



پروژه پایانی

پروژه عملی درس شامل پیاده سازی یک کامپایلر تک گذره برای نسخه ساده شده ی C است که قابلیت وجود تابع های تو در تو^۱ به آن افزوده شده است. توضیحات کلی آن در ادامه آمده است. توجه کنید که استفاده از کدهای موجود در مرجع درس یا سایر کتب کامپایلر، در صورت تسلط بر آن کد و اعلام مأخذ در مستندات همراه پروژه اشکالی ندارد ولی استفاده از کدها و برنامه های موجود در سایت ها و کدهای سایر گروه ها (در همین نیم سال یا سال های گذشته) اکیدا ممنوع است و در اکثر موارد سبب مردودی در درس خواهد شد. در این مورد تفاوتی میان گروه دهنده یا گیرنده کد وجود ندارد.

^۱nested functions

مشخصات کامپایلر

- کامپایلر تک‌گذره است و از ۴ جزء تحلیل‌گر لغوی، تحلیل‌گر نحوی، تحلیل‌گر معنایی و مولد کد میانی تشکیل شده است. شما این اجزا را در ۳ فاز پیاده‌سازی می‌کنید:

فاز	اجزاء	نمره
۱	پیاده‌سازی تحلیل‌گر لغوی (اسکنر)	۰/۵
۲	پیاده‌سازی تحلیل‌گر نحوی (اسکنر و پارسر)	۰/۵
۳	پیاده‌سازی کامپایلر نهایی (اسکنر، پارسر، تحلیل‌گر معنایی و مولد کد میانی)	۳

- ورودی کامپایلر همواره یک فایل متنی حاوی برنامه‌ای است که کامپایلر شما باید آن را ترجمه کند.
- شما باید در هر فاز، بخشی از کامپایلر مورد نظر را پیاده‌سازی کنید و با روی هم سوار کردن تدریجی این بخش‌ها در نهایت به کامپایلر نهایی دست بیابید.
- نمره‌ی در نظر گرفته شده برای هر فاز مستقل از فازهای دیگر است.

گرامر NC-Minus

گرامری که در ادامه آمده است مربوط به بخشی از زبان C است که قابلیت وجود تابع‌های تودرتو به آن افزوده شده است. در این گرامر، پایانه‌ها پررنگ تر از غیرپایانه‌ها نمایش داده شده‌اند. این گرامر صرفاً برای اطلاع شما در اختیارتان قرار گرفته و شما در فاز اول نیازی به این اطلاعات ندارید.

1. $\text{program} \rightarrow \text{declaration-list } \mathbf{EOF}$
2. $\text{declaration-list} \rightarrow \text{declaration-list } \text{declaration} \mid \epsilon$
3. $\text{declaration} \rightarrow \text{var-declaration} \mid \text{fun-declaration}$
4. $\text{var-declaration} \rightarrow \text{type-specifier } \mathbf{ID} ; \mid \text{type-specifier } \mathbf{ID} [\mathbf{NUM}] ;$
5. $\text{type-specifier} \rightarrow \mathbf{int} \mid \mathbf{void}$
6. $\text{fun-declaration} \rightarrow \text{type-specifier } \mathbf{ID} (\text{params}) \text{ compound-stmt}$
7. $\text{params} \rightarrow \text{param-list} \mid \mathbf{void}$
8. $\text{param-list} \rightarrow \text{param-list} , \text{param} \mid \text{param}$
9. $\text{param} \rightarrow \text{type-specifier } \mathbf{ID} \mid \text{type-specifier } \mathbf{ID} []$
10. $\text{compound-stmt} \rightarrow \{ \text{declaration-list statement-list} \}$
11. $\text{statement-list} \rightarrow \text{statement-list } \text{statement} \mid \epsilon$
12. $\text{statement} \rightarrow \text{expression-stmt} \mid \text{compound-stmt} \mid \text{selection-stmt} \mid \text{iteration-stmt} \mid \text{return-stmt} \mid \text{switch-stmt}$
13. $\text{expression-stmt} \rightarrow \text{expression} ; \mid \mathbf{continue} ; \mid \mathbf{break} ; \mid ;$
14. $\text{selection-stmt} \rightarrow \mathbf{if} (\text{expression}) \text{statement } \mathbf{else} \text{statement}$
15. $\text{iteration-stmt} \rightarrow \mathbf{while} (\text{expression}) \text{statement}$
16. $\text{return-stmt} \rightarrow \mathbf{return} ; \mid \mathbf{return} \text{expression} ;$
17. $\text{switch-stmt} \rightarrow \mathbf{switch} (\text{expression}) \{ \text{case-stmts default-stmt} \}$
18. $\text{case-stmts} \rightarrow \text{case-stmts } \text{case-stmt} \mid \epsilon$
19. $\text{case-stmt} \rightarrow \mathbf{case} \mathbf{NUM} : \text{statement-list}$
20. $\text{default-stmt} \rightarrow \mathbf{default} : \text{statement-list} \mid \epsilon$

- 21. $\text{expression} \rightarrow \text{var} = \text{expression} \mid \text{simple-expression}$
- 22. $\text{var} \rightarrow \mathbf{ID} \mid \mathbf{ID} [\text{expression}]$
- 23. $\text{simple-expression} \rightarrow \text{additive-expression} \text{ relop } \text{additive-expression} \mid \text{additive-expression}$
- 24. $\text{relop} \rightarrow < \mid ==$
- 25. $\text{additive-expression} \rightarrow \text{additive-expression} \text{ addop } \text{term} \mid \text{term}$
- 26. $\text{addop} \rightarrow + \mid -$
- 27. $\text{term} \rightarrow \text{term} * \text{signed-factor} \mid \text{signed-factor}$
- 28. $\text{signed-factor} \rightarrow \text{factor} \mid + \text{factor} \mid - \text{factor}$
- 29. $\text{factor} \rightarrow (\text{expression}) \mid \text{var} \mid \text{call} \mid \mathbf{NUM}$
- 30. $\text{call} \rightarrow \mathbf{ID} (\text{args})$
- 31. $\text{args} \rightarrow \text{arg-list} \mid \epsilon$
- 32. $\text{arg-list} \rightarrow \text{arg-list} , \text{expression} \mid \text{expression}$

توضیحات تحلیل‌گر لغوی

شما در این فاز باید با استفاده از روشی که در کلاس آموخته‌اید، فایل متنی که به شما داده می‌شود را tokenize کنید. مراحل پیاده‌سازی تحلیل‌گر لغوی تا رسیدن به DFA را می‌توانید به صورت دستی انجام دهید. برنامه‌ای که می‌نویسید باید یک تابع به نام `get_next_token` داشته باشد که با فراخوانی آن، تا زمانی که یک token جدید تشخیص داده شود، برنامه موجود در فایل ورودی را به صورت حرف به حرف بخواند. این روند تا تشخیص تمام token های فایل ورودی ادامه پیدا می‌کند.

توجه داشته باشید که انتظار داریم تمام token ها را در یک فایل خروجی به نام `scanner.txt` قرار دهید. در این فایل به ازای هر خط از فایل ورودی یک خط قرار دارد. در واقع ساختار این فایل بدین صورت است که در ابتدای هر خط عددی می‌آید که نشان‌دهنده شماره خط فایل ورودی است. سپس در ادامه‌ی این خط به ترتیب تمام token های موجود در خط متناظر از فایل ورودی آورده می‌شوند.

هر token نیز در قالب یک دوتایی به شکل `(token type, token string)` می‌باشد. دقت کنید انواع token هایی که باید تشخیص دهید در جدول پایین آمده است:

Token Type	Description
NUM	Any string matching: <code>[0-9]⁺</code>
ID	Any string matching: <code>[A-Za-z][A-Za-z0-9][*]</code>
KEYWORD	<code>if else void int while break continue switch default case return</code>
SYMBOL	<code> ; : , [] () { } + - * = < ==</code>
COMMENT	Any string between a <code>/*</code> and a <code>*/</code> OR any string after a <code>//</code> and before a <code>\n</code>
WHITESPACE	blank (ASCII 32), <code>\n</code> (ASCII 10), <code>\r</code> (ASCII 13), <code>\t</code> (ASCII 9), <code>\v</code> (ASCII 11), <code>\f</code> (ASCII 12)
EOF	End of file

ملاحظات تحلیل‌گر لغوی

- در فایل ورودی هنگامی شماره خط عوض می‌شود که `\n` آمده باشد.
- در صورتی که متن فعلی بخشی از یک کامنت باشد، هر گونه کاراکتری می‌تواند در آن ظاهر شود.
- در هنگام تشخیص token ها نیازی به ذخیره و گزارش WHITESPACE ها و COMMENT ها نیست و این token ها را در نظر نگیرید.

خطاپردازی در تحلیل لغوی

به منظور رفع خطاهای پیش آمده هنگام تحلیل لغوی باید از روش Panic Mode استفاده کنید. در این روش، خطاپردازی که همراه تحلیل گر لغوی طراحی می کنید، می بایست هنگام برخورد با یک خطا (کاراکتر غیر مجاز و یا خانه ای خالی در جدول) در میانه ی تشخیص یک token، تمامی کاراکترهای گرفته شده از ورودی را از ابتدای آن token کنار گذاشته و فرایند تشخیص token ها را از بعد از کاراکتر آخر، از سر گیرد.

توجه داشته باشید که همراه با فایل خروجی scanner.txt باید یک فایل دیگر با نام lexical_errors.txt نیز ارائه کنید که در هر خط از با شماره خطی که در آن خطایی رخ داده است می آید و در ادامه آن دوتایی هایی می آید که هر کدام به شکل (string, message) است که string رشته ای است که هنگام خطا کنار گذاشته شده و پس از آن برای تشخیص token بعدی بررسی می گردد. message نیز پیغام خطای لغوی موردنظر را نشان می دهد. توجه داشته باشید که در این فایل خطاها باید به ترتیب رخ دادن نوشته شده باشند.

نمونه‌ی ورودی

توجه کنید که صحت خروجی تحلیل‌گر لغوی شما و خطاپرداز مربوطه آن با ورودی دادن فایل‌های تست از پیش ساخته شده بررسی خواهد شد. این تست‌ها به گونه‌ای خواهند بود که نکات مختلفی را از منظر تحلیل لغوی و هم از منظر خطاپردازی در بر خواهند داشت. نمونه‌ای ساده از یک فایل ورودی به همراه خروجی‌های موردنظر در ادامه آورده شده است:

```
1 void main(void){
2     int a = 0;
3     // comment2
4     a = 2 + +2;
5     a = a + -3;
6     cde = a;
7     if (b /* comment1 */ == 3) {
8         a = 3;
9         cd!e = -7;
10    }
11    else
12    {
13        b = a < cde;
14        {cde = @2;
15    }}
16    return;
17 }
```

نمونه فایل ورودی

```
1. (KEYWORD, void) (ID ,main), (SYMBOL, ( ), (KEYWORD, void) (SYMBOL,)) (SYMBOL,{)
2. (KEYWORD, int) (ID, a) (SYMBOL, =) (NUM, 0) (SYMBOL, ;)
4. (ID, a) (SYMBOL, =) (NUM, 2) (SYMBOL, +) (SYMBOL, +) (NUM, 2) (SYMBOL, ;)
5. (ID, a) (SYMBOL, =) (ID, a) (SYMBOL, +) (SYMBOL, -) (NUM, 3) (SYMBOL, ;)
6. (ID, cde) (SYMBOL, =) (ID, a) (SYMBOL, ;)
7. (KEYWORD, if) (SYMBOL, ( ) (ID, b) (SYMBOL, ==) (NUM, 3) (SYMBOL, )) (SYMBOL, {)
8. (ID, a) (SYMBOL, =) (NUM, 3) (SYMBOL, ;)
9. (ID, e) (SYMBOL, =) (SYMBOL, -) (NUM, 7) (SYMBOL, ;)
10. (SYMBOL, })
11. (KEYWORD, else)
12. (SYMBOL, {)
13. (ID, b) (SYMBOL, =) (ID, a) (SYMBOL, <) (ID, cde) (SYMBOL, ;)
14. (SYMBOL, {) (ID, cde) (SYMBOL, =) (NUM, 2) (SYMBOL, ;)
15. (SYMBOL, }) (SYMBOL, })
16. (KEYWORD, return) (SYMBOL, ;)
17. (SYMBOL, })
```

نمونه فایل خروجی تحلیل‌گر لغوی

```
9. (cd!, invalid input)
14. (@, invalid input)
```

نمونه فایل خطاپرداز لغوی

انجام پروژه

- پروژه باید به صورت انفرادی و یا گروه‌های دو نفره انجام شود. در صورتی که میزان مشارکت اعضای گروه‌های دو نفره با یکدیگر برابر نباشد، فردی که مشارکت کمتری داشته، نمره‌ی کمتری نسبت به دیگری می‌گیرد.
- در صورتی که هر گونه سوالی در رابطه با تعریف پروژه دارید، آن را از طریق کوئرا مطرح نمایید.
- مهلت بارگذاری سورس کد پروژه برای فاز اول، ساعت ۲۳:۵۹ روز جمعه ۲۳ فروردین است.
- تحویل حضوری دو فاز اول و دوم همزمان بوده و زمان‌بندی دقیق آن متعاقبا اعلام خواهد شد. (توجه کنید که حضور هر دو عضو گروه‌های دو نفره در جلسه‌ی تحویل الزامی است.)

موفق باشید.