

به نام خدا

## گزارش فاز اول پروژه طراحی تحلیل گر لغوی

فاطمه شفیعی اردستانی

رضا پورمحمدی

در ابتدا برای راحتی در نوشتن قوانین متغیرهای زیر را تعریف می کنیم. `digit` به منظور نمایش ارقام ۰ تا ۹ ، `printabe` به منظور نمایش تمام کارکترهای قابل چاپ و `alpha` برای نمایش تمام حروف انگلیسی کوچک و بزرگ در نظر گرفته شده اند.

```
digit      [0-9]
printabe   [ -~]
alpha      [a-zA-Z]
```

باتوجه به پروژه تعریف شده، توکن های زیر باید شناسایی شوند:

۱-مقادیر ثابت عددی:

از آنجایی که اعداد صحیح می توانند مثبت یا منفی باشند، برنامه باید بتواند تشخیص دهد علامت قبل از عدد ثابت، عملگر است یا علامت عدد است؛ به این منظور همواره دو توکن آخری که شناسایی شده اند را ذخیره می کنیم.

سپس هنگام دیدن دنباله ای از اعداد و بعد از تشخیص اینکه توکن ثابت داریم، بررسی می کنیم که آیا توکن قبلی، عملگر + یا عملگر - بوده است. اگر چنین عملگرهایی وجود نداشتند، بررسی می کنیم که مقدار عدد ثابت، در بازه ی اعداد ۱۶ بیتی هست یا نه و با توجه به نتیجه، توکن مناسب را به آن اختصاص می دهیم.

اگر قبل از عدد، عملگر وجود داشت، بررسی می کنیم که قبل از عملگر چیست، اگر عدد یا متغیر بود، مجدداً عدد مشاهده شده را باتوجه به مقدارش با توکن مناسب، مشخص می کنیم. اما اگر قبل از علامت مثبت یا منفی، توکنی جز عدد یا متغیر داشتیم، به این معناست که علامت، متعلق به خود عدد است، بنابراین توکنی که با عنوان توکن علامت نام گذاری شده بود را، با توکن عدد ثابت، رونویسی می کنیم.

```
{digit}+ {
    if(pre2 == "TOKEN_PLUS\n" || pre2 == "TOKEN_MINUSE\n"){
        if(pre1 == "TOKEN_INTCONST\n" || pre1 == "TOKEN_INVALIDINTCONST\n" ||
pre1 == "TOKEN_IDENTIFIER\n"){
            fprintf(fp, pre1);
            pre1 = pre2;
```

```

        n = atoi(yytext);
        if (n<=32767){
            pre2 = "TOKEN_INTCONST\n";
        }
        else{
            pre2 = "TOKEN_INVALIDINTCONST\n";
        }
    }
    else{
        n = atoi(yytext);
        if (n<=32767){
            pre2 = "TOKEN_INTCONST\n";
        }
        else{
            pre2 = "TOKEN_INVALIDINTCONST\n";
        }
    }
}
else{
    fprintf(fp, pre1);
    pre1 = pre2;
    n = atoi(yytext);
    if (n<=32767){
        pre2 = "TOKEN_INTCONST\n";
    }
    else{
        pre2 = "TOKEN_INVALIDINTCONST\n";
    }
}
}
}

```

۲-مقادیر ثابت کاراکتری:

مقادیر ثابت کاراکتری بین دو دبل کوتیشن قرار دارند و ممکن است شامل هر کارکتر قابل چاپی باشند. بنابر این قانون آن به فرم زیر نوشته می شود:

```

(\"{printabe}*\\")

```

۳- نام متغیرها و توابع:

باتوجه به اینکه نام متغیر و توابع باید با حروف یا \_ شروع شوند عبارت زیر را برای تشخیص آن داریم:

```

({alpha}|_)( {alpha}|_|{digit})*

```

۴- کامنت تک خطی یا چند خطی:

در کامنت چند خطی می‌توانیم بعد از دیدن کاراکتر # هر کاراکتری به جز \n را به هر تعداد ببینیم، بنابراین عبارت منظم آن به شکل زیر است:

```
(#(.)*)
```

در کامنت چند خطی، بعد از دیدن #! تا زمان دیدن #! می‌توان همه کاراکترها را به هر تعداد دید بنابراین داریم:

```
(#!({printabe}|"\\n")*!#)
```

۵- کلمات کلیدی، عملگرها و سایر توکن‌ها:

برای تشخیص این توکن‌ها، کافی است نام آنها را در بین "" قرار دهیم.

```
"int"
```