

برای طراحی کامپایلر از کتابخانه ی regular expression (re) زبان پایتون کمک گرفته شده. [1]
ابتدا تابع compiler صدا زده میشود.

```
323 def compiler():
324     file_py=open('p.txt','r').read()
325     # print(file_py)
326     file_py_sp=file_py.split('\n')
327     # print(len(file_py_sp))
328     # give code to lexical analyzer
329     for index,code_line in enumerate(file_py_sp,start=1):
330         lexical_analysis(index,code_line,file_py)
331     show_key()
332     # give code to syntax analyzer
333     for index,code_line in enumerate(file_py_sp,start=1):
334         # if index==len(file_py_sp):
335         #     break
336         # print(index,code_line)
337         syntax_analysis(index,code_line,file_py_sp)
338
339
```

این تابع فایل p.txt را میخواند. این فایل برنامه پایتون است که در نهایت به کد زبان C تبدیل میشود. فایل را خط به خط split میکنیم و با حلقه for برای هر خط تابع lexical_analysis را صدا میزنیم. Index شماره خط است که چون در enumerate متغیر start=1 است پس از صفر شروع نمیشود و شماره خط های فایل درست است. تابع lexical_analysis در ادامه توضیح میدهم. بخش lexical_analysis:

<pre>33 def lexical_analysis(index,code,file_py): 34 global x 35 # all pattern 36 37 # pattern for symbol key 38 t_for='for' 39 t_while='while' 40 t_if='if' 41 t_elif='elif' 42 t_else='else' 43 t_break='break' 44 t_continue='continue' 45 t_return='return' 46 t_try='try' 47 t_except='except' 48 t_in='in' 49 t_range='range' 50 t_print='print' 51 t_input='input' 52 53 # pattern for variable 54 t_variable= r'/s[a-zA-Z0-9]*[a-zA-Z0-9]/s' 55 # pattern for operator 56 t_equal='==' 57 t_equal_not='!=' 58 t_more='>' 59 t_less='<' 60 t_more_equal='>=' 61 t_less_equal='<=' 62 # t_and='and' 63 # t_or='or'</pre>	<pre>63 # t_or='or' 64 # t_not='not' 65 66 t_plus_equal='+=' 67 t_mines_equal='-=' 68 t_multi_equal='*=' 69 t_dev_equal='/=' 70 71 t_plus='+' 72 t_mines='-' 73 t_multi='*' 74 t_dev='/' 75 t_bagi='%' 76 77 78 79 80 token_key=[t_for,t_while,t_if,t_elif,t_else,t_break,t_continue 81 ,t_return,t_try,t_except,t_in,t_range,t_print,t_input] 82 83 token_op=[84 t_equal, 85 t_equal_not, 86 t_more, 87 t_less, 88 t_more_equal, 89 t_less_equal, 90 # t_and, 91 # t_or, 92 # t_not, 93 t_plus_equal, 94 t_mines_equal, 95 t_multi_equal,</pre>
--	---

در این تابع ابتدا pattern هایی برای عملگر ها و متغیر ها و کلمات کلیدی مثل for, if, while ... تعریف کردیم. مثلاً برای for خود کلمه ی “for” را قرار داده ایم. حرف r قبل از string ها آنها را به raw string تبدیل میکند. اپراتور های مقایسه ای و عملگر های انتساب هم تعریف شده اند. عملگر های عادی /* ... هم

تعریف شده اند. سپس همه ی symbol key ها (مثل ... if,for) را در لیستی به نام token_key قرار داده شد و مقایسه ای ها و عملگرها و اپراتور ها در لیست token_op.

```
103
104     # get all key and append to list
105     for tk in token_key:
106         if tk in code:
107             if tk not in symbol_key:
108                 symbol_key.append(tk)
109
110     # get op in code
111     for tk in token_op:
112         if tk in code:
113             if tk not in op_list :
114                 op_list.append(tk)
115
116     # get variable in code
117     # a=12 a+=1
118     pat=r' [a-zA-Z]{1} |^[a-zA-Z][a-zA-Z0-9]*=[0-9a-zA-Z\'"]*'
119     x=re.findall(pat,file_py,re.MULTILINE)
120
121
```

هر خط کد که به تابع lexical_analysis داده میشود با متغیر code معلوم میشود. و بعد symbol_key ها و عملیات هایی که در code هستند در دو لیست ذخیره میکنیم. (op_list و symbol_key)

توضیح برای متغیر pat : عملیات انتساب

{1} در جلوی [a-zA-Z] به این معناست که یکبار حرف دیده شود. و * جلوی

[a-zA-z0-9] به همان معنی closure است. و بعد از تساوی " یا ' هم میتواند باشد.

متد show_key هم نمایش موارد بالاست. یعنی symbol_key و op_list نمایش متغیر ها و کلمات کلیدی و عملیات هایی که در فاز lexical یافت شده.

بخش syntax analysis:

در صفحه بعد کد این بخش آمده. به این تابع شماره خط هر کد index و هر خط کد everylinecode و همچنین تمام کد داده شده allcode. در ادامه نحوه ساخت دستورات مختلف توصیف شده. برای مثال در دستور for: s* به معنی این است که هرچقدر فاصله وجود داشت نادیده گرفته شود. به این دلیل که در پایتون در حلقه های تو در تو indentation نیاز است، این کار برا سادگی انجام شده. بعد کلمه for آمده و بعد دوباره یک s* که باز به معنی نادیده گرفتن فاصله است. سپس اسم متغیر حلقه با [a-zA-Z]* مشخص شده و بعد in range و سپس یک پرانتز باز میکنیم ([) و در انتها هم میبندیم. درون پرانتز دو عدد داریم که هر کدام با [0-9]* مشخص شده اند و بینشان کما است. و بعد دو نقطه. بقیه تعاریف هم مشابه هستند. در while مقایسه ها با [<=>|<|>|!=|==] مشخص شده. همه این pattern ها را در متغی token قرار میدهیم. حال در ادامه به بررسی مطابق بودن کد دریافتی با قوانینی که تعریف کردیم میپردازیم. تابع (re.findall(tk,everylinecode)) بررسی میکند که آیا هر خط کد ورودی با

Tk که همان قوانین بالا هستند که در متغیر token همه شان را گذاشتیم مطابقت دارد یا نه. اگر درست بود متغیر compile یک لیست مقدار دار میشود و none نمیشود. و سپس اگر طول all_syntax برابر یک نشد این به معنای خطا است و انرا چاپ میکنیم و خطی که ارور دارد با index مشخص میکنیم وگرنه کد درست است و چاپ میکنیم که این خط syntax درست دارد. و سپس compile_to_c() صدا میزنیم.

```

140 counter=1
141 def syntax_analysis(index, everylinecode, allcode):
142     global counter
143     # pattern
144     # for i in range(1,100):
145     pat_for=r'\s*for\s*[a-zA-Z]* in range\s*([0-9]*,[0-9]*)\s*:'
146     # while i<10:
147     pat_while=r'\s*while\s*[a-zA-Z0-9_]*\s*[!=>|<|>|<=|>=]\s*[0-9a-zA-Z\'"]*\s*:'
148     # if i==0:
149     pat_if1=r'\s*if\s*[a-zA-Z0-9_]*\s*[!=>|<|>|<=|>=]\s*[0-9a-zA-Z\'"]*\s*:'
150     # if i%2==0:
151     pat_if2=r'\s*if\s*[a-zA-Z0-9_]*\s*[%]\s*[0-9]*\s*[!=>|<|>|<=|>=]\s*[0-9]*\s*:'
152     # elif i==4:
153     pat_elif1=r'\s*elif\s*[a-zA-Z0-9_]*\s*[!=>|<|>|<=|>=]\s*[0-9a-zA-Z\'"]*\s*:'
154     # elif i%2!=0:
155     pat_elif2=r'\s*elif\s*[a-zA-Z0-9_]*\s*[%]\s*[0-9]*\s*[!=>|<|>|<=|>=]\s*[0-9]*\s*:'
156     # else:
157     pat_else=r'\s*else\s*:'
158     # print('hi')
159     pat_print1=r'\s*print\s*[(]\s*[\']|\s*[a-zA-Z0-9_]*\s*[\']\s*[(\s*)]'
160     # print(i)
161     pat_print2=r'\s*print\s*[(]\s*[a-zA-Z0-9_]*\s*[(\s*)]'
162     # a=a+1
163     pat_op=r'\s*[a-zA-Z0-9_]*\s*=\s*[a-zA-Z0-9_]*\s*[+|-|*|/] \s*[0-9a-zA-Z_\'"]*'
164     # a=1 name='mohsen'
165     pat_define_var=r'\s*^[a-zA-Z][a-zA-Z0-9_]*\s*=\s*["0-9a-zA-Z_\'"]*'
166
167     token=[pat_for,pat_while,pat_if1,pat_if2,pat_elif1,pat_elif2,pat_else,pat_print1,pat_print2,pat_op,pat_define_var]
168     all_syntax=[]
169     for tk in token:
170         compire=re.findall(tk,everylinecode,re.MULTILINE)
171         compire=''.join(compire)
172         if compire==' ' or compire==' ' or compire==None:

```

```

168     all_syntax=[]
169     for tk in token:
170         compire=re.findall(tk,everylinecode,re.MULTILINE)
171         compire=''.join(compire)
172         if compire==' ' or compire==' ' or compire==None:
173             continue
174         all_syntax.append(compire)
175     # print(all_syntax)
176     if len(all_syntax)!=1:
177         print(Fore.RED+'you have error on line : '.upper(),index,'you write : '.upper()+everylinecode)
178         print(Fore.WHITE+'')
179         quit ()
180     else:
181         print(Fore.GREEN+'code in line '.upper(),index,' is ok'.upper())
182         counter+=1
183         if counter==len(allcode):
184             compile_to_c()
185             pass
186
187
188     # print(all_line_syntax)
189     c_code="""
190     #include<stdio.h>\n
191     void main()
192     {
193     """
194
195     def compile_to_c():
196         global c_code

```

چند مثال از خطای syntax گرفتن:

در فایل py.txt مثلا به جای while مینویسیم whieele میبینیم که خطا را در خط 10ام تشخیص داد.

```
+
+
-----
ALL OF VARIABLE
1 : i
2 : index=1
3 : faktoril=5
-----
CODE IN LINE 1 IS OK
CODE IN LINE 2 IS OK
CODE IN LINE 3 IS OK
CODE IN LINE 4 IS OK
CODE IN LINE 5 IS OK
CODE IN LINE 6 IS OK
CODE IN LINE 7 IS OK
CODE IN LINE 8 IS OK
CODE IN LINE 9 IS OK
YOU HAVE ERROR ON LINE : 10 YOU WRITE : whieele index <= 5 :

p.txt
1 for i in range (1,100):
2     if i%2==0:
3         print("this is even")
4     elif i%2!=0:
5         print("this is odd")
6     else:
7         print("this is unknown")
8     index=1
9     faktoril=5
10    whieele index <= 5 :
11        faktoril=faktoril*index
12        index=index+1
13    print(faktoril)
```

یا مثلا اسم متغیر index را به index1 تغییر دهیم. در خط 8 خطا گرفت.

```
<=
*
+
-----
ALL OF VARIABLE
1 : i
2 : faktoril=5
-----
CODE IN LINE 1 IS OK
CODE IN LINE 2 IS OK
CODE IN LINE 3 IS OK
CODE IN LINE 4 IS OK
CODE IN LINE 5 IS OK
CODE IN LINE 6 IS OK
CODE IN LINE 7 IS OK
YOU HAVE ERROR ON LINE : 8 YOU WRITE : 2index=1

p.txt
1 for i in range (1,100):
2     if i%2==0:
3         print("this is even")
4     elif i%2!=0:
5         print("this is odd")
6     else:
7         print("this is unknown")
8     1index=1
9     faktoril=5
10    while index <= 5 :
11        faktoril=faktoril*index
12        index=index+1
13    print(faktoril)
```

یا تابع print را تغییر دهیم به prriint. میبینیم که خطا میگیرد در خط 7.

```
+
+
-----
ALL OF VARIABLE
1 : i
2 : index=1
3 : faktoril=5
-----
CODE IN LINE 1 IS OK
CODE IN LINE 2 IS OK
CODE IN LINE 3 IS OK
CODE IN LINE 4 IS OK
CODE IN LINE 5 IS OK
CODE IN LINE 6 IS OK
YOU HAVE ERROR ON LINE : 7 YOU WRITE : prriint("this is unknown")

p.txt
1 for i in range (1,100):
2     if i%2==0:
3         print("this is even")
4     elif i%2!=0:
5         print("this is odd")
6     else:
7         prriint("this is unknown")
8     index=1
9     faktoril=5
10    while index <= 5 :
11        faktoril=faktoril*index
12        index=index+1
13    print(faktoril)
```

بخش تبدیل به کد c:

patter های زبان c را تعریف میکنیم مثلا برای for, while, ... برای راحتی در اینجا pattern زبان پایتون فقط آن بخشهایی از کد را که نیاز داریم گذاشتیم.

```

190 c_code=""
191 #include<stdio.h>\n
192 void main()
193 {}
194 ""
195 def compile_to_c():
196     global c_code
197     file=open('p.txt','r').read().split('\n')
198
199     # define var
200     if_list=[]
201
202     # define pattern
203     # python pattern
204     pat_while_py=r'\s*[a-zA-Z0-9_]*\s*[==|>|<|>=|<=|!=]*\s*[a-zA-Z0-9_\'"]*\s*'
205     pat_print_py=r'["|\'] [a-zA-Z0-9_]*["|\']\s*\s*[()][a-zA-Z0-9_]*()'
206     pat_for_py=r'\s*[0-9]*\s*,\s*[0-9]*\s*'
207     pat_if_py=r'\s*[a-zA-Z0-9_]*\s*[==|!=|>|<|>=|<=]*\s*[a-zA-Z0-9_\'"]*\s*'
208     # pat_if1=r'\s*^if \s*[a-zA-Z0-9_]*\s*[==|!=|>|<|>=|<=]*\s*[0-9a-zA-Z\'"]*\s*;'
209     pat_op_py=r'[a-zA-Z0-9_]*\s*=\s*[a-zA-Z0-9_]*\s*[+|-|*|/] \s*[0-9a-zA-Z_\'"]*\s*'
210     # a=1 name='mohsen'
211     pat_var_py=r'\s*^[a-zA-Z][a-zA-Z0-9_]*\s*=\s*[0-9a-zA-Z_\'"]*\s*'
212     # c pattern
213     pat_for_c='for(int i={};i<{};i++)'
214     pat_while_c='while ({})'
215     pat_print_c='printf({})'
216     pat_if_c='if({})'
217     pat_else_if_c='else if({})'
218     pat_else_c='else'
219     print_='''
220
221     # help c
222     for_='''

```

مثلا while فقط به متغیر هایی که مقایسه میشوند و عملیات آن توجه میکنیم.

یا مثلا توضیح print در:

```
pat_print_py=r'["|\'] [a-zA-Z0-9_]*["|\']\s*\s*[()][a-zA-Z0-9_]*()'
```

' به معنای single quotation است و | به معنی یا. جلوی print یا یک متغیر قرار میگیرد یا یک رشته حرفی. پس در وسط یک | داریم. پس در اینجا برای راحتی خود کلمه print را نیآوریم چون نیاز به آن نداریم. البته در بخش syntax analysis آوردیم.

در اول تابع هم فایل p.txt را میخوانیم و آنرا بر اساس خط split میکنیم و حاصل در متغیر file میریزیم.

توضیح pattern های زبان C:

```
pat_for_c='for(int i={};i<{};i++)'
```

در شکل بالا از formatted string استفاده کردیم. یعنی {} یک متغیر باید در آن قرار بگیرد. که با کمک گرفتن از pattern های زبان پایتون که خودمان در بالا تعریف کردیم متغیر ها را از کد پایتون در میآورد و در ادامه این متغیر ها را توی {} قرار میدهیم.

در حلقه for بعدی هر خط file را بررسی میکنیم. باز با توجه به pattern هایی که تعریف کردیم و تابع re.findall هر خط را با pattern زبان پایتون مقایسه میکنیم و آنگاه با استفاده از تابع format patter های زبان c را پر میکنیم یعنی درون {} مقادیر را جایگذاری میکنیم. حاصل در for_ یا if_ و...

ذخیره میشود.

```
221 # help c
222 for _='
223 if ='
224 elif ='
225 else ='
226 print_=[]
227
228 for f in file[0:8]:
229     if 'for' in f:
230         x=re.findall(pat_for_py,f,re.MULTILINE)
231         x=''.join(x)
232         x=x.split(',')
233         for _=pat_for_c.format(str(x[0]),str(x[1]))+'\n{'
234         # print(for_)
235     if 'if' in f and 'elif' not in f:
236         x=re.findall(pat_if_py,f,re.MULTILINE)
237         x=''.join(x)
238         x=x.split('if')
239         if len(x)>1:
240             if _=pat_if_c.format(str(x[1]))+'\n+'{"
241             # print(if_)
242
243     if 'elif' in f:
244         x=re.findall(pat_if_py,f,re.MULTILINE)
245         x=''.join(x)
246         x=x.split('elif')
247         # print(x)
248         if len(x)>1:
249             elif _=pat_else_if_c.format(str(x[1]))+'\n+'{"
250             # print(elif_)
251
252     if 'else' in f:
253         if 'else' in f:
254             else _='else'+""\n"+"{"
255             # print(else_)
256
257         if 'print' in f :
258             x=re.findall(pat_print_py,f,re.MULTILINE)
259             x=''.join(x)
260             print_.append(pat_print_c.format(str(x))+';')
261             # print(print_)
262
263     char=""
264     list=[if_,elif_,else_]
265     if len(list)==len(print_):
266         char+=list[0].format('{'+'\n'+print_[0]+"n"}'+'\n'
267         char+=list[1].format('{'+'\n'+print_[1]+"n"}'+'\n'
268         char+=list[2].format('{'+'\n'+print_[2]+"n"}'+'\n'
269
270     # print(char)
271
272     char=for_.format('{'+'\n'+char+"n"}'+'\n')
273     # print(for_)
274
275     while _=''
276     pr _=''
277
278     for i,f in enumerate(file[7:]):
279         if i==0:
280             char+='int '+f+';'+'\n'
281         elif i==1:
282             char+='int '+f+';'+'\n'
283         elif i==2:
284             x=re.findall(pat_while_py,f,re.MULTILINE)
285             x=''.join(x)
286             x=x.split(' ')
287             # print(x)
288             if len(x)>1:
289                 while _=pat_while_c.format(str(x[1])+str(x[2])+str(x[3]))+"\n"+"{"+"\n"
290                 char+=while_
291         elif i==3:
292             char+=f+";"+'\n'
293         elif i==4:
294             char+=f+";"+'\n'+"}"+'\n'
295         elif i==5:
296             x=re.findall(pat_print_py,f,re.MULTILINE)
297             x=''.join(x)
298             if x[0]=='(' and x[-1]==')':
299                 x=x[1:-1]
300             pr _=pat_print_c.format(str(x))+';'
301             char+=pr_
302
303     # print(char)
304     c_code=c_code.format("{}+"\n"+char+"\n"+"")
305     # print(c_code)
306
307     file2=open('c.txt','w')
308     file2.write(c_code)
```

متغیر f به هر خط کد اشاره دارد.

```
279 if i==0:
280     char+='int '+f+';'+'\n'
281 elif i==1:
282     char+='int '+f+';'+'\n'
283 elif i==2:
284     x=re.findall(pat_while_py,f,re.MULTILINE)
285     x=''.join(x)
286     x=x.split(' ')
287     # print(x)
288     if len(x)>1:
289         while _=pat_while_c.format(str(x[1])+str(x[2])+str(x[3]))+"\n"+"{"+"\n"
290         char+=while_
291 elif i==3:
292     char+=f+";"+'\n'
293 elif i==4:
294     char+=f+";"+'\n'+"}"+'\n'
295 elif i==5:
296     x=re.findall(pat_print_py,f,re.MULTILINE)
297     x=''.join(x)
298     if x[0]=='(' and x[-1]==')':
299         x=x[1:-1]
300     pr _=pat_print_c.format(str(x))+';'
301     char+=pr_
302
303     # print(char)
304     c_code=c_code.format("{}+"\n"+char+"\n"+"")
305     # print(c_code)
306
307     file2=open('c.txt','w')
308     file2.write(c_code)
```

```
if 'if' in f and 'elif' not in f:
```

در این خط چون کلمه if در elif هم وجود دارد با and آنرا handle کردیم.

```
for _=pat_for_c.format(str(x[0]),str(x[1]))+'\n{'
```

متغیر x یک لیست است چون از تابع join استفاده کردیم. (مثل فاز syntax است)

بخش دوم فایل p.txt هم به همین صورت ترجمه میشود. در نهایت حاصل ترجمه را در char میریزیم.

و سپس حاصل نهایی با تابع `format` در متغیر `c_code` قرار میگیرد و یک فایل `c.txt` درست میکنیم و با تابع `c_code write` را در آن میریزیم. این کد `c` قابل اجراست و کامپایل میشود.