_hash_index and then check to see if we have a new smallest hash. Firstly, if our best hash was just popped off the front of buffer, we find the rightmost smallest hash and set that as our minimum hash, making sure to record the hash and location to a fingerprints vector. If our best hash is still in the buffer, we just check it against the newest hash. If the new one is smaller, we set that our min_hash_index and then append its information to the fingerprints vector which is returned at the end of the function.

```cpp
// Compare Fingerprints
    for (int i = 0; i < fingerprints.size(); i++)
        sort(fingerprints[i].begin(), fingerprints[i].end());

    vector<pair<float, string>> similarity_scores;

    for (int i = 0; i < fingerprints.size(); i++){
        for (int j = i+1; j < fingerprints.size(); j++){
            vector<size_t> intersection;
            set_intersection(
                fingerprints[i].begin(), fingerprints[i].end(),
                fingerprints[j].begin(), fingerprints[j].end(),
                back_inserter(intersection)
            );
            similarity_scores.push_back({(float)intersection.size() / (float)fingerprints[i].size(), filenames[i] + ' ' + filenames[j]});
            intersection.clear();
        }
    }

    sort(similarity_scores.rbegin(), similarity_scores.rend());

    for (auto score: similarity_scores){
        cout << setprecision(3) << setfill('0') << fixed << score.first << "\t" << score.second << '\n';
    }
```

```
}
```

In the main, once we have the fingerprints vectors for each of the files passed to standard input, we start our pairwise comparison. For the set_interection from stl that we are using, we must first sort all the fingerprint vectors. After this, we can get the size of the intersection of fingerprint vectors of the files, divide that by the length of the file in the outer loop, and then use this as our percentage of common fingerprints found. This percentage is printed out for that comparison after the percentages are sorted at the end.

**The results of your analysis, including identifying any submissions found by the algorithm that you believe may be plagiarism**

| | |
|---|---|
| 1.000 | bills_53.c bills_54.c |
| 1.000 | bills_09.c bills_54.c |
| 1.000 | bills_09.c bills_53.c |
| 1.000 | bills_03.c bills_54.c |
| 1.000 | bills_03.c bills_53.c |
| 1.000 | bills_03.c bills_09.c |
| 0.971 | bills_10.c bills_30.c |
| 0.960 | bills_47.c bills_48.c |
| 0.904 | bills_15.c bills_22.c |

Here, at the top of our Plagiarism Report, we have the pairs of student submission that were most similar. The first thing that must be said is that this assignment was clearly very simple, and they were all given the same basic instructions on how to write the program. This makes it so most of the submissions are very similar by default. However, students 53, 54, 9, and 3 had the same fingerprints. This is so unlikely to happen naturally, that it should be safe to say that these students plagiarized. One could assign a cutoff of perhaps 90% similarity and could reasonably say these students cheated, but the safest bet is to at least talk to the students mentioned above.