# HTMLParser Frame

*Hzy*9819

2018 年 5 月 5 日

## 1 Demand

### 1.1 Structure Mode

file, string.

### 1.2 Function

extract :

1.title(tag: h1 $\sim$ h6);

2.text(tag: p, b, code, em, font, i...);

3.url(tag: href);

...

### 1.3 Demand

1.use string, fstream, stringstream.

2*.use Regular Expression.

3*.Tree structure.

4*.Reference to http://htmlparser.sourceforge.net/

## 2 Frame

Is a good idea to use tree structure to reserve html codes.

### 2.1 Members

#### 2.1.1 Parser

```
1  class Parser {
2  private:
3      Node Head, Body;
```

```
4        Tag HTML;
5    public:
6        Parser() {}
7        ~Parser() ;
8        Parser(Node hd, Node bd, Tag T) : Head(hd), Body(bd), HTML(T) {}
9        Parser(stringstream & ss) ;
10       Parser(const char * f) {
11           using namespace std;
12           fstream in(f);
13           stringstream ss;
14           ss << in.rdbuf();
15           *this = Parser{ss};
16       }
17       Parser(std::string s) {
18           std::stringstream ss;
19           ss.str(s);
20           *this = Parser{ss};
21       }
22       void ShowHead(long long mode);
23       void ShowBody(long long mode);
24       void clear();
25   };
```

In fact we can change file and string into a stream,

so that we only need A stream's structure function .

### 2.1.2 Node

```
1    class Node {
2    private:
3            Tag Mtag ;
4            vector< Node > SubNode ;
5            string Text ;
6    public:
7            Node() {}
8            ~Node() ;
9            Node(Tag T, vector< Node > * SN, string tx) {
10                   ...
11           }
12           void SetTag(Tag T) ;
13           void AddSubNode(Node p) ;
14           void AddText(string s);
15           Tag CheckTag() ;
16
17           Node FindTag(string s);
```

2

```
18          void ShowTag(long long mode) ;
19          void ShowText() ;
20          bool CheckType(long long mode) ;
21          void ShowSubNode(long long mode);
22          void clear();
23 };
```

use bitmask express one kind of mode avoid too much port.

And the mode config will be shown in global-setting.

### 2.1.3  Tag

```
1 class Tag {
2 private:
3          string TagName;
4          long long TagType;
5          vector< Attribute > ATT;
6 public:
7          Tag() {}
8          ~Tag() ;
9          void SetTagType(string s);
10         Tag(string Typ, vector< Attribute > * att) {
11                 ...
12         }
13         string CheckTagName() ;
14         long long ShowType();
15         void ShowTagName();
16         void ShowAttribute(long long mode);
17
18         void AddAttribute(Attribute att);
19         void clear();
20 };
```

divide Tag into different types is convenience to extract different kind of information.

### 2.1.4  Attribute

```
1 struct Attribute {
2     std::string Type, Value;
3     Attribute() {}
4     Attribute(std::string tp, std::string val) : Type(tp), Value(val) {}
5     void Show(long long mode) ;
6     void clear() ;
7 };
```

maybe is no use ? Some information can be hidden in the attributions

## 2.2 Detail

### 2.2.1 structure

We use a stream to input the website(html code), it's easy to divide the code into tag and texts. So that we just need to find each $'<'$ and $'>'$ to find a tag. But it's also a matter to know a tag's scope(some tag has a $</>$ to stop it with some not), it forces us to list those special tags who need a $</>$ to stop in the global-setting.

When it comes to a tag, we also need to analyze the tag-name and select one or more tag-type(title, text, url...) to standard it. Also, attribute should not be left.

Then things become easy, we just need find tag $->$ analyze tag $->$ input text $->$ find next tag.

```cpp
Tag Build_Tag(std::string s) {
    using namespace std;
    int n = s.size() - 1, p = 1;
    string Name, typ, val;
    bool SetName = 0;
    Tag T;
    while(p < n) {
        if(s[p] == Space) {
            do {
                ++p;
            }while(p < n && s[p] == Space);
            continue;
        }
        if(!SetName) {
            while(p < n && s[p] != Space) Name += s[p++];
            Name = Tolower(Name);
            T.SetTagType(Name);
            SetName = 1;
        }
        else {
            typ = val = Empty;
            while(p < n && s[p] != Space && s[p] != EqualMark)
                typ += s[p++];
            if(s[p] == EqualMark) {
                do {
                    ++p;
                }while(p < n && s[p] == Space);
                if(s[p] != QuateMark)
                    while(p < n && s[p] != Space) val += s[p++];
                else {
                    do {
                        val += s[p++];
                    }while(p < n && s[p] != QuateMark);
```

```
34                if(p < n) val += s[p++];
35            }
36        }
37        Attribute att(typ, val);
38        T.AddAttribute(att);
39    }
40  }
41  return T;
42 }
```

```
1  std::string Build_Node(Node & n, std::stringstream & st,
2       std::string tagname, std::string Lef) {
3      using namespace std;
4      string temp = Empty, tmp = Lef;
5      size_t p1, p2;
6      bool b1 = 0, b2 = 0;
7      do {
8          temp = temp + Spacestring + tmp;
9
10         if(!b1) p1 = tmp.find(LeftMark, 0);
11         if(!b2) p2 = tmp.find(RightMark, 0);
12         if(!b1 && (p1 != string::npos))
13             b1 = 1;
14         if(!b2 && (p2 != string::npos))
15             b2 = 1;
16
17         if(b1 && b2) {
18             p1 = temp.find(LeftMark, 0);
19             p2 = temp.find(RightMark, 0);
20             if(p1 > p2) {
21                 tmp = temp.substr(0, p2);
22                 n.AddText(tmp);
23                 tmp = temp.substr(p2 + 1);
24                 temp = Empty;
25                 b1 = b2 = 0;
26                 continue;
27             }
28
29             tmp = temp.substr(0, p1);
30             n.AddText(tmp);
31
32             tmp = temp.substr(p1, p2 - p1 + 1);
33             Tag T = Build_Tag(tmp);
34             tmp = temp.substr(p2 + 1);
35             temp = Empty;
```

```
36
37              b1 = b2 = 0;
38              if(AssMark + tagname == T.CheckTagName())
39                  return tmp;
40              if(IgnoreMark == T.CheckTagName().at(0)
41              || AssMark == T.CheckTagName().at(0))
42                  continue;
43
44
45          Node sn;
46          sn.SetTag(T);
47
48
49          if(CheckBanner(T.CheckTagName())) {
50              do {
51                  temp = temp + Spacestring + tmp;
52                  if(Tolower(tmp).find(AssMark + T.CheckTagName(), 0)
53              != string::npos) {
54
55                      string tp = Empty;
56                      tp += LeftMark,
57              tp += AssMark,
58                  tp += T.CheckTagName() + RightMark;
59
60                      size_t p = Tolower(temp).find(tp, 0);
61                      tmp = temp.substr(0, p);
62                      sn.AddText(tmp);
63
64                      while(temp[p] != RightMark) p++;
65                      tmp = temp.substr(p + 1);
66                      temp = Empty;
67                      break;
68                  }
69                  if(!(st >> tmp)) break;
70              }while(1);
71          }
72          else if(CheckAss(T.CheckTagName()))
73              tmp = Build_Node(sn, st, T.CheckTagName(), tmp) ;
74          n.AddSubNode(sn);
75          continue;
76      }
77      if(!(st >> tmp)) break;
78  }while(1);
79  return temp;
80 }
```

### 2.2.2 About Bitmask

It's a huge work to write a function for extracting different kind of text, so we classify different tag to different kind of text, and use Bitmask to represent it. This method can lower code length effective and easier to use. But it also forces us to define some of constant in global-setting to represent different kind of text. And enumerate tags to classify them is also a slavery.

Here is a sample(no need to regard it as a standard):

```cpp
const LL NoneType = 0;
const LL Title = (1LL << 0);
const LL Text = (1LL << 1);
const LL Url = (1LL << 2);
const LL Img = (1LL << 3);
const LL Script = (1LL << 4);
const LL Button = (1LL << 5);
const LL Table = (1LL << 6);
const LL Video = (1LL << 7);
const LL audio = (1LL << 8);
const LL ALL = (1LL << 62) + ((1LL << 62) − 1);

const LL Att_Type = (1LL << 0);
const LL Att_Value = (1LL << 1);
...
```

### 2.2.3 About Banner

We need to find each $<>$ in the html code, but in some text(like javascript) we also can find $<>$. That may be a mistake when we deal with it using the method above. So that we do need to skip this kind of text.

```cpp
const std::string Banner = % script code ruby %;
```